

Cluster-Computing and Parallelization for the Multi-Dimensional PH-Index

Master Thesis**Author(s):**

Vancea, Bogdan Aurel

Publication date:

2015-04

Permanent link:

<https://doi.org/10.3929/ethz-a-010437712>

Rights / license:

In Copyright - Non-Commercial Use Permitted

Cluster-Computing and Parallelization for the Multi-Dimensional PH-Index

Master Thesis

Bogdan Aurel Vancea

<bvancea@student.ethz.ch>

Prof. Dr. Moira C. Norrie

Dr. Tilmann Zäschke

Christoph Zimmerli

Global Information Systems Group
Institute of Information Systems
Department of Computer Science
ETH Zurich

1st April 2015



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



Abstract

The storage and management of multi-dimensional data is an important aspect in many applications, like geo-information systems, computer vision and computer geometry. In contemporary times, when computers are able to capture and store increasing amounts of multi-dimensional data, it is important for applications to efficiently store and query this data.

This work presents a distributed version of the PH-tree, a highly efficient in-memory multi-dimensional data structure, supporting range and nearest neighbour queries. We present a distribution architecture for the PH-tree, which extends it to run on a cluster of computers. The distributed version of the PH-tree is able to use the main-memory of all of machines in the cluster to store the multi-dimensional data. Moreover, the distributed setting allows each machine to be queried independently of the other machines. The point distribution algorithm proposed in this work uses the Z-order space filling curve to assign sections of the space to the computers in the cluster. Additionally, we present an automatic data re-balancing algorithm, which attempts to maintain an equal storage load across all computers. The performance evaluation shows that the proposed distribution version of the PH-tree is able to scale with respect to the number of computer in the cluster, obtaining an almost linear increase in throughput and a similar reduction in the response time for point operations.

Another contribution of this work is the extension of the PH-tree to support concurrent write access, allowing it to take advantage of multi-core processor architectures. We present several concurrent write access strategies, with different consistency guarantees, discuss their advantages and disadvantages and evaluate the write performance. The presented copy-on-write strategy allows queries to execute on snapshots of the PH-tree, while the PH-tree itself is accessed by one writer and multiple readers. This work also presents two fine grained locking strategies, which sacrifice consistency to allow multiple reader and writer threads to access the tree in the same time. All presented concurrency strategies perform locking only in the case of the write operations, and allow read queries to execute without taking any locks, which greatly increases the performance of read operations for large numbers of reader threads.

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Objectives	2
1.3	Thesis Outline	2
2	Background	5
2.1	The PH-tree	5
2.2	Related Work	7
2.2.1	Distributed Indexes	7
2.2.2	Concurrent Data Structures	7
3	Index Distribution	9
3.1	Challenges	9
3.2	Distribution Strategies	10
3.2.1	Hashing	11
3.2.2	Spatial Splitting	12
3.2.3	Z-Order Curve Splitting	14
3.3	Architecture	15
3.4	Algorithms	16
3.4.1	Key-Mapping Design	16
3.4.2	Point Operations	18
3.4.3	Range Queries	18
3.4.4	Nearest Neighbour Queries	19
3.4.5	Entry Rebalancing	20
4	Concurrency	23
4.1	Challenges	23
4.2	PH-tree Structure	24

4.3	Concurrency Strategies	25
4.3.1	Copy-on-Write	25
4.3.2	Locking	26
4.3.3	Hybrid Approach	28
4.4	Chosen Strategies	29
5	Implementation	31
5.1	System Architecture	31
5.2	PH-tree Java API	32
5.3	Distribution	34
5.3.1	Client Design	34
5.3.2	Server Design	35
5.3.3	Configuration Server and Key-Mapping Implementation	36
5.3.4	Communication Flow	38
5.3.5	Iterators	39
5.4	Concurrency	40
6	Evaluation	43
6.1	Distribution	43
6.1.1	Experimental Setup	43
6.1.2	Experimental Results	44
6.2	Concurrency	47
6.2.1	Experimental Setup	47
6.2.2	Experimental Results	47
6.3	Results	51
7	Conclusions	53
7.1	Contributions	53
7.2	Future work	54
A	Installation and Usage Instructions	55
A.1	Compilation and Installation	55
A.1.1	Client and Server Libraries	55
A.1.2	ZooKeeper	55
A.2	Usage	56

1

Introduction

1.1 Motivation

Multi-dimensional data is widely used today, for example in domains like database management systems, geographic information systems, computer vision and computational geometry. When all of the dimensions of the data hold numerical values, this data can be viewed as a collection of points in higher dimensional spaces. Due to this nature, multi-dimensional numerical data provides the possibility of posing more complex queries based on the distance between these points in space. For example, in the context of a geo-information system, one could query for all of the points that fall inside a specific rectangle or attempt to find the nearest neighbours of an arbitrary query point.

Several point-based multi-dimensional indexing solution have been developed in the latest years, the most prominent being kD-trees [1] and quadtrees [4]. This type of data structures store the multi-dimensional data such that more complex operations, like range and nearest neighbour queries are executed efficiently. The PH-tree [11] is a new multi-dimensional data structure based on the quadtree. In addition to providing support for complex queries, the PH-tree is also space-efficient, as its space requirements are sometimes even lower than those of multi-dimensional arrays.

As technology advances and the world becomes more connected, multi-dimensional data becomes easier to acquire and store. Because of this, it is necessary that multi-dimensional data structures need to store and manage more data than would fit in the main memory of a single machine. However, traditional multi-dimensional indexes like the kD-tree and quadtree do not cover this use case as they are designed to run on a single machine.

Additionally, in the last few years the processor speed has reached the power wall and processor designers cannot increase the CPU frequency by increasing the number of transistors. Recent advances in processor design have been made by adding more cores per CPU rather

than increasing the processing frequency. Therefore, it is important that contemporary data structures be adapted to multi-core architectures by allowing them to support concurrent operations. As with the case of the increase storage requirements, traditional multi-dimensional data structures do not support concurrent write operations.

This thesis attempts to provide a solution to these two issues by extending the PH-tree to run on distributed cluster of machines and modifying it to support concurrent operations.

1.2 Objectives

We have highlighted two challenges currently faced by indexing systems: *high storage requirements* and *support for concurrent access*. This work proposes the distributed PH-tree, a version of the PH-tree that can be run on a cluster of machines, making it able to handle data sets that cannot fit in the main memory of a single machine. Moreover, the distributed PH-tree should be able to handle concurrent requests. This applies both to requests sent to different machines that are part of the cluster and concurrent requests sent by different clients to the same machine.

Specifically, the distributed PH-tree has to fulfill the following requirements:

Cluster capability The system should run across a network of machines, making use of the memory and processing resources on each machine. Furthermore, the system should attempt to balance the number of multi-dimensional entries assigned to each machine.

Cluster concurrency The system should be able to support concurrent requests to different machines in the cluster. Each machine should be able to process queries related to the entries that it stores locally.

Node concurrency Each machine should support multi-threaded read and write access to the entries that it stores.

As the thesis touches on two main subjects, distribution and concurrency, the main challenges encountered are twofold. From the distribution perspective, the challenges are the identification of suitable entry distribution and balancing strategies, devising efficient algorithms for executing queries across multiple cluster nodes, and the efficient management of a very large number of cluster nodes. For the concurrency perspective, the challenges are the identification of a suitable concurrent access strategy that can maximize the number of concurrent write operations.

1.3 Thesis Outline

This chapter gave an overview of the challenges currently faced by multi-dimensional indexing structures and briefly explained how this work seeks to address them. Additionally, this chapter also presented the main objectives of this thesis. The rest of the thesis is structured as follows:

Chapter 2 provides additional information about the PH-tree, its characteristics and supported operations. The second part of this chapter describes relevant previous work done in the areas of distributed multi-dimensional indexes and concurrent data structures.

The design of the distributed PH-tree from the point of view of a distributed system is presented in **Chapter 3**. This chapter presents the chosen data distribution strategy and touches on the possible alternatives and the consequences of this choice. Additionally, we provide an overview of how the queries spanning multiple nodes are executed by the system.

The addition of the multi-threaded read and write support for the PH-tree is presented in **Chapter 4**. Several concurrent access strategies are discussed, in terms of their advantages, disadvantages and consistency guarantees.

Chapter 5 describes the implementation-specific decisions that were taken during the development process. This chapter also presents the technologies that were used and justifies the technological choices.

The distributed PH-tree is evaluated in **Chapter 6**. The performance characteristics of the implemented distributed systems as well as those of the implemented concurrency strategy are discussed.

Chapter 7 concludes the thesis by presenting the contribution of this work in the context of distributed multi-dimensional indexing systems. We also give an outlook of the potential extensions.

2

Background

The first part of this chapter analyses the single-threaded PH-tree, a multi-dimensional data structure and the starting point of this work. It provides an overview of its specific characteristics and describes the supported operations.

The second part of this chapter presents the relevant related work in the area of distributed multi-dimensional indexes and concurrent data structures.

2.1 The PH-tree

The PH-tree [11] is a novel multi-dimensional indexing structure that focuses on fast retrieval, fast update and highly efficient storage. It combines concepts from PATRICIA-tries, quadrees and hypercubes to encode the multi-dimensional data using bit prefix sharing. Because of this, the storage requirements can sometimes be lower than those of the storage requirements of an array of objects.

An example PH-tree storing 2-dimensional 4 bit points is shown in Figure 2.1. The tree illustrated in this figure stores the points: (0, 4), (3, 4) and (3, 6). However, each point is stored as the *bit interleaving* of its values in all dimensions. Therefore, the bit interleaving of point (3, 6), represented in binary as (0011, 1010) is 01001110. By storing the points under the *bit interleaving* form, the PH-tree maps the multi-dimensional points into a 1-dimensional space. The 1-dimensional values will be referred to as *Z-values* as they form the *Z-order space filling curve*.

The PH-tree has the following important properties:

- **Z-Ordering.** The multi-dimensional entries are ordered according to the Z-order space filling curve. As a consequence, points which are close in the original multi-dimensional space will usually also be relatively close in the Z-Ordering.

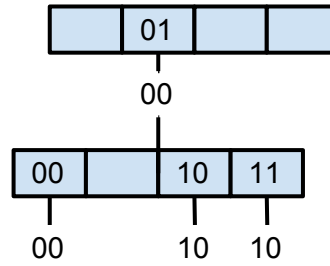


Figure 2.1: 2 dimensional 4bit PH-tree containing the points: (0000, 1000), (0011, 1000), (0011, 1010).

- **Constant maximum height.** The height of the PH-tree is equal the number of bits needed to store each value of the multi-dimensional space. Therefore, when the values in each dimension are represented in 64 bits, as long or doubles, the maximum height of the PH-tree is 64. This property is important as *the maximum height of the PH-tree is independent of the number of points it stores*.
- **Balancing is not needed.** The PH-tree does not perform balancing operation after write operation have occurred. While this fact could lead to un-balanced trees, it is usually not a problem in practice, as the height of tree is limited by a constant factor.

The PH-tree is similar to a HashMap, as it stores entries consisting of key-value pairs. In this case however, the key is a multi-dimensional data point, while the value can be an arbitrary object. The following operations are supported:

- **Point operations** operate on single multi-dimensional points. The possible queries are *get()*, *put()*, *delete()* and *updateKey()*. The *get()* query retrieves the value associated with a multi-dimensional key, and the *delete()* operation removes the value associated with a multi-dimensional key. The *put()* operation adds new key-value to the index. Finally, the *updateKey()* operations updates the key associated with an existing key-value entry.
- **Range queries** select all of the data points that fall inside a multi-dimensional hyper-rectangle. This rectangle is defined by a *lower left* and a *upper right* point.
- **Nearest neighbour queries** select the closest k points in the index to an arbitrary query point.

Currently, the amount of entries the PH-tree can store limited by the memory available on a single host. A distributed version of the PH-tree would be able to store a much larger amount of entries and would be able to make use of clusters of computing nodes. Specifically, the distributed version would be well suited for running in cloud environments, where users can easily add more computing nodes to their existing systems to handle higher storage requirements and improve performance. The challenges in creating a distributed version of the PH-tree are the implementation of range and nearest neighbour queries, which might need to be dispatched to multiple nodes of the system. The goal is to design a distribution architecture and a set of algorithms that minimize the number of nodes to which a complex query is sent.

2.2 Related Work

The following sections reviews the existing research concerning distributed indexes and concurrent data structures. The goal of this section is to understand the existing approaches that can be used for the design and implementation of the distributed PH-tree and to identify which of these approaches are best suited for the use cases targeted by this work.

2.2.1 Distributed Indexes

Distributed Hash Tables (DHT's) are a class of decentralized distributed systems which expose an API similar to that of a hash-table. These P2P systems use a structured overlay network to assign a given key to a node in the system. Requests for a certain key can be sent to any node of the system and will be forwarded to the node responsible for it. DHT's have been designed to work on internet-scale, connecting large numbers of machines across large distances. The main concerns are scalability and availability, while the *get()* and *put()* operations have "best effort" semantics. Chord [10], one of the most popular and most efficient DHT system, manages to guarantee that any request for a certain key will arrive the responsible node in at most $O(\log n)$ hops, in the context in which each node contains only $O(\log n)$ links to other nodes.

One important issues of distributed hash tables is that they only support "exact match" queries. The Prefix Hash Tree [8] adds support for range queries to a DHT by using a secondary data structure, to maintain a mapping between nodes and keys. This mapping is a trie containing the prefixes of the bit-strings of the keys, and thus, logically, leaf nodes of this trie correspond to interval regions of keys. All leaf-nods are linked using a linked-list, allowing range queries to traverse the query range starting for the first leaf that matches the range.

A similar approach is taken by the SkipIndex [12], a multi-dimensional indexing system that uses a SkipGraph to map geometrical regions to machine nodes. Each system node maintains only a partial view of the region tree and thus range and nearest neighbour queries are performed using selective multicasting.

A different approach for creating a scalable multi-dimensional index is to use a big-data platform. [6] proposes storing a block-based multi-dimensional index structure like the R-tree directly in HDFS¹ and using query models like MapReduce² as building blocks for the implementation of range and nearest neighbour queries.

2.2.2 Concurrent Data Structures

Concurrent data structures are currently a very active area of research. There are two main ways of implementing a concurrent data structures: lock-based solution and lock-free solutions. Lock-based solutions use locks to synchronize the concurrent access to either the data structure itself, called *course-grained locking*, or to parts of the data structure, which is also called *fine grained locking*. Lock-free solutions use either atomic primitives, like *compare-*

¹http://hadoop.apache.org/docs/r1.2.1/hdfs_design.html

²<http://en.wikipedia.org/wiki/MapReduce>

and-swap or software transactional memory to synchronize modification attempts on the data structure.

In the area of concurrent search trees, lock-based approaches are much more common than lock-free solutions. Fraser [5] provides a set of API's that can be used to design arbitrary concurrent data structures. While these API's can be used for any type of data structures, they are based on software transactional memory or multi-word CAS operations, currently not implemented by contemporary hardware. Brown [3] presents a general technique for non-blocking trees implemented using multi-word equivalents of the *load-link*(LL), *store-conditional* (SC) and *validate*(VL) atomic primitives.

A naive concurrent implementation of a search tree can be achieved by protecting all access using a single lock. Concurrency can be improved by synchronizing the access at the level of nodes. *Hand-over-hand locking* [7], something also called *lock-coupling*, is a fine-grained technique that states that a thread can only acquire the lock of a node if it holds the lock of the node's parent node. This technique allows multiple concurrent write operations, but the degree of concurrency is limited due to the fact that all subtrees of a locked node are not accessible, even if the other threads need to perform modifications in a different part of the tree. Another approach is to perform optimistic retries by only locking the nodes that should be modified, check if locked nodes have not been removed from the tree and then perform the update. This idea is used by [2] to implement a highly concurrent AVL tree.

However, an issue with all concurrent search trees where modifications are done to a single node is that the execution of queries which have to run on the tree for longer periods of time, like iterators or range queries, overlaps with the execution of other write-operations. Therefore, longer running queries will view the any updates that are done on the tree. In some cases, it is desirable to run the read queries on a snapshot of the tree that is isolated from other concurrent updates. This can be achieved using a strategy called *Copy-On-Write*, which allows a single writer and an arbitrary number of reads to access the same data structure concurrently. Moreover, each write operations creates a new version of the data structure that is only available to readers which have started after the update which created the new version has finished.

3

Index Distribution

This chapter describes how PH-tree is extended to be deployed and run on a number of machines part of a computer cluster. It presents the challenges and the possible design choices for extending the PH-tree to run in a distributed setting, discusses the system architecture and describes the algorithms used for the distributed queries. Concurrency related issues are addressed separately in Chapter 4.

3.1 Challenges

While extending the PH-tree to be a distributed index increases its storage capacity, there are a number of challenges to overcome in order to reach a good implementation. These challenges are:

- **Balancing the storage load.** In a cluster of n identical nodes¹, nodes which have the same hardware resources available, all nodes should store an equal amount of data. Assuming that the values associated with each key are relatively equal in size, all nodes should store a relatively equal number of key-value entries.
- **Maintaining data locality.** The PH-tree needs to provide support for range and nearest neighbour queries, types of queries which generally analyze contiguous portions of the space and return sets of points which are generally close together. If the indexed points manifest locality, points close in space reside on the same node, range and nearest neighbour queries could be sent to only a part of the nodes in the cluster, minimizing the network traffic. In the best case, range and nearest neighbour queries could be sent to only a single node. It is therefore preferable that queries are sent to as few nodes as

¹Cluster machines on which the PH-tree is deployed will be referred to as nodes, indexing hosts or indexing servers, throughout this work.

possible, minimizing both the response time and the request load to the nodes in the cluster.

- **Ensuring cluster scalability.** To support very large storage requirements, the system should handle a large number of online-indexing nodes. Furthermore, to improve scalability, adding and removing nodes to and from the system should be easy and should not require the cluster to be shut down and re-started. Online cluster re-configuration is challenging because of the data migrations that it entails in order to achieve a balanced storage load. New nodes added to the system start off as being empty and should receive some entries from the currently running active nodes. The data received by these empty nodes should preserve locality. Nodes that are being removed from the system need to first make sure that the entries currently stored are moved to different nodes in the system in a manner that preserves locality.
- **Minimizing node information.** As the PH-tree is held in-memory, it is important for the indexing nodes to be efficient in managing their memory, to maximize the number of entries they can store. Maintaining a connection between two node over a network generally requires an open socket on both participating nodes. Active connections take up system memory and CPU cycles on indexing hosts and generate traffic on the network connecting them. Reducing the number of open connections that a host has to other hosts in the cluster reduces the network traffic between the machines in the cluster and reduces the amount of memory and CPU cycles used by these machines to keep the connections alive. Therefore, it is important to devise a communication model in which each node has to be aware and communicate with a small number of other nodes in the system, to reduce the resource consumption.

Other challenges posed by distributed indexes, and distributed systems in general, are availability and security, however these issues are not tackled by this work.

By taking into account the discussed challenges, an ideal architecture of the system has the following requirements:

1. Ensures low response time by minimizing the number of network requests that have to be made.
2. Points are assigned to nodes in the cluster in a way that preserves locality.
3. New nodes can be added and removed without shutting down or stopping the cluster.
4. Each node maintains open connections to a small number of other nodes in the system.

The following sections of this chapter describe the possible strategies for distributing the points to nodes in the cluster, the strategy chosen to properly address the challenges presented in this section, the final architecture of the system and the algorithms for the query execution.

3.2 Distribution Strategies

There are two main solutions to the problem of distributing a tree data structure over a cluster of computing nodes:

- **Assigning tree nodes to machines.** In this approach, there will be only a single PH-tree containing all of the entries stored by the cluster, each node of the PH-tree being stored on an individual machine. This means that following a pointer from a parent tree node to a child tree node could require the forwarding of the request to a different machine, the one containing the child tree node. As the PH-tree depth is limited by the bit width of the data being stored, from the theoretical point of view, the number of such "forwards" is limited by a constant number. Additionally, spreading the tree nodes uniformly across the machines will lead to a balanced storage load. The drawback of this approach is that even though the number of "forwards" during a tree travels is limited by a constant number, the impact on the response time of the system is very high as each "forward" corresponds to an additional network request. For example, for indexes storing 64 bit multi-dimensional points, a point query would take 64 forwards between the machines in the cluster in the worst case. Moreover, the machines which stores the root node will receive all of the requests for any data, and will quickly become a bottleneck.
- **Assigning data points to machines.** In this approach, each machine in the system will maintain a in-memory PH-tree and data points are assigned to machines according to some mapping. This mapping will be referred to as the *key-mapping*, as it maintains a mapping between the keys of the entries stored and the machines that store them. Queries for arbitrary points will always be dispatched to the machines that hold the points, according the *key-mapping*. Therefore, if the *key-mapping* is known by the clients of the index, point queries can be resolved by a single network request. However, in such an architecture, the manner in which the points are distributed according to the *key-mapping* influences the manner in which the range and nearest neighbour queries are executed.

While both of these alternatives satisfy some of the requirements of an ideal architecture, the drawbacks of the first approach, the assignment of tree nodes to machines, outweigh the advantages it provides. The number of network requests needed by this approach when storing 32 or 64 bit data, typical for most applications, make the use of it prohibitive in a setting where queries must have a low-latency.

Therefore, this work will focus on the second distribution approach. The following sections will present possible strategies for distributing the points to PH-tree's store on different machines in the cluster.

3.2.1 Hashing

The first approach considered is using a *hashing* function to assign the multi-dimensional points to machines in the cluster. This process is illustrated in Figure 3.1. The hashing function takes the multi-dimensional point and generates a host id h based on the point values in each dimension. This approach is similar to the approach used to implement in-memory hash tables, where keys are mapped to buckets in the hash table.

The main advantage of this approach is that if the hash function is known to all of the machines in the cluster and to all of the clients, point queries can be executed in $O(1)$ network

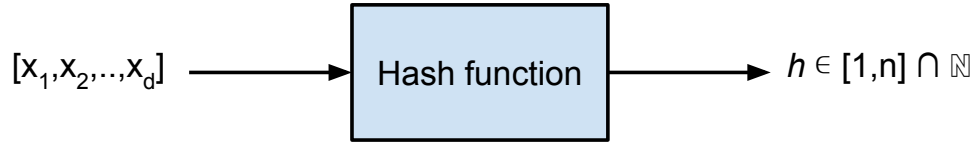


Figure 3.1: Distributing the points based on a hash function

requests. Moreover, if the hash function is uniform, the multi-dimensional points are spread uniformly across the cluster, leading to a balanced storage load on all of the machines.

However, the use of a hash function for the point distribution presents some drawbacks. First of all, a hash function works well for a fixed number of hosts n . A different cluster configuration, with a different number of machines, requires the use of a different hash function. Moreover, changing the hash function generally requires a re-hashing of all of the points already stored in the cluster², potentially leading to a situation where all of the stored points need to be moved to a different machine. This situation would cause a lot of traffic within the cluster, making it hard for the system to answer queries until the migration process is finished.

A second issue with use of hash functions is the observation that hash functions which are perceived to be good in this situation spread the multi-dimensional points uniformly across the machines in the cluster. In such a situation, it is quite likely that the hash function does not preserve point locality, and such, complex queries like range queries and nearest neighbour queries will need to be dispatched to all of the computing nodes in the cluster.

3.2.2 Spatial Splitting

Another approach that can be used for the point distribution problem is to partition the original space into a set of contiguous regions and assign each region to a host in the system. The most intuitive way to perform this split is to split the multi-dimensional space equally to obtain one region per index host. Figure 3.2 shows a potential splitting of a 2D space for 2 hosts, while Figure 3.3 presents a possible split of the same 2D space for 4 hosts.

The spatial splitting approach is not only intuitive, but it also preserves locality, as generally, points which are close together are part of the same regions and are therefore stored on the same machine. Additionally, range queries can be resolved by dispatching the range query to all of the hosts whose assigned regions intersect with the range received as input. In the best case, range queries can be dispatched to a single host. It is generally preferable to split the space into hyper-rectangles, as these regions can be stored using only two data points, the bottom right and top left corners, independently of the number of dimensions of the original space.

While the spatial splitting works well in some situations, sometimes it is unclear how to properly split the space between the machines. For example, the 2D space is easy to split into 2 or 4 equal regions, but it is more difficult to split for 3 or 5 hosts. It is always possible to split the

²This can be circumvented by using an extendible hash function, which does not perform a full re-hash when a bucket is added.

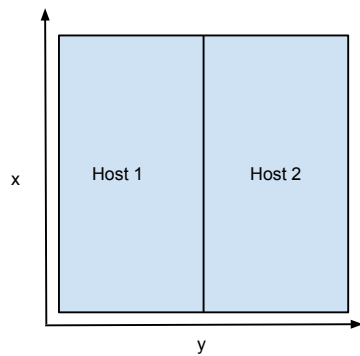


Figure 3.2: Equal geometric space splitting for 2 hosts and 2 dimensions

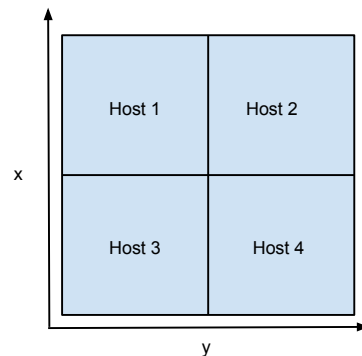


Figure 3.3: Equal geometric space splitting for 4 hosts and 2 dimensions

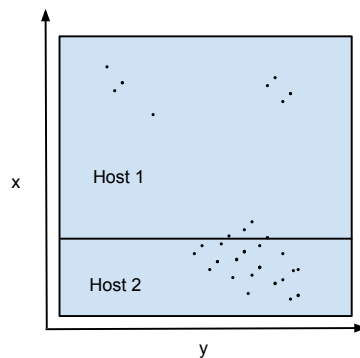


Figure 3.4: Potential partitioning of the space for a skewed input distribution.

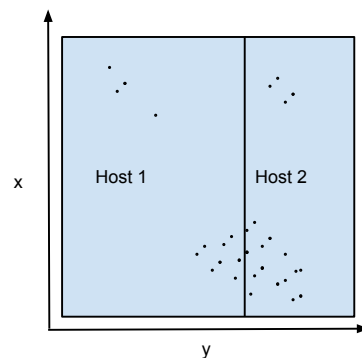


Figure 3.5: Another potential partitioning for the same skewed input distribution.

space into n equally sized "strips", however this approach diminishes the locality of the data for large values of n . Representing regions as arbitrary polygons instead of hyper-rectangles addresses the problem to a degree, but storing these polygons and computing intersection between them is more difficult than performing the same task using hyper-rectangles.

Difficulties in splitting the space can also arise even when using hyper-rectangle regions of different size. One instance is the case of skewed point distributions, presented in Figures 3.4 and 3.5. In the case of the point distributions illustrated in these figures, an equal space partitioning would not split the points equally between the hosts. However, there are many possible rectangle-based splittings of the space, and it is unclear which of these it is better. In some case, the best split could be achieved by splitting the space into more regions than machines and assign each region to a machine.

Additionally, in the case of skewed input distributions, the optimal split of the space into regions will change over time, as some regions will end up containing more and more points. In this case, the cluster needs to go through a *re-balancing* phase, when the regions are modified and some points are moved from highly populated regions to less populated neighbouring regions. An important observation is that the number of neighbouring regions for an arbitrary

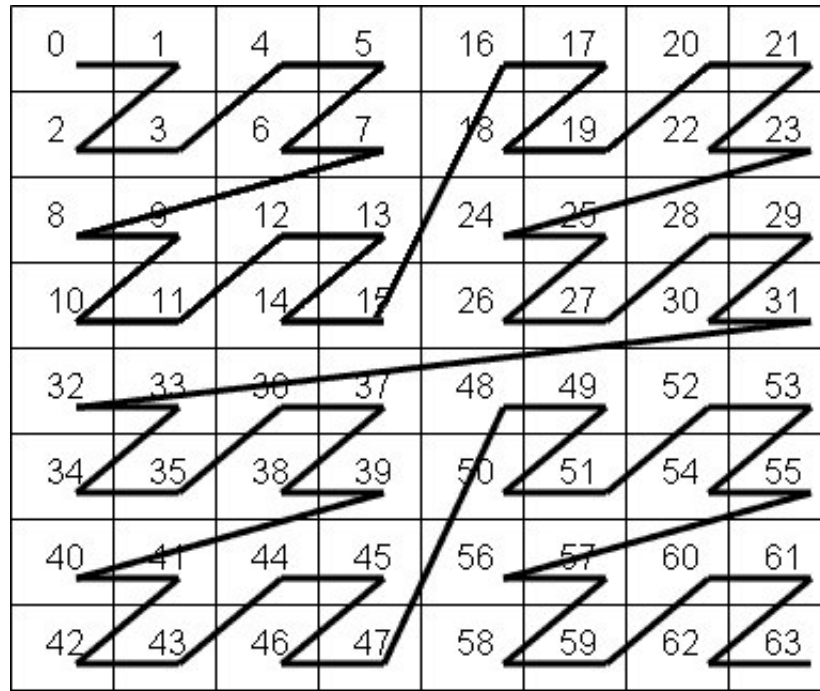


Figure 3.6: Z-order space filling curve filling the 2 dimensional space.

region is a function of the number of dimensions of the original space. This leads to an exponential number of regions participating, or at least being considered, in the re-balancing of the points from a single dense region to its neighbours.

3.2.3 Z-Order Curve Splitting

As previously mentioned, the in-memory PH-tree stores the multi-dimensional points according to the Z-ordering, by mapping the initial space to the 1-dimensional Z-order curve. This curve fills in the initial space in a contiguous and complete manner. Figure 3.6³ shows the Z-order space filling curve filling the 2D space where the values in each dimension are limited by 3 bits.

Therefore, instead of splitting the original space into regions, one could attempt to partition the Z-order curve into contiguous intervals. Figure 3.7 shows the 1 dimensional Z-order curve corresponding to the same 3 bit 2-dimensional space from Figure 3.6, split into 4 intervals of different colours: $[0, 12]$ with blue, $[13, 35]$ with yellow, $[36, 47]$ with violet and $[48, 63]$ with green. Figure 3.8 shows how the interval splitting is reflected in the original space. Each interval on the Z-order curve corresponds to either a hyper-rectangle or a set of hyper-rectangles. For each interval there are at most $O(d * w)$ of these hyper-rectangles, where d is the number of dimensions and w is the number of bits needed to represent the values in each dimension. The problem of generating the representation of the multi-dimensional space from a Z-order curve is described in [9], which provides an algorithm for generating

³Image adapted from http://www.scholarpedia.org/article/B-tree_and_UB-tree. Various other adaptations of this image have been used for exemplification purposes.

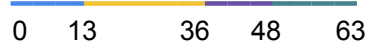


Figure 3.7: Example intervals on the 3bit Z-order curve

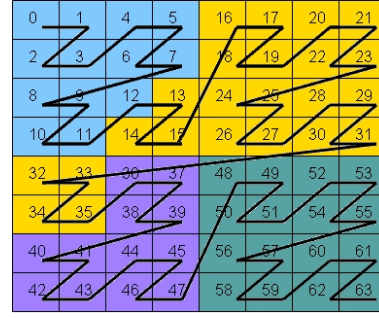


Figure 3.8: Example intervals on the 3 bit Z-order curve filling the 2D space

the hyper-rectangles associated with a Z-order curve range.

An advantage of this splitting method is the regions corresponding to the Z-order intervals manifest locality, making it well suited for a point distribution method. Moreover, even if in the original space each region has up to an exponential number of neighbouring regions, each interval on the Z-order curve has a small, constant number of neighbouring intervals (2 for the interior intervals, 1 for the two edge intervals). This makes it easier to handle the cases in which a region is densely populated and should re-balance its points to neighbouring regions.

Due to the presented advantages, the Z-order curve interval splitting is the method that was chosen as a point distribution strategy for the distributed PH-tree.

3.3 Architecture

The network architecture of the distributed PH-tree is presented in Figure 3.9. The points are distributed to the hosts in the cluster using the Z-order curve interval splitting method. The *key-mapping* is stored on the configuration server and is known to all of the nodes and the clients of the system. The configuration server uses a notification-based model and is responsible for notifying the indexing hosts and clients whenever any changes occur. This model reduces the amount of network traffic within the cluster, as nodes and clients only need to read the *key-mapping* on start-up and do not need to poll the configuration server for changes.

Clients of the index first connect to the configuration server, retrieve the *key-mapping* and register to receive notifications if the mapping changes. After receiving the mapping, clients can begin to send queries to the indexing nodes within the cluster. Clients decide which hosts need to be contacted for a query based on the type of the query and the *key-mapping*. This model is efficient as it allows the clients to decide the hosts to be queried based on local data. Furthermore, queries concerning more indexing nodes can potentially be executed by sending network requests in parallel to all of the involved nodes. A potential drawback of this model is that the configuration server could become the bottleneck in a system with many indexing hosts or many connected clients. This issue can be addressed by replicating the configuration data to a cluster of configuration servers. An alternative solution would

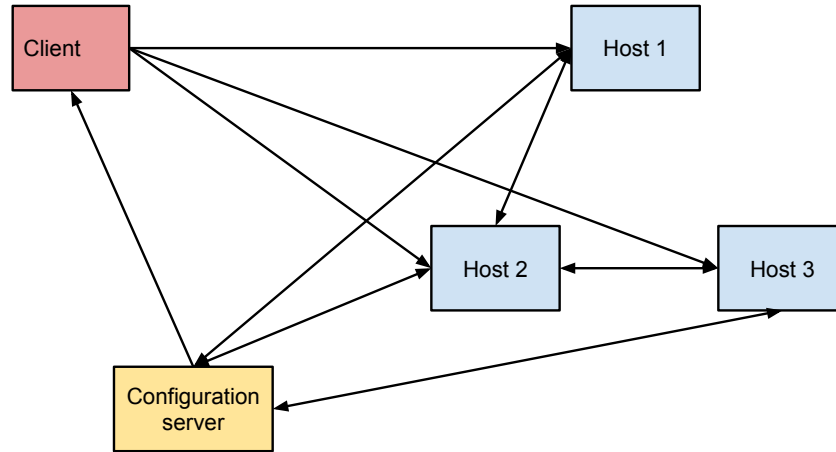


Figure 3.9: The network architecture of the distributed PH-tree

be to store only a partial representation of the mapping in each indexing node and allow queries to be "hop" from one node to the other until they reach the destination machines, as in the case of the previously discussed distributed hash tables. However, the extra latency introduced by repeated "hops" with the cluster makes this approach ill-suited for latency-sensitive applications like a distributed-index.

As the data is split according to Z-order curve intervals, each indexing host only needs to maintain two connections to other hosts in the cluster (or 1 for the hosts holding the edge intervals). This means that the resources used by a host to maintain connections to its neighbours remain constant, even if the number of the nodes in the cluster increases.

3.4 Algorithms

This section will discuss how the *key-mapping* is created from the Z-order curve intervals and how the queries are executed based on the key-mapping.

3.4.1 Key-Mapping Design

The *key-mapping* used by the distributed PH-tree assigns the points based on the Z-order curve intervals. To determine the host responsible for a certain multi-dimensional point, one has to find the Z-order interval that contains the Z-value of this point. In the case of a range query, one has to find all of the Z-order intervals for which the regions in the original space intersect the hyper-rectangle defined by the top left and bottom right points which define the range query. There are two alternative for designing the *key-mapping* based on the Z-order curve intervals:

1. Store the Z-intervals as pairs of Z-values and perform 1-dimensional interval matchings to determine the hosts that have to be contacted for certain queries. For point

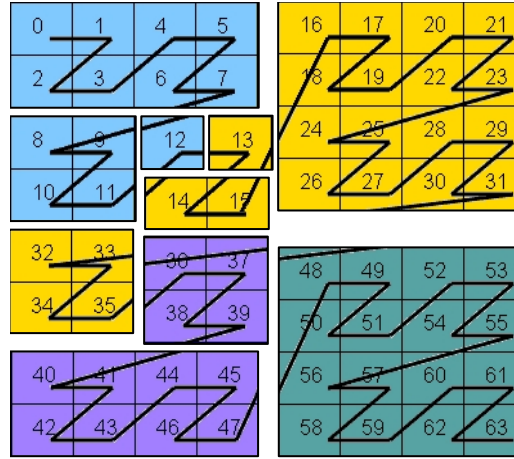


Figure 3.10: The intervals described in 3.8 split into hyper-rectangles.

operations, determining which host needs to be contacted can be done by performing a binary search of the Z-value of the query point on the Z-order intervals belonging to each host. In the case of range queries, one needs to determine the hosts storing the top left and bottom right corners using binary search, and return all of the hosts holding the intervals between the interval of the host storing the top left point and the host storing the top right point. In both cases, the complexity is $O(\log n)$, where n is the number of hosts. The correctness of this approach stems from the fact that the region represented by a d -dimensional hyper-rectangle $[X, Y]$ ⁴ is included in the set of d -dimensional regions corresponding to the Z-interval $[Z_X, Z_Y]$. This can be seen in Figure 3.11a.

2. Store the hyper-rectangles from the original space that correspond to each Z-order interval and perform intersections between the query points/rectangles from the original space and the hyper-rectangle of the interval regions. The complexity of this operation is $O(n * w * d)$, for n indexing hosts, as each region contains at most $O(w * d)$ hyper-rectangles. Figure 3.10 shows how the regions defined by the Z-order intervals can be split into hyper-rectangles.

Even if the first solution is asymptotically faster than the second one for finding the hosts to query, the first approach attempts to query a much larger area in the original space than the second one. Therefore, the first method can potentially generate many more network requests than the first one, especially in a system with a large number of online indexing hosts. Moreover, the intersection operation itself is performed in-memory on the client machine and does not generate a network request. Because of these reasons, the second approach is preferable to the first approach in a system looking to minimize the number of network requests associated to each query, like the distributed PH-tree.

An algorithm for decomposing Z-order intervals into a set of hyper-rectangles is described in [9]. This algorithm decomposes a Z-order interval into at most $O(w * d)$ hyper-rectangles in $O(w * d)$ time, where w is the bit width of values in each dimension and d is the number of dimensions. We propose storing the hyper-rectangles associated with each Z-order interval

⁴Where X is the top left corner of the hyper-rectangle and Y is the bottom right corner of the hyper rectangle

into a multi-dimensional data structure able to perform intersection between the stored hyper-rectangles and a query point or a query hyper-rectangle.

3.4.2 Point Operations

Point operations involve a single query point, like the insertion or removal of a key-value pair, obtaining the value associated with a key-value pair or changing the key associated with a key-value pair. For the point queries, the client only needs to determine which host to send the request to. This is done by determining the intersections between the query point and the hyper-rectangles corresponding to the Z-order intervals. In the case of the *get()*, *put()* and *delete()* operations, the client sends the request to the host responsible for the query point and returns the result of the operation. In the case of the *updateKey()* operations, it is possible that a different host is responsible for the new key of the value associated with the query point. In that case, the *updateKey()* operation is performed in two steps: first the client removes the old key from the host responsible for it, and then inserts the new key into the proper host.

3.4.3 Range Queries

Range queries have the form *range_query(top-left, bottom-right)*. Both of the received arguments are multi-dimensional data points and this operation returns all of the points in the index that fall inside the multi-dimensional rectangle defined by these two points. These operations are resolved in the following manner on the client. First, the client performs an intersection between the hyper-rectangle range and the hyper-rectangles associated with the Z-order intervals contained in the *key-mapping*. It then send requests to all of the hosts whose regions were intersected, receives and combines the results. It is important to note that each host returns a list of point that are already ordered via Z-ordering and because the nature of the Z-order interval split, these lists do not need to be merged, but only concatenated.

A range query can generate up to $O(n)$ network requests, as in some cases, all of the hosts might need to be contacted. However, as previously also mentioned in Section 3.4.1, the *key-mapping* performs the intersection between the query-rectangle and the regions corresponding to the Z-order intervals, instead of performing the intersection between the Z-order interval of the query range and the Z-order intervals of the machines. This guarantees that hosts which *cannot* contain points in the query range will not be queried. It is however possible that some hosts will not contain any points falling in the query hyper-rectangle, even if the regions they are responsible for intersects with the query range.

Figure 3.11a illustrates an example range query. It shows the range defined by the top left point $[(0, 2)]$ and bottom right point $[(1, 5)]$ overlaid against the 3-bit Z-order curve. The Z-order interval corresponding to this range is $[4, 19]$ and the 2 dimensional region associated to it includes the $[(0, 2), (1, 5)]$ rectangle. Figure 3.11b also displays the regions associated to each host according to the split presented in Figure 3.7. In this case, only the hosts responsible for the blue and yellow regions need to be contacted by the client.

As previously mentioned, because the Z-order range associated with the query range corresponds to a larger regions than the query hyper-rectangle, especially in high dimensional spaces,

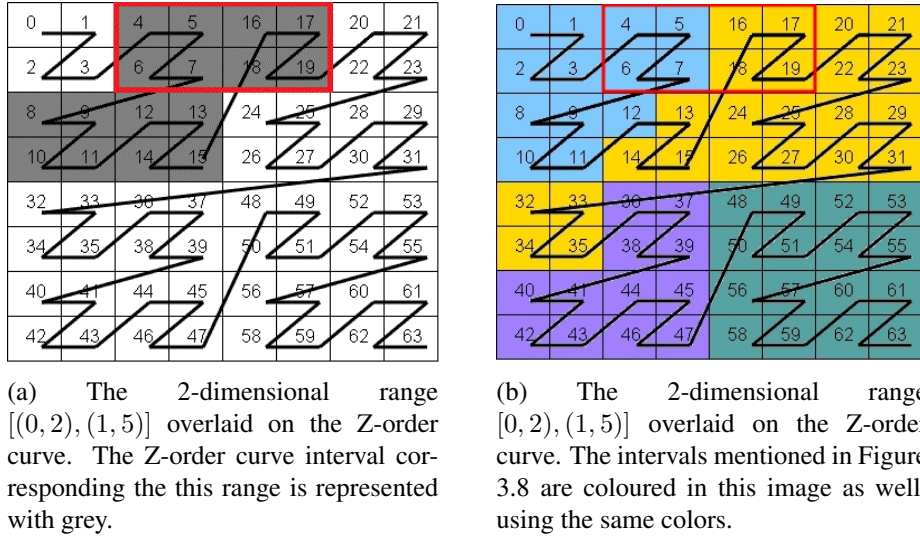


Figure 3.11: Illustration of a range query intersected with the host intervals

3.4.4 Nearest Neighbour Queries

Nearest neighbour queries return the k nearest points, according to the Euclidean distance, to an arbitrary query point, thus having the form $knn_query(q, k)$. The query point does not need to be stored in the index. The execution of a nearest neighbour query has the following steps:

1. The client first determines which hosts is responsible for the query point q and forwards the query to this host.
2. After receiving a set of candidate points from the host, the client determines the furthest neighbour fn of the query point q from this set of candidates, according to the Euclidean distance.
3. The client computes the distance $d = dist(q, fn)$ and performs a range query on the range $[q - d, q + d]$. If the hyper-rectangle range is all stored by the host that was first queried for the candidates, the retrieved candidate points are simply returned as the query result.
4. The final result of the nearest neighbour query consists of the k points from the result of the previously executed range query that are closest to q .

The additional range query is needed because the Z-order splitting, or any other geometrical splitting, does not guarantee that closest k neighbours of a point will be placed in the same regions as the point itself. This range query ensures that points which are closer to q than fn , but stored in a different region, will also be returned.

It is sometime possible that the initial query sent to the host responsible for q will return fewer than k points. This will generally only happen when the cluster contains very few points. In this case, the queried area should be iteratively increased until at least k points are returned. This can be done by computing the distance d and iterative increasing it until the range query $[q - d, q + d]$ returns at least k points.

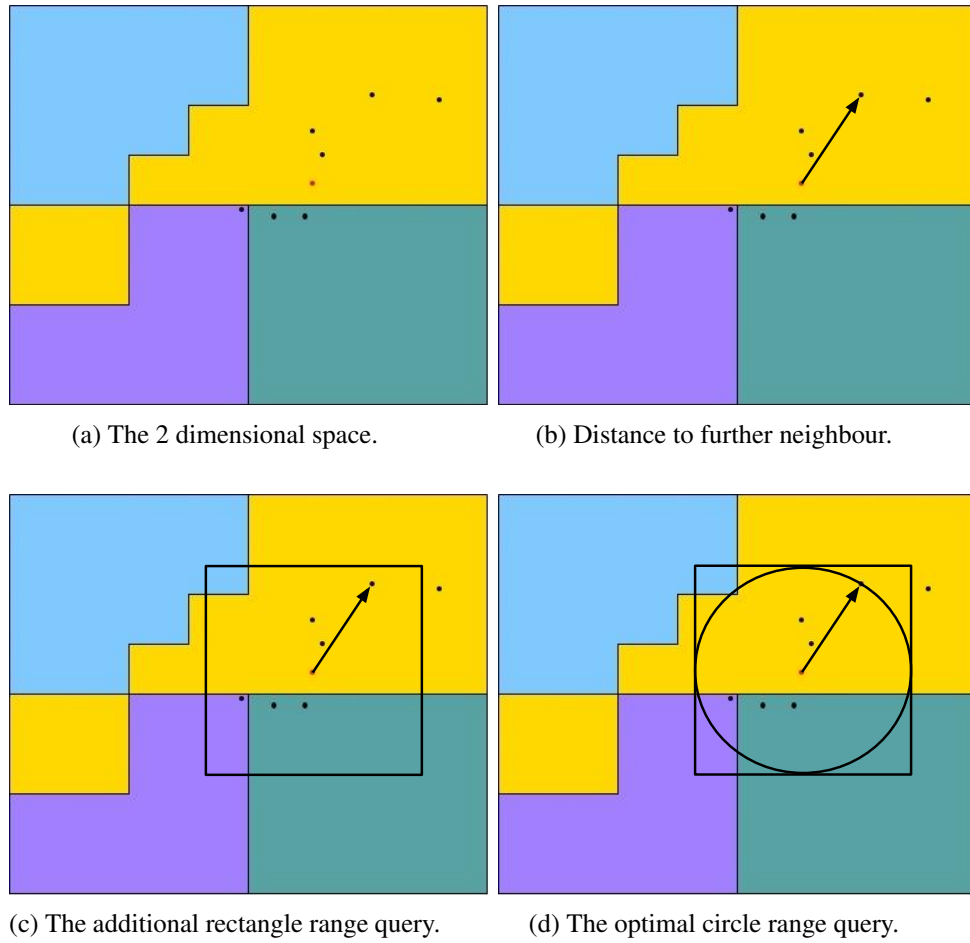


Figure 3.12: Illustration of a 3 nearest neighbour query execution in a 2 dimensional space.

While this approach for finding the k nearest neighbours of a query point is correct, it is not the most efficient one. The main issue is that the additional range query searches for other neighbours in a hyper-cubing range. The optimal strategy would be to search for these neighbours in a hyper-sphere of radius d centered at q . This hyper-sphere is inscribed in the hyper-cube search by the range query, as illustrated in Figure 3.12d. Therefore, the algorithm presented in this work can sometimes send unnecessary requests to hosts responsible for the areas outside of the hyper-sphere and still within the hypercube. This is generally not a problem in lower dimensional spaces, however the difference in volume between the hyper-cube and the hyper-sphere increases gradually with dimensionality of the space.

3.4.5 Entry Rebalancing

As previously described, data sets which are skewed towards certain areas of the multi-dimensional space will cause some machines to store more points than the others. A load balancing algorithm is needed to balance the amount of entries stored across the cluster, such that each host stores roughly the same amount of entries.

Each host is assigned a threshold t corresponding to the maximum number of entries it can store. For simplicity, we assume that all the hosts in the system are identical and have the same t . The load balancing algorithm used by the distributed PH-tree assumes that values stored with each data points have roughly the same size, and such a threshold on the number of entries a hosts can store coincides with the amount of data it can store. The presented balancing algorithm can easily be modified by simply using the amount of occupied memory as a threshold.

Furthermore, the configuration server maintains the number of entries associated to each hosts, which we call the *size* of the host. It is important to note that since each multi-dimensional point has a small size compared to the amount of main memory available on each machine, the balancing algorithm can work properly even if a host only knows an approximation of the number of entries stored by the other hosts. Therefore, to improve scalability, the hosts do not update the size stored in the configuration server after each individual write operation, but only after a certain number of write operations were performed since the last update. An alternative solution would be update the size at a fixed time period. Each host maintains a cached copy of the size information, and is notified by the configuration server when a different host updates this information.

After the number of entries stored by a certain host reaches the balancing threshold t , that host will attempt to move some of the points it stores to a neighbouring host. This process is referred to as a *re-balancing operation* and will change the Z-order intervals of both the host initiating the re-balancing and the host receiving the additional entries. Given a host h responsible for an arbitrary interval, the host responsible for the Z-order interval to the right will be called its *right neighbour*. Similarly, the host responsible for the Z-order interval to the left will be called the *h's left neighbour*.

The re-balancing algorithm proceeds as follows:

1. The *initiator* host checks the size of the hosts responsible for the intervals to the left and to the right of its Z-order interval. When the host responsible for the leftmost intervals initiates a re-balancing operation it only checks the size of its right neighbour, while the host responsible for the rightmost interval only considers re-balancing to its left neighbour.
2. The *initiator* selects the neighbouring host with the fewest entries and sends an initiate balancing message. If the neighbour is available to participate in the operation, the algorithm continues with the following step. If the neighbour is busy, the *initiator* will attempt to re-balance to the other neighbour. If that is not possible, the re-balancing operation fails and will be re-attempted at a later point.
3. The initiator sends a subset of its entries to the *receiver* host. This subset is a contiguous run of entries stored by the initiator.
4. After the entries have been moved, the *initiator* updates the *key-mapping* on the configuration server.
5. The initiator sends a commit message to the *receiver* host, notifying it that the operation was successful. After this message is sent, both hosts mark themselves as available for any other re-balancing operations.

This algorithm allows a host h to be part of only a single re-balancing operation at a moment in time. As soon as a host successfully initiates a re-balancing operation, both the *initiator* and the *receiver* mark themselves as busy and will refuse further re-balancing initiation operations until the current operation is finished. This allows up to $n/2$ balancing operations to run in parallel across the cluster, as each operation involves two neighbouring hosts.

Furthermore, clients currently balancing do not accept write requests. This is done to prevent the *initiator* host to accept updates to the section of Z-order curve that is currently re-balanced to the *receiver* host.

4

Concurrency

4.1 Challenges

The in-memory version of the PH-tree does not support concurrent access and remains consistent only if accessed by a single writer thread or multiple reader threads at a time. With the current popularity of multi-core CPU's, concurrent write access has the potential of providing a significant improvement in the throughput of the PH-tree. The addition of concurrent access presents the following challenges:

- Concurrent access strategies must guarantee that PH-tree remains consistent for any number of threads that access it and for any arbitrary interleaving of the executions of these thread. Furthermore, deadlock should not occur when the data structure is accessed by multiple threads.
- The chosen strategy should optimize the execution time for the most common operation, by minimizing for example, the amount of locks that have to be taken during for those operations. In case of an indexing data structure, we assume that read operations are more common than write operations and should be prioritized by the concurrency strategy from the response time point of view.
- Different concurrent access strategies come with different consistency guarantees. For example, consider the situation in which two processes A and B work concurrently on a shared data structure. Process A starts a read operation at time t and process B a write operation at time $t + \epsilon$. The consistency model determines if process A might see the changes made by process B, if process B finishes the operation before A.
- The PH-tree also supports an iterator over all of the entries or over all of the entries contained within a specified hyper-rectangle. The concurrency strategy should allow iterators to traverse a PH-tree, even if the tree is concurrently modified by other threads.

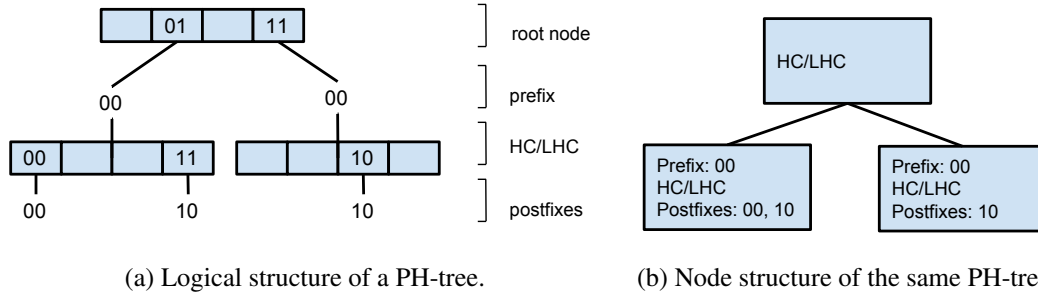


Figure 4.1: A 2 dimensional PH-tree containing three 4-bit entries: (0000, 1000), (0011, 1010), (1011, 1000). The values are omitted.

The concurrency strategy used will influence which entries of the tree are seen by the iterator.

4.2 PH-tree Structure

The PH-tree stores key-value pair entries. The keys are multi-dimensional data points and the values are optional. The entries are stored as bit strings, to exploit any prefix sharing between the bit strings of the keys and reduce the storage requirements.

Figure 4.1a illustrates a 2 dimensional PH-tree containing the entries (0000, 1000), (0011, 1010), (1011, 1000). Figure 4.1b shows the corresponding node structure for the same PH-tree. Each node contains a *prefix*, a *hypercube* (HC) and a set of *postfixes*. The *prefix* is a portion of the bit string common to all bit string portions stored in that node. A *postfix* is an ending portion of a single bit string. The *hypercube* is a data structure that stores references to *postfixes* or subnodes. Currently, the PH-tree supports multiple representation for the *hypercube*, either as an array (HC) or a sparse linear representation LHC. For example, the root node in Figure 4.1a has a *hypercube* storing references to two subnodes. The subnodes each have an associated *prefix* and a *hypercube* storing links to postfixes.

This structural representation of the PH-tree has the following important consequences:

- All information stored in a node cannot be read or written in an atomic manner, as the prefix, the hypercube and the postfixes need to be read or written in separate instructions.
- Adding a new sub-node to a node, as well as removing or replacing it can be done atomically.

The following section presents the concurrent access strategies available in the concurrent version of the PH-tree and discusses their consistency guarantees.

4.3 Concurrency Strategies

4.3.1 Copy-on-Write

Copy-on-Write is an concurrent access strategy that allows multiple threads to access the same shared data structure. When one of these threads needs to modify the shared data structure, it first makes a copy of it, performs the modifications and then updates it with its modified copy. Each write operation creates a new version of the shared data structure, making it immutable. This strategy allows one writer thread and an arbitrary number of reader threads to work on the same data in the same time. If multiple writer threads were to access the data structure, each would end up with its updated copy and the last to replace the shared data structure with the copy would overwrite the changes of the other writer.

In case of the trees, writer threads do not have to copy the full tree before applying their modifications. It is sufficient to create a copy of the node that is updated and copies of all of its ancestor nodes, up to and including the root of the tree. The number of nodes copied for each update is equal to the maximum height of the tree. In case of the PH-tree, this number is equal to the number of bits needed to represent the values on each dimensions, independent of the number of nodes or entries currently stored in the tree.

Figure 4.2 illustrates the execution of write operation. For simplicity, a general binary tree is pictured here, as the operation is executed in a similar manner on a PH-tree. In this example, node G needs to be updated. To achieve this, all of the ancestors of G are copied and the copy of G, G' is updated. Finally, the root pointer of the tree is changed from A to A' atomically. It is important that the root swap is performed in an atomic manner, such that any readers concurrently accessing the tree will either access the old or the new version of the tree.

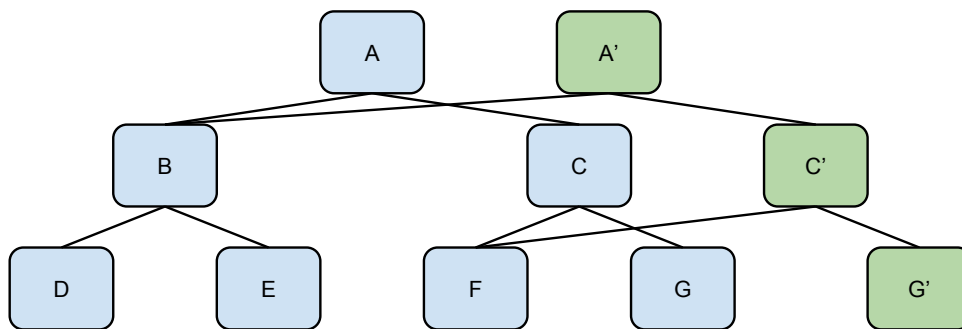


Figure 4.2: Illustration of Copy-on-Write on a binary tree where the node G has to be written

There are several possibilities to handle multiple writer threads using the Copy-on-Write strategy. One approach would be allow multiple threads to access the tree in the same time. This can be done by implementing the atomic root swap using a Test-And-Set atomic operation, essentially turning each write operation into a transaction. In case the root was changed since the thread has started the operation, the atomic Test-And-Set fails and the thread will need to restart the operation. The main drawback of this approach is that the commit fails even if the changes are made to different parts of the tree.

Iterators created on a Copy-on-Write PH-tree essentially traverse a snapshot of the tree taken

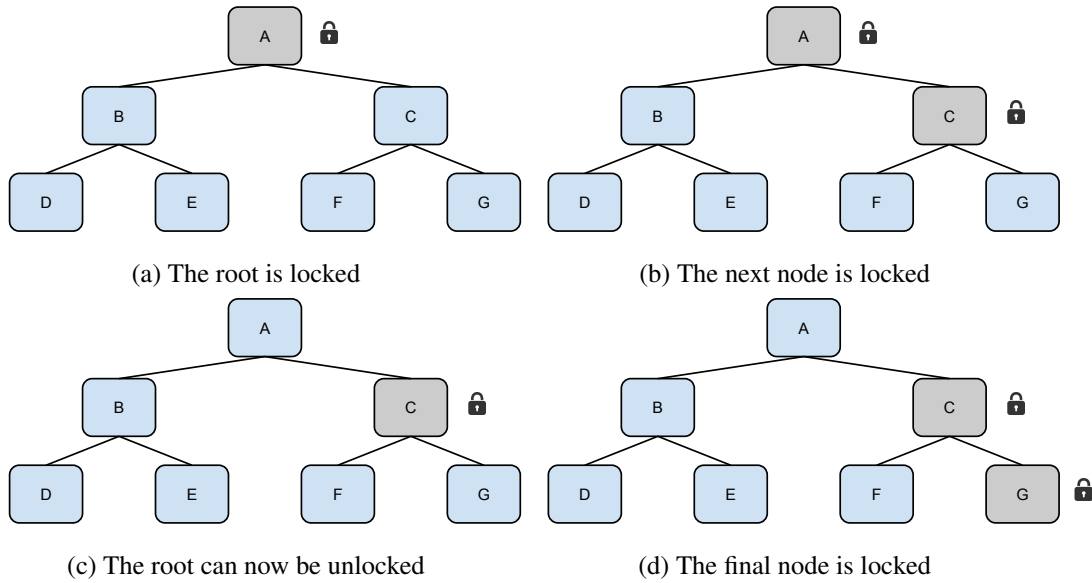


Figure 4.3: Illustration of hand-over-hand locking descent to node G

on the iterator creation time t . If an iterator is created at time t , it will view all of the entries that were contained in the tree at time t , even if these entries were removed from the tree before the iterator finishes. Moreover, iterators will never see any entries that may have been added to the PH-tree after time t .

4.3.2 Locking

A very simple way to handle concurrent writer threads is to use a global lock that each thread must take before the execution of an operation. If the lock is a mutex, the tree can only be accessed by a single thread. If the thread is a read-write lock, the tree can be accessed by either a single writer or several reader threads. However, this *coarse-grained* locking strategy does not allow one writer thread and one reader to perform the accesses in parallel. Furthermore, if the lock favors readers, writers might be permanently blocked from accessing the tree for workloads with a sustained number of read accesses. A similar situation arises for a read-write lock that favors writers, which might block reader threads indefinitely. Due to these disadvantages, we turned our attention to other locking approaches, namely *fine-grained* locking strategies.

Fine grained locking strategies perform locking at the level of tree nodes and can allow multiple writer threads to modify the tree in the same time. Two popular strategies are described in the following sections. For simplicity, in the following sections we first consider how the strategies support concurrent writers and then discuss how to incorporate concurrent read and writer access.

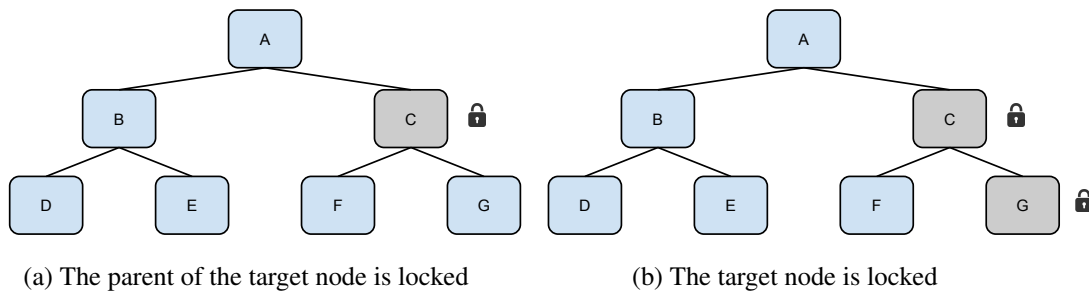


Figure 4.4: Illustration of optimistic locking descent to node G

Hand-over-Hand Locking

Hand-over-hand locking, also called *lock-coupling*, is a fine-grained locking strategy according to which a lock for a node can only be acquired if the parent lock is also held. When traversing the tree according to this strategy, at most two locks will be held by a thread at a time. An additional lock needs to be acquired before acquiring the lock of the root of the tree, to prevent conflicts on the creation of the first node of the tree.

Figure 4.3 shows the sequence of events happening when a thread needs to update a node. We consider that each node has an associated exclusive write lock. The node to be updated in this figure is node G. The thread first acquires the root, descends to node C and acquires that lock as well. The root can be unlocked before descending to node G and acquiring that lock. After the changes are done, the lock of node G is unlocked.

This strategy allows multiple threads to modify the different subtrees of the tree. The main drawback of this approach is that if the non-leaf nodes are modified, subtrees of those nodes cannot be modified of a different thread until the current thread finishes the update operation. Secondly, there will be a high contention for the locks on the higher levels of the tree, which could negatively affect performance for write-heavy workloads.

Optimistic Locking

Optimistic locking is an alternative fine grained locking strategy which seeks to minimize the number of locks which are taken during a write operation. The basic idea is to traverse the tree without acquiring any write locks until the node to be updated is reached. Once at the target node, the writer thread acquires the lock of the parent and the lock of the current node. After this happens, the write operation can continue only if both of these two nodes have not yet been removed from the tree and the current node is still a child node of the parent node. If any of these checks fail, the write operation fails and the current thread begins another attempt by starting a new traversal of the tree. Figure 4.4 shows the sequence of events for the update of a node in the tree using optimistic locking. As in the previous section, the node to be modified is G.

To speed up checking if a node is still reachable from the root node, each node has an additional removed flag, which indicates if certain node has been removed from the tree. When a thread removes a node, it sets this flag to true before releasing the lock of the node.

One advantage of this method is that in a large tree, re-tries are generally rare as threads do not generally attempt to update the same node. Moreover, a thread holding the lock of a certain node does not block other threads from performing update operations on the node's subtrees, unlike the hand-over-hand strategy.

Handling Reader Threads

It is important to mention the fact that in the case of the PH-tree, all update operation to a single node are not executed in an atomic fashion, which can cause reader threads to perform reads from nodes which are in an inconsistent state. This inconsistency can be avoided by protecting the nodes with a read-write lock and requiring the read operations to acquire the lock of a node in reading mode before performing a read. Additionally, writer threads attempt to acquire the lock of a node in write mode during the read operation. The conversion of exclusive locks to read-write locks allows both described fine-grained locking strategies to support both concurrent writers and readers.

In the case of the fine-grained locking concurrency strategies, iterators running the PH-tree operate according to a relaxed consistency semantics. Therefore, an iterator created at time t , which finished traversing the tree at time e will see all of entries that were part of the tree a time t and are still part of the tree at time e . Moreover, entries inserted or removed in the time interval $[t, e]$ may also be seen by the iterator. It is however guaranteed that all of the entries seen by the iterator are visited according to the Z-ordering.

4.3.3 Hybrid Approach

As previously mentioned, both hand-over-hand locking and optimistic locking require readers to acquire locks on the nodes they perform the reads on. This increases the latency of the read operation due to the overhead of acquiring locks. Moreover, when the tree is accessed by many reader threads, lock contention can lead to a very fast performance decrease.

We propose the modification of the fine-grained locking strategies to incorporate copy-on-write logic for the write operations. Specifically, before performing a write operation on certain node, the writer thread first creates a copy of the node, performs the update on the copy and then atomically replaces the parent node's reference to the node with the copy. As the writer thread holds the lock on the parent node, it is guaranteed that no other writer threads will overwrite this change. Furthermore, reader threads do not need to acquire locks on the nodes when reading the node information, as the current node they are accessing will not be changed by the writers.

This modification allows the read operations and iterators on the PH-tree to remain unchanged, while only the write operations need to be updated. Therefore, the response time of the read operations will not suffer from the overhead of lock contention.

4.4 Chosen Strategies

The PH-tree supports both a full Copy-on-Write concurrency strategy and the hybrid versions of the described fine grained locking strategies. The optimistic locking strategy is generally faster than the hand-over-hand locking strategy, the actual results being detailed in Section 6.2 of Chapter 6. However, optimistic locking can lead to starvation when many threads need to update the same node or nodes which are close together. Therefore, hand-over-hand locking can be chosen over optimistic locking in situations where it is important that writer threads are not starved. The choice between the Copy-on-Write strategy and the fine grained strategies comes down to the choice between the consistency guarantees provided by them. For example, snapshots of the tree can only be obtained if the tree is using the Copy-on-Write strategy.

5

Implementation

This chapter describes the implementation of the approaches described in the previous chapter in order to obtain a distributed and concurrent PH-tree. We discuss the technologies used to implement the distribution architecture and explain how the multiple concurrency strategies have been integrated in the code of the PH-tree.

5.1 System Architecture

The architecture of the distributed PH-tree is shown in Figure 5.1. The implementation consists of two artifacts: the *client library* and the *server library*.

The *client library* is deployed on the clients and contains the code necessary to connect to the configuration server, retrieve the cluster information and perform operations on the cluster. The *client library* also contains a dependency to the *Kryo* library, a Java serialization library that is used to encode the values stored on the PH-tree before the insert operation and decode the values retrieved from the server. Clients cannot modify the cluster information and can only subscribe to be notified by the configuration server when indexing hosts join or leave the cluster and when the Z-order curve intervals changes as the result of a rebalancing operation.

The *server library* contains the index server logic and needs to be deployed on the indexing hosts. The *server library* implements a multi-threaded Java server application running on the indexing hosts. The *Netty IO*¹ library is used to handle incoming client requests using a thread pool. When a server is started, it first connects to the configuration server to retrieve the cluster information and then is assigned an interval of the Z-order curve for which it is responsible. If the cluster already contains some data the new node can mark itself as a free node and will be contacted by the next node in the cluster that will begin a balancing operation. Indexing hosts update the intervals in the cluster configuration after a balancing

¹<http://netty.io/>

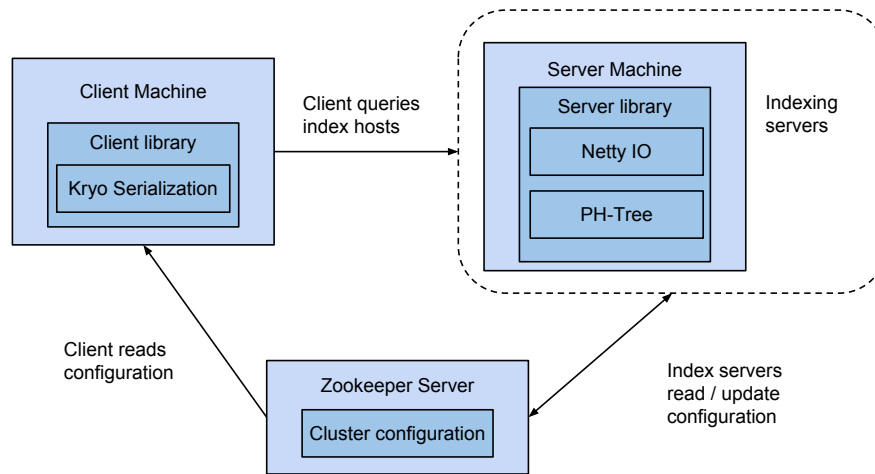


Figure 5.1: The implementation architecture of the distributed PH-tree

operation and periodically update the number of entries stored.

The *configuration server* was implemented using *ZooKeeper*², an open source server whose goal is to provide reliable distributed coordination. ZooKeeper provides a centralized service for maintaining configuration information and performing distributed synchronization. The *cluster configuration* is stored in a ZooKeeper specific format. ZooKeeper also allows the possibility of replicating the stored configuration by running a *quorum*³ of machines to improve availability.

A minimal configuration of the distributed PH-tree requires one ZooKeeper server, one indexing host running the *server library* and one client machine performing operation using the *cluster library*.

5.2 PH-tree Java API

The Java API of the PH-tree is presented in Figure 5.2. The interface `PhTreeV` is the main interface and corresponds to a PH-tree that stores entries with keys of type `long[]` and arbitrary objects as values. A set API is provided through the classes `PhTree` and `PhTreeD`. `PhTree` uses a `PhTreeV` instance with empty values to implement a set for multi-dimensional points with values of type `long`. `PhTreeD` is targeted towards application which require floating point precision for the values of each dimension. This type of tree is backed by the `PhTreeV` instance and performs pre and post processing of the multiple-dimensional points to map them from `double[]` to `long[]`.

The initial PH-tree implementation also provides support for range trees, which can store hyper-rectangles defined by two multi-dimensional points: the bottom right and top left

²<http://zookeeper.apache.org/>

³The name used by ZooKeeper for a replicated cluster of ZooKeeper servers.

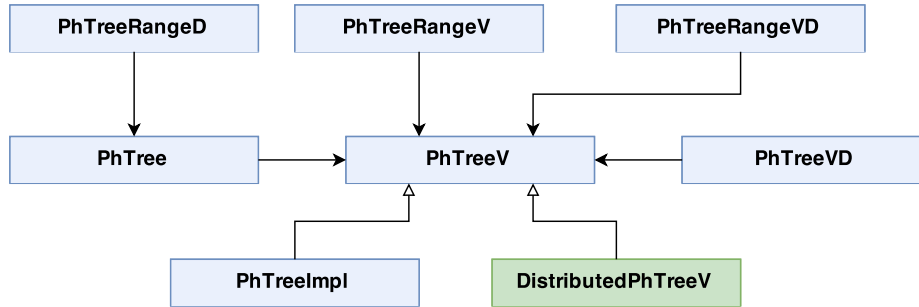


Figure 5.2: The Java API of the PH-tree. The classes coloured with blue represent the existing code, the additions in green represent the distributed implementation.

corners of the hyper-rectangle. These types of trees are implemented using special types of pre and post processors and are internally backed by point PH-trees. `PhTreeRangeV` stores hyper-rectangles defined by multi-dimensional points represented as `long[]` and allows arbitrary values. `PhTreeRangeVD` is similar but allows points to be represented as `double[]`. Both of these two implementations use an instance of `PhTreeV` to store the processed hyper-rectangles, as shown in Figure 5.2. Finally, `PhTreeRangeD` provides support for hyper-rectangle set index is backed by a `PhTree` object.

As presented in Figure 5.2, all of the types of trees discussed implement the interface `PhTreeV`. The distributed PH-tree provides another implementation for this interface. Method calls to any of the trees in the presented hierarchy are transformed into commands send by the `DistributedPhTreeV` class to the indexing hosts over the network. After the requests are serviced, `DistributedPhTreeV` combines the responses and provides the appropriate returning value.

One big advantage of this approach is that client is not aware if it is using an in-memory tree or a distributed one, except at the creation time. Therefore, applications previously using the in-memory version of the PH-tree can easily switch to the distributed implementation, with minimal code changes.

The choice of providing a separate implementation only of the `PhTreeV` class has another important consequence: the indexing hosts only need to maintain an in-memory instance of a `PhTreeV` implementation, independent of the type of tree used by the client. This simplifies the implementation and also has the advantage of off-loading the pre and post processing operations from the servers to the clients, allowing faster request processing and better throughput on the servers.

5.3 Distribution

5.3.1 Client Design

A high level diagram of the client library is presented in Figure 5.3. The `DistributedPhTreeV` class implements the `PhTreeV` interface and provides access the distributed PH-tree. Each method call on the `PhTreeV` interface is delegated to a `PhTreeIndexProxy` object. This object converts each method call into a `Request` object, similar to the *Command design pattern*⁴. The proxy class also decides which hosts need to be contacted to resolve the operation, based on a locally cached `Key-Mapping` object. This object contains the most recent cluster configuration available to the host. The `Request` object is then passed on to a `RequestDispatcher`, where it is encoded into an array of bytes and sent to the corresponding hosts via the Java TCP Socket API. The `Transport` class keeps the most recently used connections alive, to avoid having to perform the TCP handshake for each operation when repeatedly interacting with the same host.

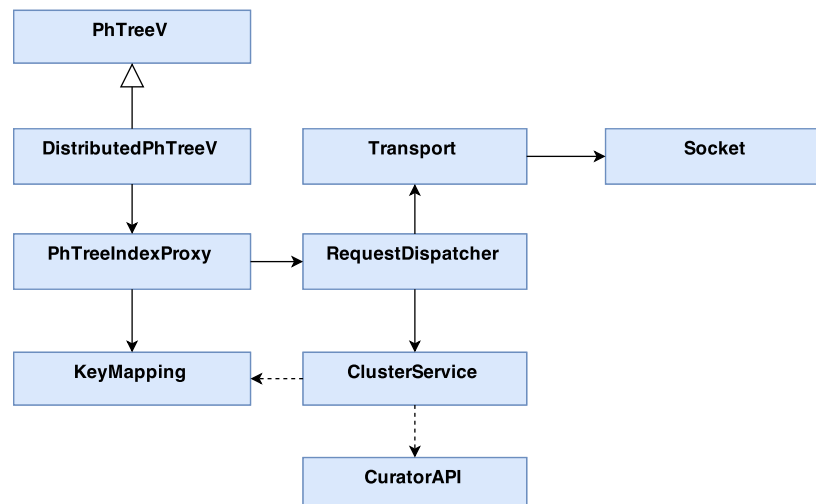


Figure 5.3: The class diagram of the client library

The calls to the `DistributedPhTreeV` methods are blocking, as the `Transport` object blocks the execution until it has received a response from the remote indexing host. The response is received as an array of bytes and is decoded into a `Response` object and returned to the `RequestDispatcher` object. If more hosts participated in the query, their corresponding `Response` objects are combined and the result of the operation is returned to the caller of the `DistributedPhTreeV` method.

⁴http://en.wikipedia.org/wiki/Command_pattern

5.3.2 Server Design

The *server library* is implemented as a multi-threaded, asynchronous, event-driven server, according to the *reactor pattern*⁵. This model uses two types of threads: *Dispatcher* threads and *Handler* threads. A *Dispatcher* thread accepts incoming connections, allocates the necessary resources for the connection and then passes the connection resources to a *Handler* thread. The *Handler* thread accepts requests from the client using the connection resource provided by the *Dispatcher* thread. For the TCP protocol, this resource is generally a reference to a *Socket* object. The *Handler* thread manages requests from more than one client by leveraging the event-based IO support provided by the operating system. The *Netty IO* Java library is used for the implementation, as it provides an easy manner to set up a pool of *Handler* threads. The chosen server architecture generally performs better than architectures which create a new thread for every request, as the *Handler* threads are pooled and are initialized before any client requests are serviced.

The high level diagram of the server library is presented in Figure 5.4. The main class of the *server library* is the `IndexMiddleware` class, which sets up the pool of handler threads and starts the dispatcher thread. The `MiddlewareChannelHandler` class is the first handler that is executed by one of the *Handler* threads upon receiving a request on a client connection. This class reads all of the bytes corresponding to a single request from a client and places them into a byte array object. This object is then passed to an instance of the `IOHandler` class, which converts the byte array into a `Request` object. The `IOHandler` then passes the `Request` object to another handler, based on the type of the request. Client requests are passed along to a `RequestHandler` object, while balancing requests, coming from other servers are passed to the `BalancingRequestHandler` object.

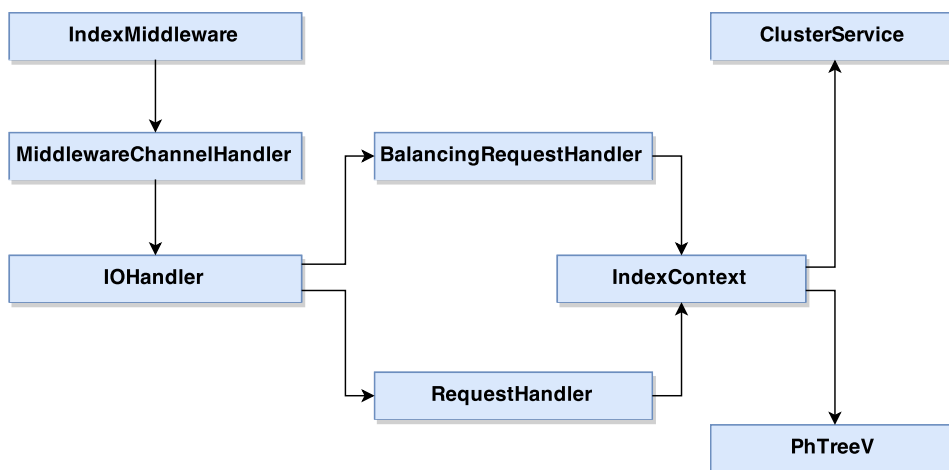


Figure 5.4: The class diagram of the server library

Information that is commonly accessed by all of the *Handler* threads is stored in the

⁵http://en.wikipedia.org/wiki/Reactor_pattern

`IndexContext` object. This object contains a reference to the `ClusterService`, used to read and write cluster configuration data, and to the in-memory PH-tree instance on that server, which is used in processing the client requests.

5.3.3 Configuration Server and Key-Mapping Implementation

The *Key-Mapping* was implemented using the ZooKeeper, an open-source coordination service. ZooKeeper provides access to a shared hierarchical namespace similar to a file system. This namespace is comprised of data registers, called *znodes*, which are identified by a path and store data in the form of a byte array. The data stored within a *znode* is guaranteed to always be read and written atomically. Additionally, each *znode* can have any number of children *znodes*. ZooKeeper provides primitives to create and delete *znodes*, check if a node exists, get and set the data stored by a node and also retrieves the list of children of a node. It is also possible to set *watches* on certain *znodes*, and receive notifications from ZooKeeper when the data of the node was changed.

Both the client and the servers access the configuration data using the `ClusterService` object, which maintains a cached version of the configuration on a `KeyMapping` object. Handlers in the *server library* module can also update the configuration through the `ClusterService` API. The `ClusterService` receives notifications from the ZooKeeper server when any data was changed and updates the `KeyMapping` object accordingly. The implementation of the `ClusterService` uses the *Curator*⁶ framework, a set of libraries that provide a simpler and more reliable API for accessing ZooKeeper.

The ZooKeeper configuration for the distributed PH-tree is illustrated in Figure 5.5. This configuration was implemented using three main *znodes*:

- The */mapping znode* contains the set of all Z-order intervals, together with the hosts responsible for each of these intervals. All of this information is serialized into a byte array and is always read and written atomically by the indexing hosts. This is done to ensure that all Z-order intervals are read by the clients in an atomic manner, and always offer them a consistent view of the system.
- The */sizes znode* stores an approximation of the number of entries stored in each indexing host, which is needed to make decisions during the *re-balancing* operation. This *znode* contains one child *znode* for each host, the child node storing the approximate size of each host. This allows all of the hosts to update their size without clashing with updates from other hosts. Even though this means that the sizes of all of the hosts cannot be read atomically, the consistency loss is not significant, as the balancing operations do not require the precise number of entries stored on the neighbours.
- Finally, the */free znode* contains a list of hosts which have registered themselves as free on start-up and will be contacted by the balancing initiators during the next balancing operations.

As the clients maintain a cached copy of the latest information in ZooKeeper, it is possible that, due to network delays, the mapping on the client may be different from the actual in-

⁶<http://curator.apache.org/>

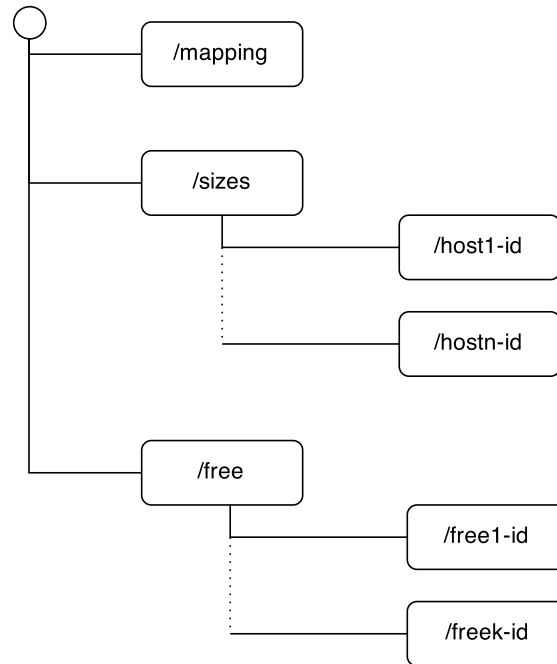


Figure 5.5: A diagram of the Key-Mapping stored in ZooKeeper

tervals on the server. Consider the following situation: a client initiates a query and after analyzing the local *Key-Mapping*, it decides it has to contact host h to retrieve the value for a certain key k . However, before the request reaches host h , this host could have already participated in a re-balancing operation and moved a part of its entries to a neighbouring host. If the entry involved in the query was moved to a different host, host h no longer has it and would respond to the client that no such host exists in the cluster. This is wrong, as the entry was simply moved to a neighbouring host. To prevent this type of situations, the *Key-Mapping* has a version number that is incremented after each re-balancing operation. Each indexing host stores the version of the mapping that was generated during the last re-balancing operation in which that host was involved. Furthermore, the clients include the version number of their local mapping in the requests sent to the hosts. Upon receiving a request, the server checks the client mapping version against the version of its last re-balancing operation. If the client version number is higher or equal, the request can be serviced without any consistency loss, as any changes that may have increase the clients version happened to other intervals in the mapping. Instead, if the client's version number is smaller than the host's version, the client is noticed that it's mapping is outdated and should retry the query. The client keeps attempting to send queries to indexing servers until it has received a response for an up to date mapping.

In a cluster of n hosts, there can be at most $n/2$ balancing operations running in parallel on the cluster, as each host can be involved in a single balancing operation at a time. Therefore, it is important that a host's updates to the mapping do not overwrite the changes made by other

hosts. The mapping update operation is performed in a *compare-and-swap* manner, using the multi-update operation support provided by ZooKeeper. Using this type of operation, it is possible to perform an update to a znode only if its version (or alternatively, another znode's version) has not changed since the previous value was read. Following a re-balancing operation, the initiator hosts attempts to write the new version of the mapping to ZooKeeper. If the mapping was changed by another host, the initiator reads the new mapping, applies its changes on the interval bounds and retries the operation until it is successful. While this approach can lead to starvation if a host is consistently blocked by other hosts to update the mapping, such a large number of concurrent re-balancing operations is not likely to happen. An alternative solution to this problem would have been to implement a distributed lock on the mapping through ZooKeeper. In that situation, the initiators would always need to perform three requests to ZooKeeper: one request to acquire the lock, the second to perform the write and the third to release the lock. Our solution optimizes the duration of the re-balancing operation for the best case, when the write will be successful, as that situation is much more common in practice.

5.3.4 Communication Flow

The distributed PH-tree uses a custom binary protocol for both the communication between a client and an indexing host and between two indexing neighbouring indexing hosts. The main communication flow is described in Figure 5.6. The `RequestDispatcher` on the client creates a `Request` object for each method call to the PH-tree API. This request object contains the request type and the arguments received by the method and it is converted into an array of bytes by the `RequestEncoder` class, according to the communication protocol. The array of bytes is sent through the network to all of the servers that participate in the operation. On the server, the request is assigned to one of request processing threads which converts the array of bytes back to a `Request` object using a `RequestDecoder`. The request is then serviced by a `RequestHandler` object, which returns a `Response` object containing the result of the request. The `Response` object is converted into an array of bytes using a `ResponseEncoder` and sent back to the client over the wire. Back on the client, the response is decoded from a byte array. In case multiple servers were contacted, all of their responses are decoded and then combined to produce the return value of the method call.

The main reasons for using a binary protocol, as opposed to a more high level protocol based on JSON⁷ or XML⁸ are the following:

- A binary protocol produces a much smaller output than a text based protocol. In the case of the custom binary protocol, all parameters are converted into arrays of bytes of the smallest possible size and combined into a large byte array. JSON and XML formats require additional metadata to store the name of the fields and to make the representation human readable.
- JSON or XML protocols follow a two-step process: first the request is converted into a JSON or XML document, then the document is converted into an array of bytes to be sent over the network. The binary protocol skips the middle step and transforms

⁷<http://json.org/>

⁸<http://en.wikipedia.org/wiki/XML>

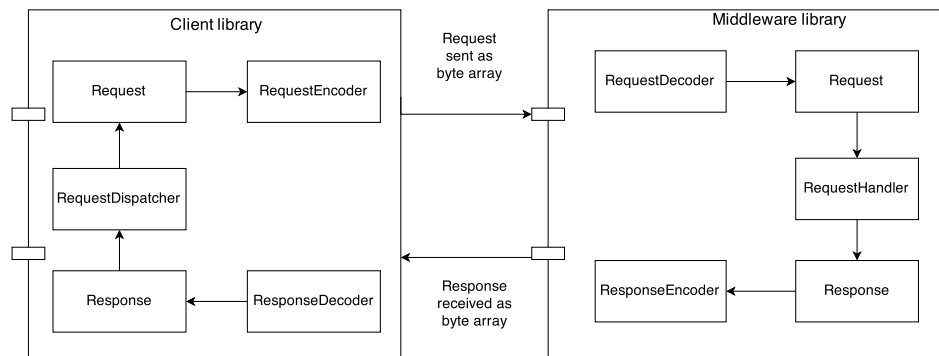


Figure 5.6: A diagram of the communication flow between the client and the server

requests directly into arrays of bytes. This makes the binary encoding faster than the text-based alternatives.

- The messages exchanged for this communication protocols are generally quite simple and thus the implementation effort for efficient custom serialization is not very high.

5.3.5 Iterators

Beside the number of queries that return a list of entries that match certain parameters, the PH-tree API also provides support for iterators. These iterators implement the Java iterator API and allow clients to iterate over all of the entries of the tree or over all of the entries that fall within a certain query hyper-rectangle.

Once an iterator is created using the client API, the client sends a request to the host holding the first key in the iterator. This host creates an in-memory iterator object that is maintained until the client has finished iterating over the entries. The client retrieves batches of entries from the iterator on the host, to avoid having to send a request each time the iterator is moved a single position. Once all entries in the iterator stored on a single host have been exhausted, the client contacts the host holding the following portion of the iterator data and repeats the procedure.

One advantage of this approach is that the client does not have to perform a network request for each entry of the iterator. Secondly, for iterators over large areas, spanning multiple hosts, it might also be unfeasible to request all of the entries in the iterator from all of the hosts and just iterate over them in-memory, on the client machine. A potential drawback of this solution is that each iterator running on a client has an associated iterator object running on a server machine. Thus, creating very large numbers of iterators could lead to large memory consumption on the servers and could potentially affect the performance in a negative manner.

5.4 Concurrency

The addition of concurrency support to the PH-tree involves the modification of the write operations to handle accesses from multiple threads. The operations which perform queries, create or execute iterators across the PH-tree were not modified, to preserve the response time of the read operations. Therefore, only the *put()*, *delete()* and *updateKey()* methods for implementations of the `PhTreeV` interface require modifications to support concurrency.

The implementation was performed by isolating the write methods into a separate interface called `PhOperations`. Thus, the PH-tree implementation delegates the execution of the write methods to an internal `PhOperations` object. This pattern allows us to easily provide multiple implementation for concurrent write access strategies by placing each strategy into a separate implementation class.

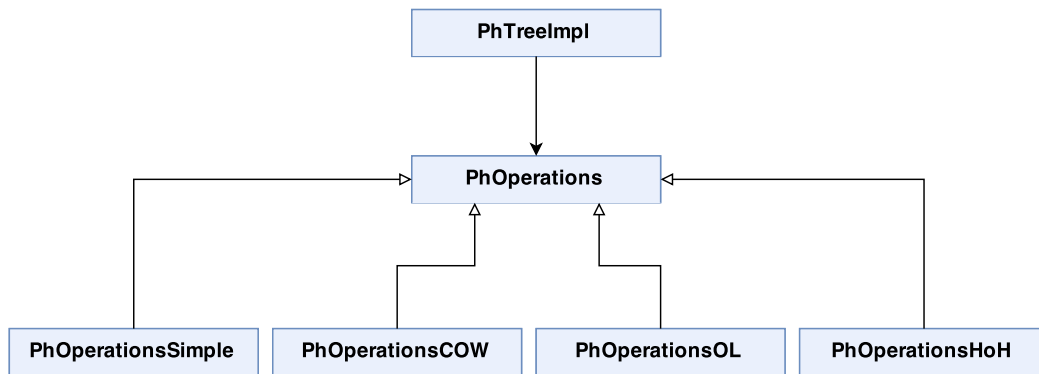


Figure 5.7: The Java API of the concurrency classes

Figure 5.7 shows the `PhOperations` interface and its implementations. First of all, `PhOperationsSimple` provides an no-concurrency implementation of the write operations, allowing applications that do not require concurrent access to by-pass the response time performance penalty added by the synchronized implementations. `PhOperationsCOW` provides a full Copy-on-Write implementation, which essentially creates a "new" version⁹ of the tree at each write operation and supports iteration on snapshots of the tree. This implementation supports a single writer thread accessing the tree concurrently with an arbitrarily large number of readers. `PhOperationsHoH` and `PhOperationsOL` are fine grained lock-based concurrency strategies, providing support for multiple reader and writer threads at the cost of strict consistency. In both of these implementations, the nodes modified are first copied and the modified copies are atomically replaced in the parent nodes. This removes the need for taking locks in the read methods. `PhOperationsHoH` implements the tree descent using hand-over-hand locking, while `PhOperationsOL` performs optimistic locking for this purpose. As previously discussed in Chapter 4, `PhOperationsOL` may need multiple descents on the tree until the locking attempt was successful and should generally perform better than `PhOperationsHoH` for workloads in which the writers are not accessing

⁹Only a single path from the root to the update node is copied at each write operations, the rest of the tree is shared with the previous version.

the same area of the tree at the same time.

Another big advantage is that the `PhOperations` reference can be swapped at runtime, which changes the concurrency strategy used by the PH-tree. This allows a tree to use optimistic locking while it is "loaded" with entries from multiple threads, and then switch to copy-on-write to allow iterators or range queries to execute on snapshots of the tree.

The implementation uses the locks provided by the Java API for the write locks. An instance of `ReentrantLock` is added to each node object, and the locks are taken by the writer threads according to the concurrency strategy used. Furthermore, the `ReentrantLock` object is used to synchronize the writer threads in the copy-in-write implementation of the tree. An alternative to using lock objects on each node is to perform synchronization using a separate concurrent data structure, in an attempt to reduce the amount of memory associated with each node. For example, to "lock" a node, a writer would insert the node reference into a concurrent hash set. Similarly, "unlocking" a node is done by removing the node from the set. Therefore, a node is "locked" if it is stored in the set and "unlocked" otherwise. This approach would reduce the space overhead added by the lock object, however its performance depends greatly on the performance of the concurrent hash set. Furthermore, if a thread wants to acquire a lock that is already taken it has to repeatedly poll the data structure, entering a busy-wait state. This means that the thread will keep using the CPU and generate a lot of traffic on the memory bus until the lock is released. We have finally chosen the solution that assigns a lock object to each node, as it avoids busy-waiting when writer threads have to wait on a lock.

6

Evaluation

This chapter contains the performance evaluation of the distributed PH-tree. Section 6.1 presents the performance of the PH-tree from the distributed point of view, for indexing clusters of different sizes. The concurrent access strategies presented in Chapter 4 are evaluated in Section 6.2.

6.1 Distribution

6.1.1 Experimental Setup

The experiments were executed in order to measure the insertion response time and throughput for distributed clusters containing different numbers of server hosts. The machines used were provided by the Google Computer Engine¹ cloud platform. Each server process ran on a virtual instance containing 2 virtual CPU's and 7.5 GB of RAM memory². Google Compute Engine does not provide any information regarding the physical hardware used by the virtual machines, however the virtual machine emulated the Intel Core i7-3800 processor family, with each of the 2 CPU cores having a 2.6GHz frequency. The same type of machines have been used to run the ZooKeeper server, the indexing server processes and generate client requests.

Both the throughput and the response time were measured on the client machines. In the measurements, we consider that an operation performed on the distributed index starts in the moment when the method corresponding to this operation is called in the client code and ends when the same methods returns. Therefore, the duration of an operation includes the request encoding on the client side, the time needed by the request to travel over the network to the

¹<https://cloud.google.com/compute/>

²The code used by Google Compute engine for this type of virtual machine is `n1-standard-2`

server machine, the processing time on the server machine, the time needed by the server response to travel over the network back to the client and the response decoding.

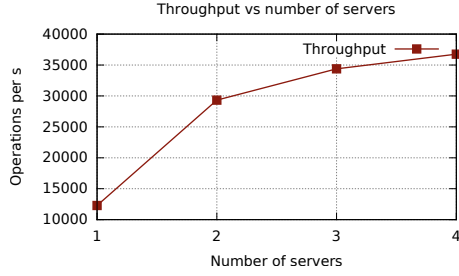
The experiments are run the help of a shell script, which starts the necessary processes on the appropriate machines, tracks the execution of the experiment, collects the client logs and stops the processes at the end of the experiment. A typical experiment runs in the following manner. First, the shell script connects to virtual machine assigned for the ZooKeeper server and starts the ZooKeeper server process. Next, the script connects to all server machines part of the cluster configuration for the current experiment and starts an indexing server process on each server machine. Finally, the shell script connects to a set of client machines and starts a Java program that will generate the appropriate workload for the current experiment. Each client process logs the duration of each request together with the time stamp at which the request was completed. After all clients have finished executing, the log files are gathered from the client machines and merged into a master log.

Because not all the clients start executing at precisely the same time, each experiment will be formed of three phases: *the build-up* phase, when some clients perform requests and others have not started yet, the *stable request load* phase, when all clients are performing requests and the *cool down* phase, when some clients still perform requests and others have finished. These phases are identified by analyzing the trace of the throughput of the system for the whole duration of the experiment. The time intervals corresponding to the *build up* and *cool down* phases are removed from the master log and the stats presented are computed only on the data from the *stable request load* phase.

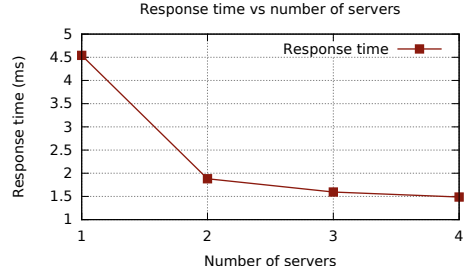
6.1.2 Experimental Results

All distributed experiments used a 2 dimensional PH-tree with 64 bit long values in each dimension. The client requests were generated from a number of 7 client virtual machines and the number of servers is increased from 1 to 4 in increments of 1. We evaluated the performance for four different workloads: a insertion workload, a point query workload, a range query workload and nearest neighbour workload. The values in each dimension have been sampled from the Gaussian distribution using the `nextGaussian()` method of the `Random` Java class. Because this method returns a double precision value sampled from the Gaussian distribution with mean 0.0 and standard deviation 1.0, the sampled value has been scaled to a higher space by multiplying it with $2^{64} - 1$. Because each type of experiment uses the exact same workload, generated by the same number of clients, for clusters of different sizes, it is expected that environments with more servers will lead to a higher overall throughput and a lower response time.

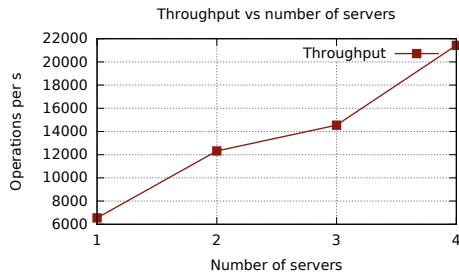
During the insertion workload, each client machine created requests using 8 separate threads, each thread inserting 50000 2-dimensional points. At the end of each run of this experiment, the PH-tree contains 2.8 million entries. Figures 6.1a and 6.1b present evolution of the throughput, respectively of the response of the PH-tree for the insert operation. Figure 6.1a shows that the throughput of the system increases together with the number of indexing servers. We see that increasing the number of indexing hosts from one to two leads to a linear increase in throughput, while the addition of the third and fourth server increases the throughput even more. Figure 6.1b presents the evolution of the response time for the insert operation, showing that the response time decrease almost linearly when increasing the



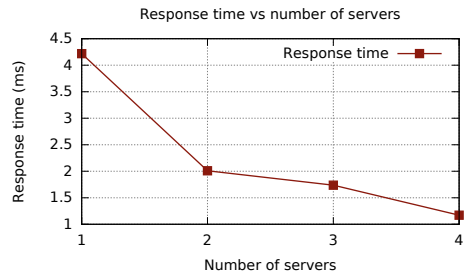
(a) Write throughput



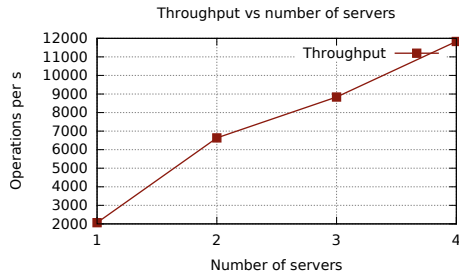
(b) Write response time



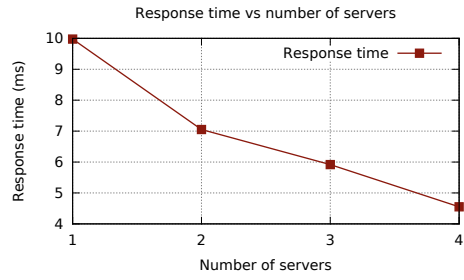
(c) Read throughput



(d) Read response time



(e) Range query throughput



(f) Range query response time

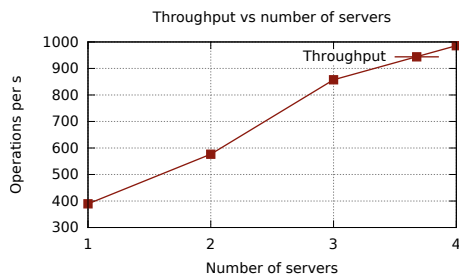
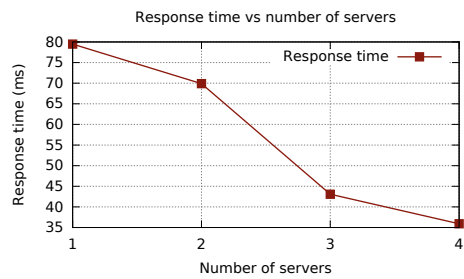
(g) KNN with $k = 100$ throughput(h) KNN with $k = 100$ response time

Figure 6.1: Evolution of the throughput and response time of the distributed PH-tree for clusters of different size

number of hosts.

The point query workload is based on the insertion workload, using the same number of client machines and threads. Each thread inserts 50000 points, and also retrieves the each point with a point query after the insertion operation was performed. For this workload, we measured only the duration of the `get()` operation. Figure 6.1c presents the evolution of the throughput, while Figure 6.1d presents the evolution of the response time. The evolution of the throughput and response time is similar to the insertion workload, as increasing the number of indexing servers increase the throughput of the system and reduces the response time.

For the range query workload, each of the 8 threads running on each machine performs 10000 range queries on a 2-dimensional PH-tree generated using the insertion workload described previously. Each range query is performed on a hyper-rectangle centered over an existing query point. The top left corner of the query hyper-rectangle is obtained by subtracting a random displacement from each dimension of the query center point. The top right corner is obtained by adding the same randomly generated displacements to each dimension of the center point. The displacements are sampled from the uniform distribution in the integer range $[0, 2^{21} - 1]$ using the `nextInt()` method of the `Random` class. With this scaling, the number of results is kept under 1000 for each query. Furthermore, more than 70% of all range queries need to query a single host to retrieve the response. Figures 6.1e and 6.1f present the evolution of the throughput and response time for the described workload. In such a workload, when most of the range queries are resolved by sending the range request to a single indexing host, the throughput and response time evolve similar to workloads involving point operations, i.e. each new added indexing server leads to an increase in throughput and a decrease in response time.

It is important to add that in the case of range queries, the measured performance depends on the number of hosts that need to be contacted to resolve each range query. As range queries need to send requests to all of the indexing servers in the worst case, some other workloads might lead to each query being forwarded to all hosts in the cluster. In that case, even if the client library attempts to send the requests in parallel, the response time of the range operation will be dependant on the number of indexing hosts n .

The nearest neighbour workload, similar to the range query workload, is executed on a 2-dimensional PH-tree obtained using the insertion workload and containing 2.8 million points. Each of the 8 threads running on the client machine perform 5000 nearest neighbour queries with $k = 100$. For this experiment, each query points considered in the nearest neighbour query is chosen from the points already stored in the tree. In the case of this workload, around 90% of the range queries are resolved after the host containing the query point is queried, while the rest of the queries need to be expanded to multiple hosts. Figure 6.1g shows the that the throughput of the system increases as new servers are added. Moreover, additional servers also decrease the average response time, as shown in Figure 6.1b. As is the case with range queries, the response time of the nearest neighbour query is not independent of the number of indexing hosts. As k increases, the section of the space containing the k neighbours increases as well, and a larger number of hosts increases the potential number of hosts whose regions intersect with the section containing the k neighbours.

6.2 Concurrency

6.2.1 Experimental Setup

The aim of the concurrency tests is to evaluate the throughput and response time of each implemented concurrency strategy. Furthermore, it is important to view how each concurrency strategy scales with respect to the number of concurrent threads accessing the PH-tree. Additionally, we identify the impact the concurrency strategies have when the PH-tree is only accessed by a single thread in the same time.

The concurrent experiments have been run on a Google Compute Engine VM³ instance with 16 virtual CPU cores and 60GB of RAM. Similar to the machines used for the distribution computation, information regarding the physical hardware used by the machine is not available. However, the emulated hardware family is Intel Xeon E5-1600, and each virtual CPU has a frequency of 2.5 GHz. All performed experiments were performed with 3 repetitions and the response time and throughput have been averaged across the repetitions.

6.2.2 Experimental Results

Write performance

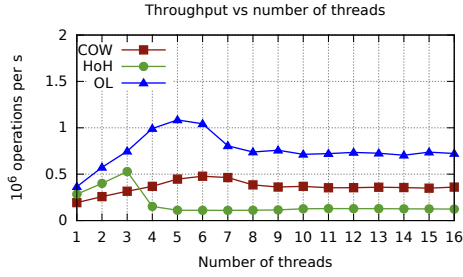
The first type of tests evaluate the throughput and response time of the concurrency strategies for an increasing number of threads. In each run, an array of 5 million multi-dimensional elements is inserted into an empty PH-tree. The values are sampled from the Gaussian distribution. The array is split equally into non-overlapping intervals between the inserting threads and each thread will insert a certain section of it. As the number of inserted entries is equal across the runs, it is expected that more inserting threads will lead to a higher throughput.

The experiments were performed for multi-dimensional entries with an increasing number of dimensions. Figure 6.2 shows the response time and throughput for the insertion operation, with the number of threads varying from 1 to 16. The number of writer threads is limited to the number of available CPU cores, to limit the effects of context switching on the experiment execution. The figure shows the results of the same experiment on PH-tree having dimensions 1, 2, 3, 6 and 10.

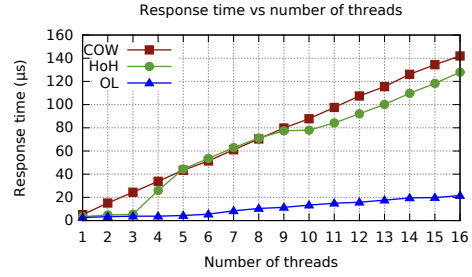
For all of the tests, the Copy-on-Write strategy maintains a relatively steady throughput, caused by the fact that this strategy only allows a single writer thread and multiple readers at a time. The contention on the writer lock also causes the response time to increase as the number of writer threads increases.

For up to 6 dimensions, both the Hand-over-Hand and Optimistic locking strategies generally perform better than the Copy-on-Writer strategy for a small number of writer threads, having a higher throughput and a lower response time. For both of these strategies, increasing the number of threads leads to an increase in throughput, up to a point where the lock contention becomes so high that it causes the throughput to drop. Optimistic locking generally outperforms Hand-over-Hand locking, having a better throughput and a lower response time in all cases. Additionally, Optimistic locking generates less lock contention and manages to main-

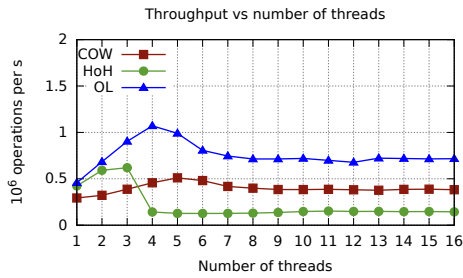
³The code of the virtual machine on the Google Compute Engine platform is `n1-standard-16`.



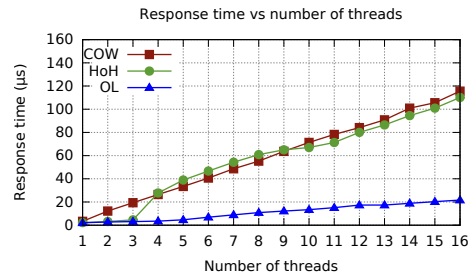
(a) 1D PH-tree throughput



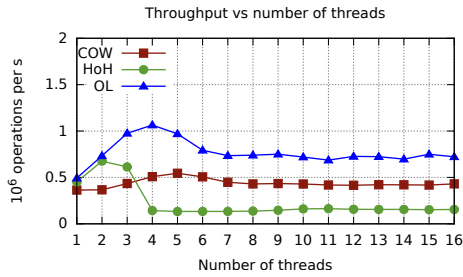
(b) 1D PH-tree response time



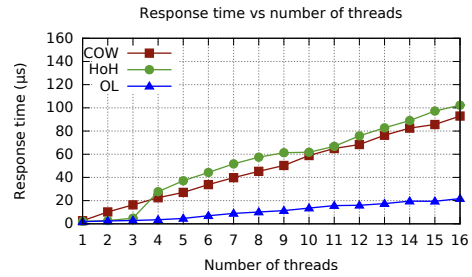
(c) 2D PH-tree throughput



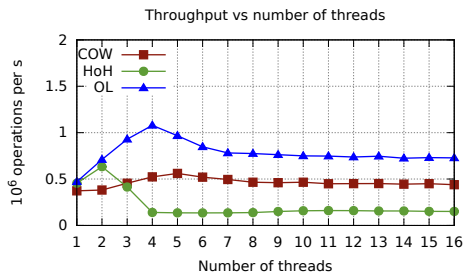
(d) 2D PH-tree response time



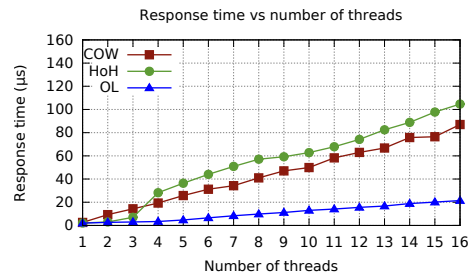
(e) 3D PH-tree throughput



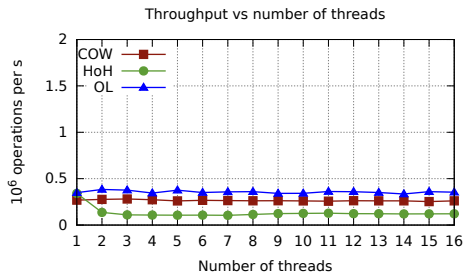
(f) 3D PH-tree response time



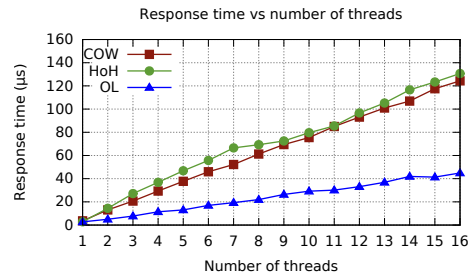
(g) 6D PH-tree throughput



(h) 6D PH-tree response time



(i) 10D PH-tree throughput



(j) 10D PH-tree response time

Figure 6.2: Evolution of write throughput and response time for the PH-tree

tain the increase in throughput for a higher number of writer threads than Hand-over-Hand locking. For both fine grained locking strategies, an increase in the number of threads leads to an increase in the lock contention, which leads to an increase in response time.

Another important factor in the multi-threaded performance of the PH-tree is the contention on the memory bus. As the number of dimensions of the PH-tree increases, each thread needs to both read and write an increasing amount of data from and to the main memory. Figure 6.2i shows that for 10 dimensional trees, the difference in throughput between the concurrency approaches is much smaller than in the 1, 2, 3 or 6 dimensional experiments.

Single Thread Write Performance

The second type of tests measures the penalty induced by the concurrency strategies when the tree is accessed by a single writer thread. For these experiments, 1 million 3-dimensional entries were inserted into an empty PH-tree by a single thread. The hardware used for this tests is an Intel i5-760 with 4 CPU cores, each running at a 2.8 GHz frequency. The values for each entry were generated uniformly at random. Figure 6.3 illustrates the time needed for the insertion workload for each concurrency strategy and the single-threaded version of the tree. For the single-threaded version, the insertion finished in 1.1 seconds. The version will be used as a baseline for the other measurements.

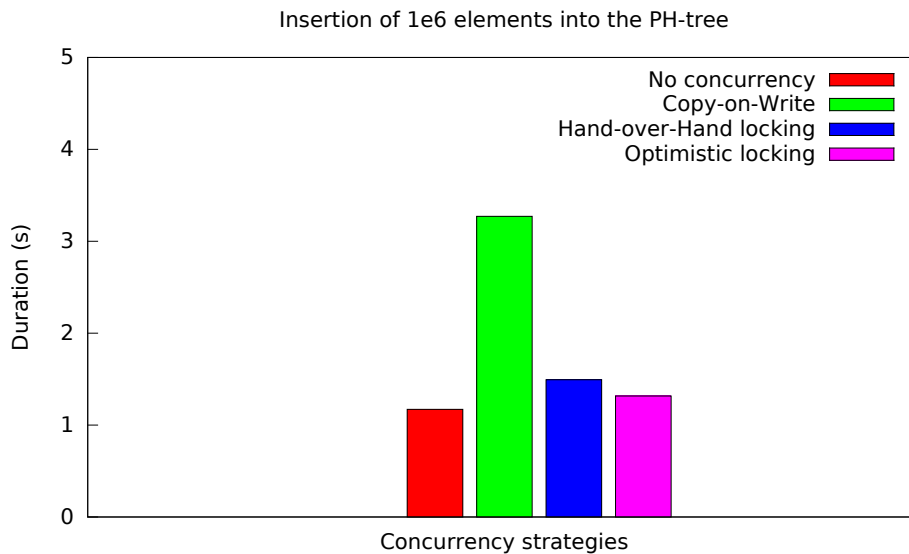
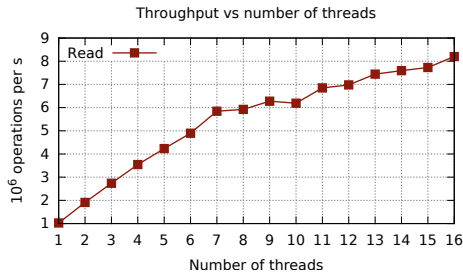
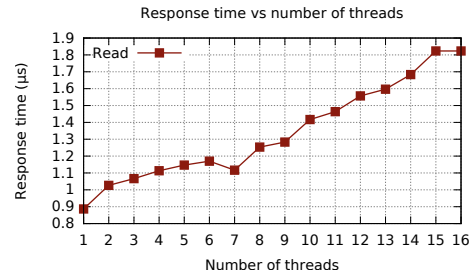


Figure 6.3: Time needed by 1 thread to insert 1^6 entries into the PH-tree

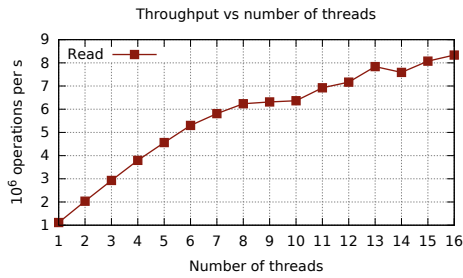
In the case of the Copy-on-Write version, the insertion workload finished in 3.2 seconds, taking almost three times longer than the baseline. This is due to the fact that each insert copies a full-path from the root of the tree to the modified node. Hand-over-Hand locking finished in 1.5 seconds, and was approximately 30% slower than the baseline version. In the case of Hand-over-Hand locking, all nodes in the path from the root to the modified node are locked and copies are created only for the nodes which are updated. This makes the Hand-over-Hand strategy faster than the full Copy-on-Write strategy, but slower than the



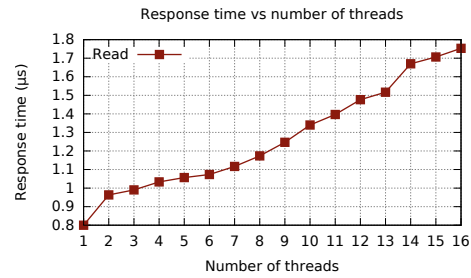
(a) 3D PH-tree throughput



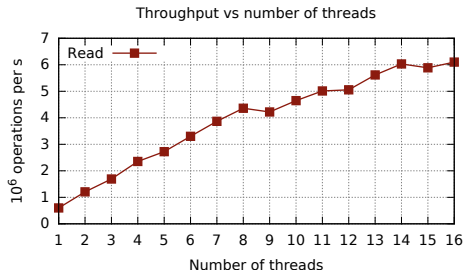
(b) 3D PH-tree response time



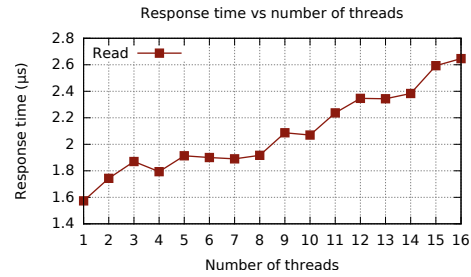
(c) 6D PH-tree throughput



(d) 6D PH-tree response time



(e) 10D PH-tree throughput



(f) 10D PH-tree response time

Figure 6.4: Evolution of read throughput and response time for the PH-tree

single-threaded implementation, which does not take any locks or perform any node copies. Finally, the Optimistic locking strategy performs the best from the implemented concurrency strategies, finishing the insertion in 1.3 second. This strategy is faster than the Hand-over-Hand strategy, as only two nodes are locked, as opposed to a full path, but is still slower than the single-threaded implementation.

Read performance

Finally, the last type of tests evaluate the scalability of the PH-tree read operations with respect to the number of concurrent reader threads. In this test, the PH-tree is loaded with 5 million entries and then all entries inserted in the tree are retrieved with the *get()* point query. The entries are split into non-overlapping partitions and each partition is assigned to a reader thread. Each thread then retrieves all of the keys in its partition. As previously mentioned, the concurrent access strategies did not modify the read operations and all implemented concurrent strategies perform the read operations in the same manner.

These tests have been run on the same 16 CPU machine described in 6.2.1. The number of reader threads is varied from 1 to 16, to reduce the effects of thread context switching. The test was repeated for PH-trees of dimensions 1, 2, 3, 6 and 10, to observe the effect of the higher dimensional data in the overall performance of the PH-tree. Figure 6.4 presents the evolution of the throughput and response time for the 3, 6 and 10 dimensions. The plots for dimensions 1 and 2 have been omitted as the performance of the PH-tree in those cases is very similar to the 3 dimensional case.

The results show that an increased number of reader threads increases the overall throughput of the PH-tree. The increase in throughput is almost linear for up to 6 reader threads, after which each additional reader thread adds a smaller increase to the throughput. The throughput stabilizes after the number of readers goes over 10. As in the case of insertion performance, the number of dimensions of the data stored in the PH-tree influences both the throughput and response time. As the number of dimensions increases, the throughput decreases and the response time decreases. This is due to the fact that the PH-tree nodes themselves are bigger and each thread needs to transfer more data over the memory bus.

6.3 Results

The distributed experiments showed that the distributed PH-tree scales with an increasing number of server machines. In all experiments of this type, the addition of server machines to the cluster increased the overall throughput of the system and decreased the response time for each individual request. This increase in throughput and decrease in response time will always happen in workloads consisting mostly of point operations. In the case of range and nearest neighbour queries, the performance is dependent on the number of hosts that need to be contacted for each individual query, a number that generally increases if the number of hosts is increased. However, we found that if the majority of the range and nearest neighbour queries are resolved on a single host, the system displays similar throughput and response time characteristics as in the case of point operation workloads.

The concurrent experiments compared the performance of the three implemented concurrency strategies: *Copy-on-Write*, *Hand-over-Hand locking* and *Optimistic locking*. *Optimistic locking* performed better for the write operations, having a higher throughput and a lower response time than both other strategies. *Hand-over-Hand locking* performs better than *Copy-on-Write* for smaller dimensions and up to 3 or writer threads, depending on the operation. *Hand-over-Hand locking* performs worse in settings with a higher number of writers due to an increased lock and memory bus contention. Generally, *Optimistic locking* provides the best performance for both read and write operations.

7

Conclusions

7.1 Contributions

This work presents a complete solution for implementing a distributed and parallel multi-dimensional index based on the PH-tree. We propose a distribution architecture for the PH-tree and devise a set of efficient distributed algorithms for the PH-tree operations. Furthermore, the in-memory PH-tree was extended to support concurrent write operations, allowing each machine in the cluster to service more requests in an concurrent manner. The main contributions of this work are the following:

The distribution architecture partitions the multi-dimensional space in a manner that allows fast point operation and efficient range and nearest neighbour queries. The proposed balancing algorithm ensures that all hosts will eventually hold a comparable number of entries, even in the case of very skewed input distributions. The presented distribution solution can be easily adapted for other types of indexes, as long as the key domain for of this indexes can be mapped to a 1 dimensional line, allowing the construction of efficient and scalable range-queryable, distributed key-values stores.

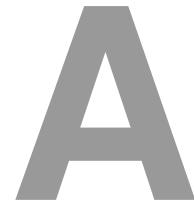
The in-memory PH-tree was updated to support concurrent access according to two consistency paradigms. The Copy-on-Write approach offers full consistency guarantees, providing snapshot semantics to queries and isolating readers from any further modifications to the tree. For situations where full consistency is not needed, the Optimistic locking approach offers relaxed consistency semantics and provides improved write throughput by allowing concurrent write operations. Hand-over-Hand locking provides the same consistency as Optimistic Locking, and also guarantees that starvation will not occur in the case of the writers. Furthermore, none of presented approaches require the readers to acquire any locks and thus add no overhead to the execution of the read operations.

7.2 Future work

One aspect of the distributed PH-tree that could be improved is the execution of the distributed nearest neighbour queries. As previously mentioned in section 3.4, the nearest neighbour algorithm first obtains the k nearest neighbours from the host responsible for the query point and then proceeds to run a range query on the cluster to check for any closer points on different hosts. The additional check is done by exploring a hyper-cubic section of the initial space. A more efficient approach would be to check a hyper-spheric region centered in the query point. This improvement could reduce the number of hosts queried on the second step of the nearest neighbour queries and could potentially improve the execution time of such queries.

From the point of view of the balancing algorithm, hosts participating in re-balancing operations do not accept write requests for the duration of these operations, to simplify the update of the mapping between the points and the hosts. The balancing algorithm could be modified to allow the hosts to accept write requests to a separate buffer and integrate the buffer to the stored PH-tree after the balancing operation has finished.

The concurrency strategies presented by this work rely on locking and copy-on-write to allow concurrent access. Further work could focus on investigating the performance of concurrent access strategies which use atomic operations instead of locks as a synchronization mechanism. Such a solution would remove the overhead added by the lock contention and potentially provide better performance in some situations.



Installation and Usage Instructions

The distributed PH-tree was implemented in Java and uses Maven¹ for compilation, packaging and dependency management. Compiling the project requires the Java Development Kit 8 and Maven 3 to be installed on the local computer. Moreover, the source code uses the PH-tree library as a compilation and run-time dependency, and therefore it is required that the PH-tree Maven artifact be installed in the local computer repository prior to compilation.

A.1 Compilation and Installation

A.1.1 Client and Server Libraries

As previously mentioned, the compilation is performed using Maven, which can be downloaded from <http://maven.apache.org/download.cgi>. To compile the project, open a console to the directory where the source code is stored and type the command: `maven package`. This command will compile both the client and server libraries and execute all the unit tests. When compiling from the command line, it is necessary that the Maven binary be present in the PATH variable of the command line environment. Alternatively, the project can be compiled from an IDE, as the majority of the contemporary IDE's provide wizards for importing and managing Maven projects.

A.1.2 ZooKeeper

The cluster of machines running the distributed PH-tree needs at least one ZooKeeper server to manage the cluster configuration information. An up-to-date guide for downloading and installing a basic ZooKeeper server is available at the following location: <https://>

¹<https://maven.apache.org/>

```
// connect to a cluster of machines storing
// its configuration in the ZooKeeper at
// the following address
String zkHost = "localhost";
int zkPort = 2181;
PHFactory factory = new ZKPHFactory(zkHost, zkPort);

//create a 2D PH-tree with 64 bit width
PhTree tree = factory.createPHTreeSet(2, 64);
```

Figure A.1: Example code snippet for connecting to the PH-tree from the Java API

zookeeper.apache.org/doc/trunk/zookeeperStarted.html. A short summary of operations needed to install and start a ZooKeeper server is presented below.

- Download the latest stable ZooKeeper package and copy it to the desired installation directory. This directory will be referred to as `ZOO_HOME`.
- In the `ZOO_HOME/conf` directory, rename the `zoo_sample.cfg` file to `zoo.cfg`. The initial file contains a sample configuration file for ZooKeeper, and renaming it to `zoo.cfg` will make ZooKeeper use this configuration when started.
- Start the ZooKeeper server using the `zkServer.sh` script in the directory `ZOO_HOME`.

A.2 Usage

To set up the distributed PH-tree, first start the ZooKeeper server on an assigned computer. Indexing servers can be started using the `middleware-1.0-SNAPSHOT.jar` Java jar. The following code snippet starts a server on the machine with `ip server-ip` on port 5000 that connects to a ZooKeeper server running on `zk-ip`:

```
java -jar middleware-1.0-SNAPSHOT.jar server-ip:5000 zk-ip:2181
```

Additional servers can be added by executing this command on a different machine and changing `server-ip` accordingly.

Clients can access the distributed PH-tree using the API provided in the client library module. An example for connecting to a running cluster of machines is listed in Figure A.1. The `PHFactory` class creates references to PH-tree classes, which can be used to perform operations on the distributed PH-tree.

List of Figures

2.1	2 dimensional 4bit PH-tree containing the points: (0000, 1000), (0011, 1000), (0011, 1010).	6
3.1	Distributing the points based on a hash function	12
3.2	Equal geometric space splitting for 2 hosts and 2 dimensions	13
3.3	Equal geometric space splitting for 4 hosts and 2 dimensions	13
3.4	Potential partitioning of the space for a skewed input distribution.	13
3.5	Another potential partitioning for the same skewed input distribution.	13
3.6	Z-order space filling curve filling the 2 dimensional space.	14
3.7	Example intervals on the 3bit Z-order curve	15
3.8	Example intervals on the 3 bit Z-order curve filling the 2D space	15
3.9	The network architecture of the distributed PH-tree	16
3.10	The intervals described in 3.8 split into hyper-rectangles.	17
3.11	Illustration of a range query intersected with the host intervals	19
3.12	Illustration of a 3 nearest neighbour query execution in a 2 dimensional space.	20
4.1	A 2 dimensional PH-tree containing three 4-bit entries: (0000, 1000), (0011, 1010), (1011, 1000). The values are omitted.	24
4.2	Illustration of Copy-on-Write on a binary tree where the node G has to be written	25
4.3	Illustration of hand-over-hand locking descent to node G	26
4.4	Illustration of optimistic locking descent to node G	27
5.1	The implementation architecture of the distributed PH-tree	32
5.2	The Java API of the PH-tree. The classes coloured with blue represent the existing code, the additions in green represent the distributed implementation.	33
5.3	The class diagram of the client library	34

5.4	The class diagram of the server library	35
5.5	A diagram of the Key-Mapping stored in ZooKeeper	37
5.6	A diagram of the communication flow between the client and the server	39
5.7	The Java API of the concurrency classes	40
6.1	Evolution of the throughput and response time of the distributed PH-tree for clusters of different size	45
6.2	Evolution of write throughput and response time for the PH-tree	48
6.3	Time needed by 1 thread to insert 1^6 entries into the PH-tree	49
6.4	Evolution of read throughput and response time for the PH-tree	50
A.1	Example code snippet for connecting to the PH-tree from the Java API . . .	56

Bibliography

- [1] Jon Louis Bentley. Multidimensional binary search trees used for associative searching. *Commun. ACM*, 18(9):509–517, September 1975.
- [2] Nathan G. Bronson, Jared Casper, Hassan Chafi, and Kunle Olukotun. A practical concurrent binary search tree. *SIGPLAN Not.*, 45(5):257–268, January 2010.
- [3] Trevor Brown, Faith Ellen, and Eric Ruppert. A general technique for non-blocking trees. In *Proceedings of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, PPoPP ’14, pages 329–342, New York, NY, USA, 2014. ACM.
- [4] Raphael A. Finkel and Jon Louis Bentley. Quad trees: A data structure for retrieval on composite keys. *Acta Inf.*, 4:1–9, 1974.
- [5] Keir Fraser and Tim Harris. Concurrent programming without locks. *ACM Trans. Comput. Syst.*, 25(2), May 2007.
- [6] Haojun Liao, Jizhong Han, and Jinyun Fang. Multi-dimensional index on hadoop distributed file system. In *Networking, Architecture and Storage (NAS), 2010 IEEE Fifth International Conference on*, pages 240–249, July 2010.
- [7] M. Moir and N. Shavit. Concurrent data structures. In *Handbook of Data Structures and Applications*, D. Metha and S. Sahni Editors, pages 47–14 – 47–30, 2007. Chapman and Hall/CRC Press.
- [8] Sriram Ramabhadran, Sylvia Ratnasamy, Joseph M. Hellerstein, and Scott Shenker. Prefix hash tree: An indexing data structure over distributed hash tables. Technical report, 2004.
- [9] Toms Skopal, Michal Krtek, Jaroslav Pokorný, and Vclav Snel. A new range query algorithm for Universal B-trees. *Information Systems*, 31:489–511, 2006.
- [10] Ion Stoica, Robert Morris, David Karger, M. Frans Kaashoek, and Hari Balakrishnan. Chord: A scalable peer-to-peer lookup service for internet applications. In *Proceedings of the 2001 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications*, SIGCOMM ’01, pages 149–160, New York, NY, USA, 2001. ACM.
- [11] Tilmann Zäschke, Christoph Zimmerli, and Moira C. Norrie. The ph-tree: A space-efficient storage structure and multi-dimensional index. In *Proceedings of the 2014*

- ACM SIGMOD International Conference on Management of Data*, SIGMOD '14, pages 397–408, New York, NY, USA, 2014. ACM.
- [12] Chi Zhang and Arvind Krishnamurthy. Skipindex: Towards a scalable peer-to-peer index service for high dimensional data. Technical report, Princeton Univ, 2004.



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

Title of work (in block letters):

CLUSTER COMPUTING AND PARALLELISATION FOR THE MULTI-DIMENSIONAL PH-INDEX

Authored by (in block letters):

For papers written by groups the names of all authors are required.

Name(s):

VANCEA

First name(s):

BOGDAN AUREL

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the 'Citation etiquette' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

Place, date

ZURICH, 1.04.2015

Signature(s)

For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.