

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

Master Thesis

Bogdan Aurel Vancea

<bvancea@student.ethz.ch>

Prof. Dr. Moira C. Norrie
Tilman Zaeschke
Christoph Zimmerli

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

21st November 2014



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



Abstract

Here comes the abstract.

Contents

1	Introduction	1
1.1	Multi-dimensional Indexes	1
1.1.1	Background information	1
2	Design and Implementation	3
2.1	Algorithms	3
2.1.1	Data partitioning	3
2.1.2	Data balancing	3
2.1.3	Basic Index operation	4
2.1.4	Iterators	4
2.1.5	Range search	4
2.2	Implementation	4
3	Performance Analysis	5
3.1	Benchmark	5
4	Conclusions	7

1

Introduction

1.1 Multi-dimensional Indexes

1.1.1 Background information

Need to add an introduction here.

This is an example of how to cite a scientific publication [1] from your bibliography (BibTeX¹ file). And this example shows how you create links within your documents, e.g. link to section 1.1.

¹<http://en.wikipedia.org/wiki/BibTeX>

2

Design and Implementation

2.1 Algorithms

2.1.1 Data partitioning

Describe here how the key-value pairs are partitioned across the hosts.

2.1.2 Data balancing

Describe how the key-value pairs are balanced across the index nodes. The cluster should be able to properly balance the amount of keys that are stored in each server.

A simple load balancing strategy would be the following:

1. Upon reaching a certain threshold t of keys stored, a host decides it has to split its zone.
2. This host first sends a broadcast to all nodes to request the number of keys they all store. It then chooses the hosts with the fewest keys and considers that this key is the split receiver. The splitter decides to split its zone in half (or in such manner that around half of the keys held in the initial zone are moved to the receiver).
3. The splitting host sends the keys to the receiver. (What does the receiver do with them? Duplicates could appear in case of KNN or range queries.) .
4. After the receiver received all keys, the mapping is updated on the Zookeeper.
5. The splitter then deletes all of the keys sent from its own index (maybe this can be done really fast by removing a bunch of nodes from the index).

Ideas that might not work:

- Storing information like the number of keys held by each node in Zookeeper. This would mean that the ZooKeeper would need to be notified on each insert, which would cause severe scalability and availability issues. What could work is that the ZK is notified when each host reaches a certain key threshold (i.e, every 1000 keys inserted, a request is send to ZK).

2.1.3 Basic Index operation

Describe the point operations here. These are operations that affect a single key-value pair, like get, put, delete, contains, etc.

2.1.4 Iterators

Describe how iterators are handle when dealing with a cluster of index servers. Describe the current algorithm used and the alternatives.

2.1.5 Range search

Describe how the range search is performed. Describe how the number of hows that need to be queried is reduced and what the alternatives are.

2.2 Implementation

Describe the technologies used, the reasons for which these technologies were chosen and any alternatives.

Currently used frameworks:

ZooKeeper ZooKeeper is used for stored cluster metadata and membership. Currently, the only alternative would have been to implement such a distributed storage manager manually. Using ZooKeeper saved a lot of development time.

Netty Netty is a Java IO library and it is used to implement the server request handling component. Alternatives would have been Java NIO library.

Kryo Kryo is very fast serialization library for Java and it is used to serialize the values that have to be stored on the server. These objects need to be transformed into a representation that can be sent over the network. Kryo is faster than the Java serialization, does not require the implementation of the Serializable interface and transforms the objects into byte arrays. This should make the representation smaller than simply transforming the object to a string.

3

Performance Analysis

3.1 Benchmark

4

Conclussions

List of Figures

List of Tables

Acknowledgements

Bibliography

- [1] Alfonso Murolo. Designing wordpress themes by example. Master's thesis, ETH Zurich, 2013.