

PROJEKTBERICHT

# HAMAUBE

KAI MARTINEN, FABIAN KOHLER, ...

JUNE 12, 2018



UNIVERSITÄT HAMBURG

DEPARTMENT OF COMPUTER SCIENCE

CHAIR OF DISTRIBUTED SYSTEMS AND INFORMATION  
SYSTEMS

BETREUT DURCH STEFFEN FRIEDRICH

# Abstract

Abstract in English

# Kurzfassung

Kurzfassung auf Deutsch

# Contents

<b>Abstract</b>	<b>I</b>
<b>Kurzfassung</b>	<b>II</b>
<b>List of Tables</b>	<b>VI</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Initial goal and contributions . . . . .	1
1.2 Thesis outline . . . . .	1
<b>2 Preliminaries</b>	<b>2</b>
2.1 Topic 1 . . . . .	2
2.1.1 Subtopic1 . . . . .	2
2.1.2 Subtopic2 . . . . .	2
2.2 Topic 2 . . . . .	2
2.2.1 Subtopic1 . . . . .	2
<b>3 Conclusion</b>	<b>3</b>
3.1 Future work . . . . .	3
<b>A Glossary</b>	<b>5</b>
<b>B Appendix C</b>	<b>10</b>
B.1 Section 1 . . . . .	10
B.2 Section 2 . . . . .	10

<b>C</b>	<b>Appendix C</b>	<b>11</b>
C.1	Einführung . . . . .	11
C.2	Datenmodell . . . . .	11
C.3	Refinements . . . . .	15
C.4	Performance Auswertung . . . . .	16
<b>D</b>	<b>Appendix C</b>	<b>17</b>
D.1	Section 1 . . . . .	17
D.2	Section 2 . . . . .	17
<b>E</b>	<b>Appendix C</b>	<b>18</b>
E.1	Section 1 . . . . .	18
E.2	Section 2 . . . . .	18
<b>F</b>	<b>Appendix C</b>	<b>19</b>
F.1	Section 1 . . . . .	19
F.2	Section 2 . . . . .	19
<b>G</b>	<b>Appendix D</b>	<b>20</b>
G.1	Section 1 . . . . .	20

# List of Figures

C.1	.....	12
C.2	.....	13
C.3	.....	16

# List of Tables







# Chapter 1

## Introduction

Introduction.

You can reference the only entry in the .bib file like this: [1]

### 1.1 Initial goal and contributions

### 1.2 Thesis outline

# Chapter 2

## Preliminaries

### 2.1 Topic 1

#### 2.1.1 Subtopic1

#### 2.1.2 Subtopic2

### 2.2 Topic 2

#### 2.2.1 Subtopic1

# Chapter 3

## Conclusion

Write here you conclusions

### 3.1 Future work



# Appendix **A**

## Glossary

Just comment `\input{AppendixA-Glossary.tex}` in `Masterthesis.tex` if you don't need it!

### Symbols

---

\$ US. dollars.

### A

---

A Meaning of A.

### B

---

### C

---

### D

---

---

**E**

---

**F**

---

**G**

---

**H**

---

**I**

---

**J**

---

**M**

---

**N**

---

---

P

---

Q

---

R

---

S

---

T

---

U

---

V

---

W

---

X

---



---

---



# Appendix **B**

## Appendix C

Some text.

### **B.1 Section 1**

### **B.2 Section 2**

## Appendix C

Some text.

### C.1 Einführung

Täglich kommen mehrere Petabytes an Daten von über 60 Google Anwendungen zusammen. Dafür verantwortlich sind mehr als 1000 Computer die untereinander vernetzt sind. Um diese Daten verwalten zu können wurde Bigtable ins Leben gerufen. Das Ziel der Datenbank war es in vielen Bereichen anwenden zu können. Dazu sollte es Skalierbar sein sowie eine hohe Performance und Verfügbarkeit besitzen.

### C.2 Datenmodell

Bigtable ist eine verteiltes, persistentes multidimensionale sortierte Map. Diese Map ist indexiert über eine row key, column key und einem timestamp. Jeder Wert in dieser Map ist ein Array mit Bytes. Das folgende Datenmodell soll eine Speicherung von Webseiten veranschaulichen.

Das Datenmodell besteht aus zwei Familien dem Inhalt und den Ankerpunkten. Die erste Familie beinhaltet den Inhalt der Webseite, mit drei unterschiedlichen Zeitstempeln ( $t_3, t_5, t_6$ ). Drei unterschiedliche Zeitstempeln bedeutet, dass die Website `www.cnn.com` in drei unterschiedlichen Versionen abgespeichert wurde. Die Anker Familie beinhaltet jeweils nur eine Version. Den Anker mit „CNN.com“ mit dem Zeitstempel  $t_8$  und dem Anker „CNN“

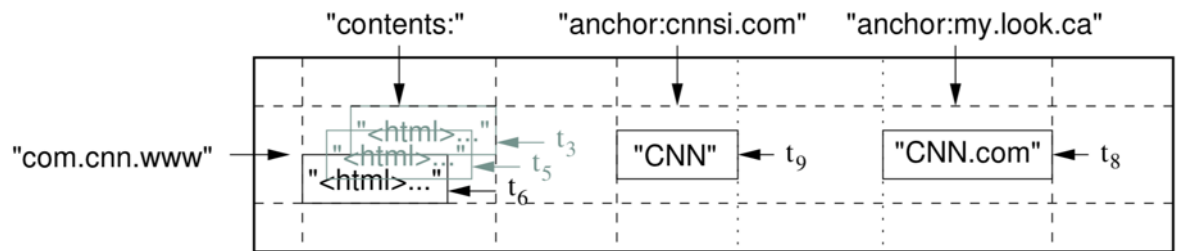


Figure C.1: Tabellen Beispiel

mit dem Zeitstempel t9.

**Rows** Bigtable speichert Daten in lexikographischer Reihenfolge und sortiert diese nach Zeilen. Eine row range beinhaltet alle gleichnamigen URL's, so dass alle mit der gleichen Domain zusammen abgespeichert werden. Das vereinfacht die Analyse und das Hosting der gleichnamigen Domains und macht dies zudem effizienter.

**Column families** Verschieden column keys werden in eine gemeinsame Gruppe gespeichert. Das nennt man column families. Alle Daten, welche in der gleichen Gruppe gespeichert werden sind vom gleichen Typ. Bevor Daten in einer Gruppe gespeichert werden können, muss diese column family als erstes erstellt werden.

**Timestamp** Jede Zelle in Bigtable kann mehrere Versionen der gleichen Daten beinhalten. Dieser Versionen werden indexiert durch den Zeitstempel (timestamp). Der Zeitstempel ist bis auf die Micro Sekunde genau. Durch die „two per-column-family“ kann der Benutzer festlegen, wie viele Versionen der gleichen Daten gespeichert werden sollen. Alle weiteren Versionen werden automatisch gelöscht.

**API** Die API von Bigtable ermöglicht das Erstellen und Löschen von Tabellen und Spaltennamen, sowie das Ändern von Tabellen, Cluster und Metadaten einer Spaltenfamilie. Das folgende Codebeispiel wurde in C++ geschrieben und verändert den Inhalt der Tabelle Webservice.

Der Code verändert die Spalte „com.cnn.www“ und fügt einen neuen Anker hinzu. Im nächsten Schritt wird ein vorhandener Anker „anchor:www.abc.com“

```
// Open the table
Table *T = OpenOrDie("/bigtable/web/webtable");

// Write a new anchor and delete an old anchor
RowMutation rl(T, "com.cnn.www");
rl.Set("anchor:www.c-span.org", "CNN");
rl.Delete("anchor:www.abc.com");
Operation op;
Apply(&op, &rl);
```

Figure C.2: Zugriffs Beispiel mit der BigTable API

gelöscht und festgeschrieben.

**Building Blocks** Bigtable ist auf mehreren anderen Teilen der Google-Infrastruktur aufgebaut. Zum Beispiel benutzt Bigtable das verteilte Google File System (GFS). Welches für die Speicherung von Logs und Daten verantwortlich ist. Ein Bigtable Cluster operiert typischerweise auf einem verteilten Pool von Computern. Auf diesem Pool laufen eine breite sparte von verschiedenen Anwendungen. Bigtable basiert auf einem Cluster Management System für Zeit-Planung-Jobs, Ressourcenmanagements auf geteilten Computer, agieren mit Computerfehlern und dem Anzeigen des Computer Status. Das SSTable Format stellt eine persistente, geordnete unveränderliche Abbildung von Schlüsseln zu Werten zur Verfügung, wobei sowohl Schlüssel als auch Werte willkürliche Bytefolgen sind. Operationen werden bereitgestellt, um den Wert zu suchen, der einem bestimmten Schlüssel zugeordnet ist.

**Implementation** Bigtable besteht aus drei Haupt Komponenten, einer Bibliothek, einem Master Server und vielen weiteren tablet servern. Die Bibliothek ist in jeden Client verlinkt, somit kann der Client auf alle Funktionen zugreifen. Der Master Server wird zufällig ausgewählt. Dieser teilt den tablet servern die tablets zu. Außerdem ist für die Verteilung der Lasten zuständig und ist für die garbage collection. Sobald eine Tabelle zu groß wird, wird diese von einem tablet server gesplittet. So wird sichergestellt, dass eine Tabelle nie größer als 100-200 MB ist.

**Tablet location** Um Daten zu speichern, verwendet Google bei Bigtable eine „three-level hierarchy“. Das erste Level ist eine Datei, welches auch das Chubby file genannt wird, dort wird der Speicherort des root tablets hinterlegt. Das root tablet beinhaltet alle Speicherorte aller tablets in einer METADATA Tabelle. Das spezielle an dieser Tabelle ist, dass egal wie groß sie wird, diese niemals geteilt wird. Somit wird sichergestellt, dass die „three-level hierarchy“ eingehalten wird. Die METADATA Tabellen speichern die Orte aller anderen tablets in einer Tabelle ab.

**Tablet Zuweisung** Jedes tablet ist zu einem Zeitpunkt immer nur einem tablet server zugeordnet. Der Master server verfolgt die lebenden tablet servern und die aktuell zugeordneten tablets zu den tablet servern inklusive aller unzugeordneten tablet servern. Beim Start einer Bigtable führt der Master Server folgende Schritte aus:

1. Wählt einen einzigartigen Master Lock in Chubby
2. Scannt die Server Verzeichnisse um die lebenden tablet server zu finden
3. Kommuniziert mit den vorhandenen tablet server um bereits zugeordnete tablets zu identifizieren
4. Master scannt METADATA Tabelle um die vorhandene Zugehörigkeiten zu lernen

**Tablet serving** Ein persistenter Zustand eines tablets wird in GFS gespeichert. Alle Updates werden auf „well-formed“ geprüft und anschließend in einem commit-log gespeichert. Die neusten Updates werden in eine memtable gespeichert, ältere updates werden in die SSTable geschrieben. Wenn Daten aus dem tablet server abgefragt werden muss ein merge zwischen den neuen Daten in der memtable sowie den älteren Daten in der SSTable erstellt werden.

**Compaction** Je mehr Daten gespeichert werden, umso größer wird die memtable. Damit diese tabelle nicht zu groß wird, gibt es ein „minor compaction“. Diese Funktion friert eine memtable ein sobald diese eine bestimmte Größe erreicht hat und erstellt eine neue memtable. Die gefrorene

mentable wird zu einer SSTable konvertiert. Je mehr Daten gespeichert werden, desto unordentlicher wird die Ansammlung von SSTable. Um die SSTable zu sortieren wird periodisch ein „merging compaction“ ausgeführt. Dies strukturiert die SSTable neu und es werden Ressourcen, durch die Löschung von Daten, freigegeben. Außerdem werden gelöschte Daten endgültig gelöscht, das ist wichtig für Services, welche sensible Daten beinhalten.

## C.3 Refinments

Um die hohe Performance, Verfügbarkeit und Zuverlässigkeit beizubehalten, werden einige Verbesserungen (refinments) benötigt.

**Lokale Gruppen** Gruppierung erspart Zugriffszeit. Zum Beispiel bei dem Datenmodell Webseite. Die „page metadata“ und „content“ der Webseite werden in einer anderen Gruppe gespeichert. So muss eine Anwendung, welche nur die Metadaten möchte, nicht durch den kompletten Inhalt einer Seite iterieren. Zudem gibt es Tuningparameter, welche bestimmen ob Daten in den Arbeitsspeicher geladen werden um die Zugriffszeit zu minimieren.

**Kompression** Ein Benutzer kann selbst bestimmen ob SSTable komprimiert wird und falls ja, in welchem Ausmaß. Jeder SSTable Block kann einzeln ausgewählt werden. Für die Komprimierung kommen die Verfahren Bentley und McIlroy's zum Einsatz. Diese können mit 100-200Mb/s kodiert und mit 400-1000 MB/s enkodiert werden.

**Caching für Lesezugriffe** Für das Caching von Lesezugriffen gibt es zwei Verfahren. Der Scan Cache (high-level), speichert key-value Paare und liefert eine SSTable zurück. Das Block Cache (low-level) Verfahren speichert SSTable Blocks, die von der GFS gelesen werden.

**Bloom Filter** Benutzer legt selbst fest ob ein Filter zum Einsatz kommt. Der Vorteil eines Filters liegt darin, dass wenn Daten gesucht werden, nicht jede SSTable nach den bestimmten Daten durchsucht werden muss. Ein Bloom Filter erlaubt es, nach einer bestimmten Art von row/column Paaren zu fragen, ob diese in einer SSTable gespeichert sind.



**Beschleunigte tablet Wiederherstellung** Wenn der Master ein tablet von einem Server zu einem anderen Server verschiebt, führt der Ursprungs Server erst ein „minor compaction“ aus um die Ladezeit für den neuen tablet server zu verkürzen.

**Unveränderlichkeit ausnutzen** Es können nur Daten verändert, welche in der memtable stehen. Daten in der SSTable können nicht verändert werden. Das macht man sich zu nutzen indem man keine Synchronisation braucht, wenn auf die Daten zugegriffen wird. Memtable sind die einzigen Daten auf die man schreiben kann und gleichzeitig lesen. Damit es zu keinen konflikten kommt, setzt Bigtable hier auf „Copy-on-write“.

## C.4 Performance Auswertung

Experiment	# of Tablet Servers			
	1	50	250	500
random reads	1212	593	479	241
random reads (mem)	10811	8511	8000	6250
random writes	8850	3745	3425	2000
sequential reads	4425	2463	2625	2469
sequential writes	8547	3623	2451	1905
scans	15385	10526	9524	7843

Figure C.3: Percormance Übersicht

Die Performance wird hauptsächlich durch die verwendete CPU (2ghz) begrenzt. Zudem kann man erkennen, dass bei einem tablet server der Durchsatz bei ca. 75MB/s liegt ( $1000 \text{ bytes} * 64 \text{ KB Block size} = 75 \text{ MB/s}$ ). Damit der Durchsatz bei einem Single tablet server erhöht wird, wird die die SSTable größe in der regel von 64KB auf 8kb gesenkt. Zudem wird erkannt, dass der Durchsatz nicht Linear ansteigt. Bei einer Erhöhung der tablet server von eins auf 500 liegt die Erhöhung des Durchsatzes bei gerade mal dem 300 fachen ( $10811 / (500 * 6250) = 350$ ). Diese Begrenzung liegt wie bei einem tablet server an der CPU der tablet servern.

# Appendix **D**

## Appendix C

Some text.

### **D.1   Section 1**

### **D.2   Section 2**

# Appendix **E**

## Appendix C

Some text.

### **E.1 Section 1**

### **E.2 Section 2**

# Appendix **F**

## Appendix C

Some text.

### **F.1 Section 1**

### **F.2 Section 2**

# Appendix **G**

## Appendix D

### G.1 Section 1



# Bibliography

- [1] Marcelo Arenas, Claudio Gutierrez, and Jorge Pérez. On the semantics of sparql. In *Semantic Web Information Management*, pages 281–307. 2009.