# Bigtable: A Distributed Storage System for Structured Data

# A Comparison of Approaches to Large Scale Data Analysis

Jillian Preece

March 7, 2017

# Main Idea of Bigtable paper

- Bigtable- A storage system for managing structured data
  - Designed for huge amounts of data

- Google is largest user of Bigtable
  - Web indexing, Google Earth, Google Finance, Google Analytics

- Main idea- How Bigtable works, how it is implemented, and use cases

- Advantages- Easy to examine data change over a period of time (stores data for a customizable length of time), self-managing, uses building blocks (including mapreduce)

- Broken up into column families, columns, and rows

# Implementation of Bigtable

- Library – links into every client

- Master server – assigns each tablet to a specific tablet server, balances the tablet server load, handles schema changes within the database

- Tablet servers – each server can hold anywhere from 10 to thousands of tablets, handles read and write commands to each tablet that it has, splits tablets that have grown too large into smaller tablets

- Tablet – large table broken up at its' row boundaries
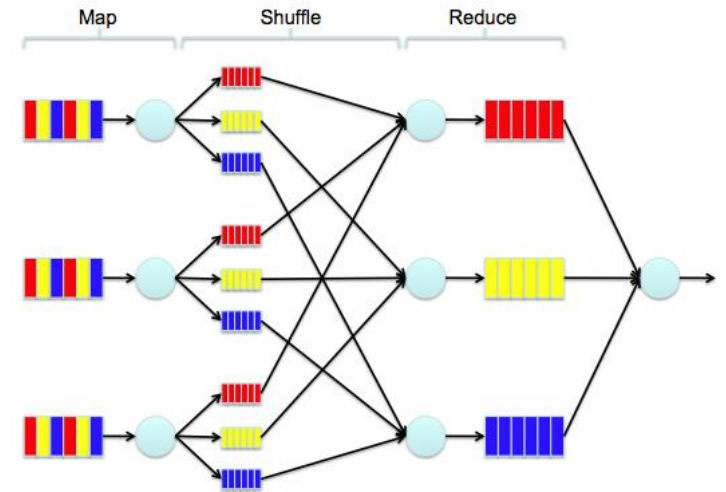
# Analysis of Implementation

- The implementation of Bigtable seems fairly well organized and efficient
- Each server is self-managing
  - Tablet servers handle the sizes of tablets and break them up as necessary
  - Master server handles changes in the database as a whole
- Data is split into many levels within each tablet, which allows for easy storage and accessibility of data
- Everything seems to be done in the fastest and best way possible
  - Not surprising, as it was created by Google

# Main Ideas of Approaches to Large Scale Data Analysis

- Main Idea- Comparison of MapReduce to Parallel DBMS

- MapReduce
  - Sorts the data into different sections, then filters through this data
  - Difficult to implement, but cost effective
  - 2-3 times slower in query execution when compared to Parallel

- Parallel DBMS
  - Has existed for over 20 years
  - Easy to execute the code, but it is expensive to do so
  - Uploading data takes longer

# Implementation of Approaches to Large Scale Data Analysis

- MapReduce is typically implemented by:
  - Importing the MapReduce function into the database
  - Applying it to code

- Essentially splitting up larger databases into smaller subsets and then checking each of those for the desired data

# Analysis of Implementation

- Based on the amount of time that it took for MapReduce to complete the tasks of
  - Load times
  - Cluster data
  - Selection
  - Aggregation …

- ….it is clear that MapReduce took substantially more time to complete each task compared to other systems

- Popular because it is cheap
  - Other data analysis techniques should be considered just based on how long it takes to complete a task

# Comparison of the two papers

**Bigtable**

- Discusses the system that Google uses currently to store information for several of it's projects

- Implementation is efficient, finding and sorting data is made simple through the use of column families, row keys, and timestamps

**Approaches to Large Scale Data Analysis**

- Compares MapReduce and Parallel DBMS
  - Both do essentially the same thing, one does it faster and more expensive and other is slow but cheap

- Implementation is inefficient, it is complicated and takes much longer than necessary to perform tasks

# Main idea of Stonebraker talk

- During the 1980s and 1990s, people believed in the methodology of one size fits all as far as DBMSs go

- However…one type of database engine does not meet the needs of all different kinds of databases – one size does not fit all

- There are a huge number of different database engines that serve different purposes

- Traditional row stores are basically obsolete and are the epitome of 'one size fits none', they are not good in many different types of markets
  - Data Warehouse, OLTP, NoSQL, etc.

# Advantages and Disadvantages of the Main Idea of Bigtable in the Context of Comparison Paper and Stonebraker

**Advantages**

- Cost efficient data sorting and location by both MapReduce and Bigtable

- Not a 'one size fits all' database engine

- Bigtable is adaptable based on the self-managing basis of their servers

**Disadvantages**

- DBMS researchers may focus on changing Bigtable rather than creating new DBMSs

- Because Bigtable was created by Google, it doesn't allow for innovation of smaller companies/people