

Scalable Batch Processing in the Cloud

Keijo Heljanko

Department of Information and Computer Science
School of Science
Aalto University
`keijo.heljanko@aalto.fi`

28.6-2011

Business Drivers of Clouds

- ▶ Large data centers allow for economics of scale
 - ▶ Cheaper hardware purchases
 - ▶ Cheaper cooling of hardware
 - ▶ Example: Google paid 40 MEur for a Summa paper mill site in Hamina, Finland: Data center cooled with sea water from the Baltic Sea
 - ▶ Cheaper electricity
 - ▶ Cheaper network capacity
 - ▶ Smaller number of administrators / computer
- ▶ Unreliable commodity hardware is used
- ▶ Reliability obtained by replication of hardware components and a combined with a **fault tolerant software stack**

Warehouse Scale Computing

- ▶ Thus the smallest unit of computation in Google scale is:
[Warehouse full of computers](#)

- ▶ Luiz André Barroso, Urs Hölzle: *The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines* Morgan & Claypool Publishers 2009

<http://dx.doi.org/10.2200/S00193ED1V01Y200905CAC006>

- ▶ The book says:
“... we must treat the datacenter itself as one massive warehouse-scale computer (WSC).”

Cloud Computing

A collection of technologies aimed to provide elastic “pay as you go” computing

- ▶ **Virtualization of computing resources**: Amazon EC2, Eucalyptus, OpenNebula, Open Stack Compute, ...
- ▶ **Scalable datastore**: Amazon S3, Google Bigtable, HBase, Amazon Dynamo, Apache Cassandra, ...
- ▶ **Scalable Web Applications hosting**: Google App Engine, Microsoft Azure, Heroku, ...
- ▶ **Scalable File Storage**: GFS, HDFS, ...
- ▶ **Scalable batch processing**: Google MapReduce / Apache Hadoop, PACT, Microsoft Dryad, Google Pregel, ...

Big Data

- ▶ As of May 2009, the amount of digital content in the world is estimated to be 500 Exabytes (500 million TB)
- ▶ EMC sponsored study by IDC in 2007 estimates the amount of information created in 2010 to be 988 EB
- ▶ Worldwide estimated hard disk sales in 2010:
≈ 675 million units
- ▶ Data comes from: Video, digital images, sensor data, biological data, Internet sites, social media, . . .
- ▶ The problem of such large data massed, termed **Big Data** calls for new approaches to storage and processing of data

Example: Simple Word Search

- ▶ Example: Suppose you need to search for a word in a 2TB worth of text that is found only once in the text mass using a compute cluster
- ▶ Assuming 100MB/s read speed, in the worst case reading all data from a single 2TB disk takes ≈ 5.5 hours
- ▶ If 100 hard disks can be used in parallel, the same task takes less than four minutes
- ▶ Scaling using **hard disk parallelism** is one of the design goals of scalable batch processing in the cloud

Scaling Up vs Scaling Out

- ▶ **Scaling up:** When the need for parallelism arises, a single powerful computer is added with more CPU cores, more memory, and more hard disks
- ▶ **Scaling out:** When the need for parallelism arises, the task is divided between a large number of less powerful machines with (relatively) slow CPUs, moderate memory amounts, moderate hard disk counts

Pros and Cons of Scaling Up vs Scaling Out

- ▶ **Scaling up is more expensive than scaling out.** Big high-end systems have much higher pricing for a given: CPU power, memory, and hard disk space
- ▶ **Scaling out is more challenging for fault tolerance.** A large number of loosely coupled systems means more failures in hardware and in networking. **Solution: Software fault tolerance**
- ▶ **Scaling out is more challenging for software development** due to larger number of components, larger number of failures both in nodes and networking connecting them, and increased latencies. **Solution: Scalable software development frameworks**

Google MapReduce

- ▶ A scalable batch processing framework developed at Google for computing the Web index
- ▶ The MapReduce framework takes care of all issues related to parallelization, synchronization, load balancing, and fault tolerance. All these details are hidden from the programmer
- ▶ When deciding whether MapReduce is the correct fit for an algorithm, one has to remember the fixed data-flow pattern of MapReduce. The algorithm has to be efficiently mapped to this data-flow pattern in order to efficiently use the underlying computing hardware

MapReduce and Functional Programming

- ▶ Based on the **functional programming in the large**:
 - ▶ User is only allowed to write **side-effect free functions** “**Map**” and “**Reduce**”
 - ▶ **Re-execution is used for fault tolerance**. Side effects in functions would make this impossible
 - ▶ The functions themselves are usually written in a standard imperative programming language such as Java or C++

Why No Side-Effects?

- ▶ Side-effect free programs will produce the same output irregardless of the number of computing nodes used
- ▶ Running the code on one machine for debugging purposes produces the same results as running the same code in parallel
- ▶ It is easy to introduce side-effect to MapReduce programs as the framework does not enforce a strict programming methodology. However, the **behavior of such programs is undefined** by the framework, and should therefore be avoided.

Map and Reduce Functions

- ▶ The framework only allows a user to write two functions: a “**Map**” function and a “**Reduce**” function
- ▶ The **Map**-function is fed blocks of data (blocksize 64-128MB), and it produces (key, value) -pairs
- ▶ The framework groups all values with the same key to a (key, (... , list of values, ...)) format, and these are then fed to the **Reduce** function

MapReduce Diagram

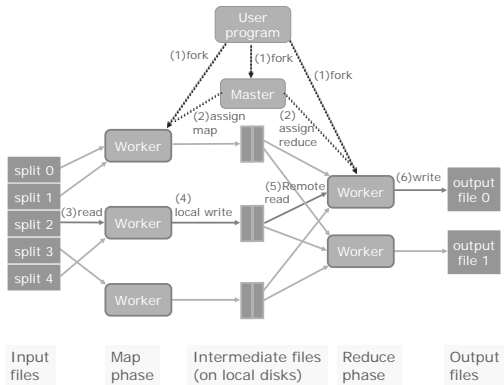


Figure: J. Dean and S. Ghemawat: *MapReduce: Simplified Data Processing on Large Clusters*, OSDI 2004

Example: Word Count

- ▶ Classic word count example from the Hadoop MapReduce tutorial:

`http://hadoop.apache.org/common/docs/current/mapred_tutorial.html`

- ▶ Consider doing a word count of the following file using MapReduce:

```
Hello World Bye World  
Hello Hadoop Goodbye Hadoop
```

Example: Word Count (cnt.)

- ▶ Consider a Map function that reads in words one a time, and outputs (word, 1) for each parsed input word
- ▶ The Map function output is:

(Hello, 1)

(World, 1)

(Bye, 1)

(World, 1)

(Hello, 1)

(Hadoop, 1)

(Goodbye, 1)

(Hadoop, 1)

Example: Word Count (cnt.)

- ▶ The Shuffle phase between Map and Reduce phase creates a list of values associated with each key
- ▶ The Reduce function input is:

(Bye, (1))

(Goodbye, (1))

(Hadoop, (1, 1))

(Hello, (1, 1))

(World, (1, 1))

Example: Word Count (cnt.)

- ▶ Consider a reduce function that sums the numbers in the list for each key and outputs (`word`, `count`) pairs. The output of the Reduce function is the output of the MapReduce job:

(Bye, 1)

(Goodbye, 1)

(Hadoop, 2)

(Hello, 2)

(World, 2)

Phases of MapReduce

1. A **Master** (In Hadoop terminology: Job Tracker) is started that coordinates the execution of a MapReduce job. Note: **Master is a single point of failure**
2. The master creates a predefined number of **M Map workers**, and assigns each one an input split to work on. It also later starts a predefined number of **R reduce workers**
3. Input is assigned to a free Map worker 64-128MB split at a time, and each **user defined Map function** is fed (key, value) pairs as input and also produces (key, value) pairs

Phases of MapReduce(cnt.)

4. Periodically the Map workers flush their (key, value) pairs to the local hard disks, partitioning by their key to R partitions (default: use hashing), one per reduce worker
5. When all the input splits have been processed, a Shuffle phase starts where $M \times R$ file transfers are used to send all of the mapper outputs to the reducer handling each key partition. After reducer receives the input files, reducer sorts the pairs by the key
6. User defined Reduce functions iterate over the (key, (... , list of values, ...)) lists, generating output (key, value) pairs files, one per reducer

Google MapReduce (cnt.)

- ▶ The user just supplies the Map and Reduce functions, nothing more
- ▶ The framework can be used to implement a **distributed sorting algorithm** by using a custom partitioning function
- ▶ The framework does **automatic parallelization and fault tolerance** by using a centralized Job tracker (Master) and a distributed filesystem that stores all data redundantly on compute nodes
- ▶ Uses **functional programming paradigm** to guarantee correctness of parallelization and to implement fault-tolerance by re-execution

Apache Hadoop

- ▶ An Open Source implementation of the MapReduce framework, originally developed by Doug Cutting and heavily used by e.g., Yahoo! and Facebook
- ▶ “Moving Computation is Cheaper than Moving Data” - Ship code to data, not data to code.
- ▶ Map and Reduce workers are also storage nodes for the underlying distributed filesystem: Job allocation is first tried to a node having a copy of the data, and if that fails, then to a node in the same rack (to maximize network bandwidth)
- ▶ Project Web page: <http://hadoop.apache.org/>

Apache Hadoop (cnt.)

- ▶ Builds reliable systems out of unreliable commodity hardware by replicating most components (exceptions: Master/Job Tracker and NameNode in Hadoop Distributed File System)
- ▶ Tuned for large (gigabytes of data) files
- ▶ Designed for very large 1 PB+ data sets
- ▶ Designed for streaming data accesses in batch processing, designed for high bandwidth instead of low latency
- ▶ For scalability **NOT a POSIX filesystem**
- ▶ Written in Java, runs as a set of userspace daemons

Hadoop Distributed Filesystem

- ▶ **A distributed replicated filesystem:** All data is replicated by default on three different Data Nodes
- ▶ Inspired by the Google Filesystem
- ▶ Each node is usually a Linux compute node with a small number of hard disks (4-12)
- ▶ A single NameNode that maintains the file locations, many DataNodes (1000+)

Hadoop Distributed Filesystem (cnt.)

- ▶ Any piece of data is available if at least one datanode replica is up and running
- ▶ Rack optimized: by default one replica written locally, second in the same rack, and a third replica in another rack (to combat against rack failures, e.g., rack switch or rack power feed)
- ▶ Uses large block size, 128 MB is a common default - designed for batch processing
- ▶ For scalability: **Write once, read many filesystem**

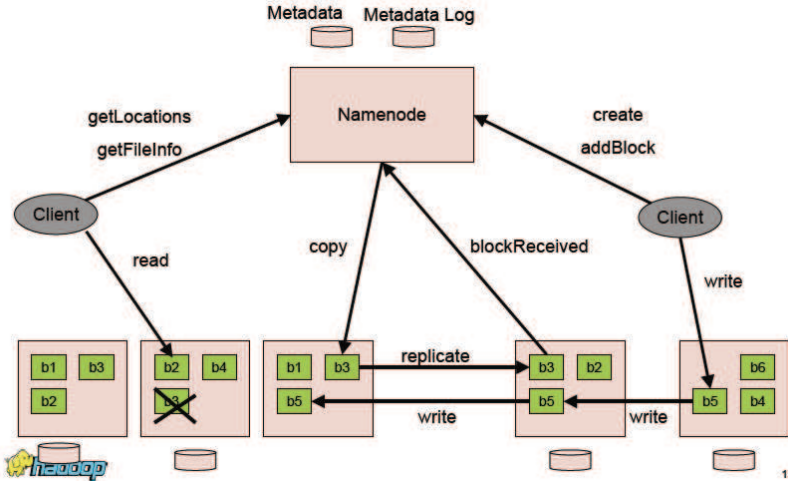
Implications of Write Once

- ▶ All applications need to be re-engineered to only do sequential writes. Example systems working on top of HDFS:
 - ▶ HBase (Hadoop Database), a database system with only sequential writes, Google Bigtable clone
 - ▶ MapReduce batch processing system
 - ▶ Apache Pig and Hive data mining tools
 - ▶ Mahout machine learning libraries
 - ▶ Lucene and Solr full text search

HDFS Architecture

- ▶ From: HDFS Under The Hood by Sanjay Radia of Yahoo

<http://assets.en.oreilly.com/1/event/12/HDFS%20Under%20the%20Hood%20Presentation%201.pdf>



HDFS Architecture

- ▶ NameNode is a single computer that maintains the namespace (meta-data) of the filesystem. Implementation detail: Keeps all meta-data in memory, writes logs, and does periodic snapshots to the disk
- ▶ All data accesses are done directly to the DataNodes
- ▶ Replica writes are done in a daisy chained (pipelined) fashion to maximize network utilization

HDFS Scalability Limits

- ▶ 20PB+ deployed HDFS installations (10 000+ hard disks)
- ▶ 4000+ DataNodes
- ▶ Single NameNode scalability limits: The HDFS is NameNode scalability limited for write only workloads to around HDFS 10 000 clients, [K. V. Shvachko: HDFS scalability: the limits to growth:](http://www.usenix.org/publications/login/2010-04/openpdfs/shvachko.pdf)
<http://www.usenix.org/publications/login/2010-04/openpdfs/shvachko.pdf>
- ▶ Currently a distributed NameNode design is being implemented to address scalability (not fault tolerance) issues

Hadoop Hardware

- ▶ Reasonable CPU speed, reasonable RAM amounts for each node
- ▶ 4-12 hard disks per node seem to be the current suggestion
- ▶ CPU speeds are growing faster than hard disk speeds, so newer installations are moving to more hard disks / node
- ▶ Gigabit Ethernet networking seems to be dominant

Hadoop Network Bandwidth Consideration

- ▶ Hadoop is fairly network latency insensitive
- ▶ Mapper reads can often be read from the local disk or the same rack (intra-rack network bandwidth is cheaper than inter-rack bandwidth)
- ▶ For jobs with only small Map output, very little network bandwidth is used
- ▶ For jobs with large Map output, Hadoop likes large inter-node network bandwidth (Shuffle phase),
- ▶ To save network bandwidth, Mappers should produce a minimum amount of output

Terasort Benchmark

- ▶ Hadoop won the Terasort benchmark in 2009 by sorting 100 TB in 173 minutes using: 3452 nodes x (2 Quadcore Xeons, 8 GB memory, 4 SATA disks/node)

Two Large Hadoop Installations

- ▶ Yahoo! (2009): 4000 nodes, 16 PB raw disk, 64TB RAM, 32K cores
- ▶ Facebook (2010): 2000 nodes, 21 PB storage, 64TB RAM, 22.4K cores
 - ▶ 12 TB (compressed) data added per day, 800TB (compressed) data scanned per day
 - ▶ A. Thusoo, Z. Shao, S. Anthony, D. Borthakur, N. Jain, J. Sen Sarma, R. Murthy, H. Liu: *Data warehousing and analytics infrastructure at Facebook*. SIGMOD Conference 2010: 1013-1020.
<http://doi.acm.org/10.1145/1807167.1807278>

Cloud Software Project

- ▶ ICT SHOK Program Cloud Software (2010-2013)
 - ▶ A large Finnish consortium, see:
<http://www.cloudsoftwareprogram.org/>
 - ▶ Case study at Aalto: CSC Genome Browser Cloud Backend
 - ▶ Co-authors: [Matti Niemenmaa](#), [André Schumacher](#) (Aalto University, Department of Information and Computer Science), [Aleksi Kallio](#), [Eija Korpelainen](#), [Taavi Hupponen](#), and [Petri Klemelä](#) (CSC — IT Center for Science)

CSC Genome Browser

- ▶ CSC provides tools and infrastructure for bioinformatics
- ▶ Bioinformatics is the largest customer group of CSC (in user numbers)
- ▶ Next-Generation Sequencing (NGS) produces large data sets (TB+)
- ▶ Cloud computing can be harnessed for analyzing these data sets
- ▶ 1000 Genomes project (<http://www.1000genomes.org>): Freely available 50 TB data set of human genomes

CSC Genome Browser

- ▶ Finland has limited resources to produce data
- ▶ Analysis software systems is area where Finland has potential for global impact

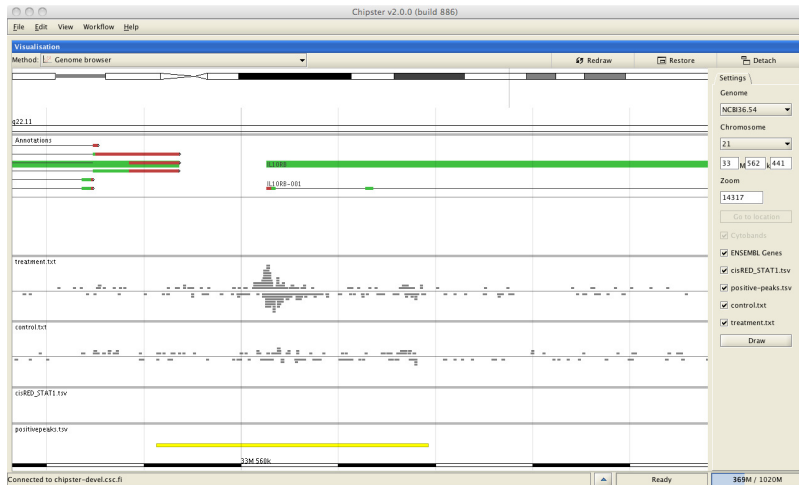
CSC Genome Browser

- ▶ Cloud computing technologies will enable scalable NGS data analysis
- ▶ There exists prior research into sequence alignment and assembly in the cloud
- ▶ Visualization is needed in order to understand large data masses
- ▶ Interactive visualization for 100GB+ datasets can only be achieved with preprocessing in the cloud

Genome Browser Requirements

- ▶ Interactive browsing with zooming in and out, “Google Earth”-style
- ▶ Single datasets 100GB-1TB+ with interactive visualization at different zoom levels
 - ▶ Preprocessing used to compute summary data for the higher zoom levels
 - ▶ Dataset too large to compute the summary data in real time using the real dataset
 - ▶ Scalable cloud data processing paradigm map-reduce implemented in Hadoop used to compute the summary data in preprocessing (currently upto $15 \times 12 = 180$ cores)

Genome Browser GUI

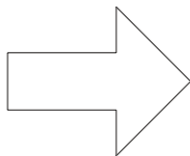


Read aggregation: problem

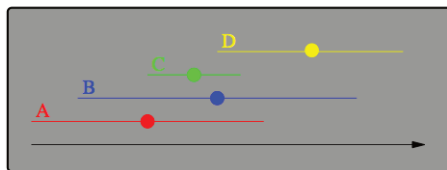
- ▶ A common problem in DNA sequence analysis is the *alignment* of a large number of smaller subsequences to a common *reference sequence*
- ▶ Once the alignment has been determined, one needs to analyze how well the reference sequence is *covered* by the subsequences (a.k.a. *reads*)
- ▶ For interactive visualization the large data set has to be summarized
- ▶ Cloud computing enables interactive visualization of sequencing data

Read aggregation: problem (cnt.)

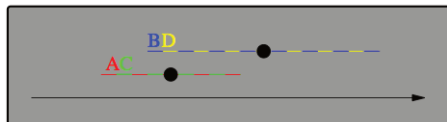
A	0	10	N	GTAC
B	2	14	N	AACG
C	5	9	N	TTAG
D	8	16	N	GCAG



N	3	9	2
N	5	15	2

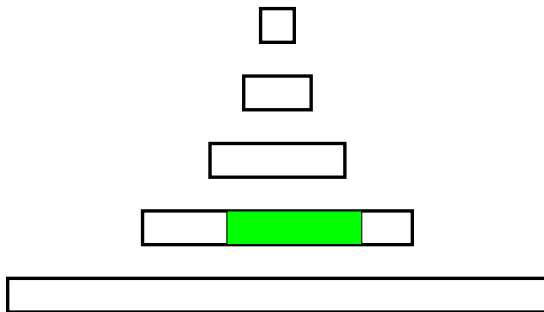


Input (BAM)

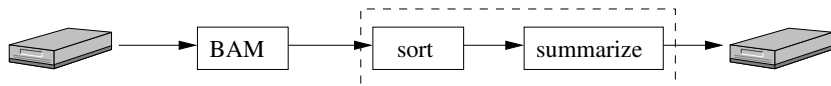


Output (Summary)

Summary Files



Read aggregation via Hadoop



- ▶ Center point of each read computed (*Map*)
- ▶ Reads sorted according to that center point
- ▶ For fixed *summary-block size* B , every B reads are merged into a single aggregated read (*Reduce*)
- ▶ Note: in the previous example we had $B = 2$

Triton Cloud Testbed

- ▶ Cloud computing testbed at Aalto University
- ▶ 112 AMD Opteron 2.6 GHz compute nodes with 12 cores, 32-64GB memory each, totalling 1344 cores
- ▶ Infiniband and 1 GBit Ethernet networks
- ▶ 30 TB+ local disk
- ▶ 40 TB+ fileserver work space



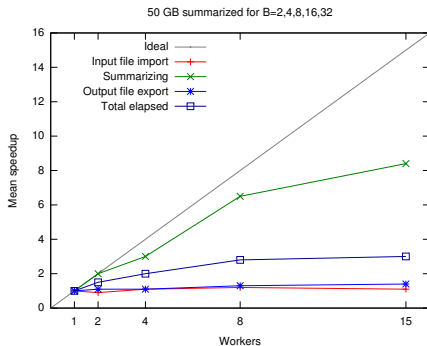
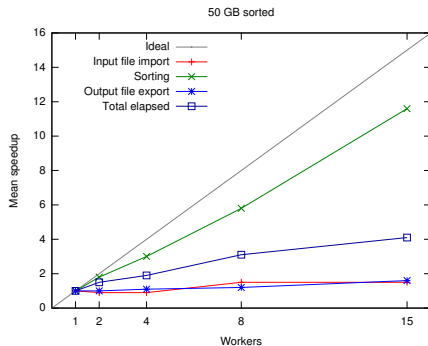
Triton Cloud Testbed Fileserver



Experiments

- ▶ One 50 GB (compressed) BAM input file from *1000 Genomes*
- ▶ Run on the Triton cluster 1-15 compute nodes with 12 cores each
- ▶ Four repetitions for increasing number of worker nodes
- ▶ Two types of runs: sorting according to starting position (“sorted”) and read aggregation using five summary-block sizes at once (“summarized”)

Mean speedup



Results

- ▶ An updated CSC Genome browser GUI
- ▶ Chipster: Tools and framework for producing visualizable data
- ▶ An implementation of preprocessing for visualization using Hadoop
- ▶ Scalability studies for running Hadoop on 50 GB+ datasets on the Triton cloud testbed
- ▶ Software released:
<http://sourceforge.net/projects/hadoop-bam/>

Future plans

- ▶ Further implementation of data analysis for the cloud
- ▶ Continued development of the genome browser GUI and bioinformatics tools
- ▶ Prestudy comparing the Genome Analysis Toolkit

Current Research Topics

- ▶ Aalto and CSC both have datacenters which can be used as testbeds for cloud computing technologies
- ▶ Focus on cloud based data analysis for “Big Data”
- ▶ MapReduce (Hadoop) scalable batch processing technologies in the cloud
- ▶ Scalable datastores: HBase (Hadoop Database) has been evaluated, also other cloud based datastores such as Cassandra are of interest