

Yannick Deuster

06.02.2014

Main Questions

What is Hadoop?

How does it work?

Who uses Hadoop?

Is it any good?

How do you get started?

What is Hadoop?



WIKIPEDIA The Free Encyclopedia "Apache Hadoop is an open-source software framework for storage and large scale processing of data-sets on clusters of commodity hardware."

Java Framework

"Big Data"

Large clusters

Cheap hardware

History

- Doug Cutting and Mike Cafarella in 2002
- Nutch
- Google File System (2003)
- MapReduce (2004)
- Cutting hired by Yahoo!
- Hadoop now independent from Nutch
- Latest release: October 15th, 2013



- Hadoop
 Common
- Utilities that support other modules

- Hadoop Distributed File System
- YARN
 - Manages and coordinates computing resources
- Hadoop MapReduce
- YARN based system for parallel processing

Extensions



How does it work?



How does it work?

- NameNode and DataNodes part of HDFS
- JobTracker and TaskTrackers for MapReduce



HDFS

• Designed for very large amounts of data

• High throughput

• Load balancing

• Fault tolerant

HDFS



MapReduce

- Programming model by Google
- Parallel computing
- Two phases
 - Map
 - Reduce



Hadoop MapReduce

• One Mapper for each data block (64MB)

• JobTracker moves computation close to data

• Mappers sort output

• Streaming API

MapReduce

mapper.py

```
1 #!/usr/bin/env python
2
 3 import sys
 4
   for line in sys.stdin:
 5
 6
           # remove leading whitespaces
7
8
9
           line = line.strip()
           words = line.split()
10
11
           for word in words:
12
                    print '%s\t%s' % (word, 1)
```

MapReduce

```
reducer.py
```

```
1 #!/usr/bin/env python
3 from operator import itemgetter
 4 import sys
6 current_word = None
 7 current_count = 0
 8 word = None
 Q
10 for line in sys.stdin:
       line = line.strip()
11
12
13
      word, count = line.split('\t', 1)
14
15
      try:
16
           count = int(count)
17
       except ValueError:
18
           continue
19
20
       if current word == word:
21
           current_count += count
22
       else:
23
           if current_word:
24
                   print '%s\t%s' % (current_word, current_count)
25
           current_count = count
26
           current_word = word
27
28 if current_word == word:
       print '%s\t%s' % (current_word, current_count)(word, 1)
29
```

DEMO

Who uses Hadoop? Coose III. Age











Google

	e e e e e e e e e e e e e e e e e e e

"Research purposes"

",University Initiative to Addres Internet-Scale Computing Challenges"



Builds Amazons product search indices

532 nodes cluster

search optimization



690 nodes cluster

7500+ daily Hadoop jobs

scheduler "Luigi"

Listening behavior in Sweden





http://files.meetup.com/5139282/SHUG%201%20-%20Hadoop%20at%20Spotify.pdf

Who uses Hadoop? YAHOO!

- more than 40.000 computers running Hadoop
- biggest cluster: 4500 nodes
- support reasearch for "Ad Systems" and web search

- Two clusters
 - 1. 1100 nodes with 8800 cores
 - 2. 300 nodes with 2400 cores
- June 13, 2012: 100 PB of data
- grows by ~0.5 PB each day
- store internal data
- analytics and machine learning



MySQL, MySQL, MySQL

twitter

- We all start there.
- But MySQL is not built for analysis.
- select count(*) from users? Maybe.
- select count(*) from tweets? Uh...
- Imagine joining them.
- And grouping.
- Then sorting.

http://www.slideshare.net/kevinweil/hadoop-and-pig-attwitter-oscon-2010-4824988



MapReduce Workflow



 Challenge: how many tweets per user, given tweets table?

twitter >>

- Input: key=row, value=tweet info
- Map: output key=user_id, value=1
- Shuffle: sort by user_id
- Reduce: for each user_id, sum
- Output: user_id, tweet count
- With 2x machines, runs 2x faster

http://www.slideshare.net/kevinweil/hadoop-and-pig-attwitter-oscon-2010-4824988

Main Questions

What is Hadoop? ✓

How does it work? \checkmark

Who uses Hadoop? ✓

Is it any good?

How do you get started?

Is it any good?

Advantages	Disadvantages	
 Build for cheap hardware 	 Difficult to set up 	
 Extremely Scalable 	 JobTracker and NameNode single points of failure 	
 Fault tolerant 	 HDFS cannot be mounted directly 	
 High throughput file system 	 Security disabled by default 	

Scientific Computing

- Main uses: Search, Text analysis, Machine learning
- No "real" math libraries
- MapReduce math algorithms not well studied
- Weak: Single computations
- Excels: Big Data analysis (Weather, DNA)

Hadoop vs. Disco

- Disco is easier to set up
- Hadoop "heavyweight" in comparison
- Disco fast on small data
- Hadoop superior on large data

Wordcount: 1 byte file		Wordcou	Wordcount: English Wikipedia	
Disco	359 ms	Disco	~22 mins	
Hadoop	12324 ms	Hadoop	~18 mins	

http://www.erlang-factory.com/upload/presentations/778/ef2013-disco.pdf

Awards

- MediaGuardian Innovator of the Year 2011
 - beat WikiLeaks and iPad

- Winner "Terabyte Sort Benchmark" 2013
 - 102,5 TB of data in 4328 seconds (~72mins)
 - 1,42 TB/min
 - 2100 nodes, 4200x2,3Ghz

How do you get started?

- What Hadoop needs
 - JRE 1.6 or higher
 - SSH access
- Download from http://hadoop.apache.org
- Archive: <u>http://archive.apache.org/dist/hadoop/</u>

How do you get started?

Official Documentation

http://hadoop.apache.org/docs/

Helpful Tutorials

http://www.michael-noll.com/tutorials/

A good book

Hadoop - The Definitive Guide