### Apache Hadoop

### Large scale data processing



Speaker: Isabel Drost



### **Isabel Drost**

### Nighttime:

### Came to nutch in 2004. Co-Founder Apache Mahout. Organizer of Berlin Hadoop Get Together. Daytime:

Software developer @ Berlin

### Hello Information Retrieval course!

# Agenda

• Motivation.

• A short tour of Map Reduce.

• Introduction to Hadoop.

• Hadoop ecosystem.



January 8, 2008 by Pink Sherbet Photography http://www.flickr.com/photos/pinksherbet/2177961471/

Illin

# Massive data as in:

### Cannot be stored on single machine. Takes too long to process in serial.

Idea: Use multiple machines.

### Challenges.

#### ............ ................ ...... 0000000 \*\*\*\*\*\*\*\* UIIII ....... end to fail ngle mach nes PowerEdge 111111 S Hard dis 2650 ..... 0 ar S

11111

0 po 2650

# More machines – increased failure probability.

January 11, 2007, skreuzer http://www.flickr.com/photos/skreuzer/354316053/

10

63

## Requirements

- Built-in backup.
- Built-in failover.

# **Typical developer**



- Has never dealt with large (petabytes) amount of data.
- Has no thorough understanding of parallel programming.
- Has no time to make software production ready.

### Tynical developer

Failure resistant: What if service X is unavailable? Failover built in: Hardware failure does happen. Documented logging: Understand message w/o code. Monitoring: Which parameters indicate system's health? Automated deployment: How long to bring up machines? Backup: Where do backups go to, how to do restore? Scaling: What if load or amount of data double, triple? Many, many more.



Heoroughunde.ling ofparallertramming.

with

 Has no time o make software production ready.

# Requirements

- Built-in backup.
- Built-in failover.

- Easy to use.
- Parallel on rails.

http://www.flickr.com/photos/jaaronfarr/3384940437/ March 25, 2009 by jaaron

February 29, 2008 by Thomas Claveirole http://www.flickr.com/photos/thomasclaveirole/2300932656/

http://www.flickr.com/photos /jaaronfarr/35357564827 March 25, 2009 by jaaron

May 1, 2007 by danny angus http://www.flickr.com/photos/killerbees/479864437/

**Open source developers** 

Developers interested in large scale applications

**Open source developers** 

Developers interested in large scale applications

Java developers

Open source developers

## Requirements

- Built-in backup.
- Built-in failover.

- Easy to use.
- Parallel on rails.

• Java based.



http://www.flickr.com/photos/cspowers/282944734/ by cspowers on October 29, 2006

# Requirements

- Built-in backup.
- Built-in failover.

- Easy to use.
- Parallel on rails.

- Easy to administrate. Java based.
- Single system.

### We need a solution that:

Is easy to use.

Scales well beyond one node.

Java based implementation.

Easy distributed programming.

Well known in industry and research.

Scales well beyond 1000 nodes.



- 2008:
  - 70 hours runtime
  - 300 TB shuffling
  - 200 TB output
- In 2009
  - 73 hours
  - 490 TB shuffling
  - 280 TB output
  - 55%+ hardware
  - 2k CPUs (40% faster cpus)

- 2008
  - 2000 nodes
  - 6 PB raw disk
  - 16 TB RAM
  - 16k CPUs
- In 2009
  - 4000 nodes
  - 16 PB disk
  - 64 TB RAM
  - 32k CPUs (40% faster cpus)

### Example use cases

- Distributed Grep.
- Distributed Sort.
- Link-graph traversal.
- Term-Vector per host.
- Web access log stats.

- Inverted index.
- Doc clustering.
- Machine learning.
- Machine translation.

### Some history.

Feb '03 first Map Reduce library @ Google

Oct '03 GFS Paper

Dec '04 Map Reduce paper

Dec '05 Doug reports that nutch uses map reduce

Feb '06 Hadoop moves out of nutch

Apr '07 Y! running Hadoop on 1000 node cluster

Jan '08 Hadoop made an Apache Top Level Project

### Hadoop assumptions

### Assumptions:

Data to process does not fit on one node. Each node is commodity hardware. Failure happens.







### Ideas:

Distribute filesystem. Built in replication. Automatic failover in case of failure.



Moving data is expensive. Moving computation is cheap. Distributed computation is easy.



Move computation to data. Write software that is easy to distribute.



### Assumptions:

Systems run on spinning hard disks. Disk seek >> disk scan.



### Ideas:

Improve support for large files.
File system API makes scanning easy.

### Hadoop by example

?xml version="1.0" encoding="UTF-8"?>

<copml version="1.0" >

<head>

<text></text>

</head>

⊲body>

dutline htmlUrl="http://eventseer.net" title="EventSeer - A Digital Library of Call for Papers" useC alDefault" version="RSS" type="rss" xmlUrl="http://eventseer.net/feeds/main/rss.xml" id="312053548" tex tseer.net" />

<outline isOpen="false" id="669809145" text="Silent" >

<outline htmlUrl="http://www.theserverside.com" title="TheServerSide.com: Patterns" useCustomFetchIn ersion="RSS" type="rss" xmlUrl="http://www.theserverside.com/rss/theserverside-j2eepatterns-rss2.xml" i taining up-to-date news, discussions, patterns, resources, and media" />

doutline htmlUrl="http://chadwa.wordpress.com" title="Chad's Search Blog" useCustomFetchInterval="fa
S" type="rss" xmlUrl="http://chadwa.wordpress.com/feed/" id="545368194" text="Chad's Search Blog" descr
" />

dutline htmlUrl="http://www.find23.net/Site/Blog/Blog.html" title="My Blog" useCustomFetchInterval=
"RSS" type="rss" xmlUrl="http://www.find23.net/Site/Blog/rss.xml" id="1620106192" text="My Blog" description")

doutline htmlUrl="http://emotion.inrialpes.fr/~dangauthier/blog" title="Yet Another Machine Learning
eMode="globalDefault" version="RSS" type="rss" xmlUrl="http://emotion.inrialpes.fr/~dangauthier/blog/fe
g" />

<outline htmlUrl="http://ml.typepad.com/machine\_learning\_thoughts/" title="Machine Learning Thoughts
="globalDefault" version="RSS" type="rss" xmlUrl="http://ml.typepad.com/machine\_learning\_thoughts/rss.xmlurl="http://ml.typepad.com/machine\_

doutline htmlUrl="http://yaroslavvb.blogspot.com/" title="Machine Learning, etc" useCustomFetchInter ion="RSS" type="rss" xmlUrl="http://yaroslavvb.blogspot.com/feeds/posts/default" id="805998569" text="M doutline htmlUrl="http://ptufts.blogspot.com/" title="Pinhead's Progress" useCustomFetchInterval="fa S" type="rss" xmlUrl="http://ptufts.blogspot.com/" title="Misc Research Stuff" useCustomFetchInterval="fa doutline htmlUrl="http://ptufts.blogspot.com/" title="Misc Research Stuff" useCustomFetchInterval on="RSS" type="rss" xmlUrl="http://resnotebook.blogspot.com/" title="Misc Research Stuff" useCustomFetchInterval on="RSS" type="rss" xmlUrl="http://resnotebook.blogspot.com/" title="Absolutely Regular" useCustomFetchInterval doutline htmlUrl="http://absolutely-regular.blogspot.com/" title="Absolutely Regular" useCustomFetch version="RSS" type="rss" xmlUrl="http://absolutely-regular.blogspot.com/" title="Absolutely Regular" useCustomFetch version="RSS" type="rss" xmlUrl="http://absolutely-regular.blogspot.com/" title="Absolutely Regular" useCustomFetch version="RSS" type="rss" xmlUrl="http://absolutely-regular.blogspot.com/feeds/posts/default" id="17850! doutline htmlUrl="http://atomai.blogspot.com/" title="Data Mining, Analytics and Artificial Intellige Mode="globalDefault" version="RSS" type="rss" xmlUrl="http://atomai.blogspot.com/feeds/posts/default" in nt in data mining, artificial intelligence, analytics, intelligent agents, semiconductors, distributing siness Objects, Oracle, Intel, AMD, or Pentaho. Heuristic, Six Sigma, or CMM. Contractor or in-house. H ail\_com" /> isabel@h1349259:~\$ more data/feeds.opml | grep -o "http://[0-9A-Za-z\-\_\.]\*" | s

- ort | uniq --count | sort | tail
  - 3 http://agbs.kyb.tuebingen.mpg.de
  - 3 http://irgupf.com
  - 3 http://jeffsutherland.com
  - 4 http://ml.typepad.com
  - 4 http://weblogs.java.net
  - 4 http://www.gridvm.org
  - 4 http://yaroslavvb.blogspot.com
  - 5 http://feeds.feedburner.com
  - 6 http://blogsearch.google.com
  - 10 http://arxiv.org

#### pattern="http://[0-9A-Za-z\-\_\.]\*"

grep -o "\$pattern" feeds.opml | sort | uniq --count



```
pattern="http://[0-9A-Za-z\-_\.]*"
```

```
grep -o "$pattern" feeds.opml
M A P
```

| uniq --count | R E D U C E




Local to data.



Local to data. Outputs a lot less data. Output can cheaply move.



Local to data. Outputs a lot less data. Output can cheaply move.



Local to data. Outputs a lot less data. Output can cheaply move. Shuffle sorts input by key. Reduces output significantly.

```
private IntWritable one = new IntWritable(1);
private Text hostname = new Text();
```

```
public void map(LongWritable key, Text value,
OutputCollector<Text, IntWritable> output,
Reporter reporter) throws IOException {
   String line = value.toString();
   StringTokenizer tokenizer = new StringTokenizer(line);
   while (tokenizer.hasMoreTokens()) {
      hostname.set(getHostname(tokenizer.nextToken()));
      output.collect(hostname, one);
   }
}
```

```
public void reduce(Text key, Iterator<IntWritable>
values, OutputCollector<Text, IntWritable> output,
Reporter reporter) throws IOException {
    int sum = 0;
    while (values.hasNext()) {
        sum += values.next().get();
    }
    output.collect(key, new IntWritable(sum));
}
```



## Petabyte sorting benchmark

Bytes	Nodes	Maps	Reduces	Replication	Time
500,000,000,000	1406	8000	2600	1	59 seconds
1,000,000,000,000	1460	8000	2700	1	62 seconds
100,000,000,000,000	3452	190,000	10,000	2	173 minutes
1,000,000,000,000,000	3658	80,000	20,000	2	975 minutes

Per node: 2 quad core Xeons @ 2.5ghz, 4 SATA disks, 8G RAM (upgraded to

16GB before petabyte sort), 1 gigabit ethernet.

Per Rack: 40 nodes, 8 gigabit ethernet uplinks.





#### **Petabyte Task Timeline**



Minutes

# What was left out

- Combiners compact map output.
- Language choice: Java vs. Dumbo vs. PIG ...
- Size of input files does matter.
- Facilities for chaining jobs.
- Logging facilities.
- Monitoring.
- Job tuning (number of mappers and reducers)

### Hadoop ecosystem.

### Higher level languages.

# Cascading









Suppose you have user data in one file, website data in another, and you need to find the top 5 most visited pages by users aged 18 - 25.





----------ировал ибранийский стала ибранова. Бако та с ули сталайский ислаг ибранова. Бако ули настро ( рано съ справутско с. съ. ст. перекото и съ. съ. ст. станата съ. ст. станата спорт на р станата ст. ст. ст. ст. ст. ст. п. ст. ст. ст. ст. ст. ст. ст. ст. на ст. ст. ст. ст. ст. ст. ст. ст. ст. на станица и последна у станита и стани од станита и представата и станита и последна и представата и станита и последна и представата и станита и представата и представата и станита и представата и представата и станита и представата и представата и представата и представата и станита и представата и представата и представата и представата и станита и представата и станита и представата и пр party and to your contract where we are a first or the party one appropriate to the loss COLORADO IN 2 March 100 100 Party and a state of the state of the state of the рато сталоров од, власти и посталоров од, ворато сталоров и сталоров и посталоров сталоров и сталоров и посталоров сталоров и сталоров и сталоров и посталоров сталоров и сталоров и сталоров и A 494 A An approximate the manufacture of the second se second sec

NAMES OF ADDRESS OF A D

на со оказарање на отоке со оказа на развијата се на стран на развијата се на стран на развијата се на стран на се оказарање на се на се на се пора се оказарање на . .

Party and a state state and so the state state of

parts see app 32 G. наралите слатина спротокото из. Паралите правлят слатина парада на 1 нариски практики слова настрана ( нариски практики слова настрана) на слована и при самар ( на слована и при самар ( на слована) и при самар ( на слов par parts in the constant value of a market.
 we have a new parts.
 we have a new parts.

рано и и о на велото на продока. Прока водото, спристо, телотория с. телото с

party one many - то са ната на продата на правита на правита на правита на правита на продата на правита на п на правита на на правита на правит на правита н на правита на пр на правита на на правита н на правита н на правита на правита на правит where reserves the state of the

алектору, не свутелериет

рато и с. о учал салоточа услова просоходо. По просока партите постарата с. техно, салоточа о. 22.1

the second secon аланаа (радинаалан, раску) мараландарынаа, калаландарын ( мараландарынаалан, раску)

bare an example and the second second second

на лан. - н рати становани наукатани оду-каранитани оду-каранитани каранани, окатан аранатаранан каранан аранат

water between the second to the second secon

рано протива разврати стала получи и при стала стала получи стала получи получи стала получи стала получи получи стала получи стала получи стала получи получи стала получи стала получи стала получи получи стала получи стала получи стала получи стала получи получи стала получи стала получи стала получи стала получи получи стала получи стала получи стала получи стала получи стала получи получи стала получи стала получи стала получи стала получи стала получи стала получи получи стала п

 A set of the set of na ay 1 али салание - во серено разко салание - во серено разко салание разначите са сърза -разко салание сумание са сърза -разко салание суманието са сърза разко саланието суманието са сърза разко саланието са сърза -на сърза - сърза - сърза -на сърза - сърза - сърза - сърза -на сърза - сърза - сърза - сърза -на сърза - сърза - сърза - сърза - сърза -на сърза - сърза - сърза - сърза - сърза -на сърза - сърза - сърза - сърза - сърза -сърза - сърза - сърза - сърза - сърза - сърза -сърза - сърза - сърза - сърза - сърза - сърза -сърза - сърза - сърза - сърза - сърза - сърза -сърза - сърза - сърза - сърза - сърза - сърза -сърза - сърза - сърза - сърза - сърза - сърза -сърза - сърза - сърза - сърза - сърза - сърза -сърза - сърза -сърза - сърза - съ -----BALL MADE разантира и служарания разродно служарния разродно служа разродно служарния разродно служа разродно служарния разродно служарния разродно служарния разродно служарния разродно служарния разродно служа разродно служарния разродно служарния разродно служа разродно служа разродно служа разродно служа разродно служа разродно служарни CARD COMMON COLUMN AND A DESCRIPTION OF A DESCRIPTIO учария водот водот на россите со сало на селото со селото со селото со селото горина, на селото со селото селото со селото горина, на селото се wages and you say in the андан ал арагана араг APR - 2000 and a product of the second second second In the second second 1.000

Example from PIG presentation at Apache Con EU 2009



```
Users = load 'users' as (name, age);
Fltrd = filter Users by
        age >= 18 and age <= 25;
Pages = load 'pages' as (user, url);
Jnd = join Fltrd by name, Pages by user;
Grpd = group Jnd by url;
Smmd = foreach Grpd generate group,
        COUNT(Jnd) as clicks;
Srtd = order Smmd by clicks desc;
Top5 = limit Srtd 5;
store Top5 into 'top5sites';
```

Example from PIG presentation at Apache Con EU 2009





```
{publisher: 'Scholastic',
author: 'J. K. Rowling',
title: 'Deathly Hallows',
year: 2007},
{publisher: 'Scholastic',
author: 'J. K. Rowling',
title: 'Chamber of Secrets',
year: 1999,
reviews: [
   {rating: 10, user: 'joe', review: 'The best ...'},
   {rating: 6, user: 'mary', review: 'Average ...'}]},
{publisher: 'Scholastic',
author: 'J. K. Rowling',
title: 'Sorcerers Stone',
year: 1998},
{publisher: 'Scholastic',
author: 'R. L. Stine',
title: 'Monster Blood IV',
year: 1997,
reviews: [
   {rating: 8, user: 'rob', review: 'High on my list...'},
  {rating: 2, user: 'mike', review: 'Not worth the paper ...',
   discussion:
      [{user: 'ben', text: 'This is too harsh...'},
       {user: 'jill', text: 'I agree ...'}]}]},
{publisher: 'Grosset',
author: 'Carolyn Keene',
title: 'The Secret of Kane',
year: 1930}
```

1

Example from JAQL documentation.



```
// Query 2. Find the authors and titles of books that have received
// a review.
for( $b in hdfsRead('books') )
  if( exists($b.reviews) )
    [{ $b.author, $b.title }];
// result...
[
    {author: 'J. K. Rowling', title: 'Chamber of Secrets'},
    {author: 'R. L. Stine', title: 'Monster Blood IV'}
];
```

Example from JAQL documentation.

### (Distributed) storage.



#### Project Voldemort A distributed database

#### **About Dynomite**

Dynomite is an eventually consistent d Amazon's Dynamo paper. Dynomite cu









### Libraries built on top.





















### Alternative approaches.











#### Get involved!

#### Solving hard problems?

Solving hard problems? Communicating your solution?

Solving hard problems? Communicating your solution? Working with excellent teams?


#### Do you love:

### Solving hard problems? Communicating your solution? Working with excellent teams?

Picture by: July 9, 2006 by trackrecord, http://www.flickr.com/photos/trackrecord/185514449



Source control system.

Continuous integration.

Test-f rst development.

Issue-tracker.



Create readable patches.

Communicate and discuss solutions.

Review others code.

Work in large, distributed teams.



# How?

- First time users:
  - Documentation in wiki.
- Experimenting:
  - Write examples.

- Found a bug:
  - Go to JIRA, file a bug.
  - Describe the bug.
  - Create a test to show.
  - Provide a patch.

- Evaluating:
  - Test performance.
  - Provide comparison.
- Participate on-list.
  - Answer questions.
  - Discuss vour use-

## Recipe to Apache

- Download the release and use it.
- Subscribe to the mailing-list.
- Questions:
  - Documentation: Wiki.
  - Discussions: Mailing list.
  - Current status: JIRA.
  - History: JIRA for patches, mailing-list for votes.
- Checkout the code and built it.



\*-user@lucene.apache.org\*-dev@lucene.apache.org

Love for solving hard problems. Interest in production ready code. Interest in parallel systems.

July 9, 2006 by trackrecord http://www.flickr.com/photos/trackrecord/18551449 Documentation, code, examples.

Contact Ross Gardler for more information on Apache at universities worldwide.

Message view	
From	Grant Ingersoll <gsing@apache.org></gsing@apache.org>
Subject	Re: Lucene Branding: the TLP, and "Lucene Java"
Date	Wed, 11 Apr 2007 01:13:36 GMT

No, you are not the only one... Many a sleepless night spent on it... :-)

I usually try to refer to it as Lucene Java, but old minist die hard and often times I just call it Lucene. I think the name has a good brand at this point and is very strongly associated w/ the Java library. I seem to recall when they were forming the TLP, that the original proposal was search.a.o, but then changed b/c the ASF didn't like generic names (or at least that is how I recall it.) And, of course, with Hadoop and the potential for Tika/Lius, it isn't just search anymore. I have often thought about an Apache "Text" project, that could eventually hold a whole family of text based tools like Lucene, Tika, Hadoop, Solr, etc. plus things like part of speech taggers, clustering/classification algorithms, UIMA, etc. all under one roof. But that is just my two cents and I don't know if it fits with what other people have in mind. There are a lot of OSS tools out there for these things, but none bring together a whole suite under a brand like Apache.

-Grant

## Why go for Apache?

#### Jumpstart your project with proven code.

January 8, 2008 by dreizehn28 http://www.flickr.com/photos/1328/2176949559

## Discuss ideas and problems online.

November 16, 2005 [phil h] http://www.flickr.com/photos/hi-phi/64055296









### Become part of the community.







