

Data driven applications

based on Apache projects

Image by Studio Tempura - http://www.flickr.com/photos/zero101/3335373387



Isabel Drost

Daytime: Co-Founder Berlin Buzzwords. Software developer at Nokia Gate 5 GmbH.

Nighttime:

Co-Founder Apache Mahout. Member Apache Software Foundation. Founder of Berlin Hadoop Get Together.

Hello GoTo Con visitors.

Apache Hadoop

- Know it?
 - Use it?
- How many nodes?
 - Contributed patches?



Apache Lucene/ Solr



Apache Mahout





Amazon EMR

By: mondays child - http://www.flickr.com/photos/normanlowery/4729806370/

0

Agenda

• Introducing example • Types of data.

Data driven design.
Toolset.

A highly representative example



BUZZWORD:

ALCHALM.



Image by Zaprittsky: http://www.flickr.com/photos/zaprittsky/4520788183/

EOPLE

WITH ASTISTS

POCEMA

Ways to uncover requirements?





On-site customer

(Product owner in Scrum)





(Switch the oracle against metrics.)

Data-based approach

("Get out of the building")



Learn from past users' behaviour



Users seldom give explicit feedback.

Take a snapshot of user interactions.

Analyze what they do.

Derive needs for new functionality from that.

Your search did not match any documents...

Searching for:



date when none is indexed. project pig instead of the animal. typing AMS but mean Amsterdam. typing confrence but mean conference.

Discover content without effort by ...

Showing content that I might like based on:



my personal interests and tags. talks I already indicated as "want to attend. abstracts and speaker info pages I read. my prior search history and twitter messages. Observe how your application is used









By Lab2112, http://www.flickr.com/photos/lab2112/462388595/

Section .

111

OWEISED

1.33

11111

(awa

ama

ANI CA

(and







Intel 63 0 munting 00 uffer(): 86 Idev Che For Your 88 86 tected ickung . unlacked iphone ... hacking King filesystems... IIIe King Held King 80 81 88 80 81 88 king 101 skast plus volume. d a case-sensitive d a case-sensitive d Extents Overflow HFS Kung volum C en o D A Catalog hecking Catalog Indis king Extende 10 BIE VOLUM 00 volu Intsh. . ing ø HFS plus LittleBear AR182.UserBundle a case-sensi 88 **Overflow** Extents 16. IK. and Attributes Catalog hierarchy . HI IN EC B data 83 88 80 valume 11 catalog. file. file. HINB Flow file. catalog. FILe.

B





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

Called Call

E STITTE

Ser .

Illing

en du mon

e A

CERVE WITH PRIDE

aller

w

- CH

"a million Letters

Rissint Harrister Market and the state of th

Homb CORM HEB HA, OH31

D Sharon Pruitt

2003

1e -

Hugisto -








Discover content without effort by ...

Showing content that I might like based on:



my personal interests and tags. talks I already indicated as "want to attend. abstracts and speaker info pages I read. my prior search history and twitter messages.





33%

Image by zigazou76: http://www.flickr.com/photos/zigazou76/3622235298

2





Image by Garrett Rooney: http://www.flickr.com/photos/rooneg/276527837/





Analysis tool-set



Structured relational data

- Attendee database.
- Speaker information.
- User transactions.



Extracting information from log files

- Access logs.
- Health check results.
- Response time logs.



What about...

- Learning what users look for from search logs?
 - Do they query by room name?
 - ... by speaker name?
 - ... by tag?
- Learning what content to show from clicks?
 - Show different site per user?
 - Generate links between talks automatically?

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

Illin

Peter Norvig (paraphrase): Rather than argue about whether this algorithm is better than that algorithm, all you have to do is get ten times more training data. And now all of a sudden, the worst algorithm ... is performing better than the best algorithm on less training data.

013

UTH

Worry about the data first before you worry about the algorithm.

USA FIRST-CLASS FOREVE

Single machines tend to fail redae Harebook 11111

-

.....

201

20

...

11

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.

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/3385756482/ March 25, 2009 by jaaron

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





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

Easy distributed programming.

Well known in industry and research.

Scales well beyond 1000 nodes.



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

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.

Assumptions:

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



Ideas:

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.





(Graphics: Thanks to Thilo.)



?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" />

dutline 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" />

dutline 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_learning_thoughts/rss.xmlurl="http://ml.typepad.com/machine_learning_thoughts/rss.xmlurl="http://ml.typepad.com/machine_learning_thoughts/rss.xmlurl="http://ml.typepad.com/machine_learning_thoughts/" title="Machine_learning_thoughts/rss.xmlurl="http://ml.typepad.com/machine_learning_thoughts/rss.xmlu

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 soutline htmlUrl="http://resnotebook.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	sort	uniqcount
MAP	SHUFFLE	IREDUCE











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

```
public void map(K key, V value, Context context) {
   String line = value.toString();
   StringTokenizer tokenizer = new StringTokenizer(line);
   while (tokenizer.hasMoreTokens()) {
      hostname.set(getHostname(tokenizer.nextToken()));
      context.write(hostname, one);
   }
}
```

```
public void reduce(K2 key, Iterable<V2> values,
    OutputCollector<K3, V2> output) {
    int sum = 0;
    while (values.hasNext()) {
        sum += values.next().get();
        }
        output.collect(key, new IntWritable(sum));
    }
```

Anatomy of a map/reduce job



Slide inspired by: "Hadoop - The definitive guide", Tom White, O'Reilly

Anatomy of a map/reduce job



Slide inspired by: "Hadoop – The definitive guide", Tom White, O'Reilly
Anatomy of a map/reduce job



Slide inspired by: "Hadoop – The definitive guide", Tom White, O'Reilly

Anatomy of a map/reduce job



Slide inspired by: "Hadoop – The definitive guide", Tom White, O'Reilly

Anatomy of a map/reduce job



Slide inspired by: "Hadoop – The definitive guide", Tom White, O'Reilly

Properties of good map reduce citizens:

do not crash on increasing amount of input data.

do not overload external resources.



self contained.







Welcome to Apache HCatalog!

Hadoop Eclipse Plug-in

PigPen

PigPen is an eclipse plugin that helps users ci

About MRUnit

MRUnit is a Java library that helps devel







Adds to Hadoop:

Random access to your data

Near-Realtime read/write

Host very large tables (billions of rows, millions of columns)











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.



Example from PIG presentation at Apache Con EU 2009



----------- NO ировал ибранийские орологие и орологие Бало и стороди сталийские истористорие Бало улистика и сталийские и орологие и 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 станита и представата и представата 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 представата и станита и представата и и представата и представ party and to your contract where we are a first of pare de oppageere d'en de COLORADO IN 2 March 100 100 Party and a state of the state of the state of the рато час коло род. чод. в лика и поло род. чод. в лика и поло в лика. в рако и пора и поло и поло во р на род. и пора и поло во рако и поло и поло и поло и поло на род. и пора и поло и поло и поло и поло и поло и поло на род. и поло чито и поло и поло и поло и поло и поло на род. и поло чито и поло и поло и поло и поло и поло на род. и поло чито и поло на род. И поло чито и поло и A 494 A An approximate the manufacture of the second se second sec

NAMES OF A DESCRIPTION OF A

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 (настрана) и практики (настрана) и практи 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 recently a sub-are a recent by the matrix of the recent of the the matrix of the sub-

алектору, во свутелно разл

рато и с. о учал салоточа услова просоходо. По просока партите постарата с. текато, салоточато. 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

ware been in the second to the second second

 A state of the sta na ay 1 A MARKA - BA MAIN -----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 TARKS AND THE REAL PROPERTY AND андан ал арагана араг APR - 2000 and a product of the second second second 1.000

In the second second

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































Trigger the workflow:

Manually On data arrival Based on a schedule

Avoid job conflicts.





Show users content they like





Show users content they like



Monitor media and social media



Show users content they like



Monitor media and social media



Deliver ads targeted to users

















facebook / scribe Source Commits Network Pull F Switch Branches (4) + Switch Tags (0) Brand Cloudera / flume Source Commits Network Switch Branches (12) + Switch Tags (14) +

<u>chukwz</u>











(Thanks to Thilo for helping set up the cluster, Thanks to packet and masq for two of the three machines.)



```
package { 'openssh-server':
   ensure => installed,
file { '/etc/ssh/sshd config':
   source => 'puppet://modules/sshd/
sshd config',
   owner => 'root',
   group => 'root',
   mode => '640',
   notify => Service['sshd'], # sshd
              will restart whenever you
              edit this file.
   require => Package['openssh-server'],
service { 'sshd':
   ensure => running,
   enable => true,
   hasstatus => true,
   hasrestart => true,
```



- Managing:
 - local files
 - packages
 - Services
- Configuring:
 - cron
 - users and groups







 Introducing Puppet
 How to Buy Documentation



- Started as distributed synchronization.
- Configuration management.
- Provides for distributed naming.
- Group services.





Jumpstart your project with proven code.

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

There is help



O'REILLY" Lars George

There is help

- Cloudera
- Hortonworks
- MapR
- JTeam

- 101tec
- Datameer
- Sematext
- •

Discuss ideas and problems online.

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

WARNING



http://flickr.com/photos/iamdietrich/4075363845/ By: Iam Dietrich

10.2/0.4

May 1, 2007 by danny angus http://www.flickr.com/photos/killerbees/479864437/ *user@*.apache.org
dev@.apache.org



Image by: Patrick McEvoy

Interest in solving hard problems. Being part of lively community. Engineering best practices.

Bug reports, patches, features. Documentation, code, examples.