# Hadoop & Pig

Dr. Karina Hauser
Senior Lecturer
Management & Entrepreneurship
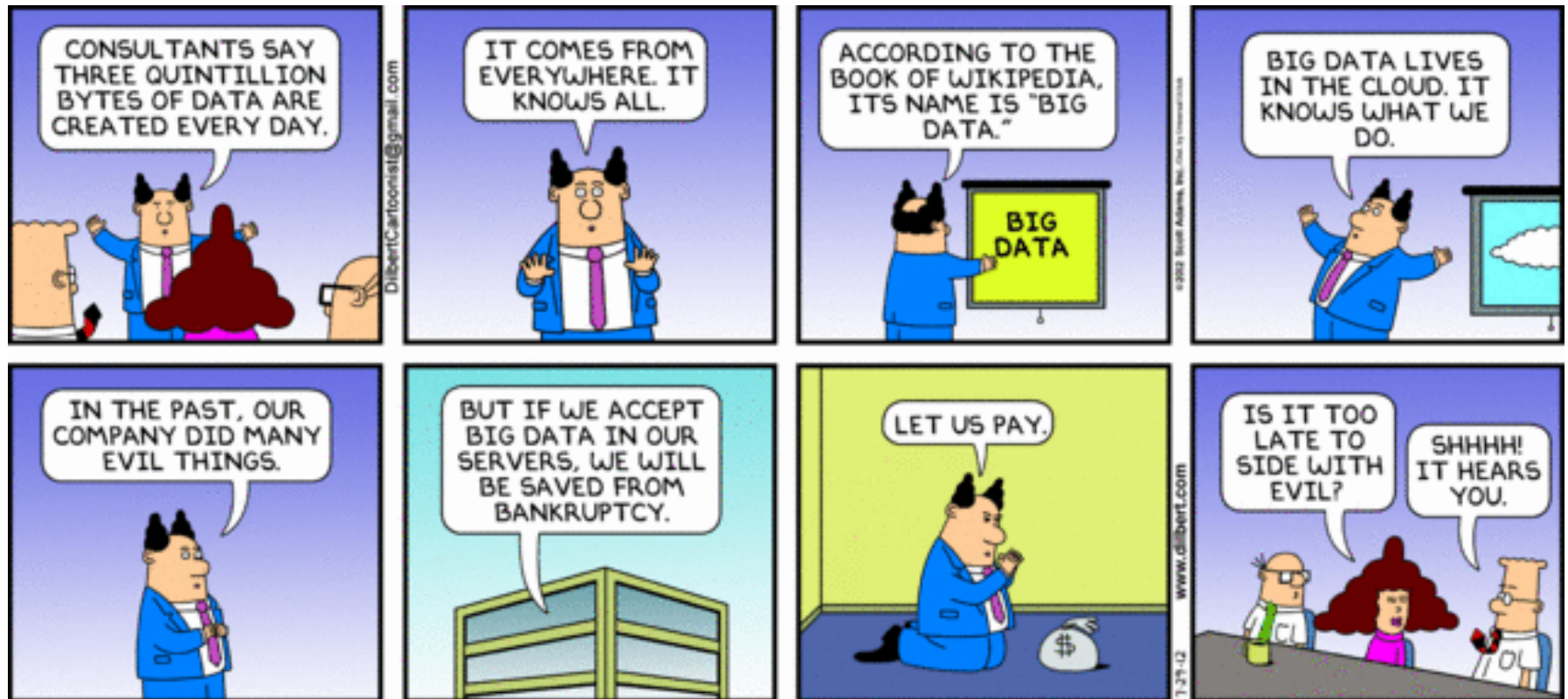
# Outline

- Introduction

- (Setup)

- Hadoop, HDFS and MapReduce

- Pig

# Introduction

- **What is Hadoop and where did it come from?**

# Big Data

# Big Data Sources

- **Every day 2.5 quintillion bytes or 2.5 exabytes ($10^{18}$) are generated, that number is estimated to double every 40 month**
- **Astronomy**
  - Sloan Digital Sky Survey (SDSS) began collecting astronomical data in 2000; 200 GB ($10^9$) per night
  - Large Synoptic Survey Telescope (LSST) (~2020); estimated ~20 TB ($10^{12}$) per night
- **Business**
  - Twitter: 12 terabytes ($10^{12}$) of Tweets every day
  - Walmart: 2.5 petabytes ($10^{15}$) of data every hour from its customer transactions

# Big Data - Big Business

- **IDC predicts big data technology and services will grow worldwide from $3.2 billion in 2010 to $16.9 billion in 2015. This represents a compound annual growth rate of 40 percent — about seven times that of the overall information and communications technology market.**

Leeds School of Business
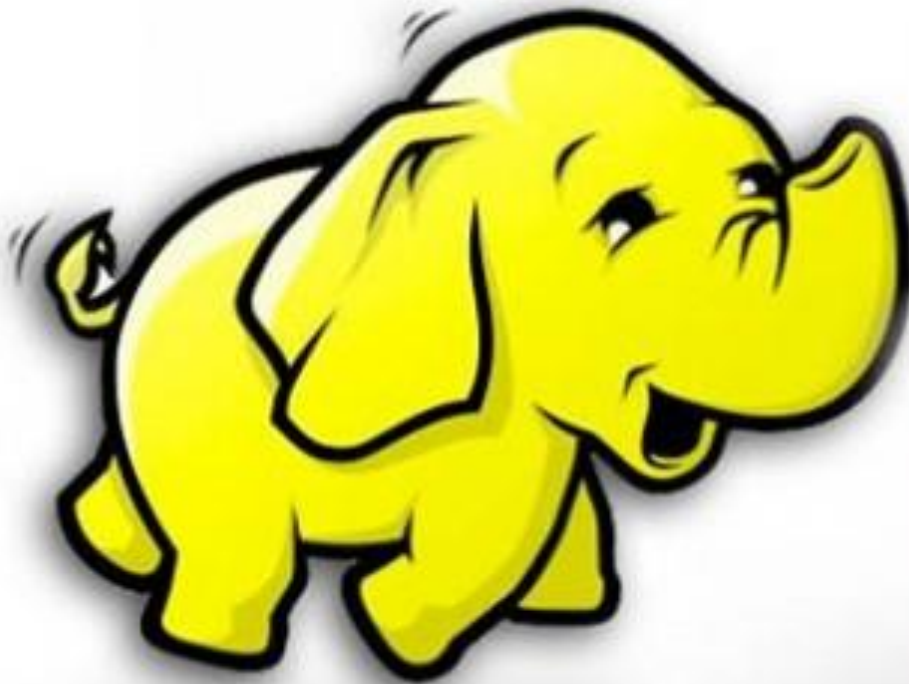UNIVERSITY OF COLORADO BOULDER

# Big Data in the News

- **The Economist Intelligence Unit study showed that nine out of 10 surveyed business leaders believe data is now the fourth factor of production, as fundamental to business as land, labor and capital.**

Leeds School of Business
UNIVERSITY OF COLORADO BOULDER

# What is Big Data ?

- **"When the data itself becomes part of the problem"**
- **Three (to five) dimensions:**
  - Volume
  - Variety
  - Velocity
  - (Veracity)
  - (Value)

# The Solution

# Short History

- **Created by Doug Cutting, named after his son's toy elephant**
- **2002 - Nutch, search engine, scalability problems**
- **2004 - Google papers on GFS and MapReduce**
- **2006 - Yahoo hires Doug to improve Hadoop**
- **2008 - Hadoop becomes Apache Top Level Project**

# Hadoop Today

- **Moving from "Internet" companies**
  - Yahoo
  - Google
- **to business and science applications**
  - Customer relationship management
  - Bioinformatics
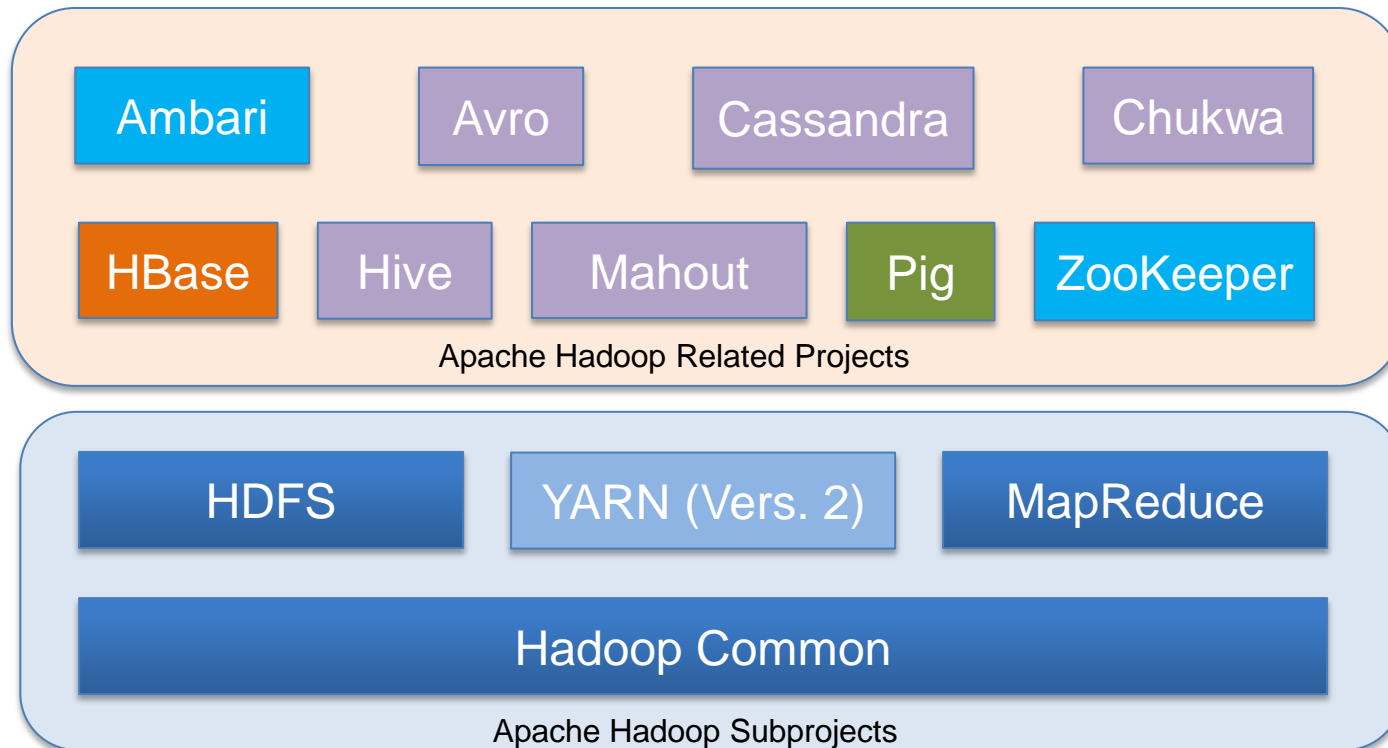  - Astrophysics

# Hadoop Definition

- **"Framework that allows for the**
  - distributed processing of
  - large data sets
  - across clusters of computers
  - using a simple programming model"
- **Open-source software, maintained by "[The Apache Software Foundation](#)"**
- **[http://hadoop.apache.org/](http://hadoop.apache.org/)**
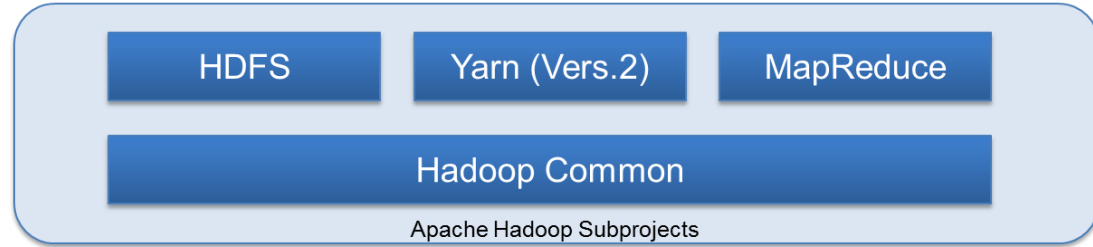
# Apache Hadoop Projects

# Hadoop Cluster



- **Commodity hardware**
- **Individual disk space on each node**
- **Hadoop framework handles:**
  - Data "backups" (through replication)
  - Hardware failure
  - Parallelization of code (through MapReduce paradigm)

# HDFS



HDFS | Yarn (Vers.2) | MapReduce
Hadoop Common
Apache Hadoop Subprojects

- **Write-once, read-many**

- **Each file is stored as sequence of same-sized blocks (default size 64MB)**

- **Blocks are replicated across different nodes**

- **Highly reliable:**
  – Redundant data storage
  – Heartbeat messages to detect connectivity problems
  → Automatic failover

**Leeds** School of Business
UNIVERSITY OF COLORADO BOULDER

Pipelining

# MapReduce



- **Programing model designed for**
  - batch processing of large volumes of data
  - in parallel
  - by dividing the work into a set of independent tasks
- **Not limited to Hadoop**

# MapReduce WordCount Example



The overall MapReduce word count process

# Problems suited for MapReduce

- **Iterate over a large number of records**
- **Extract something of interest from each**
- **Shuffle and sort intermediate results**
- **Aggregate intermediate results**
- **Generate final output**

# Hadoop 1.x Components

# Hadoop 1.x Components

- ## **Name Node**
  - Stores metadata (filenames, replications factors …)
    - ***fsimage: latest checkpoint of namespace***
    - ***edits: log of changes to namespace***
  - Checks data node availability (Heartbeat)
  - If possible: Separate machine
- ## **Data Node**
  - Stores data
  - Replicates blocks
  - Computation

# Hadoop 1.x Components

- **Secondary/Checkpoint Name Node**
  - Periodically creates checkpoints of namespace
    - *Downloads fsimage and edits, creates new namespace and updates Name Node*
  - Separate machine
- **Job Tracker**
  - Schedules and manages jobs
- **Task Tracker**
  - Executes MapReduce jobs on individual data node

# Hadoop 2.0

# Setup

- **Three options:**
  - Standalone (single Java process)
  - **Pseudo-Distributed (separate Java processes)**
  - Fully-Distributed
- **Prerequisites (on virtual machine):**
  - Ubuntu server 12.04
  - Ubuntu desktop (for monitoring)
  - Oracle (Sun) Java 1.7.0_25
    - *http://www.webupd8.org/2012/01/install-oracle-java-jdk-7-in-ubuntu-via.html*
  - SSH

# Setup Files

- **hosts (ip address)**
- **.bashrc (Java dir, home dir)**
- **hadoop configuration files in /usr/local/hadoop/conf**
  - hadoop-env.sh (Java dir)
  - hdfs-site.xml (replication factor)
  - mapred-site.xml (host/port for jobtracker)
  - master/slave (ips for multi-node cluster)

Leeds School of Business
UNIVERSITY OF COLORADO BOULDER

# Login

- **User: rmacc**
- **Password: rmacc**
- **Start terminal (Ctrl+Alt+F1)**
- **Login as hduser**
  - User: hduser
  - Password: hduser
- **Change directory to hadoop**
  - with cd $HADOOP_PREFIX
  - or cd /usr/local/hadoop

# Starting Hadoop Daemons

- **All:**

    $ bin/start-all.sh

- Log output in *logs* directory

# Check Daemons

- **jps** →
  - 1367 Jps
  - 8695 DataNode
  - 8609 NameNode
  - 6318 SecondaryNameNode
  - 2600 JobTracker
  - 2830 TaskTracker

# Web Interfaces

- http://localhost:50030 → Cluster status and jobs

- http://localhost:50070 → HDFS

# HDFS

- **Hadoop Distributed File System**

# Ubuntu (Linux)

# Hadoop

bin/hadoop  dfs  -copyFromLocal mis6110/Files/file.txt  mis6110/KDD/file.txt

bin/hadoop  dfs  -copyToLocal  mis6110/Movies/file.txt  mis6110/Files/file.txt

local file system
/usr/local/hadoop/...

(HDFS)
/user/hduser/...

mis6110

mis6110

Files          Scripts

KDD          Movies

| | | | |
|---|---|---|---|
| Create Directory | *mkdir* | | *bin/hadoop dfs -mkdir* |
| Show content of directory | *ls* | | *bin/hadoop dfs -ls* |
| Create Directory | *sh mkdir* | | *fs -mkdir* |
| Show content of directory | *sh ls* | | *fs -ls* |

# HDFS Shell

- **bin/hadoop dfsadmin –help → all admin commands**

- **bin/hadoop dfs –help → all commands**

- **Most commands similar to unix**

  - dfs –copyFromLocal

  - dfs  -ls

- Shell commands: http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/FileSystemShell.html

# Importing Data

- **KDD Example:** Legcare sales data from [KDD Cup 2000](#)
  *"We wish to thank Blue Martini Software for contributing the KDD Cup 2000 data"*

  - Cleaned for easier/faster use

  - Copy file KDDCupCleaned.txt to hdfs

# Copying File to HDFS

- **Create new directory:**

  $ bin/hadoop  dfs  -mkdir rmacc

- **Copy files:**

  $ bin/hadoop  dfs  -copyFromLocal ../rmacc/Beowulf.txt rmacc/

  $ bin/hadoop  dfs  -copyFromLocal ../rmacc/KDD* rmacc/

- **Check:**

  *$ bin/hadoop dfs   -lsr  rmacc*

  - *Localhost:50070 →Browse the filesystem*

# MapReduce WordCount Example



The overall MapReduce word count process

From: http://www.rabidgremlin.com/data20/#(3)

# Java Code for WordCount Example

```java
1.      package org.myorg;
2.
3.      import java.io.IOException;
4.      import java.util.*;
5.
6.      import org.apache.hadoop.fs.Path;
7.      import org.apache.hadoop.conf.*;
8.      import org.apache.hadoop.io.*;
9.      import org.apache.hadoop.mapred.*;
10.     import org.apache.hadoop.util.*;
11.
12.     public class WordCount {
13.
14.        public static class Map extends MapReduceBase implements Mapper<LongWritable, Text, Text, IntWritable> {
15.          private final static IntWritable one = new IntWritable(1);
16.          private Text word = new Text();
17.
18.          public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
19.            String line = value.toString();
20.            StringTokenizer tokenizer = new StringTokenizer(line);
21.            while (tokenizer.hasMoreTokens()) {
22.              word.set(tokenizer.nextToken());
23.              output.collect(word, one);
24.            }25.            }
26.        }
27.
28.        public static class Reduce extends MapReduceBase implements Reducer<Text, IntWritable, Text, IntWritable> {
```

# Java Code for WordCount Example

```
29.          public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
30.            int sum = 0;
31.            while (values.hasNext()) {
32.              sum += values.next().get();
33.            }
34.            output.collect(key, new IntWritable(sum));
35.          }
36.        }
37.
38.        public static void main(String[] args) throws Exception {
39.          JobConf conf = new JobConf(WordCount.class);
40.          conf.setJobName("wordcount");
41.
42.          conf.setOutputKeyClass(Text.class);
43.          conf.setOutputValueClass(IntWritable.class);
44.
45.          conf.setMapperClass(Map.class);
46.          conf.setCombinerClass(Reduce.class);
47.          conf.setReducerClass(Reduce.class);
48.
49.          conf.setInputFormat(TextInputFormat.class);
50.          conf.setOutputFormat(TextOutputFormat.class);
51.
52.          FileInputFormat.setInputPaths(conf, new Path(args[0]));
53.          FileOutputFormat.setOutputPath(conf, new Path(args[1]));
54.
55.          JobClient.runJob(conf);
57.        }
58.
```

**Leeds** School of Business
UNIVERSITY OF COLORADO BOULDER

# There Must be an Easier Way

# Pig (Latin)

# Pig

- **High-level data processing language (Pig Latin)**
- **Resides on user machine, not cluster**
- **Pig Latin compiled into efficient MapReduce jobs**

# Test Pig Installation

- **(Hadoop has to be running)**
  - $ cd $PIG_HOME
  - $ bin/pig  -help

# How to Run Pig

- **Grunt interactive shell**
  - Two modes:
    - *Local, standalone (pig –x local)*
    - *Hadoop, distributed (pig –x mapreduce)*
- **Scripts (.pig)**
- **Embedded in Java or Python**
- **PigPen, Eclipse plugin**

# Pig Statements in Grunt

- **LOAD → Transform data → DUMP or STORE**

- **Example:**

  - *grunt> A = load 'student' using PigStorage()*
    *AS (name:chararray, age:int, gpa:float);*

  - *grunt> B = foreach A generate name;*

  - *grunt> dump B*

  – *"A " is called a "relation" or "outer bag"*

# Pig Example: KDD Data

| Key | Date | Time | Unit Price | Order LineID | Qty | Order Status | Tax | Amount | Weekday | Hour | City | State | Customer ID |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2/27/2000 | 08\:06\:35 | 15 | 1 | 15 | Shipped | 1.3 | 16.3 | Sunday | 8 | Westport | CT | 62 |
| 2 | 3/30/2000 | 10\:00\:18 | 9 | 1 | 9 | Shipped | 0 | 9 | Thursday | 10 | Westport | CT | 62 |
| 3 | 1/28/2000 | 14\:43\:34 | 12 | 1 | 12 | Shipped | 1.02 | 13.02 | Friday | 14 | San Francisco | CA | 96 |
| 4 | 1/29/2000 | 10\:22\:37 | 12 | 1 | 12 | Shipped | 0.87 | 12.87 | Saturday | 10 | Novato | CA | 132 |
| 5 | 2/1/2000 | 08\:44\:48 | 6.5 | 1 | 6.5 | Shipped | 0.55 | 7.05 | Tuesday | 8 | Cupertino | CA | 168 |
| 6 | 2/29/2000 | 10\:31\:42 | 15 | 1 | 15 | Shipped | 1.24 | 16.24 | Tuesday | 10 | San Ramon | CA | 184 |
| 7 | 2/29/2000 | 10\:31\:42 | 14 | 1 | 14 | Shipped | 1.16 | 15.16 | Tuesday | 10 | San Ramon | CA | 184 |
| 8 | 2/29/2000 | 10\:31\:42 | 6.5 | 2 | 6.5 | Shipped | 1.07 | 14.07 | Tuesday | 10 | San Ramon | CA | 184 |
| 9 | 3/8/2000 | 16\:48\:47 | 11 | 3 | 11 | Shipped | 2.72 | 35.72 | Wednesday | 16 | San Ramon | CA | 184 |
| 10 | 1/30/2000 | 14\:13\:57 | 10 | 1 | 10 | Shipped | 0 | 10 | Sunday | 14 | Scarsdale | NY | 224 |
| 11 | 1/30/2000 | 14\:13\:57 | 13.5 | 1 | 13.5 | Shipped | 0 | 13.5 | Sunday | 14 | Scarsdale | NY | 224 |
| 12 | 2/26/2000 | 03\:42\:17 | 12.7 | 1 | 12.7 | Shipped | 0 | 12.7 | Saturday | 3 | Novato | CA | 236 |
| 13 | 3/30/2000 | 11\:51\:44 | 6.5 | 1 | 4.88 | Shipped | 0 | 4.88 | Thursday | 11 | Novato | CA | 236 |
| 14 | 3/30/2000 | 11\:51\:44 | 6.5 | 1 | 4.88 | Shipped | 0 | 4.88 | Thursday | 11 | Novato | CA | 236 |
| 15 | 3/30/2000 | 11\:51\:44 | 7 | 1 | 7 | Shipped | 0 | 7 | Thursday | 11 | Novato | CA | 236 |
| 16 | 3/30/2000 | 11\:51\:44 | 12 | 1 | 12 | Shipped | 0 | 12 | Thursday | 11 | Novato | CA | 236 |

# Pig LOAD Function

- **"Pigs eat anything"**
- **LOAD 'data' [USING function] [AS schema];**
- **USING**
  - PigStorage → structured text file (default)
  - TextLoader → unstructured UTF-8 data
  - Other and User Defined Functions
- **AS**
  - (Field1[:type], Field2:[type], ... FieldX[type])
  - Bytearray default type

# Pig Example: Loading Data

> pig –x mapreduce

> a = LOAD 'rmacc/KDDCupCleaned.txt';


> All Statements end with semicolon !!!

# Pig Debugging Statements

| Debug Operator | Description |
| --- | --- |
| DUMP | Display results |
| DESCRIBE | Display schema of relation |
| EXPLAIN | Display execution plan |
| ILLUSTRATE | Display step-by-step execution |
| *Full list* | http://pig.apache.org/docs/r0.11.1/test.html#diagnostic-ops |

- **Describe and illustrate only work if schema is provided**

# Pig Example: Loading Data

> a = LOAD 'rmacc/KDDrmacc' , AS (key,date,time, qty:float);

- **Columns can be accessed by**
  - $0   for second column (first contains key)
  - or name   (key, date, time....)

# Pig Data Types

- **Simple:**
  - int, long, float, double, chararray, bytearray, boolean
- **Complex:**
  - tuple - a set of fields  (10, 5, alpha)
  - bag - a collection of tuples {(10,5,alpha) (8,2,beta)}
  - map - a set of key value pairs [key#value]

# Pig Example: Reduce # of Fields

> a = LOAD 'rmacc/KDDCupCleaned.txt' AS (key:int,date,time,up,ol,qty,os,tax,amount:float, wd,hour,city,state,ci);

> b = FOREACH a GENERATE key, date, time, amount;

> STORE b INTO 'rmacc/KDDCupShort'.txt USING PigStorage(',') ;

# Pig Example:Group per Date

> a = LOAD 'rmacc/KDDrmacc.txt' AS (key,date,time,qty:float);

> groupday = GROUP  a  BY date;

> illustrate groupday  →

  > *group* is new key for each bag (day)

  > tuples within data within each bag

# Pig Example: Sum per Date

> a = LOAD  'rmacc/KDDrmacc.txt'  AS (key,date,time,qty:float);

> groupday = GROUP  a  BY date;

> sumday = FOREACH groupday GENERATE group, SUM(a.qty); sr

> STORE sumday INTO 'rmacc/sumday';

**Leeds** School of Business
UNIVERSITY OF COLORADO BOULDER

# Evaluation Functions

(Case Sensitive)

| Function | Description |
|---|---|
| AVG | Calculates average |
| CONCAT | Concatenates two expressions of identical type |
| COUNT | Counts the number of elements in a bag |
| COUNT_STAR | Like count by includes NULL values in count |
| DIFF | Compares two fields in a tuple |
| IsEmpty | Checks if a bag or map is empty |
| MAX | Calculates maximum |
| MIN | Calculates minimum |
| SIZE | Computes the number of elements  (characters) |
| SUM | Calculates sum |
| TOKENIZE | Splits a string and outputs a bag of words |
| *List with examples* | http://pig.apache.org/docs/r0.11.1/func.html#eval-functions |

**Leeds** School of Business
UNIVERSITY OF COLORADO BOULDER

# Other Functions

- **Math Functions**
- **String Functions**
- **Datetime Functions**
- **Tuple, Bag, Map Functions**

- **User Defined Functions**

# Pig Example: Filter Purchases> $100

> a = LOAD 'rmacc/sumday' AS (date,sum:float);

> bigpur = FILTER a BY sum>1000;

> STORE  bigpur  INTO  'rmacc/bigpur';

# Relational Operators

| Operators | Description |
|---|---|
| LOAD | Loads data from the file system |
| GROUP | Groups the data in one or more relations |
| FOREACH | Generates data transformations based on columns of data |
| FILTER | Selects tuples from a relation based on some condition |
| *Full list with examples* | http://pig.apache.org/docs/r0.11.1/basic.html#Relational+Operators |

# Other Operators

- **Arithmetic Operators**
- **Boolean Operators**
- **Cast Operators**
- **Comparison Operators**
- **Type Construction Operators**
- **Dereference Operators**
- **Disambiguate Operator**
- **[Flatten Operator](#)**
- **Null Operators**
- **Sign Operators**

# Pig WordCount Code

> b = LOAD 'rmacc/Beowulf.txt' AS be;

> beowords = FOREACH b GENERATE flatten(TOKENIZE(beo)) as bw;

> wg = GROUP beowords BY bw;

> wc = FOREACH wg GENERATE group, COUNT(beowords) as bc;

> sumord =ORDER wc BY bc;

> STORE sumord INTO 'rmacc/beowulfwc'

# Hadoop Summary

- **Accessible – runs on commodity hardware**
- **Robust – handles hardware failures**
- **Scalable – by adding more nodes**
- **Simple – allows users to quickly write efficient parallel code**
- **Pig - easy to learn**

# Conclusions

- **Individual components "easy" to setup → integration more complicated**

- **Resources**

  - Apache Hadoop (download and docu)

    - ***http://hadoop.apache.org***

  - Online Searches

  - Books for overview, not technical details

- **"Evolving project" → constantly changing, documentation can't keep up with development**

# Questions ?