# Big Data Engineering With MapReduce and Hive

Data Science Dojo

# Agenda

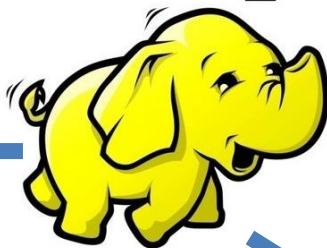Hadoop Implementations

# Hadoop

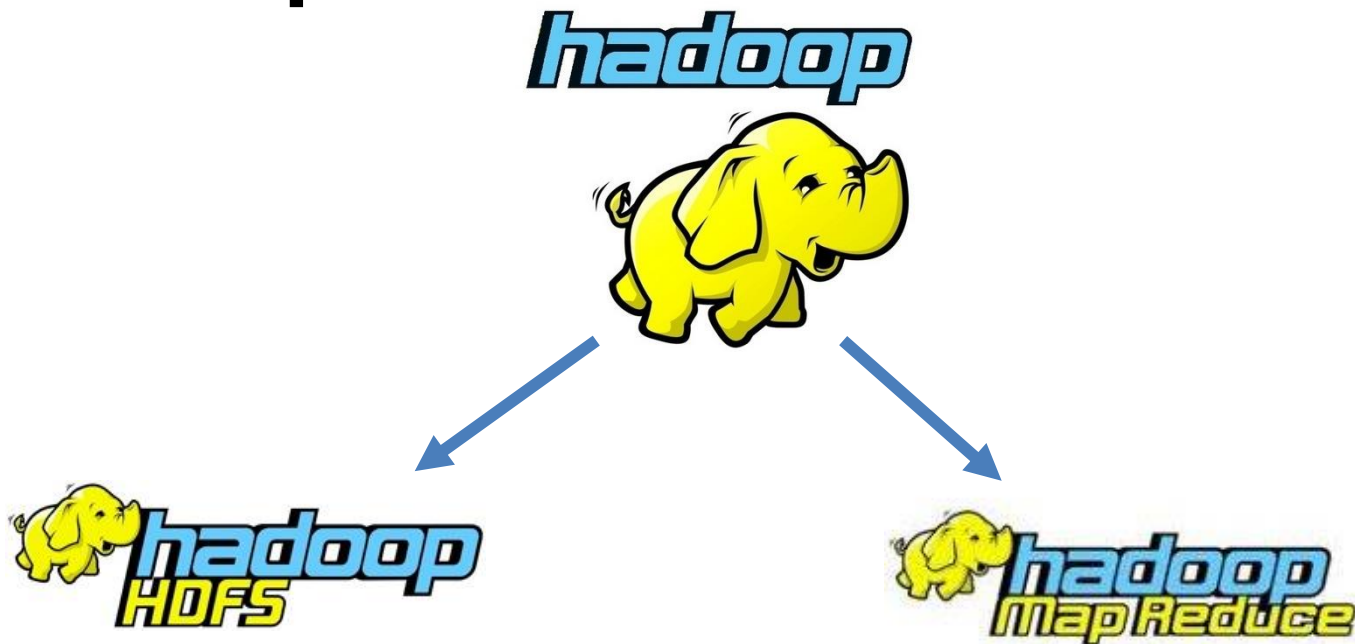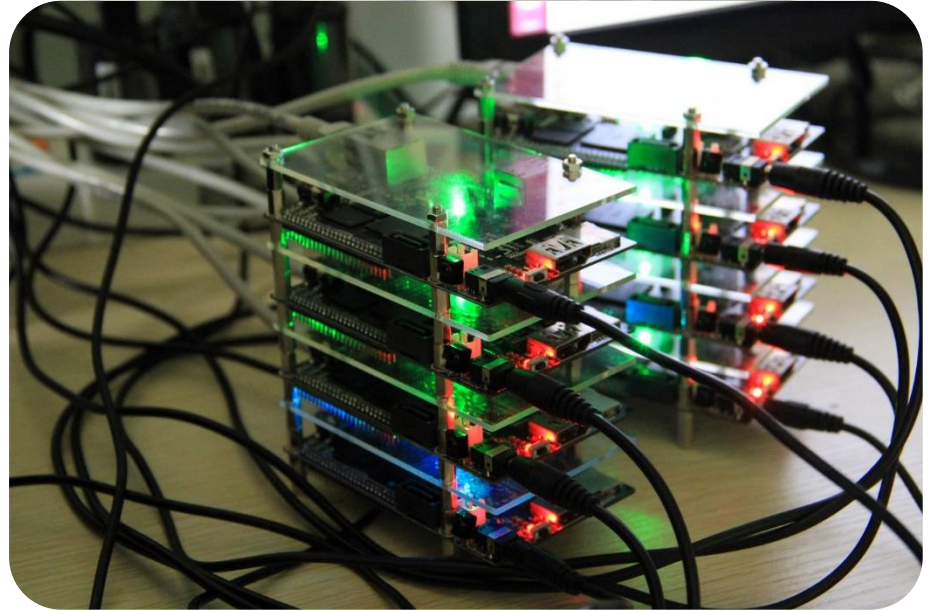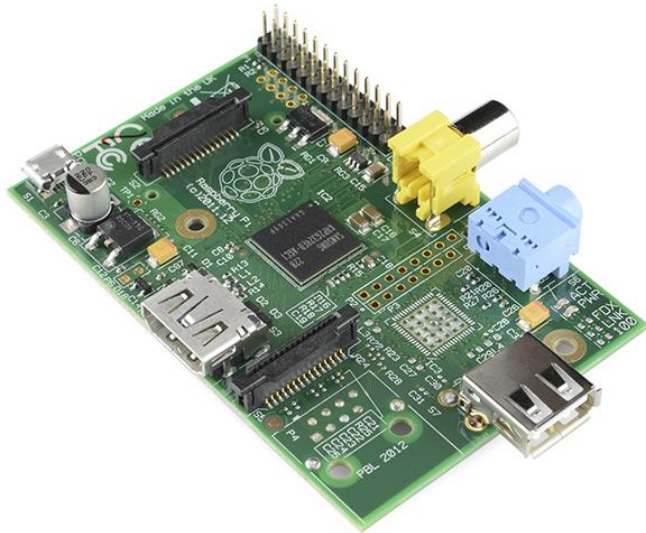# Turn Back The Clock, The Mainframe

# Distributed Computing

# Cloud Computing

# If dogs were servers...

Head Node
(Named Node)

Data Nodes

# HDFS & MapReduce

60gb of Tweets

60gb →

1 Computer

Processing: 30 hours

datasciencedojo
unleash the data scientist in you

# HDFS & MapReduce

60gb of Tweets

30gb

30gb

2 Computers

Processing: 15 hours

# HDFS & MapReduce

20Gb

60 Gb of Tweets

20Gb

20Gb

Processing: 10 hours

3 Computers

datasciencedojo
unleash the data scientist in you

# Most Cases, Linear Scaling Of Processing Power

| Number of Computers | Processing Time (hours) |
|:---:|:---:|
| 1 | 30 |
| 2 | 15 |
| 3 | 10 |
| 4 | 7.5 |
| 5 | 6 |
| 6 | 5 |
| 7 | 4.26 |
| 8 | 3.75 |
| 9 | 3.33 |

datasciencedojo
unleash the data scientist in you

# HDFS Redundancy

# Limitations with MapReduce

- ~70 lines of code to do anything
- Slow
- Troubleshooting multiple computers
- Good devs are scarce
- Expensive certifications

```java
1  package org.apache.hadoop.examples;
2
3  import java.io.IOException;
4  import java.util.StringTokenizer;
5
6  import org.apache.hadoop.conf.Configuration;
7  import org.apache.hadoop.fs.Path;
8  import org.apache.hadoop.io.IntWritable;
9  import org.apache.hadoop.io.Text;
10 import org.apache.hadoop.mapreduce.Job;
11 import org.apache.hadoop.mapreduce.Mapper;
12 import org.apache.hadoop.mapreduce.Reducer;
13 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
14 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
15 import org.apache.hadoop.util.GenericOptionsParser;
16
17 public class WordCount {
18
19   public static class TokenizerMapper
20       extends Mapper<Object, Text, Text, IntWritable>{
21
22     private final static IntWritable one = new IntWritable(1);
23     private Text word = new Text();
24
25     public void map(Object key, Text value, Context context
26                    ) throws IOException, InterruptedException {
27       StringTokenizer itr = new StringTokenizer(value.toString());
28       while (itr.hasMoreTokens()) {
29         word.set(itr.nextToken());
30         context.write(word, one);
31       }
32     }
33   }
```
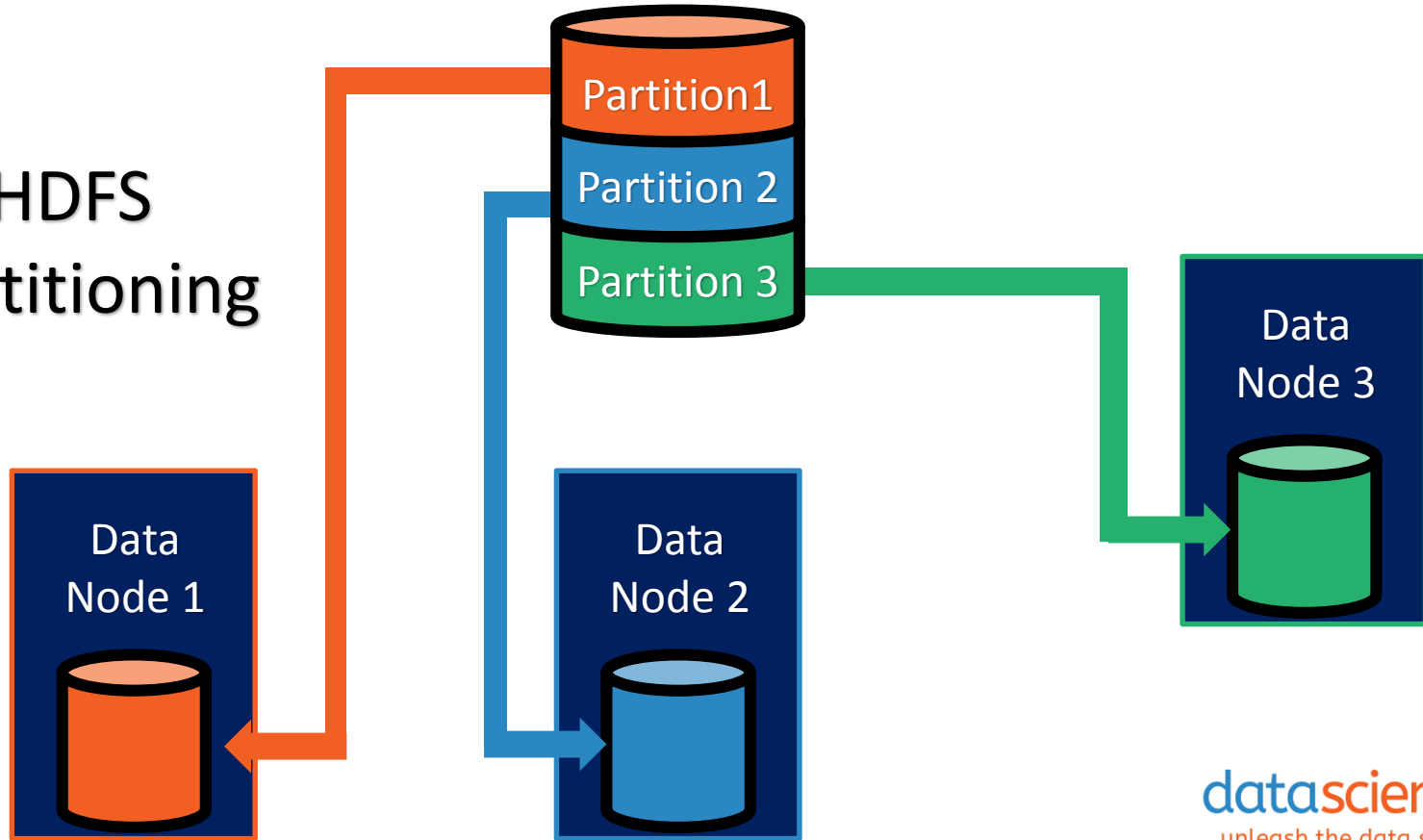
**Ambari:** Cluster provisioning, management, and monitoring

**Avro** (Microsoft .NET Library for Avro): Data serialization for the Microsoft .NET environment

**HBase:** Non-relational database for very large tables

**HDFS:** Hadoop Distributed File System

**Hive:** SQL-like querying

**Mahout:** Machine learning

**MapReduce and YARN:** Distributed processing and resource management

**Oozie:** Workflow management

**Pig:** Simpler scripting for MapReduce transformations
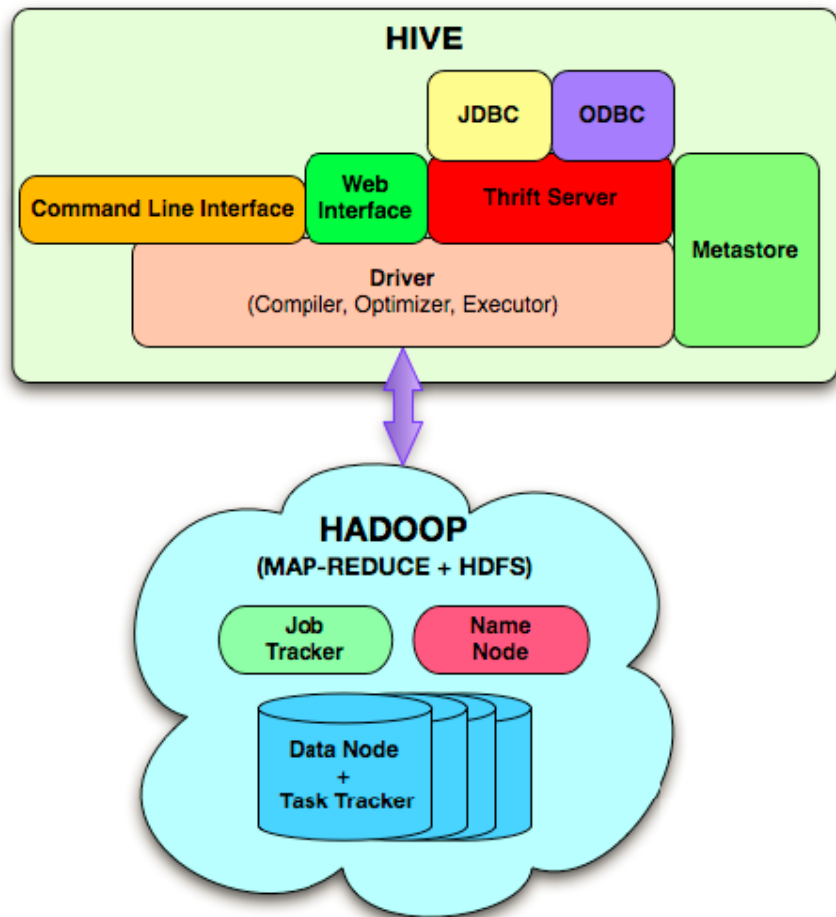
**Sqoop:** Data import and export

**Storm:** Real-time processing of fast, large data streams

**Zookeeper:** Coordinates processes in distributed systems

# Hive Jobs

HiveQL Statement → Translation & Conversion → MapReduce Job

datasciencedojo
unleash the data scientist in you

# Hive Architecture

**Data File** = Unstructured Data

**Data File** + **Metadata File/DB** = Structured Data

datasciencedojo
unleash the data scientist in you

# Semi Structured Data

Self Describing
Flat Files
- XML
- JSON
- CSV
- TSV

[
  {
    "created_at":"Thu May 07 18:06:23 +0000 2015",
    "id":5963755540631646210,
    "id_str":"5963755540631646210",
    "text":"Expert usable tips differently the press
    "source":"<a href=\"http://twitterfeed.com\" rel
    "truncated":0,
    "in_reply_to_status_id":null,
    "in_reply_to_status_id_str":null,
    "in_reply_to_user_id":null,
    "in_reply_to_user_id_str":null,
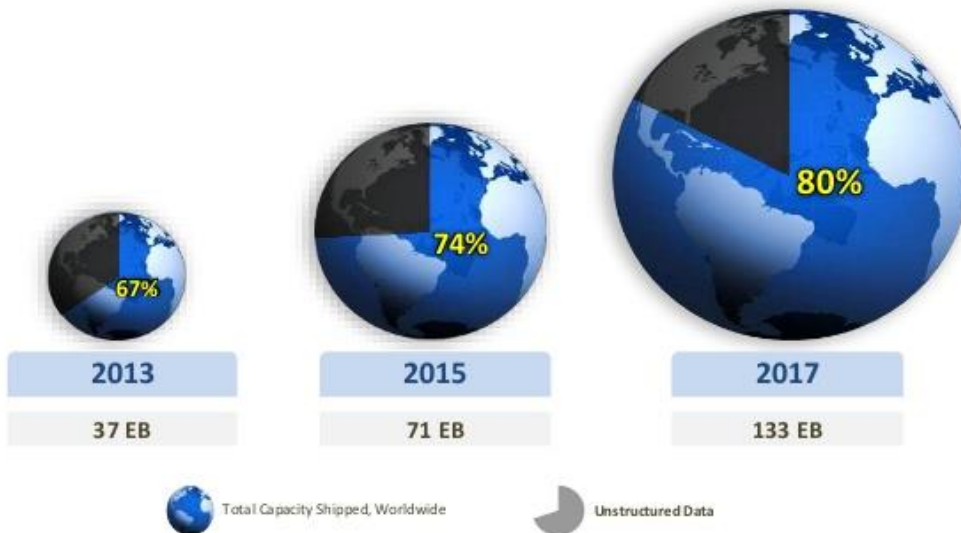
datasciencedojo
unleash the data scientist in you

# Why Hive?



- SQL spoken here (HiveQL)
- ODBC driver
- BI Integration
- Supports only Structured Data

# Limitations



80%
of organized data
is unstructured

Structured vs. Unstructured Data Growth

| 2013 | 2015 | 2017 |
|------|------|------|
| 67% | 74% | 80% |
| 37 EB | 71 EB | 133 EB |

Total Capacity Shipped, Worldwide    Unstructured Data
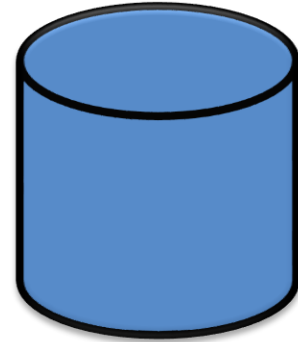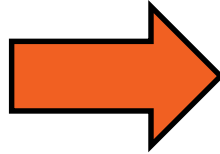
Source: IDC

datasciencedojo
unleash the data scientist in you

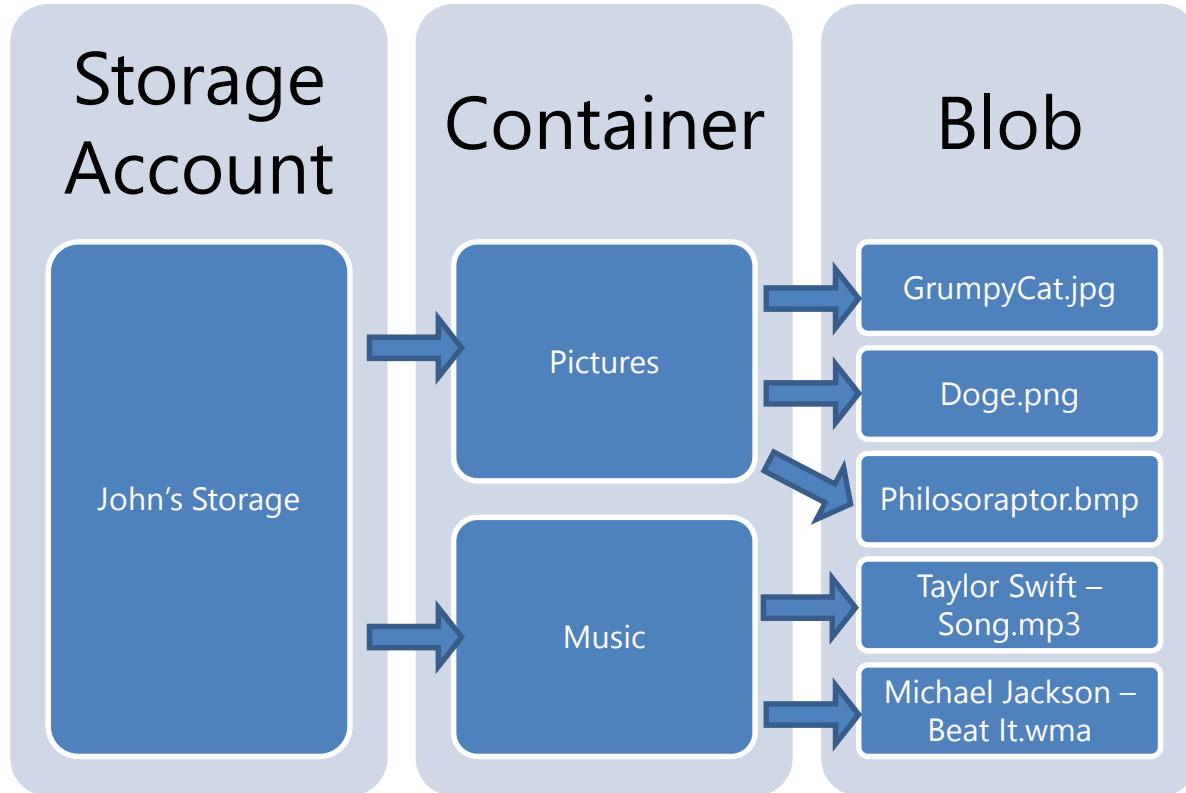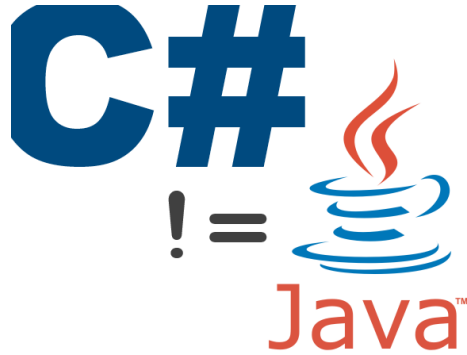# Azure Blob Storage



HDInsight

Blob Storage

# Azure Blob Storage
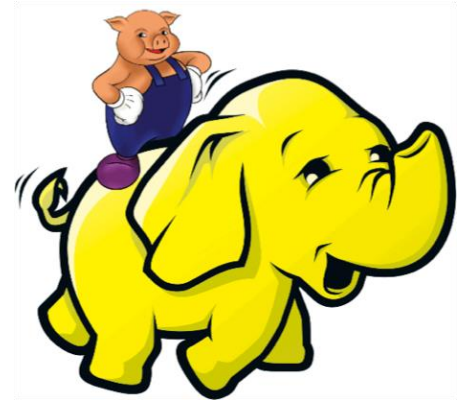
# When to Use Each
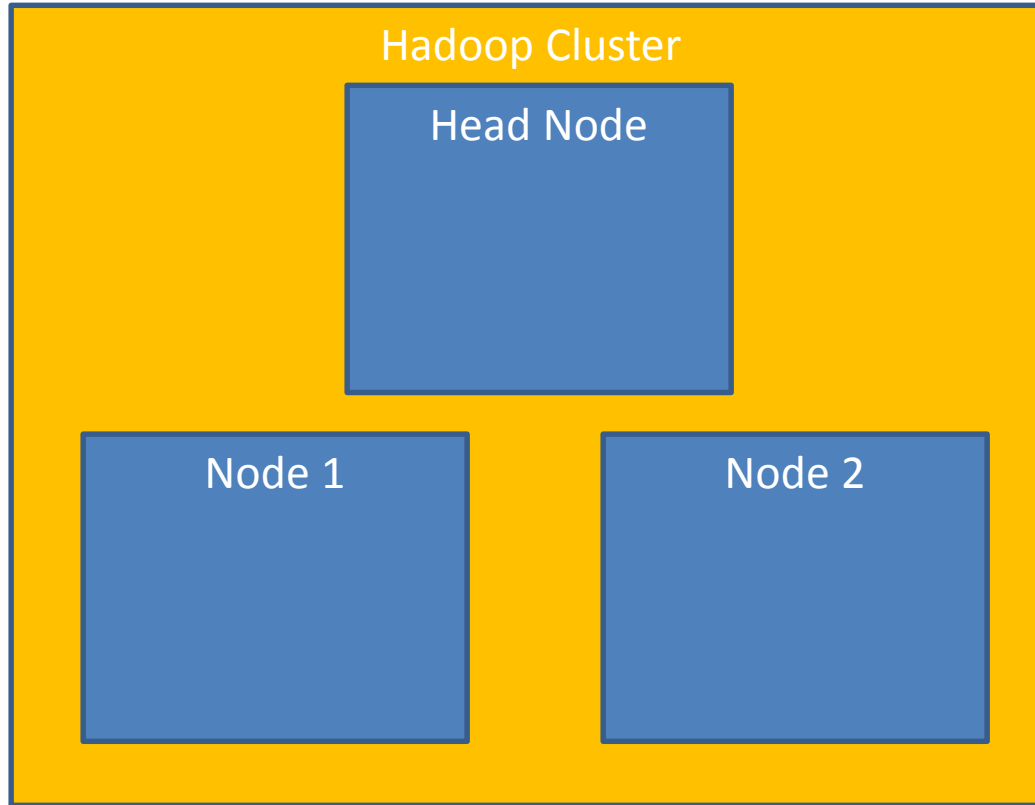


C#
Java
MapReduce

VS

Hive

VS

Pig

# MapReduce, via Playing Cards



Let's count the number of spades, clubs, hearts, and diamonds in a stack of cards, the way map reduce would.

- Each card represents a row of data

- Each suite represents an attribute of the data

# Using a 2 Data Node Cluster
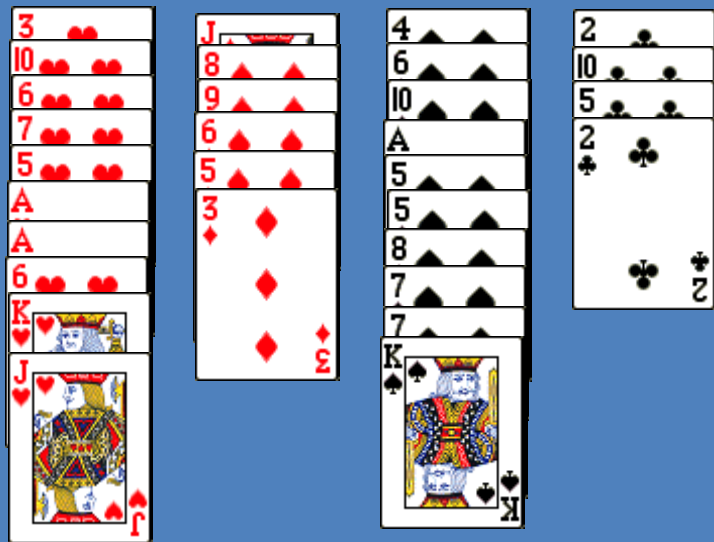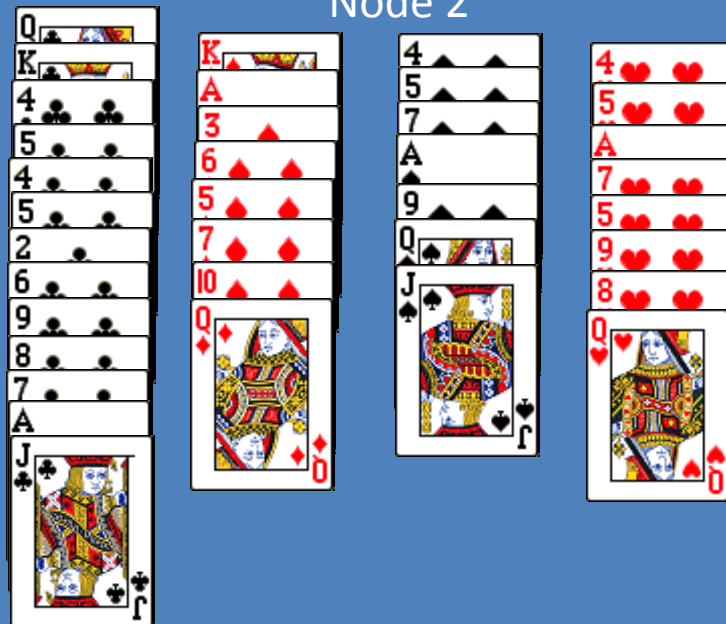
# Mapping: Each Node's HDFS
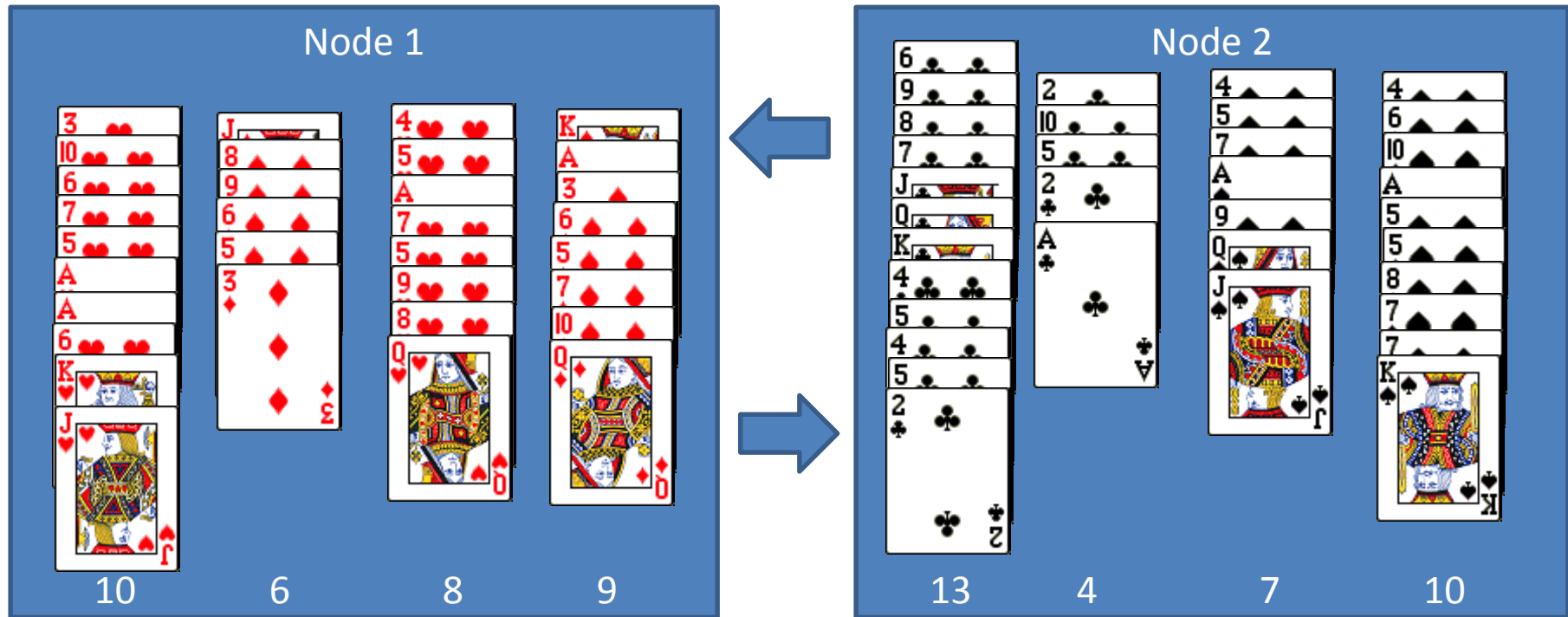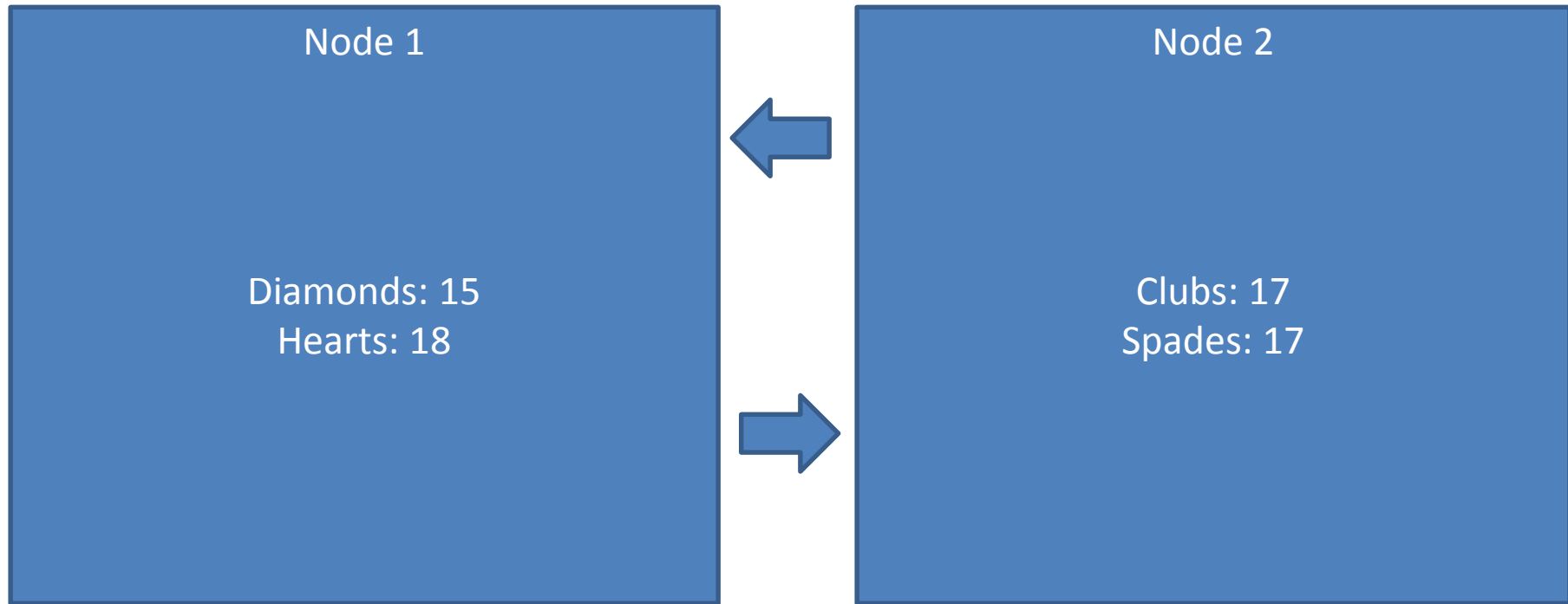
# Mapping: Node Sorting

# Mapping: Node Shuffle, Data Transfer

# Mapping: Node Shuffle, Data Transfer

Node 1

Diamonds: 15
Hearts: 18

Node 2

Clubs: 17
Spades: 17

datasciencedojo
unleash the data scientist in you

# Database = Normalization

select * from **courses**;

| id | name | teacher_id |
|---|---|---|
| 10001 | Computer Science 142 | 1234 |
| 10002 | Computer Science 143 | 5678 |
| 10003 | Computer Science 154 | 9012 |
| 10004 | Informatics 100 | 1234 |

select * from **grades**;

| student_ids | course_id | grade |
|---|---|---|
| 123 | 10001 | B- |
| 123 | 10002 | C |
| 456 | 10001 | B+ |
| 888 | 10002 | A+ |
| 888 | 10003 | A+ |
| 404 | 10004 | D+ |
| 456 | 10002 | D- |
| 404 | 10002 | B |
| 789 | 10003 | D+ |

select * from **students**;

| id | name | email | password |
|---|---|---|---|
| 123 | Bart | bart@fox.com | bartman |
| 404 | Ralph | ralph@fox.com | catfood |
| 456 | Milhouse | milhouse@fox.com | fallout |
| 789 | Nelson | muntz@fox.com | haha! |
| 888 | Lisa | lisa@gmail.com | vegan |

# Normalization, joining

```sql
select
    g.student_id,
    c.name as course,
    g.grade as grade
from grades as g
join courses c
    on g.course_id = c.id
join students s
    on g.student_id = s.id
;
```

| student | course | grade |
|---------|--------|-------|
| Bart | Computer Science 142 | B- |
| Milhouse | Computer Science 142 | B+ |
| Bart | Computer Science 143 | C |
| Lisa | Computer Science 143 | A+ |
| Milhouse | Computer Science 143 | D- |
| Ralph | Computer Science 143 | B |
| Lisa | Computer Science 154 | A+ |
| Nelson | Computer Science 154 | D+ |
| Ralph | Informatics 100 | D+ |

# Data Warehouse = Denormalization

| student | course | grade |
|---------|--------|-------|
| Bart | Computer Science 142 | B- |
| Milhouse | Computer Science 142 | B+ |
| Bart | Computer Science 143 | C |
| Lisa | Computer Science 143 | A+ |
| Milhouse | Computer Science 143 | D- |
| Ralph | Computer Science 143 | B |
| Lisa | Computer Science 154 | A+ |
| Nelson | Computer Science 154 | D+ |
| Ralph | Informatics 100 | D+ |

Tables:
- Students Table
- Courses Table
- Roster Table

# Costs, Storage vs Processing

Storage

| US – N. Virginia | US – N. California | EU – Ireland |
|---|---|---|

| Standard On-Demand Instances | Linux/UNIX Usage | Windows Usage |
|---|---|---|
| Small (Default) | $0.085 per hour | $0.12 per hour |
| Large | $0.34 per hour | $0.48 per hour |
| Extra Large | $0.68 per hour | $0.96 per hour |

Processing

| US – Standard | US – |
|---|---|

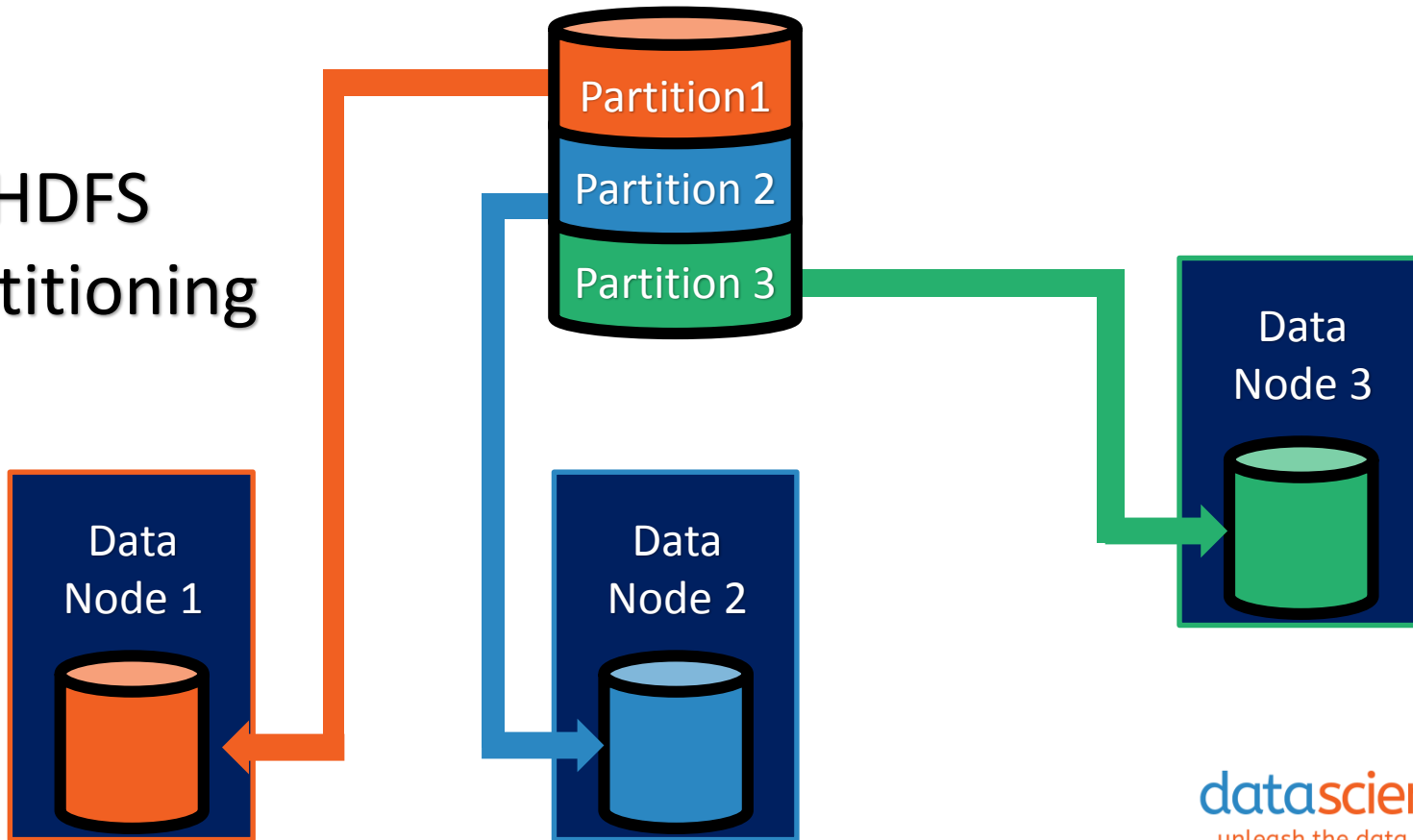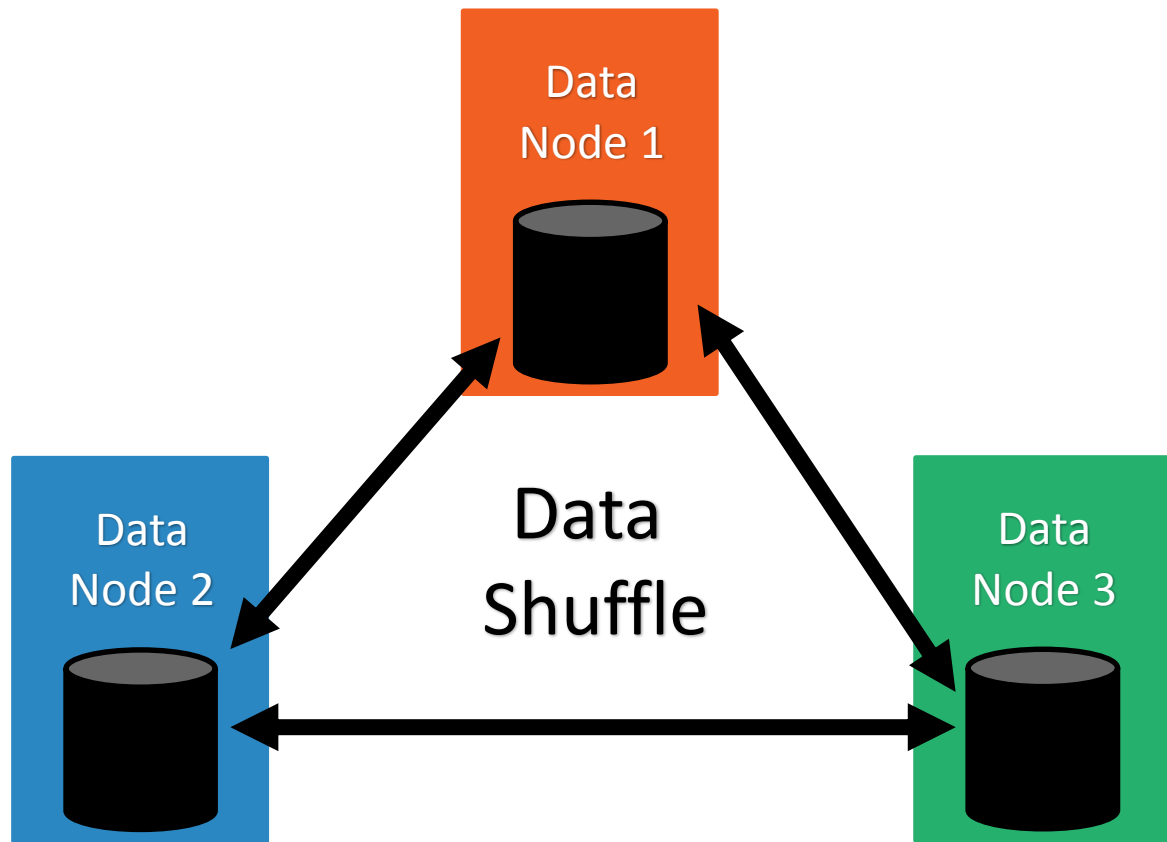| Storage | |
|---|---|
| **Tier** | **Pricing** |
| First 50 TB / Month of Storage Used | $0.150 per GB |
| Next 50 TB / Month of Storage Used | $0.140 per GB |
| Next 400 TB / | $0.130 per GB |


DENORMALIZE ALL THE THINGS

- Distributed Machine Learning
- Installed into Hadoop & Spark
- R-like language Implementation

# Classification

| | PC | Mahout | Spark |
|---|---|---|---|
| **Logistic Regression - trained via SGD** | | **<10m** | |
| **Naive Bayes** | | | |
| **Random Forest** | | | |
| **Hidden Markov Models** | | | |
| **Multilayer Perceptron** | | | |

# Recommendation Engines

|  | PC | Mahout | Spark |
|---|---|---|---|
| **User-Based Collaborative Filtering** | 🟩 | 🟧 | 🟩 |
| **Item-Based Collaborative Filtering** | 🟩 | 🟩 | 🟩 |
| **Matrix Factorization with ALS** | 🟩 | 🟩 | |
| **Matrix Factorization with ALS on Implicit Feedback** | 🟩 | 🟩 | |
| **Weighted Matrix Factorization, SVD++** | 🟩 | | |

# Clustering

|  | PC | Mahout | Spark |
|---|---|---|---|
| k-Means Clustering | 🟩 | 🟧 | 🟩 |
| Fuzzy k-Means | 🟩 | 🟩 | |
| Streaming k-Means | 🟩 | 🟩 | |
| Spectral Clustering | | 🟩 | |

datasciencedojo
unleash the data scientist in you

Spark

Running time (s)

120 ─ 110
 90
 60
 30
  0 ─ 0.9

■ Hadoop
■ Spark

In-Memory: 100x times faster than Hadoop

datasciencedojo
unleash the data scientist in you

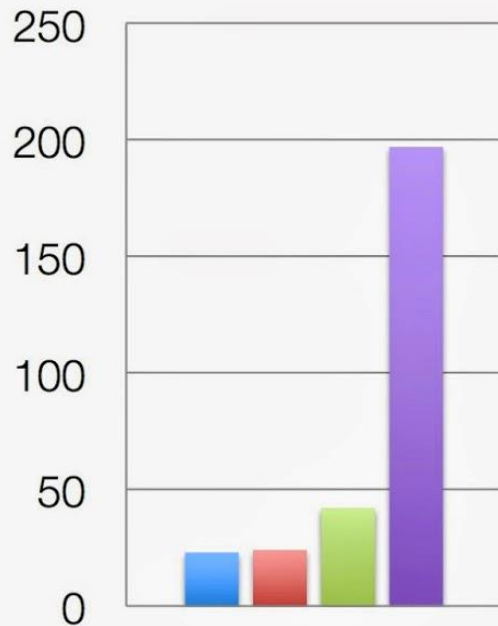3x faster on 10x few machines

Datona GraySort Benchmark: Sort 100 TB of data

Previous World Record:
- Method: Hadoop
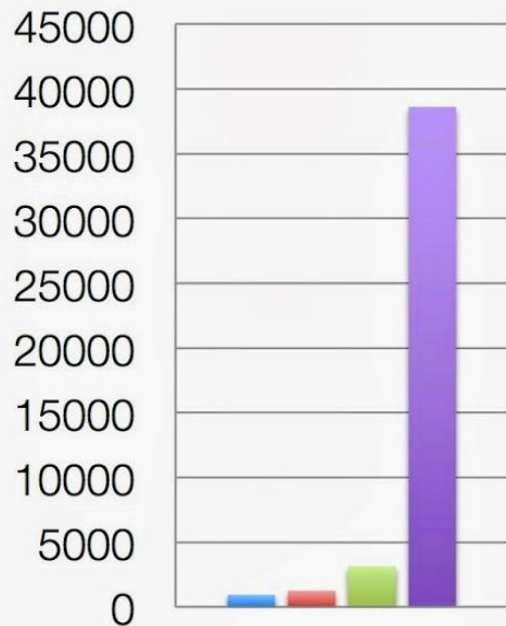- Yahoo!
- 72 Minutes
- 2100 Nodes

2014:
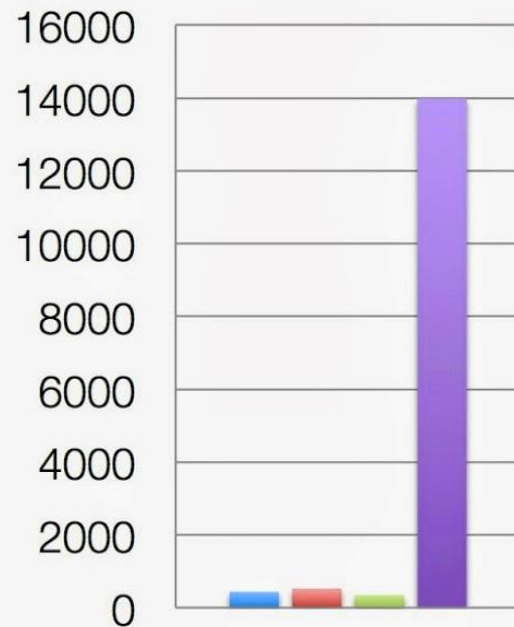- Method: Spark
- Databricks
- 23 Minutes
- 206 Nodes

Source: https://databricks.com/blog/2014/10/10/spark-petabyte-sort.html

Activity in last 30 days

Source: Xiangrui Meng, Data Bricks

# Technology adoption life cycle

# QUESTIONS

datasciencedojo
unleash the data scientist in you