# Big Data Engineering with Distributed Systems

Data Science Dojo

# Agenda

- Introduction:
  - Data engineering for data scientists
  - The "5 Vs" of Big Data
- A key problem – machine learning at scale
- Distributed computing with Apache Hadoop & Hive
- Hadoop in the Azure cloud
- Machine learning at scale with Apache Mahout
- Distributed computing v2.0 – Apache Spark

datasciencedojo
data science for everyone

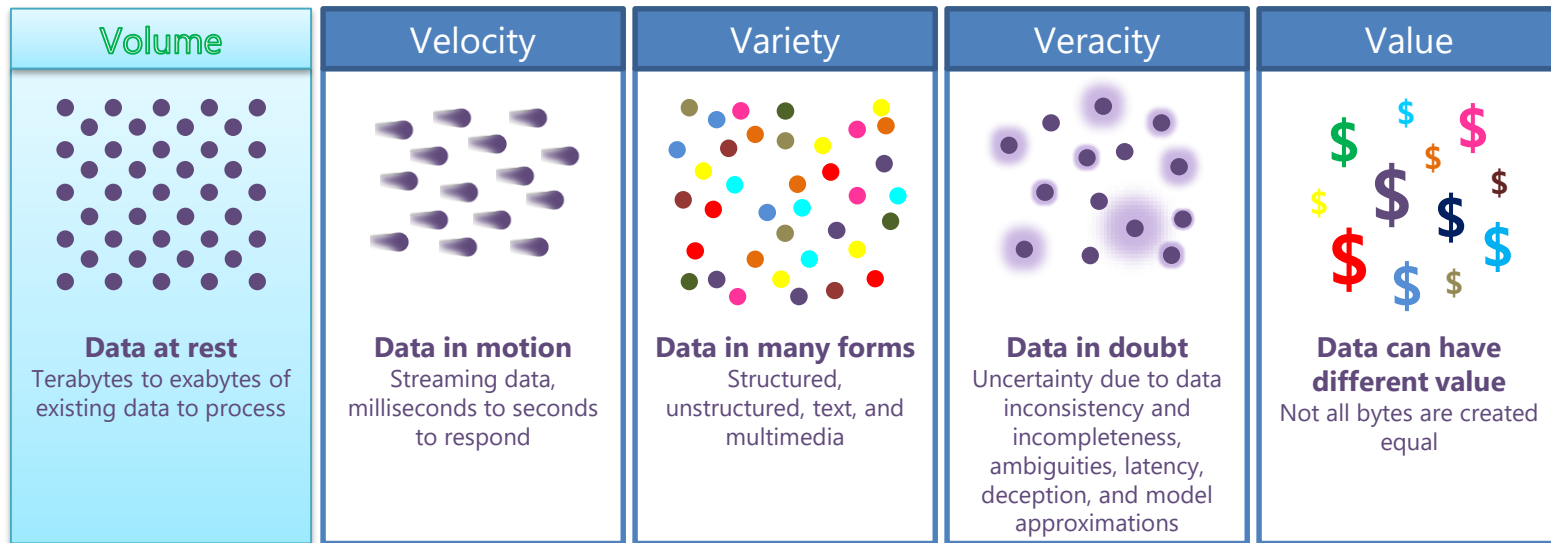# Data Engineering for Data Scientists



Driving a car

vs



Servicing a car

**Goals:**
- Teach you about data engineering topics/concepts

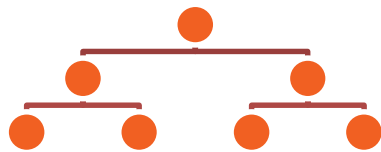**Non goals:**
- Managing or administering a Hadoop cluster

datasciencedojo
data science for everyone

# 5 Vs of Big Data

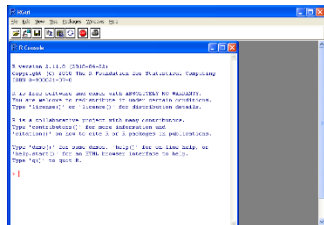| Volume | Velocity | Variety | Veracity | Value |
|--------|----------|---------|----------|-------|
| **Data at rest** Terabytes to exabytes of existing data to process | **Data in motion** Streaming data, milliseconds to seconds to respond | **Data in many forms** Structured, unstructured, text, and multimedia | **Data in doubt** Uncertainty due to data inconsistency and incompleteness, ambiguities, latency, deception, and model approximations | **Data can have different value** Not all bytes are created equal |

- **Goal:** As data scientists we want cost-effective access to the raw materials for our data products!

# MACHINE LEARNING AT SCALE

# OSS R Limits
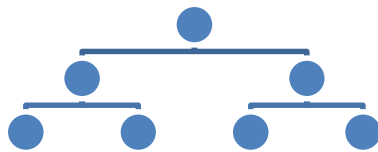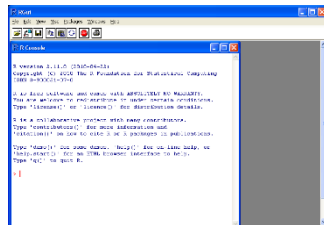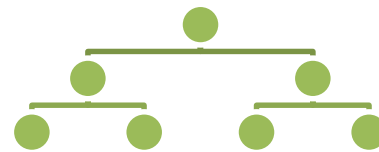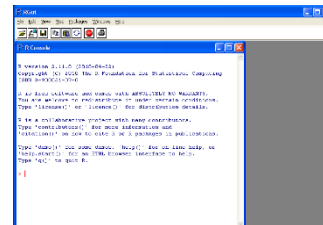
- Single core
- Single threaded
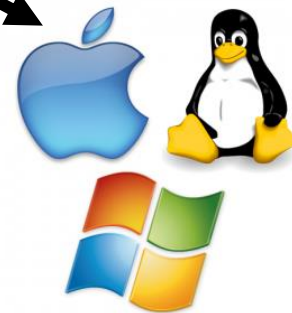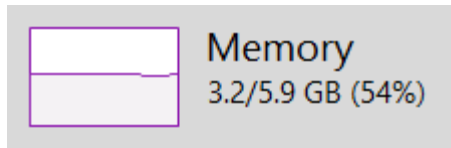
# OSS R Limits

- Single core
- Single threaded
- All in memory (RAM)
- Vectors & Matrices capped at 4,294,967,295 elements (rows) if 32-bit version; 2^32 - 1

# OSS R Limits: RAM

- All in memory (RAM)

$Max\ Data\ Limit = (\ Total\ RAM\ Access\ x\ 80\%) - Normal\ RAM\ Usage$

Laptop Example:

Memory
3.2/5.9 GB (54%)

$Max\ Data\ Limit = (\ 5.9\ gb\ x\ 80\%)\ - 3.2\text{gb}$
$Max\ Data\ Limit = \sim1.52gb$

*R data frames actually bloats data files by 3x
$R\ Data\ Limit = \sim1.52gb\ \div 3 = \sim506.7mb$

datasciencedojo
data science for everyone

# OSS R Limits: RAM

| INSTANCE | CORES | RAM | DISK SIZES [1] | PRICE |
|----------|-------|-----|----------------|-------|
| M64MS | 64 | 1,750.00 GiB | **2,000** GB | $10.34/hr |
| M128S | 128 | 2,000.00 GiB | **4,000** GB | $13.34/hr |

Azure's VM with largest RAM[*]:

$$Max\ Data\ Limit = (\ 2000gb\ x\ 80\%\ ) - 1gb$$
$$Max\ Data\ Limit = \sim 1600gb$$
$$R\ Data\ Limit = \sim 1600gb \div 3 = \sim 533.33\ gb$$

**24x7x52 Annual Cost: $116,938.44!**

*Data collected 06/07/2017

9

datasciencedojo
data science for everyone

# Machine Learning Scaling

**Programs**

- Excel

**Programming**

- R
- Python
- SAS

**Cloud**

- Azure ML
- AWS ML
- Big ML
- Cloud Virtual Machines

**This only gets us so far!**

**Distributed**

- Hadoop
- Spark
- H20
- Microsoft R Server

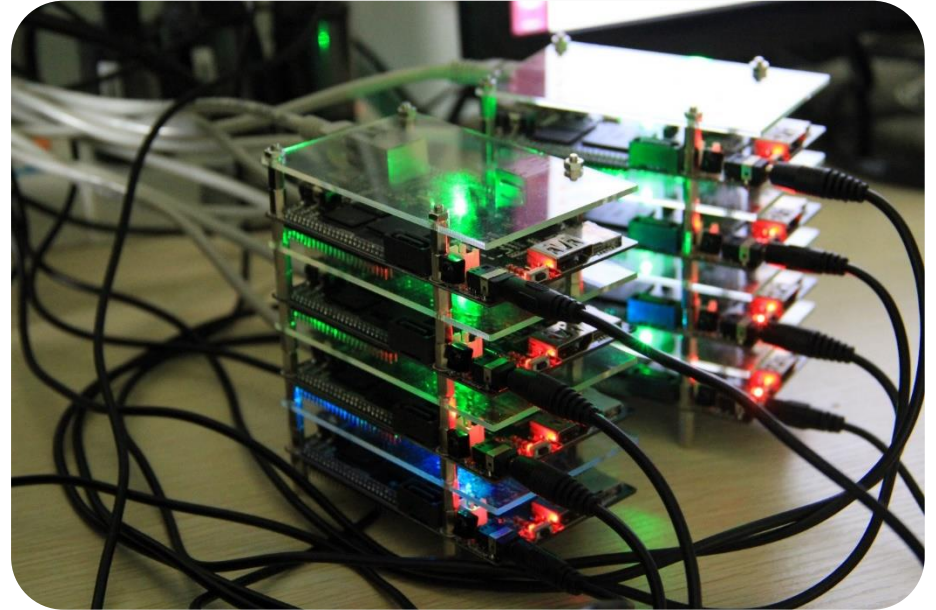**Big data scale!**
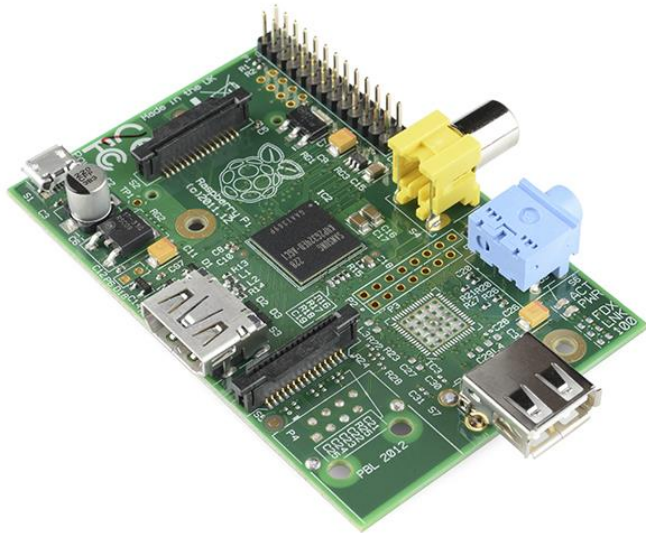
datasciencedojo
data science for everyone

# DISTRIBUTED COMPUTING WITH APACHE HADOOP
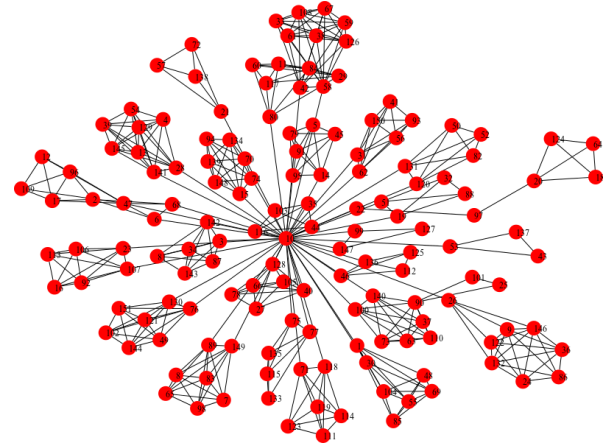
# Turn Back The Clock, The Mainframe

- "Big Iron"
- Backbone of computing for decades.
- Still widely used.
- "Scale-up" model of shared computing.
- Core platform is cost effective, ecosystem is not (e.g., software licensing).
- The original VM host!

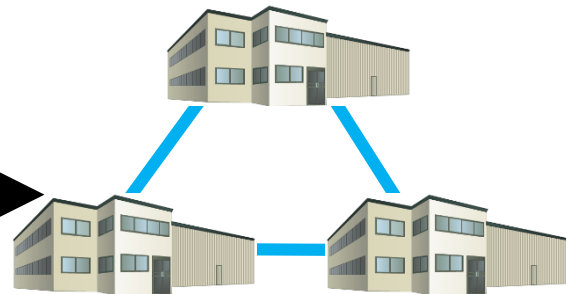# Distributed Computing

# Cloud Computing



- Conceptually – a combination of mainframe and distributed computing.
- VM hosts are now the "Big Iron".
- Many VMs work together to distribute workloads.
- Some workloads on dedicated HW (e.g., SAP HANA).

datasciencedojo
data science for everyone

# Scaling Computational Power



Old Scaling:
- Vertical Scaling, Scaling UP
- High performance computers

New Scaling:
- Horizontal Scaling, Scaling OUT
- Commodity hardware, distributed

datasciencedojo
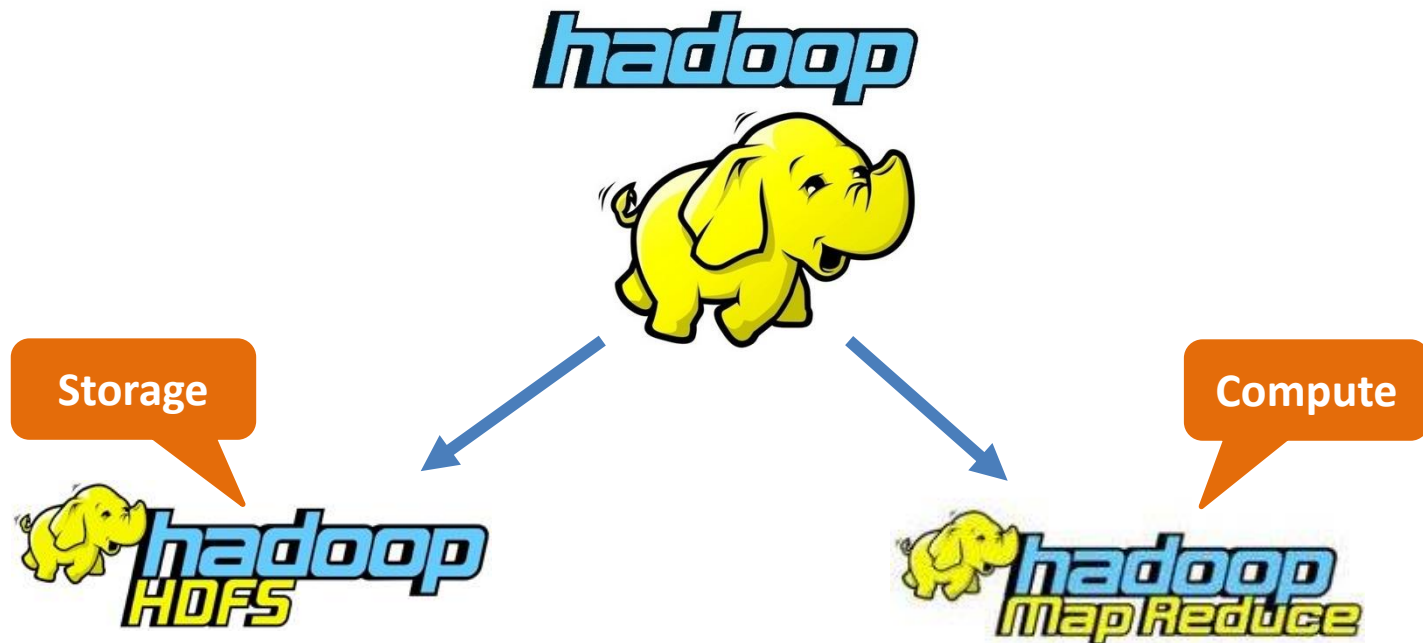data science for everyone

# What is Hadoop?

- OSS Platform for distributed computing over Internet-scale data.

- Originally built at Yahoo!

- Implementation of ideas (e.g., MapReduce) published by Google.

- The de facto standard big data platform.

- Named after a stuffed animal belonging to Doug Cutting's son.
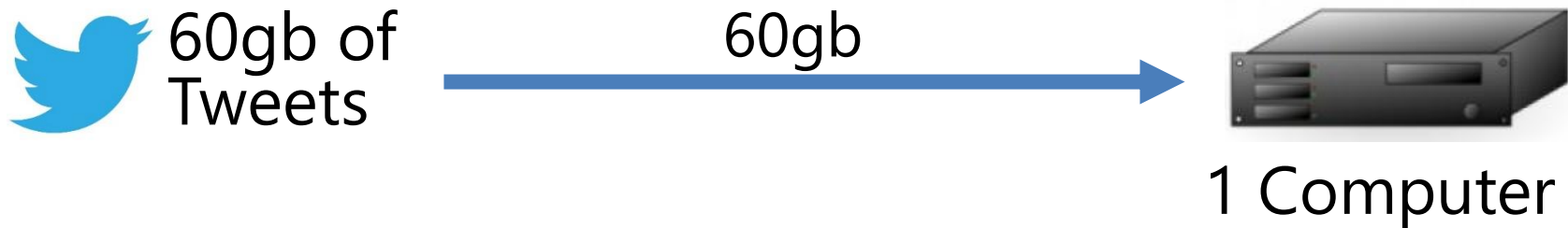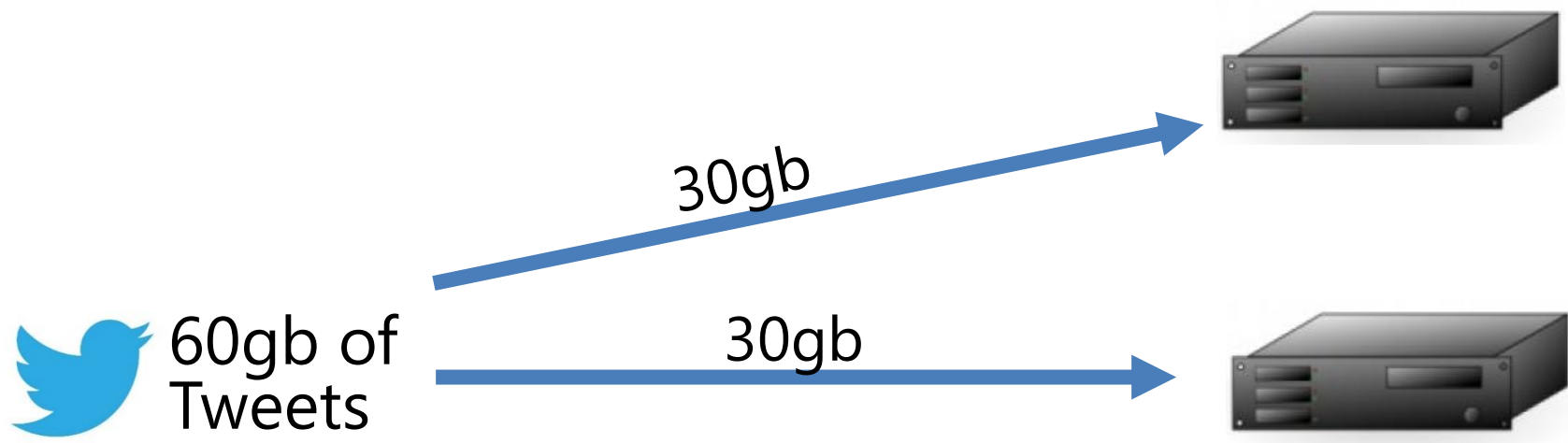
# Hadoop at Base



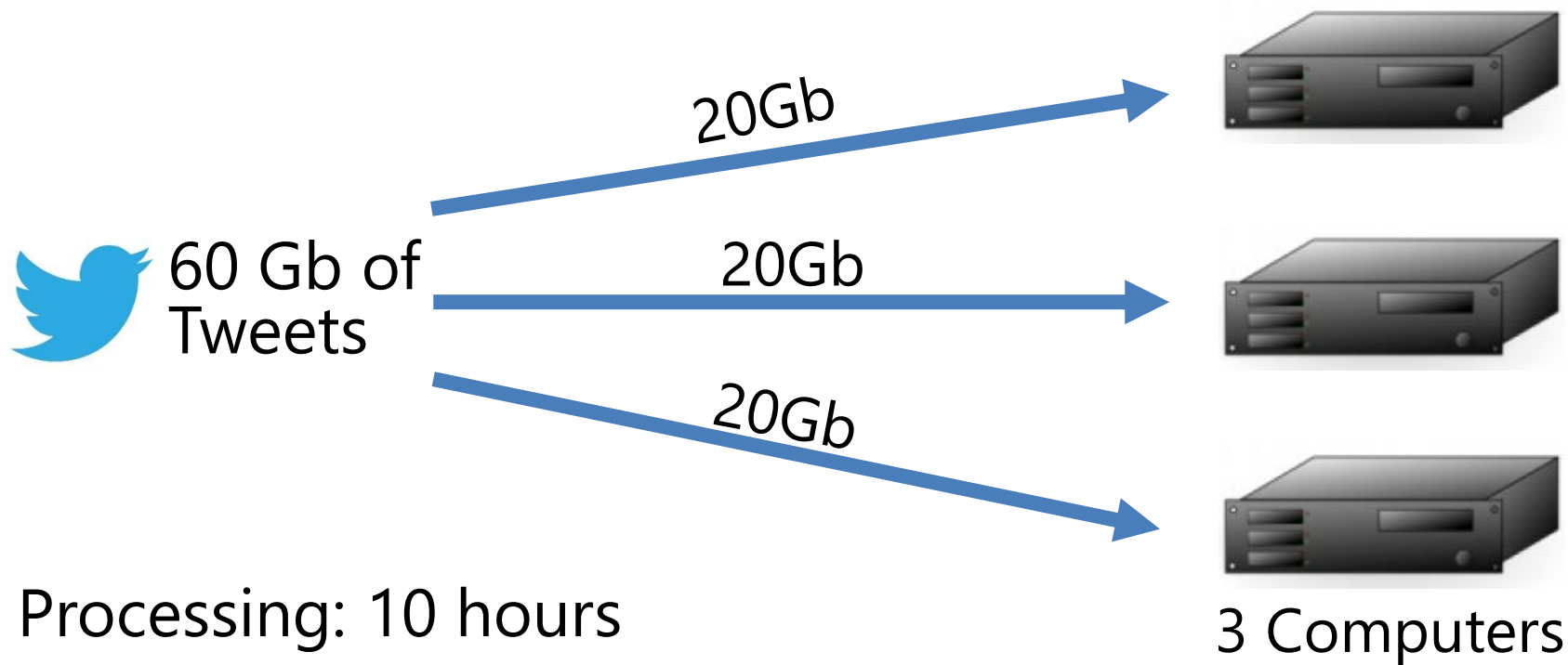Distributed batch processing engine for big data.

# HDFS & MapReduce

60gb of Tweets

60gb

1 Computer

Processing: 30 hours

# HDFS & MapReduce

60gb of Tweets

30gb

30gb

2 Computers

Processing: 15 hours

19

datasciencedojo
data science for everyone

# HDFS & MapReduce



60 Gb of
Tweets

20Gb

20Gb

20Gb

Processing: 10 hours

3 Computers

# 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
data science for everyone

# If dogs were servers...

Head Node
(Named Node)

Data  Nodes
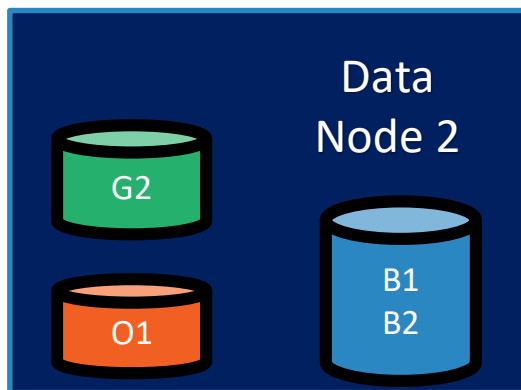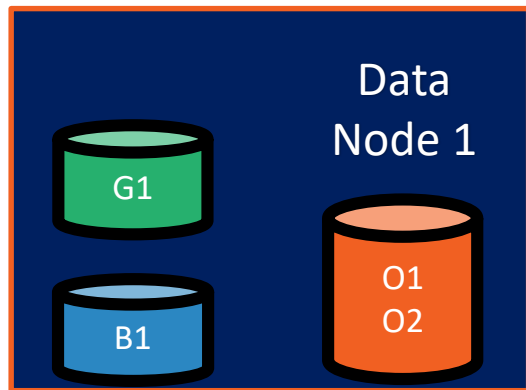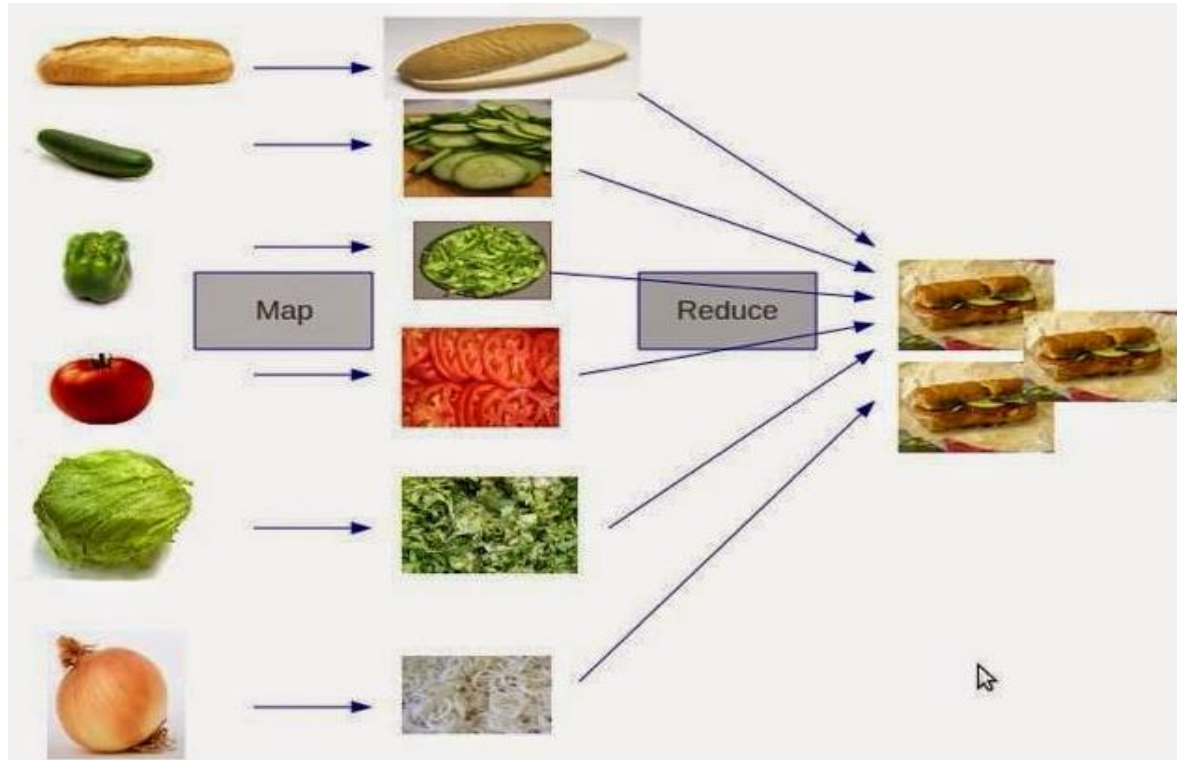
# HDFS Redundancy

# MapReduce – Sandwich Analogy

# 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      }
```

26

data science for everyone

# DISTRIBUTED COMPUTING WITH APACHE HIVE

# What is Hive?

- Abstraction built on top of MapReduce & HDFS.

- Makes Hadoop look like an RDBMS (e.g., coding in SQL).

- Developed by Facebook to democratize Hadoop.

- Applies structure to data at runtime ("schema on read").

# Schema on Read
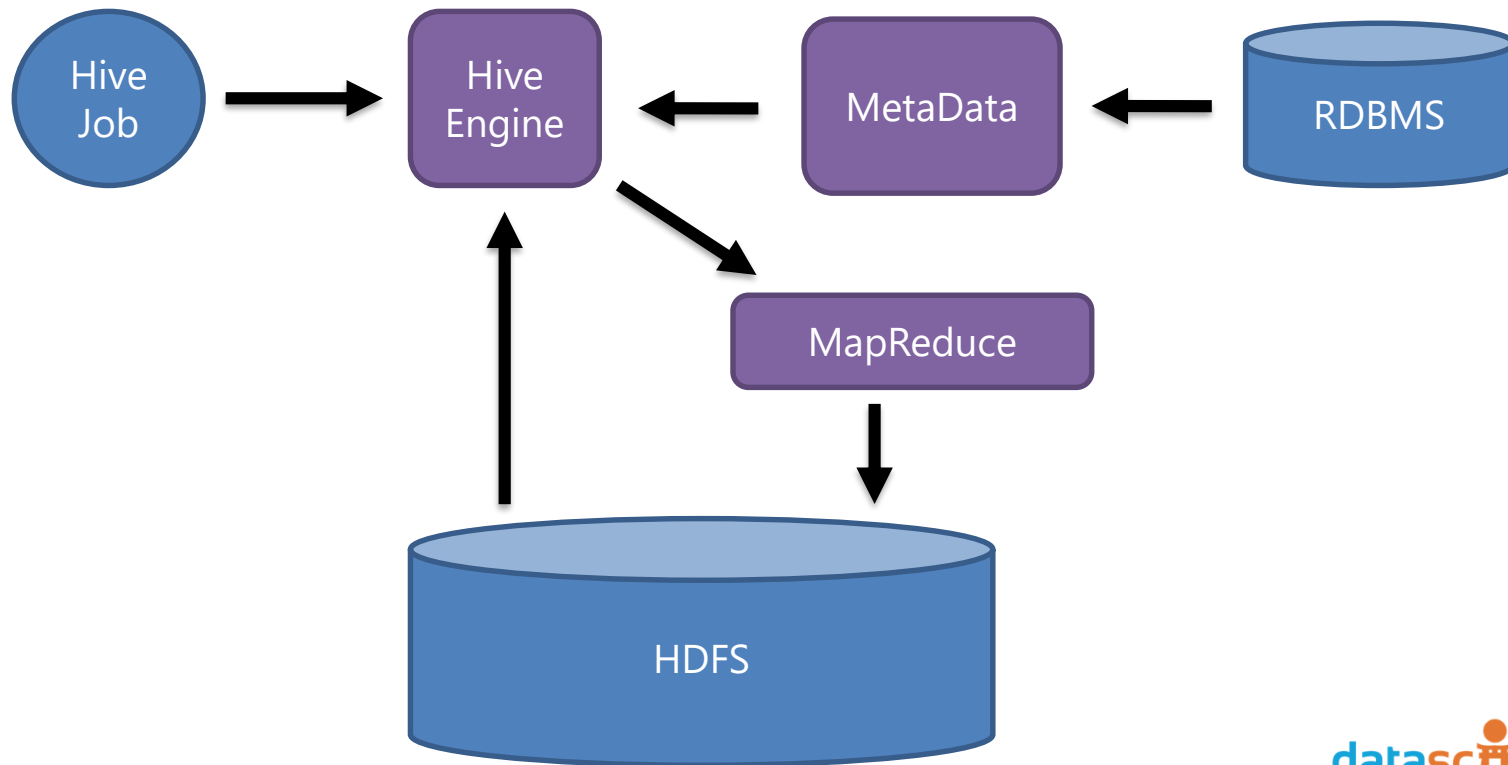
**Data File** = Unstructured Data

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

datasciencedojo
data science for everyone

# Hive Architecture

# Hive Jobs

HiveQL Statement → Translation & Conversion → MapReduce Job

# Word Count Revisited

```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      }
```

**vs.**

```sql
SELECT word,
       COUNT(*) AS word_count
FROM words
GROUP BY word
```

32

data**science**dojo
data science for everyone

# SQL Don'ts in HIVE

**SELECT * FROM ANYTHING:** This brings back everything. Everything doesn't fit on a single computer.

**JOIN:** Join will take hours or days to perform and eat up all cluster bandwidth for everyone else trying to use it in the queue.

**ORDER BY:** Sorting is very computationally expensive.

**Sub Queries:** A sub query essentially creates a secondary table, which will be huge in HIVE.

**Interactivity:** SQL in DBMS is interactive because its almost instantaneous.

datasciencedojo
data science for everyone

# Execution Engine: Tez

## The Stinger Initiative

2011, the world got together and declared MapReduce to be terrible.

- 44 companies
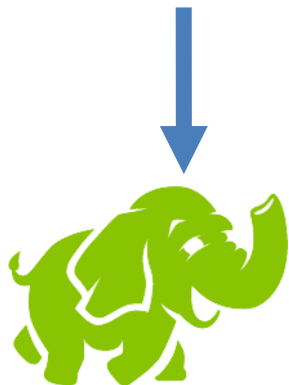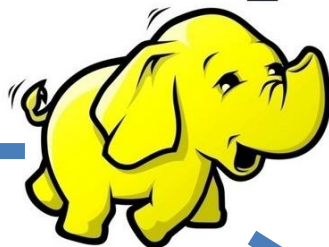- 145 developers
- 392k lines of Java code

## Hadoop 2.0 with Yarn & Tez

- Tez dropped Hive query times by **90%, 100x performance**
- Utilizes Apache Yarn
  - Yarn: resource manager for multi-cluster computing
- Introduced partial in-memory, local head nodes
- Rewrote HiveQL as an actual language, instead of translation

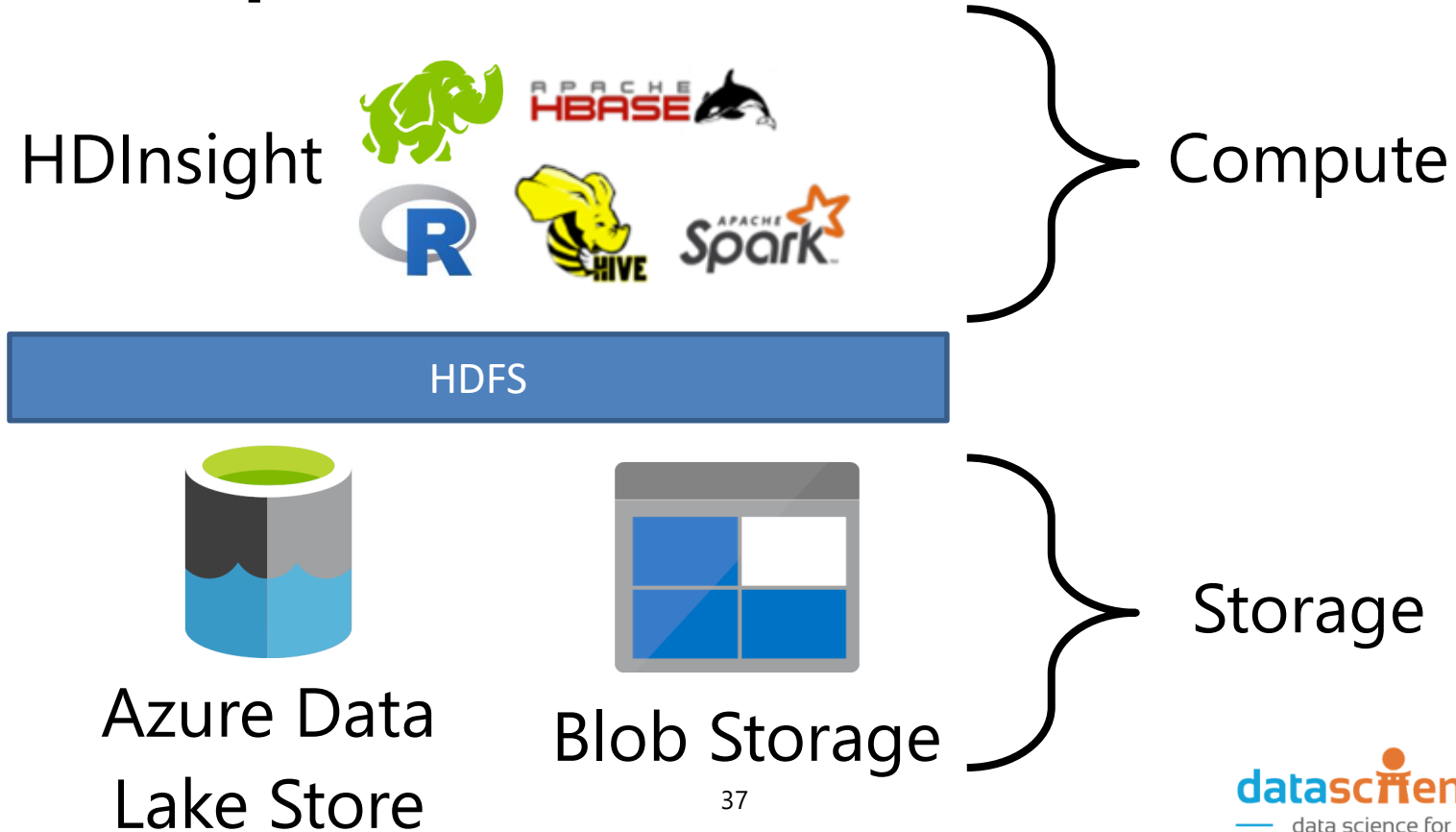# HADOOP IN THE AZURE CLOUD

# Hadoop Implementations



Hortonworks

HDInsight

MAPR

cloudera hadoop

Amazon Elastic MapReduce

datasciencedojo
data science for everyone

# Hadoop in Azure

HDInsight



Compute

HDFS

Azure Data Lake Store

Blob Storage

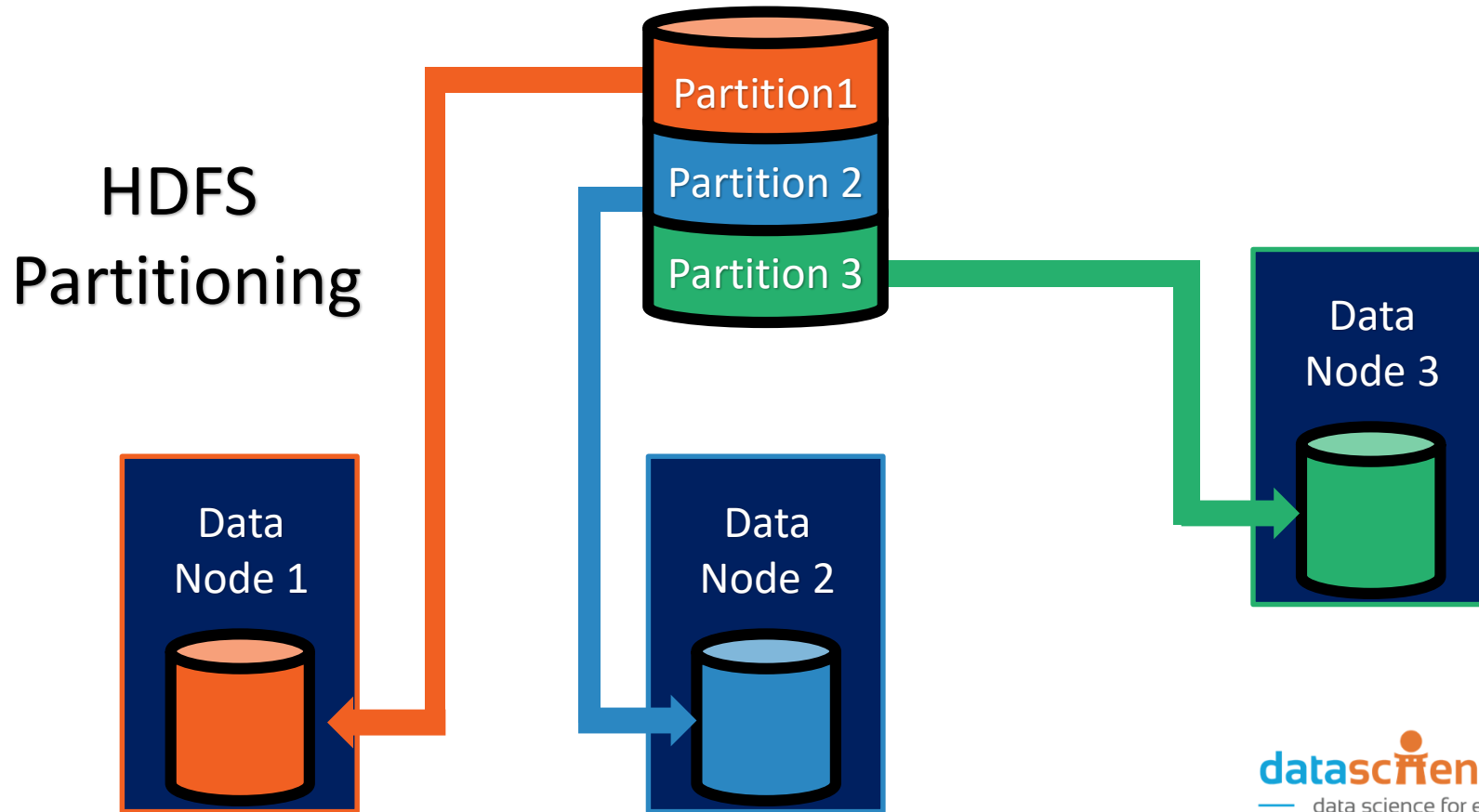Storage

# MACHINE LEARNING AT SCALE - REVISITED

datasciencedojo
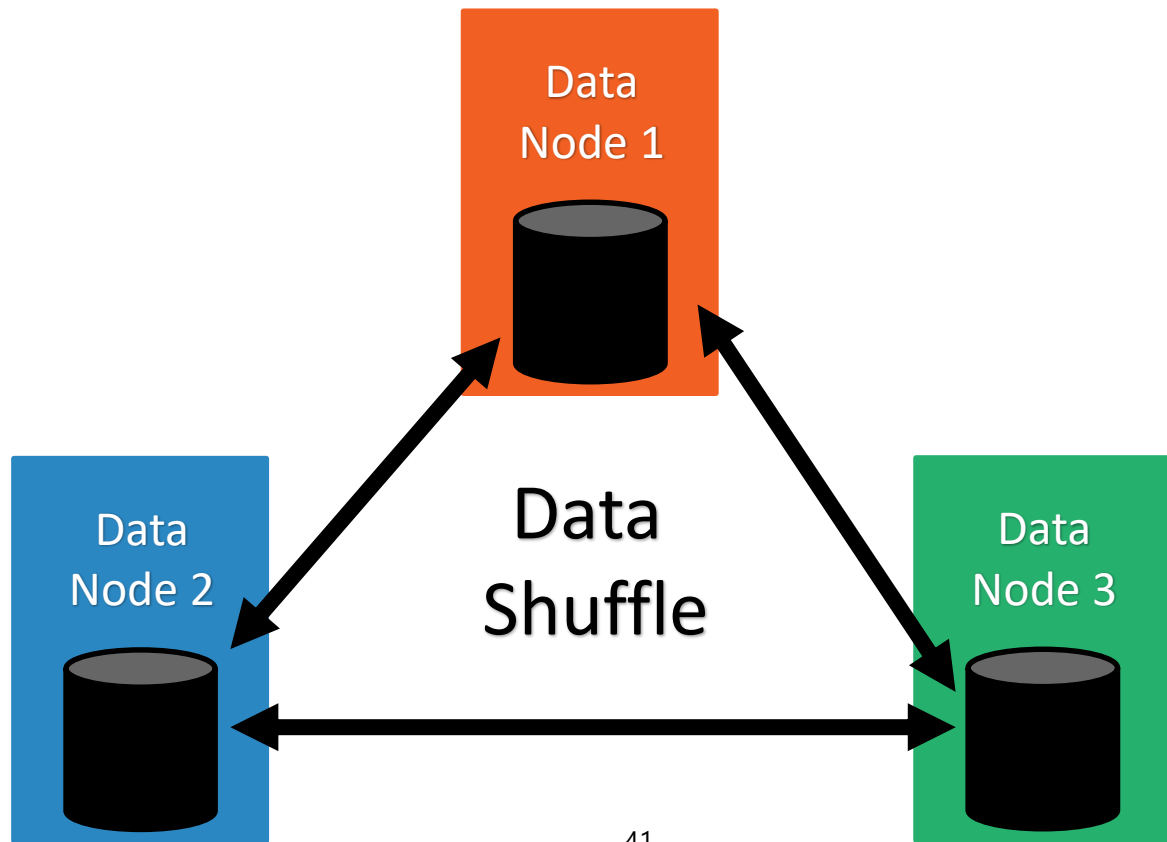data science for everyone

# What is Mahout?

- Distributed Machine Learning platform.

- Built on top of MapReduce and HDFS.

- Script-based and command line interfaces.

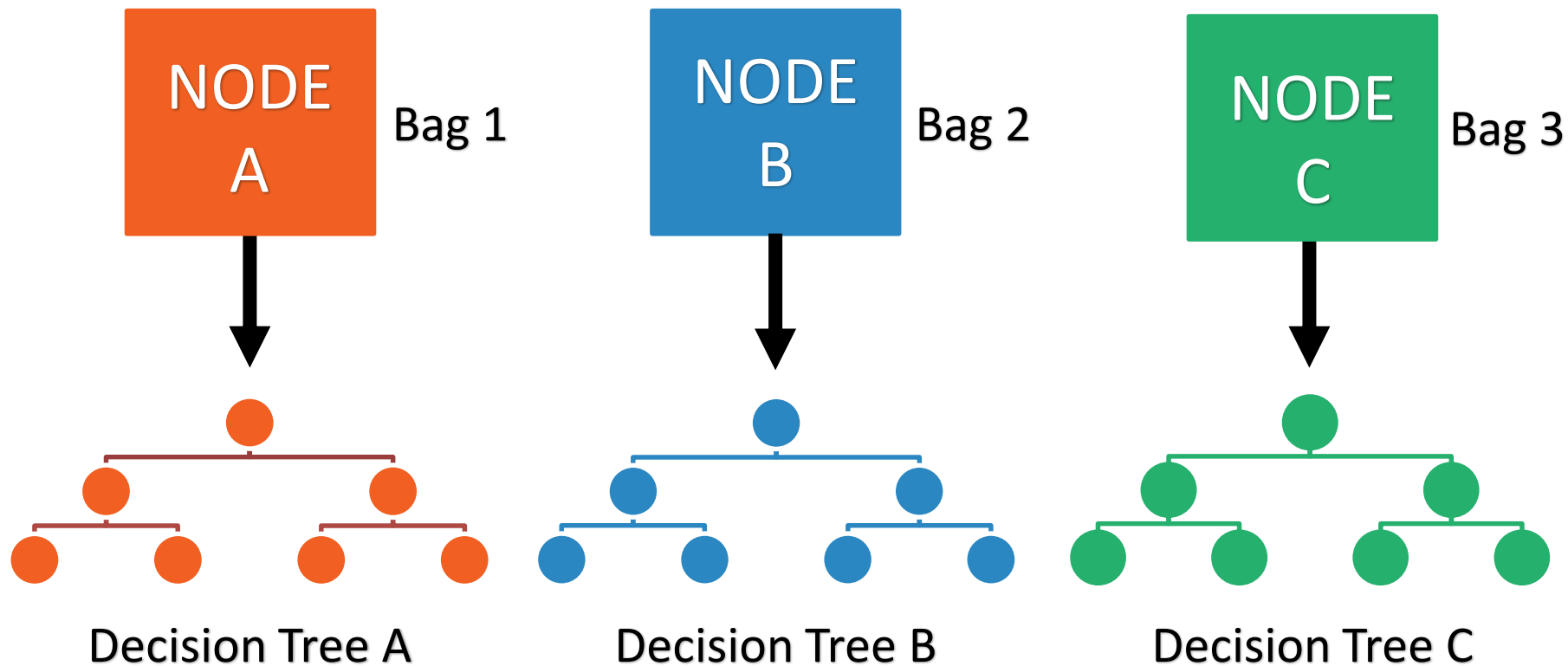- R-like language implementation.
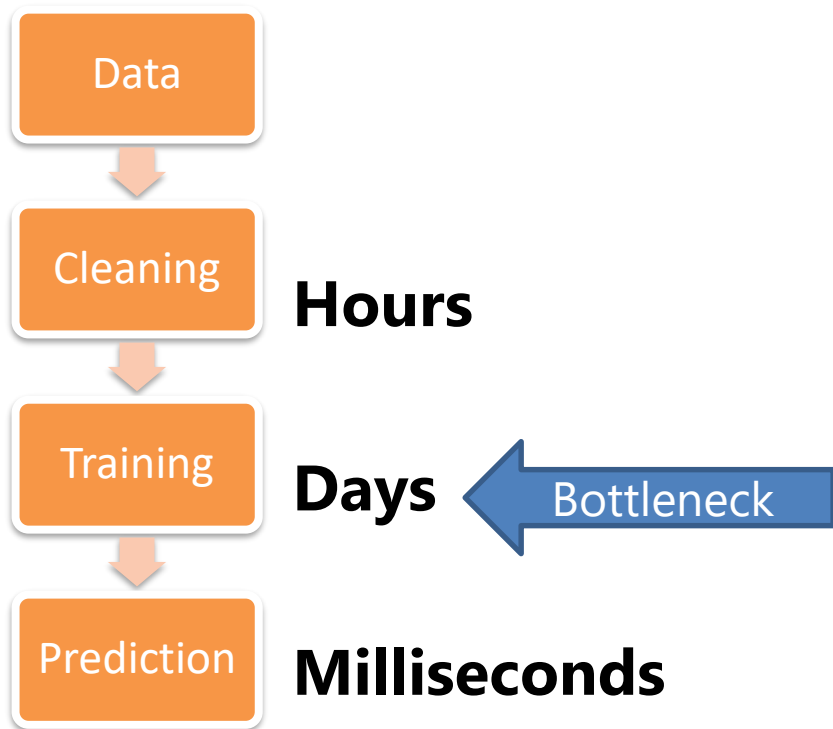
# Distributed Random Forest

# Distributed Random Forest



NODE A    Bag 1

NODE B    Bag 2

NODE C    Bag 3

Decision Tree A     Decision Tree B     Decision Tree C

# Processing Times - Machine Learning



**Data**

↓

**Cleaning** — **Hours**

↓

**Training** — **Days** ← Bottleneck

↓

**Prediction** — **Milliseconds**

- Large scale systems are only needed for training
- Phones can use models outputted by mahout to predict new data
- After a model is trained, save the model to any IO file type and reload it where you want
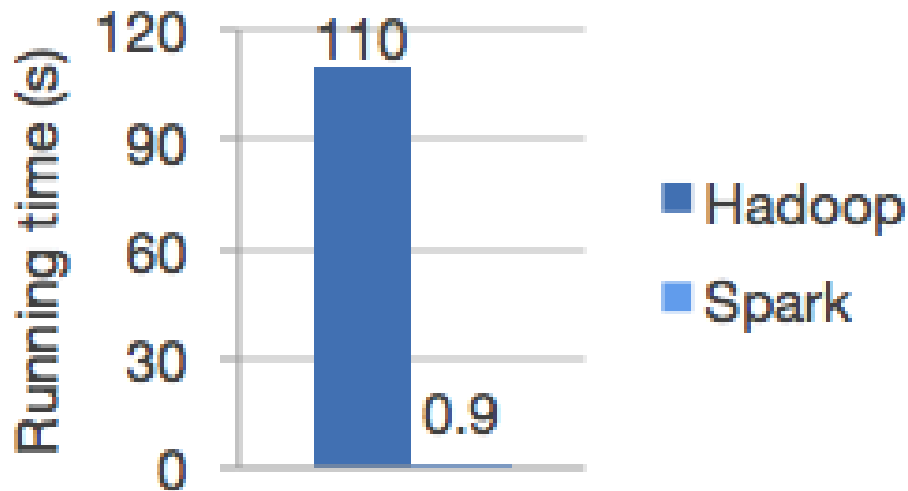
# DISTRIBUTED COMPUTING V2.0 – APACHE SPARK

# What is Spark?

- "A fast and general engine for large-scale data processing."

- Designed to incorporate the goodness of Hadoop and address Hadoop's shortcomings.

- Can complement Hadoop via integration with both HDFS and Hive.

# Why Spark? Improved Perf!

- Up to 10x faster than Hadoop working with data from disk.

- Up to 100x faster working with data stored in memory!

# Big Data, Faster!

**3x faster on 10x fewer machines!**

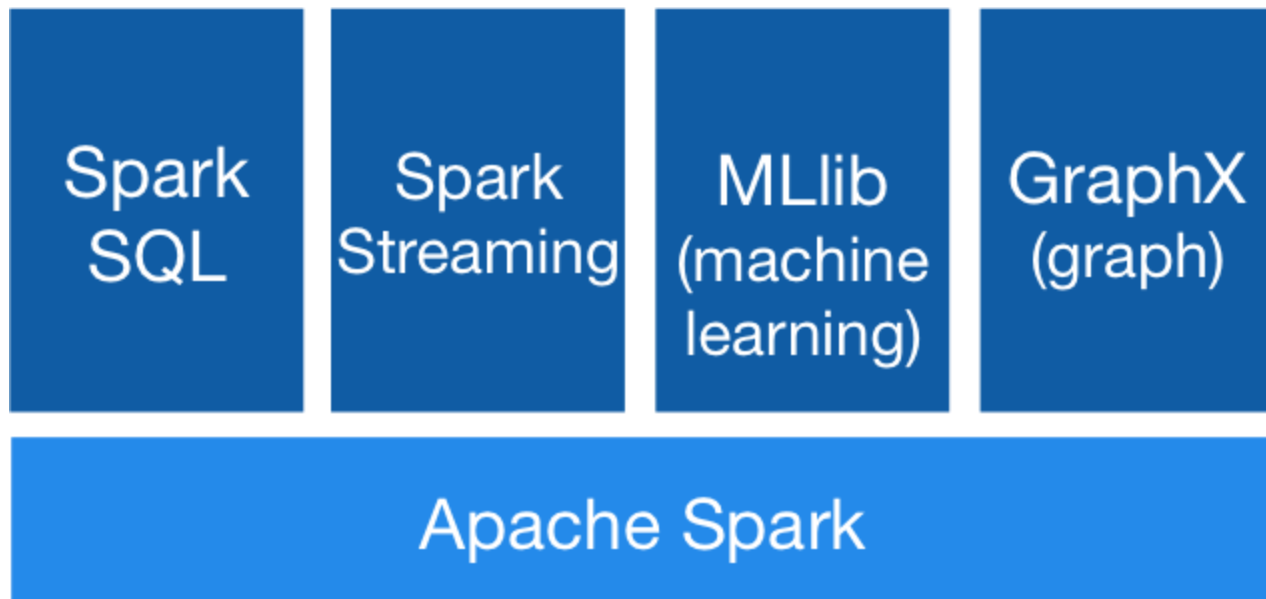Daytona GraySort Contest: Sort 100 TB of data!

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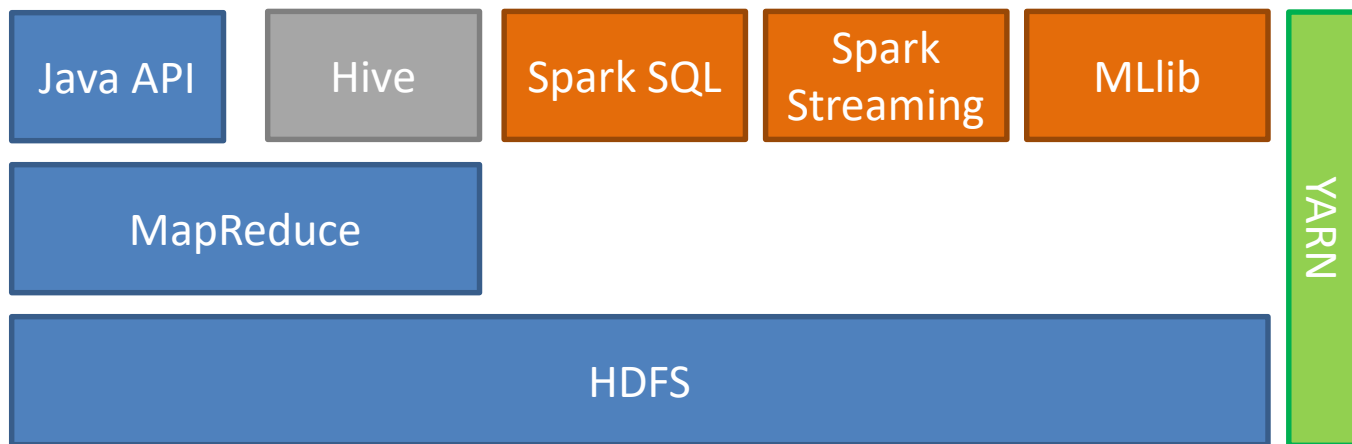- Yahoo!
- 72 Minutes
- 2100 Nodes

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

datasciencedojo
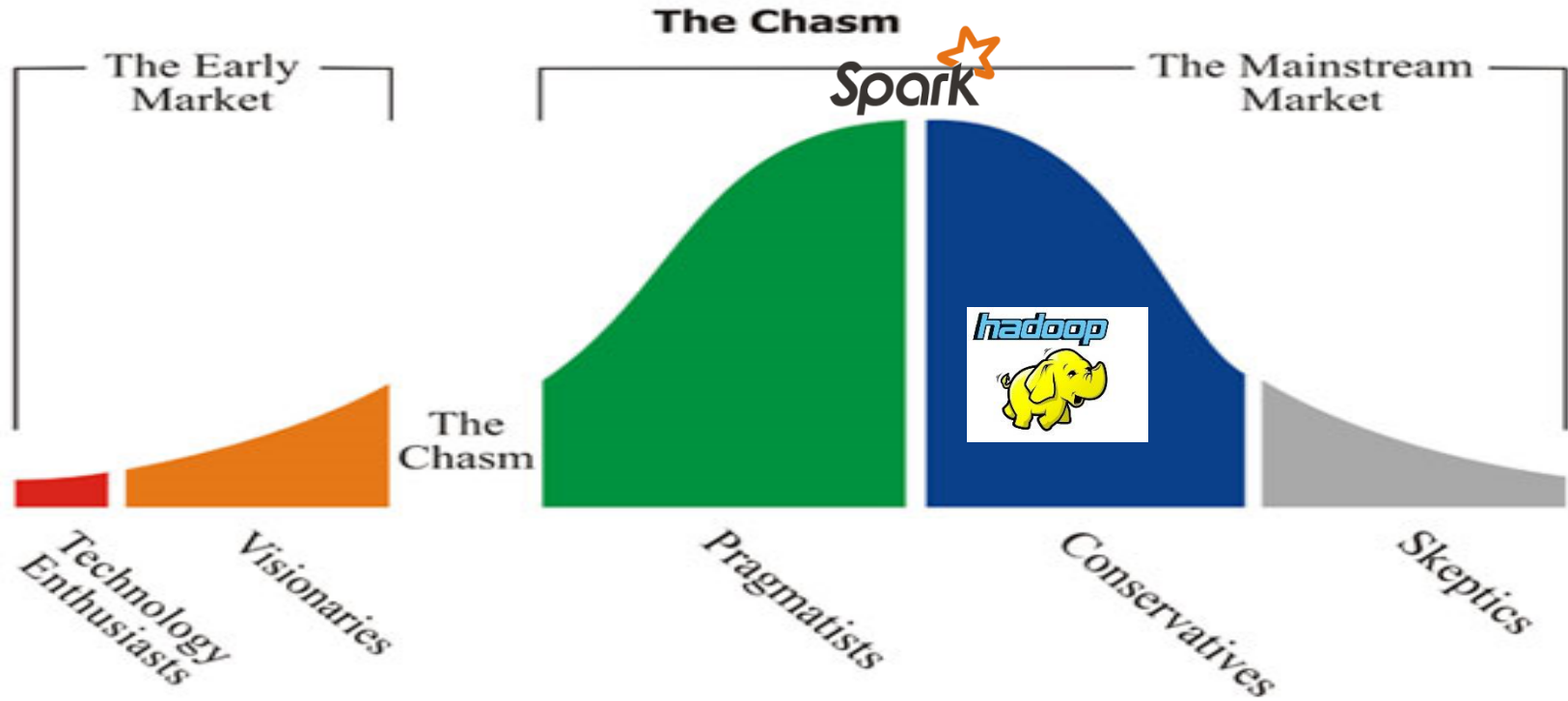data science for everyone

# Conceptual Architecture

# Spark and Hadoop



- Spark can be deployed on a Hadoop cluster and share cluster resources via YARN.
- Spark, however, does not require Hadoop!

datasciencedojo
data science for everyone

# Why is Spark Faster?

- First, Spark processing implements *lazy execution*:
  - Data operations are either *transformations* or *actions*.
  - Transformations are not executed immediately, but are stored.
  - When an action is issued, Spark evaluates all stored transformations and optimizes processing before executing.

- Second, Spark performs most processing in-memory:
  - RAM is far faster than using disk storage – even SSD drives.
  - More RAM in the cluster allows Spark to processes data faster.

datasciencedojo
data science for everyone

# Technology adoption life cycle



**The Chasm**

The Early Market

The Mainstream Market

The Chasm

Technology Enthusiasts

Visionaries

Pragmatists

Conservatives

Skeptics

Source: http://carlosmartinezt.com/2010/06/technology-adoption-life-cycle/

datasciencedojo
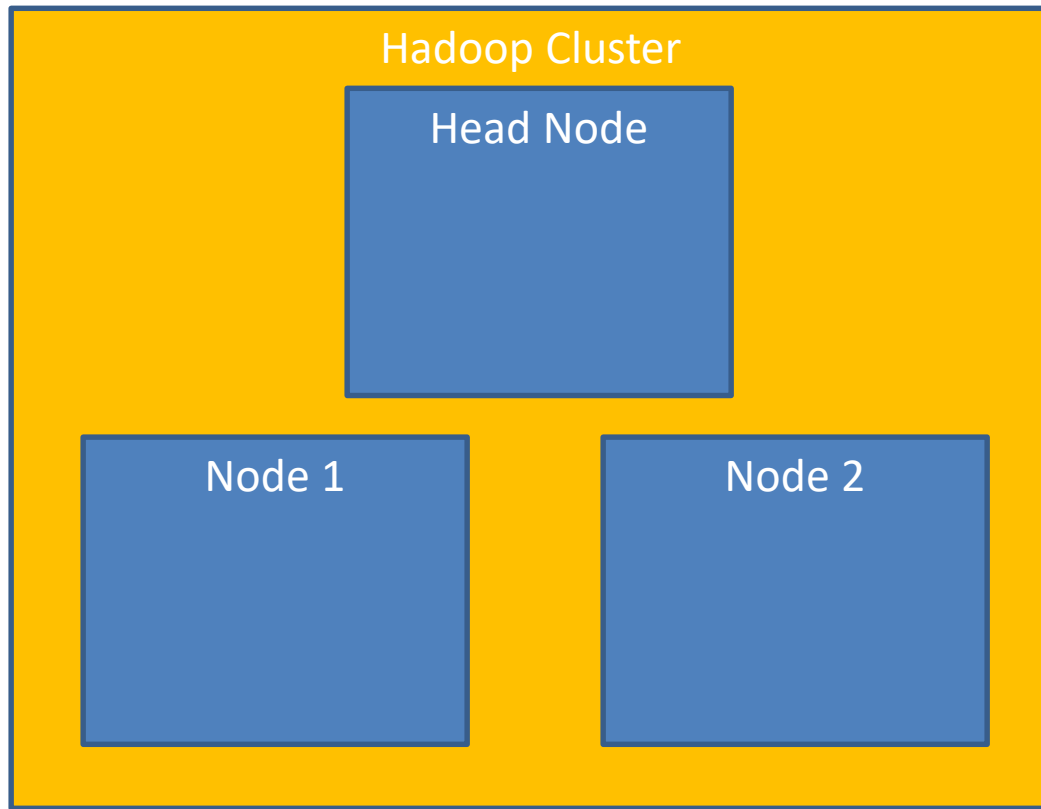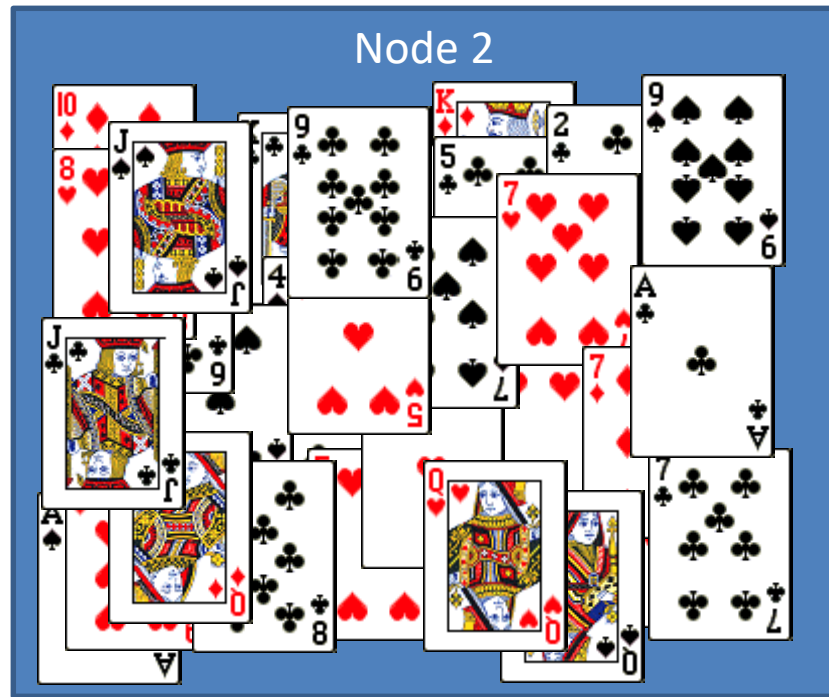data science for everyone

# QUESTIONS

# APPENDIX
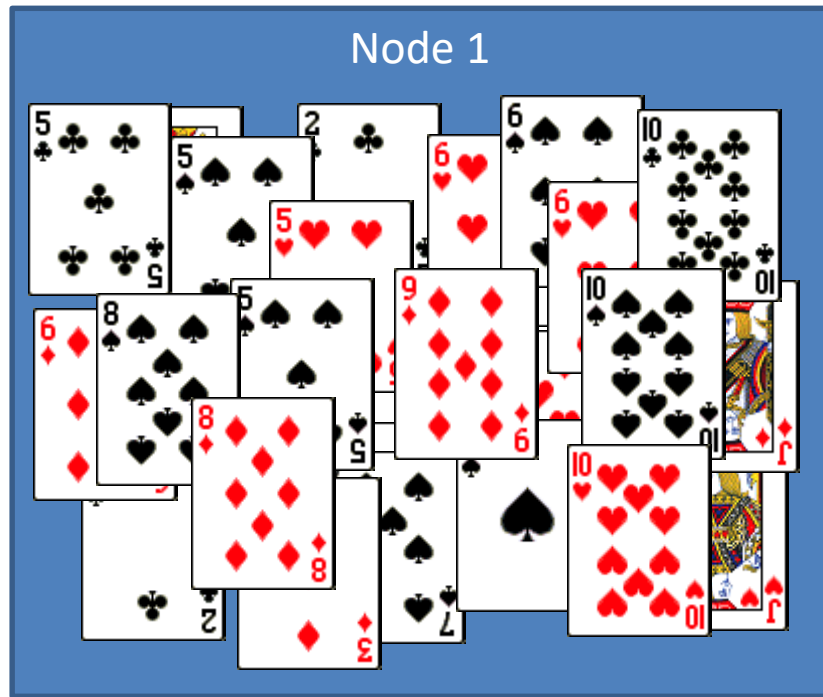
# MapReduce, via Playing Cards

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
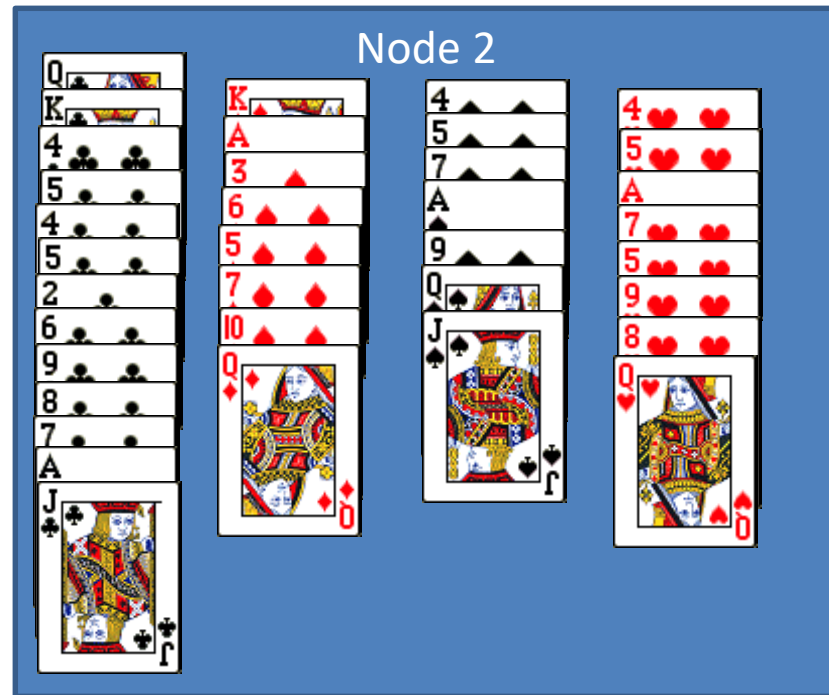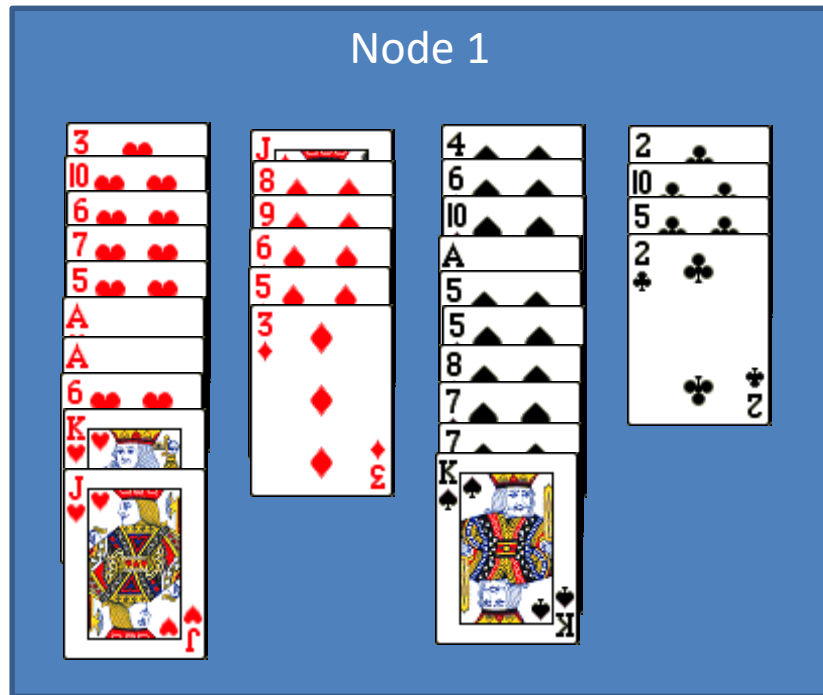- Each suit & number represents an attribute of the data
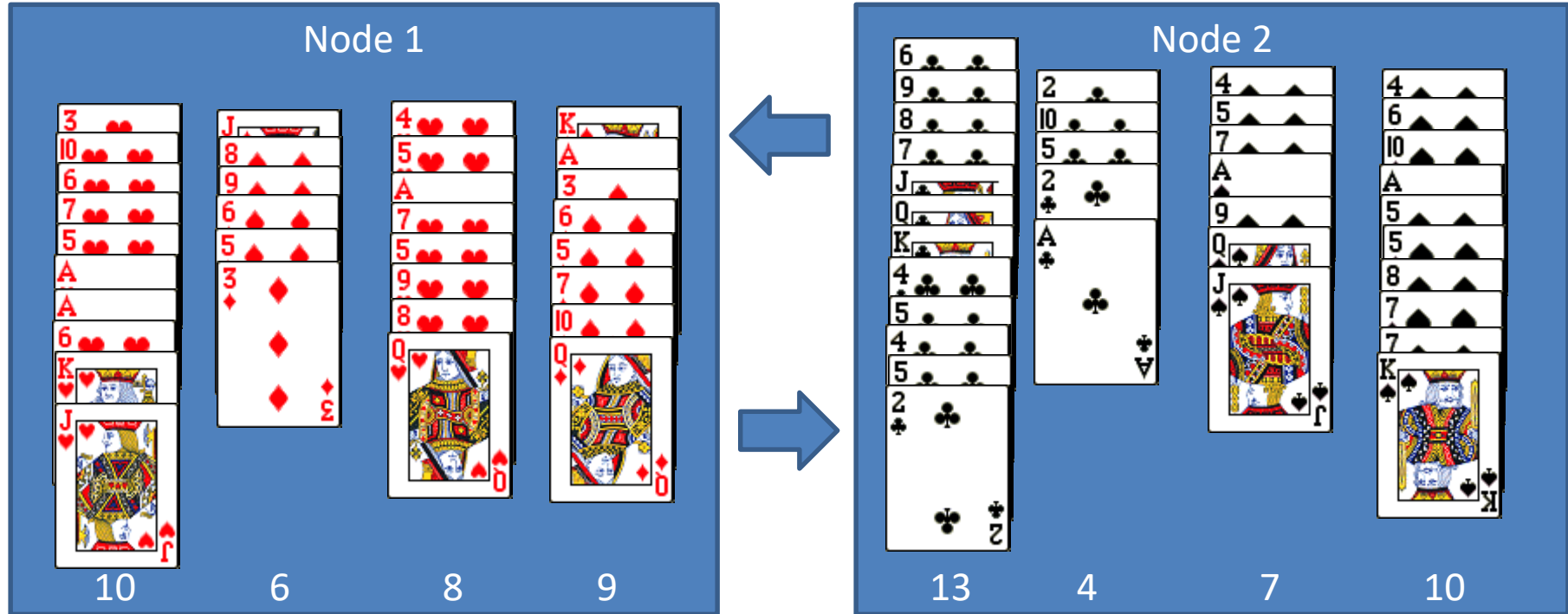
# Using a 2 Data Node Cluster

# 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