# Big Data Engineering With MapReduce and Hive

Data Science Dojo

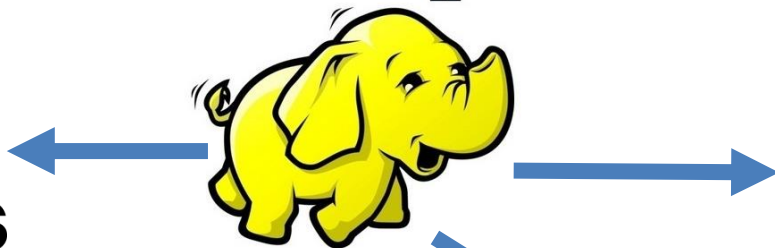datasciencedojo
unleash the data scientist in you

# Agenda

# Hadoop Implementations

# Hadoop

# Turn Back The Clock, The Mainframe

# Distributed Computing

# Cloud Computing
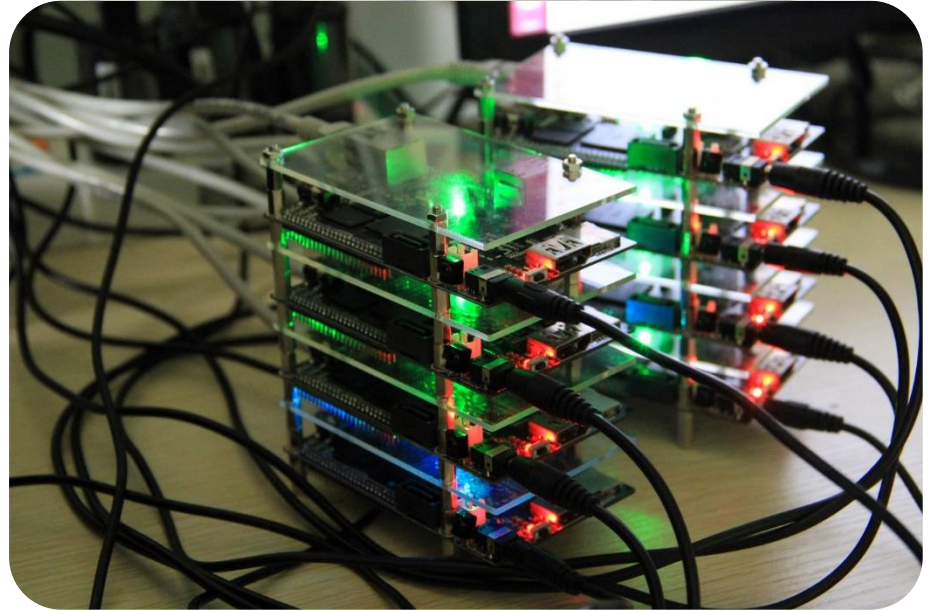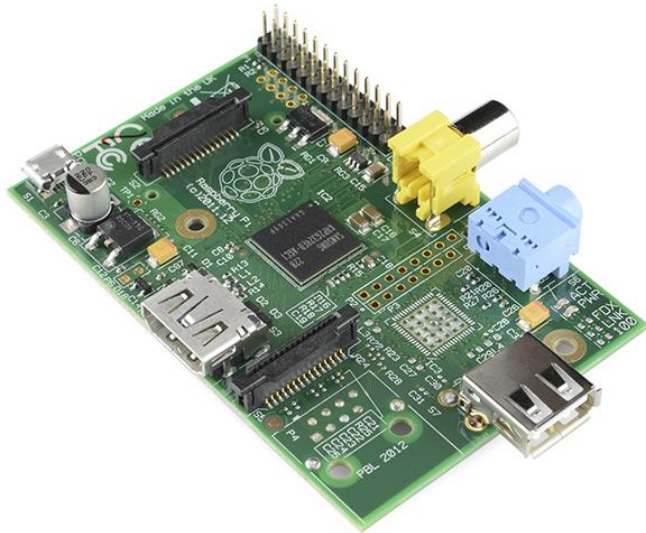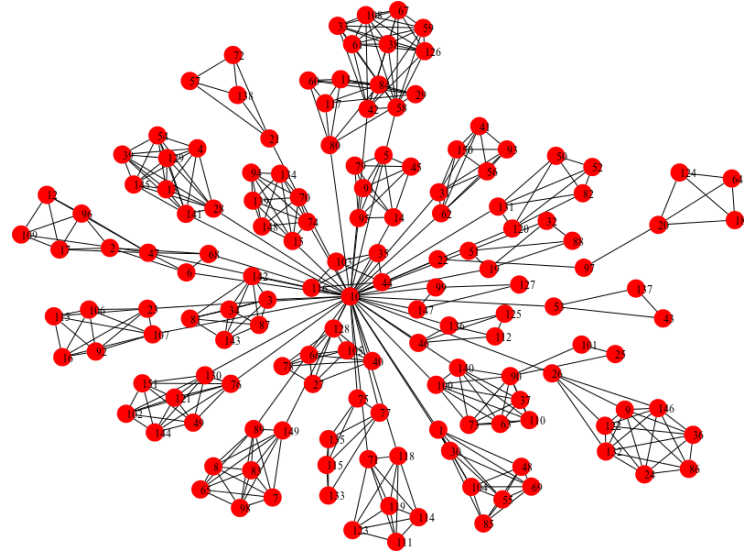
# If dogs were servers...

Head Node
(Named Node)

Data  Nodes

datasciencedojo
unleash the data scientist in you

# 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

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

# Limitations with MapReduce

- ~200 lines of code to do anything
- Slow
- Troubleshooting multiple computers
- Good devs are scarce
- Expensive certifications

```
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

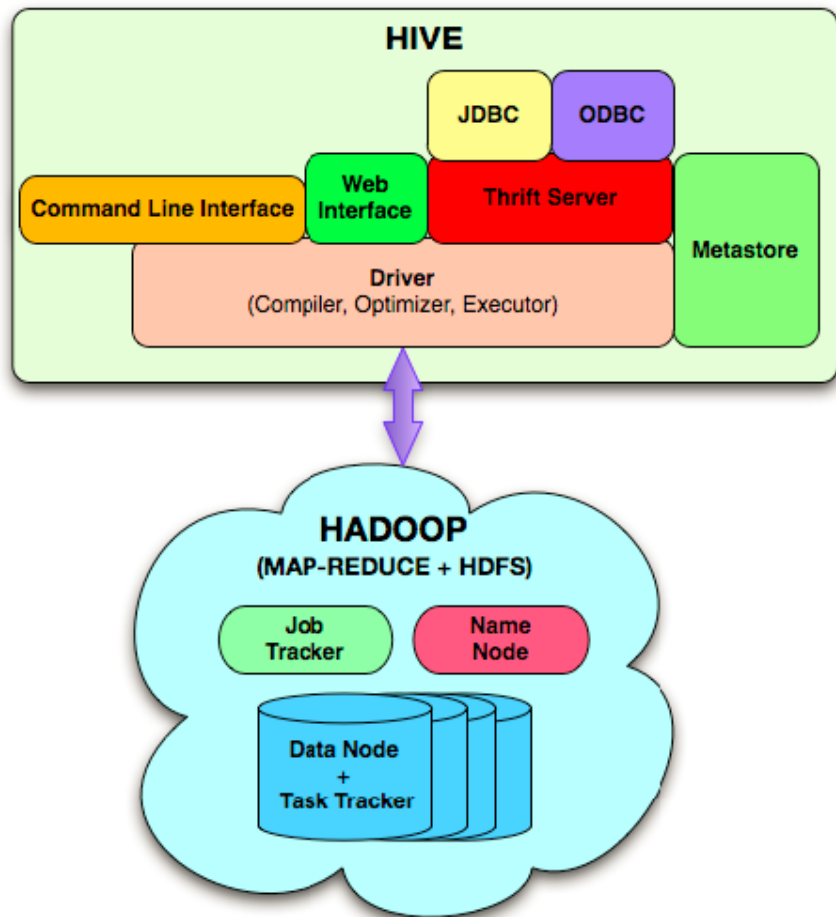**Pig:** Simpler scripting for MapReduce transformations

**Sqoop:** Data import and export

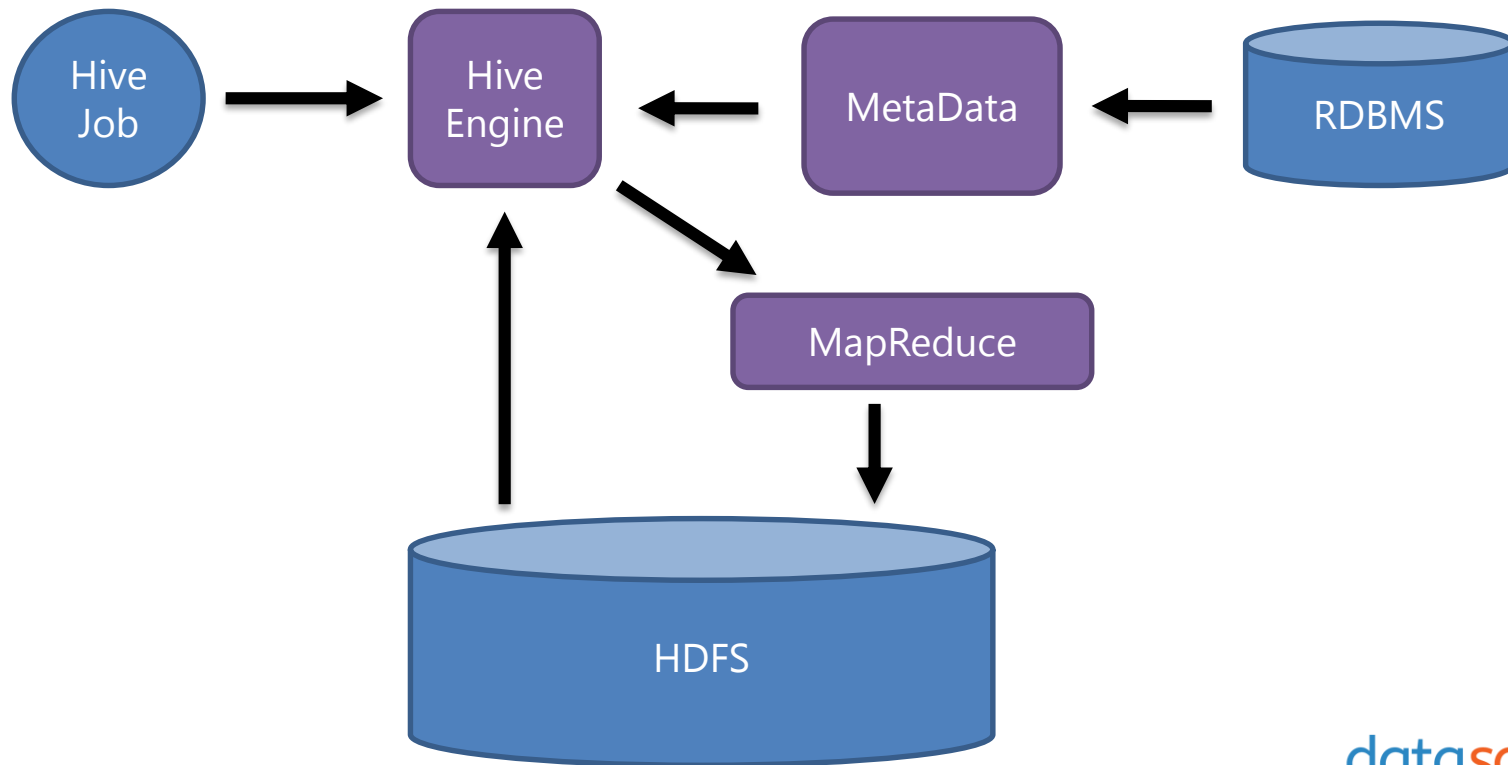**Storm:** Real-time processing of fast, large data streams

**Zookeeper:** Coordinates processes in distributed systems

# Hive Jobs

HiveQL Statement → Translation & Conversion → MapReduce Job

# 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
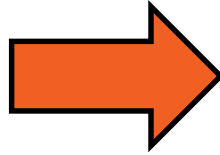
# Limitations



80% of organized data is unstructured

## Structured vs. Unstructured Data Growth

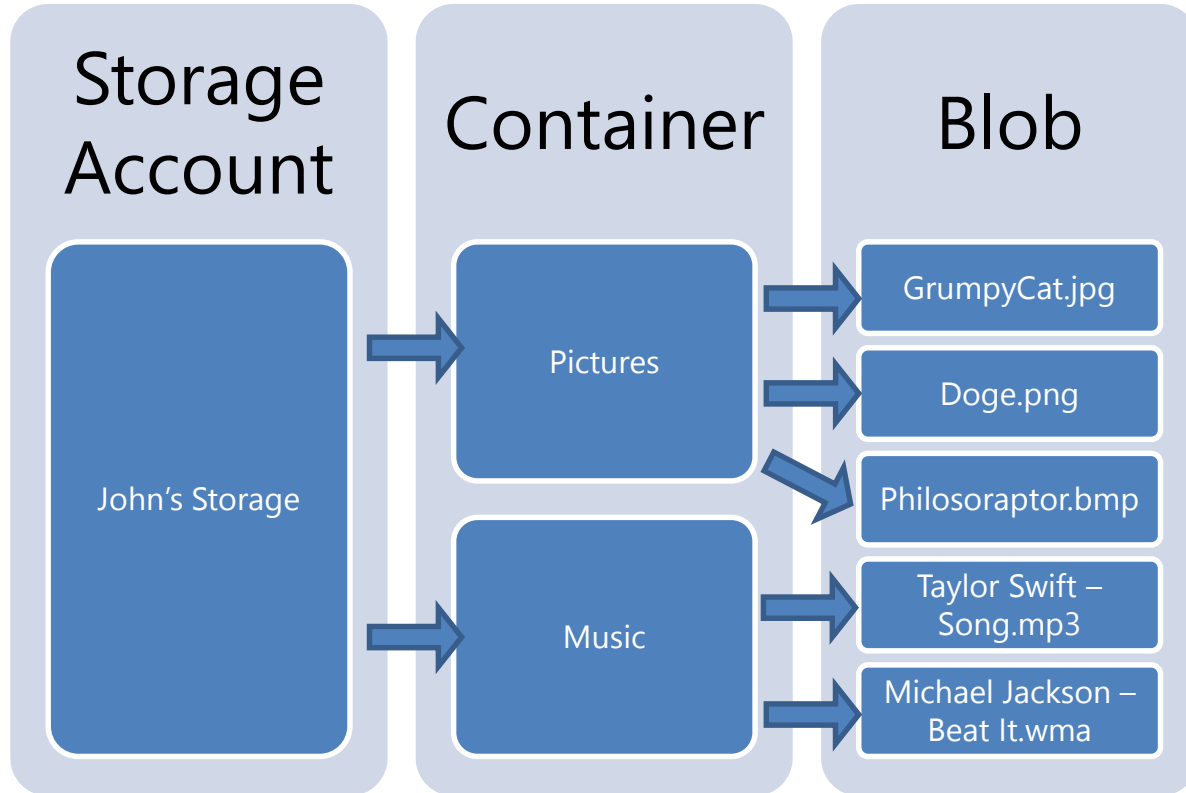| 2013 | 2015 | 2017 |
|------|------|------|
| 67% | 74% | 80% |
| 37 EB | 71 EB | 133 EB |

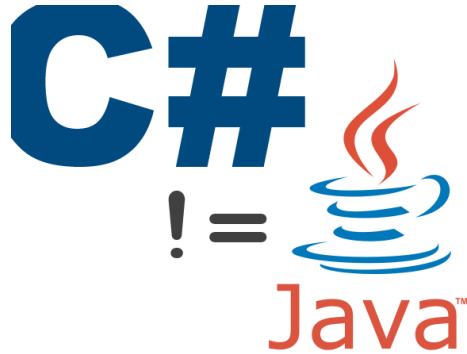Total Capacity Shipped, Worldwide    Unstructured Data
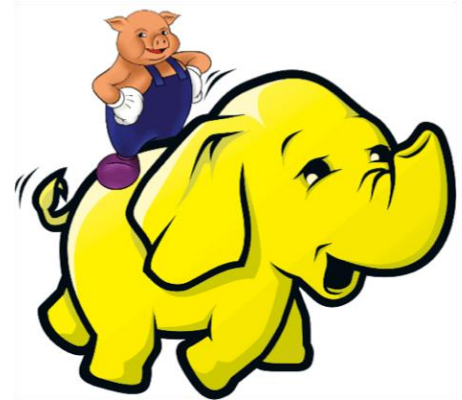
Source: IDC

# Azure Blob Storage

# When to Use Each



C#
Java
MapReduce
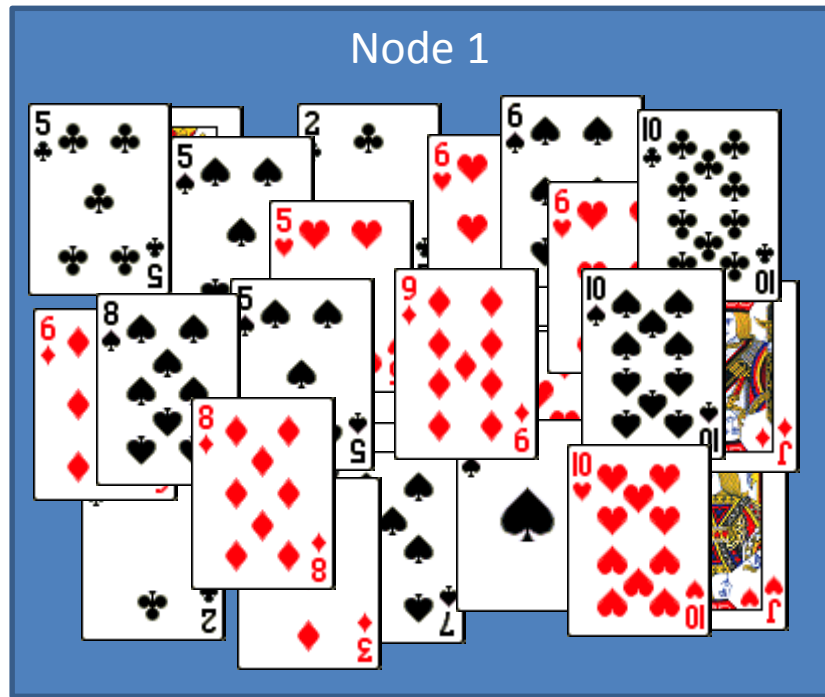
VS

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
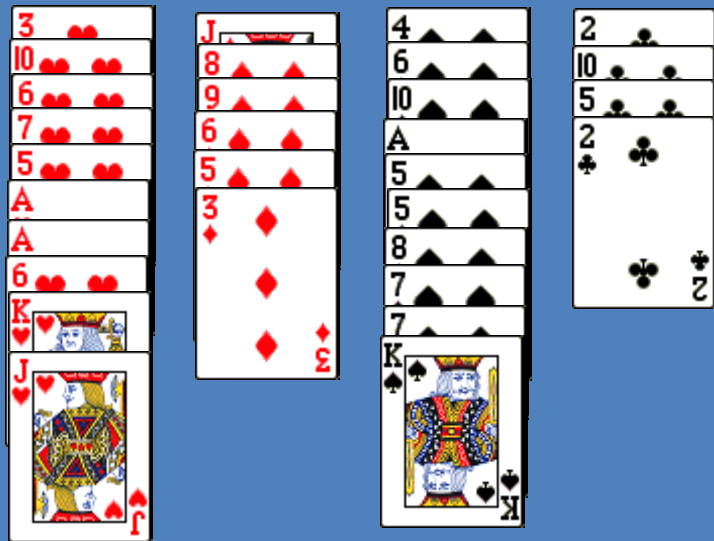- Each suite represents an attribute of the data

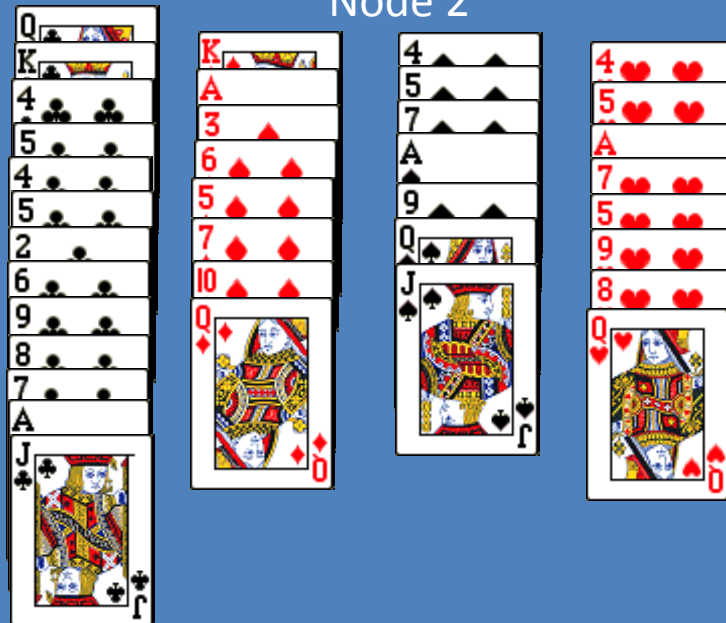# Using a 2 Data Node Cluster

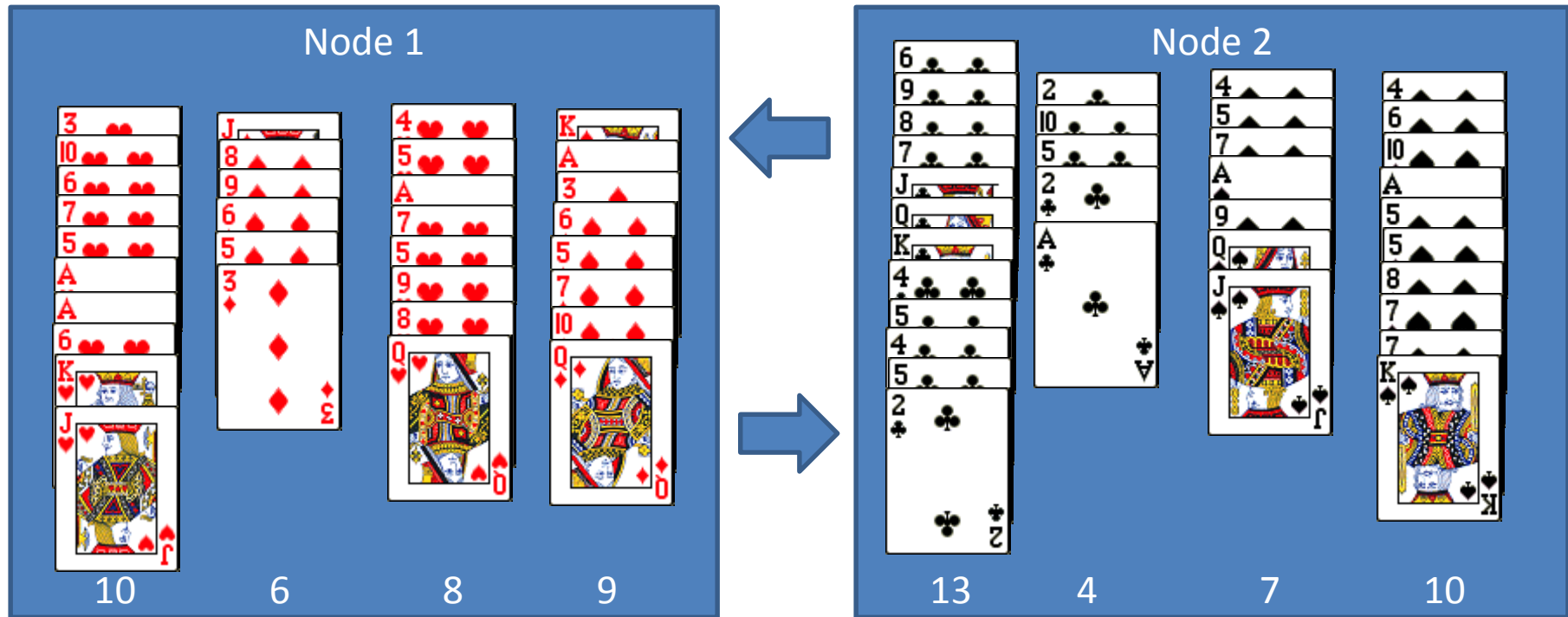# Mapping: Each Node's HDFS
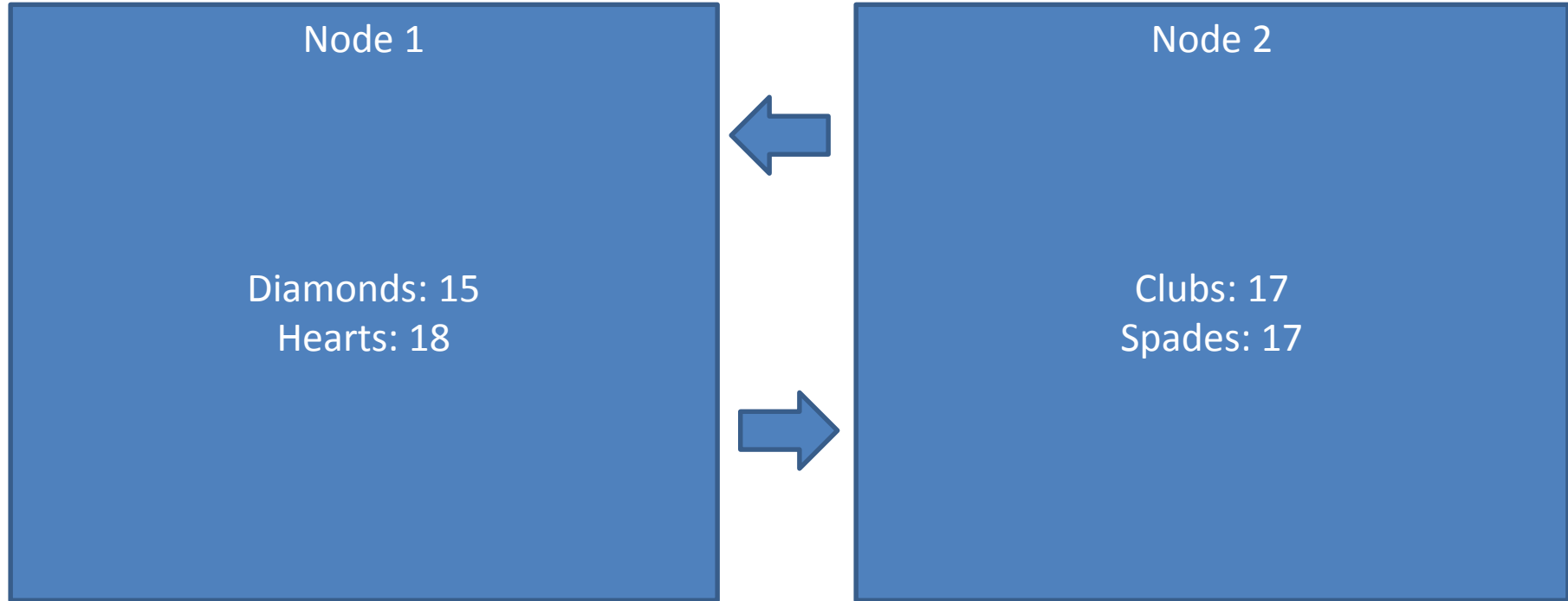
# Mapping: Node Sorting

# Mapping: Node Shuffle, Data Transfer

# Mapping: Node Shuffle, Data Transfer

# QUESTIONS