…An Introduction

# WHY?

- Volume - so much

- Velocity - so fast

- Variety - so many types

Traditional

**Expensive
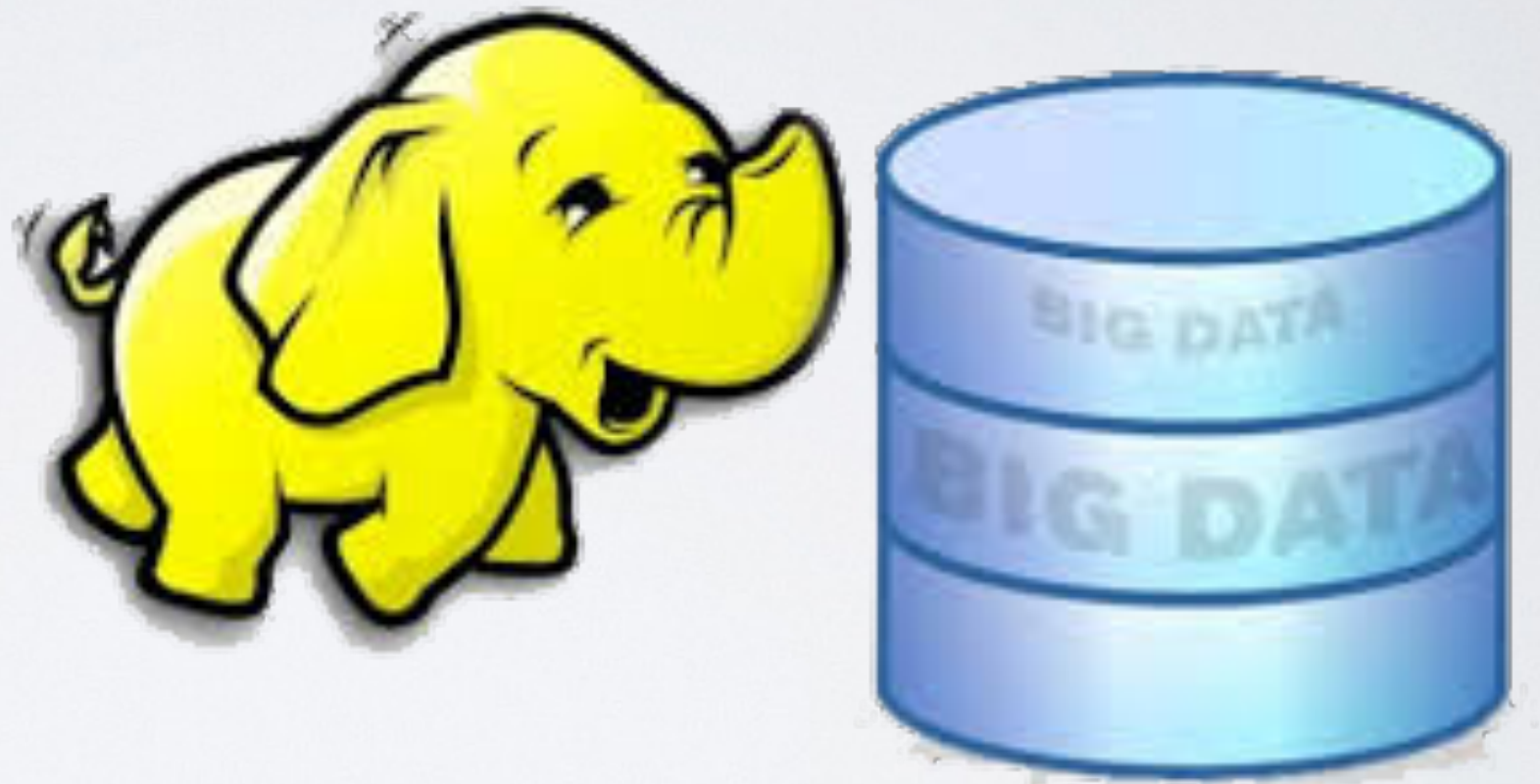Difficult
Complex**

VS.

**Cheap
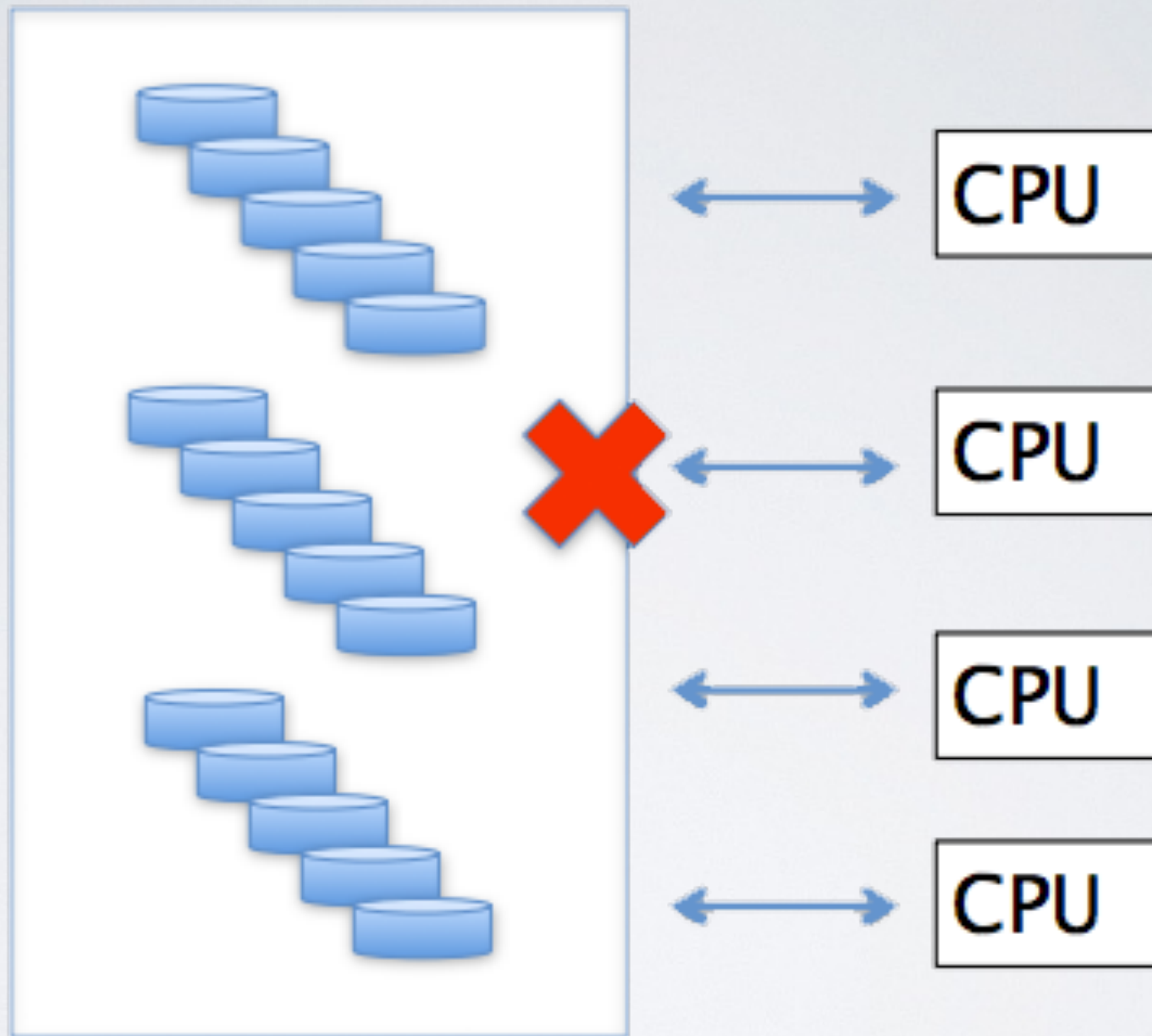Simple
Easy**

Hadoop

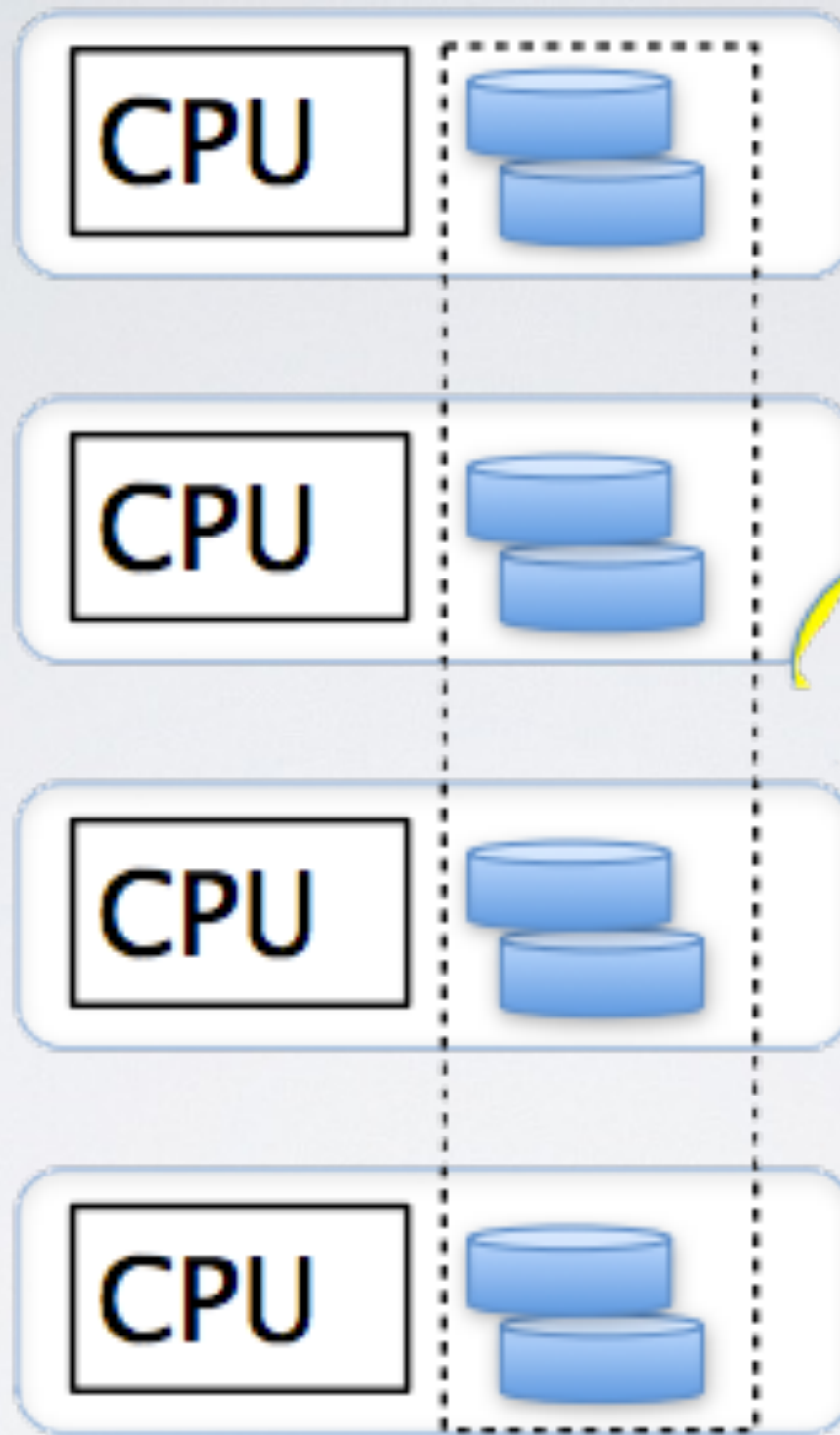The old way just wasn't cutting it

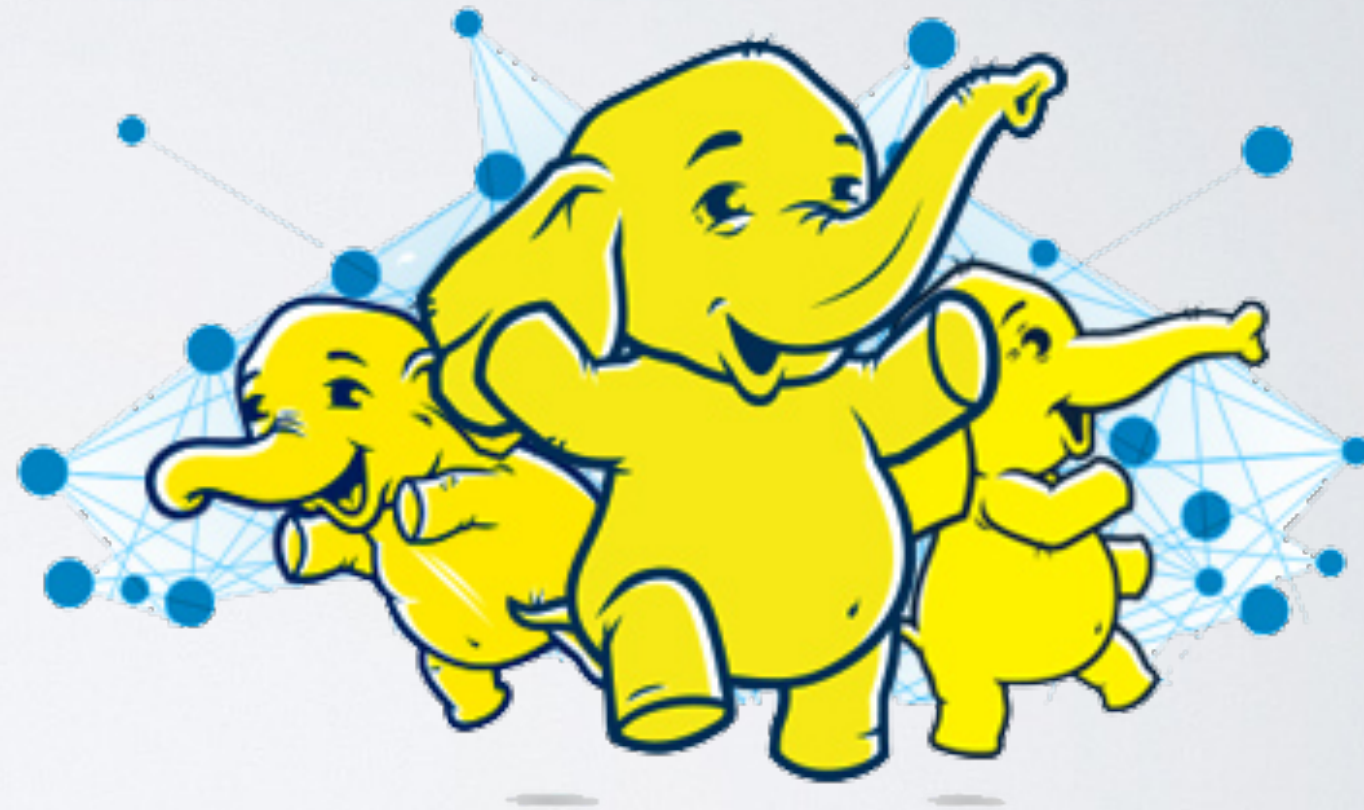# SO HOW DO WE DO IT?

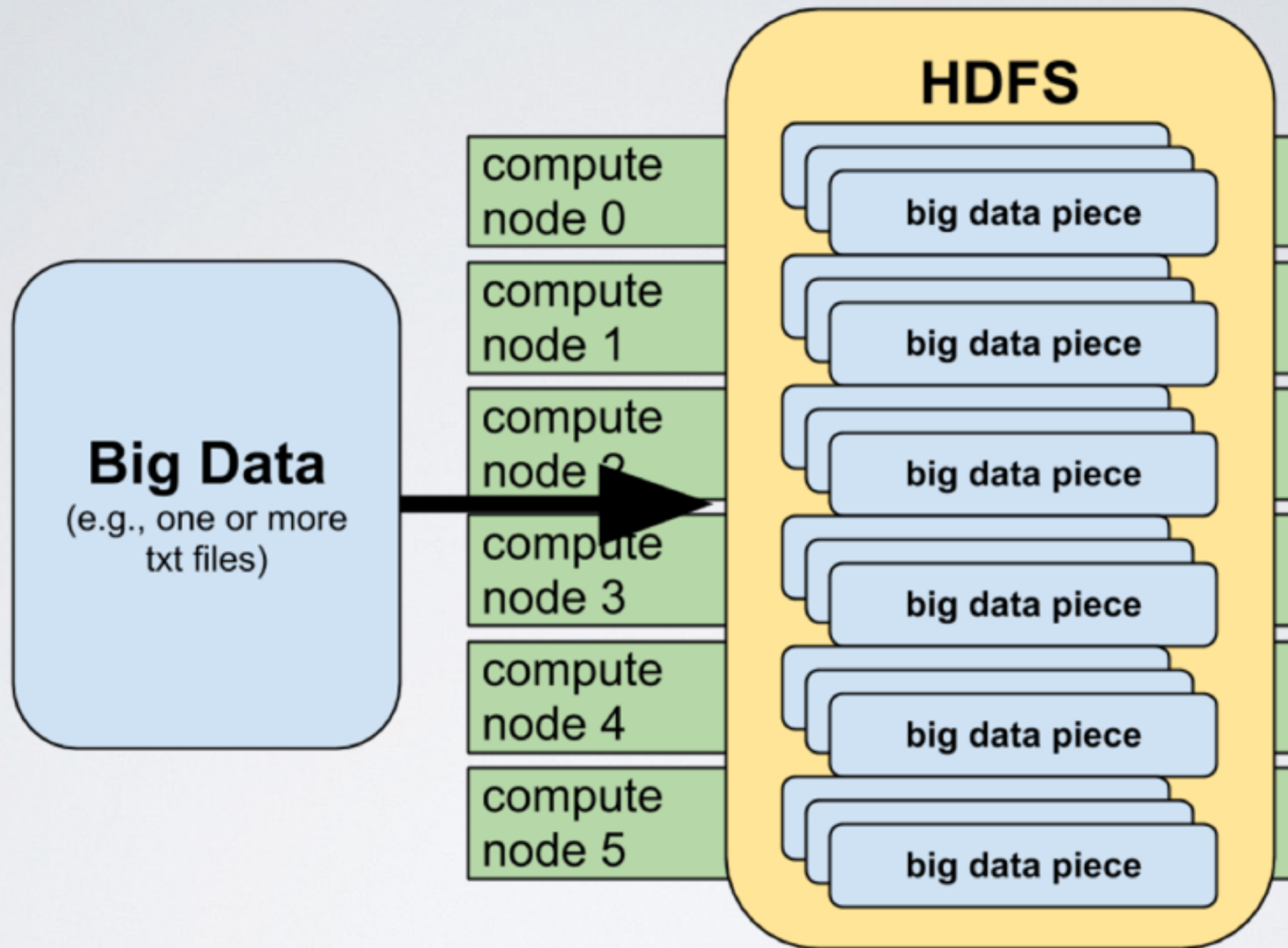- HDFS

- Map Reduce
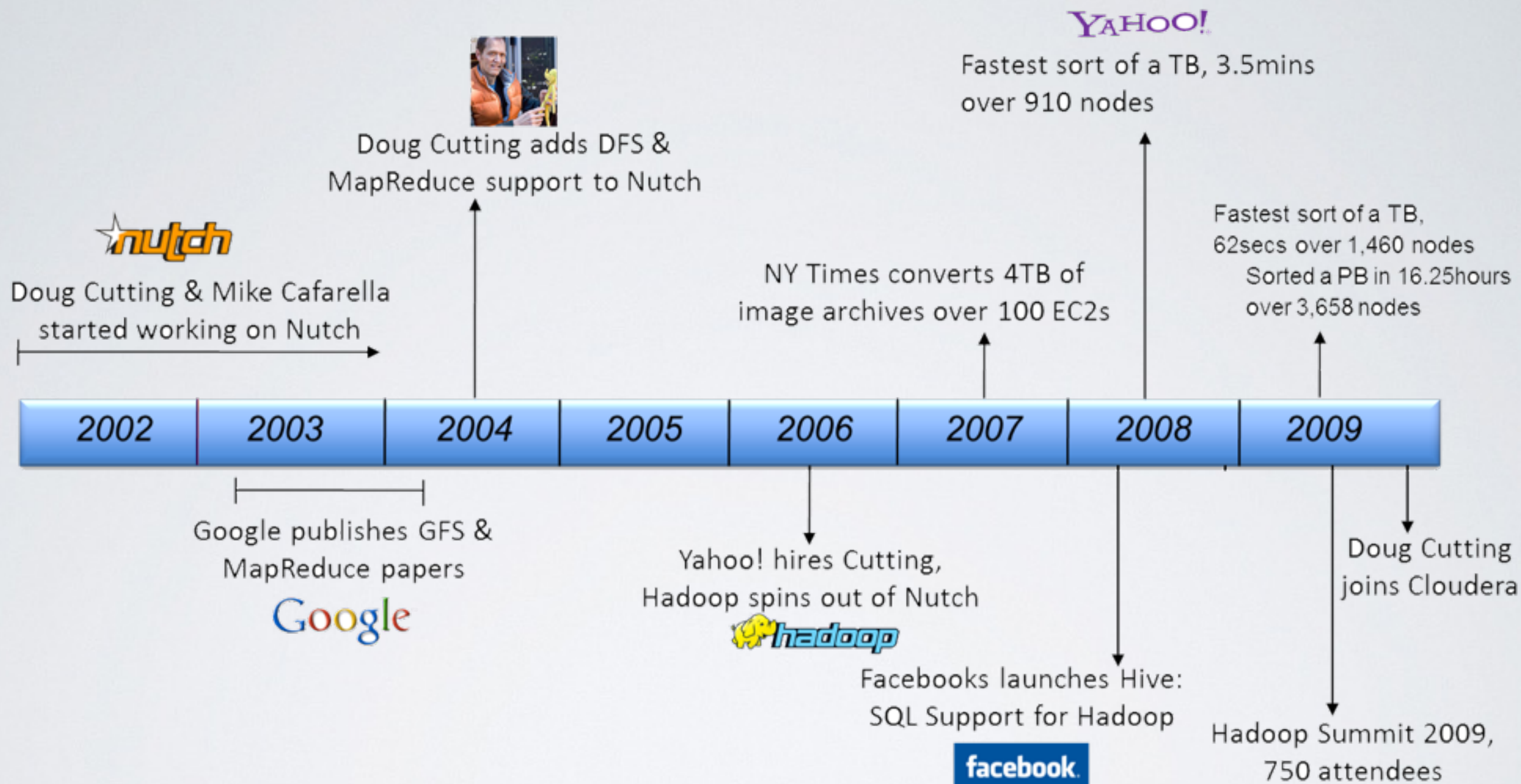
It just doesn't scale

That's better

# ENTER HADOOP

- Open source Java

- Framework for big data

- Scalable

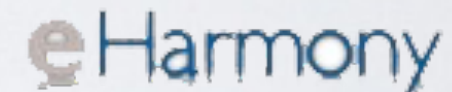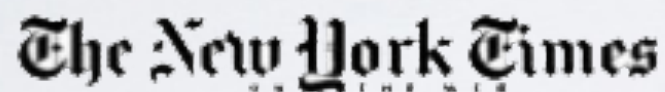- Fault tolerant

- Commodity hardware

Distribute the data and the processing

A brief history

Some hadoop users

# HDFS

- Master/slave architecture

- Fault tolerant

- Commodity hardware

- Streaming access to FS data

# HDFS - BLOCKS

- Splits data into blocks

- Easier to distribute work

# HDFS - REPLICATION

- Blocks replicated across nodes

- Fault tolerant

- 3x replications by default

Client ask NN for file

NN returns DNs that has it

Client ask DN for data

H D F S

**Namenode (NN)**

**Datanode 1**   **Datanode 2**   **Datanode N**

## Namenode - Master
- FS metadata
- Block replication
- R/W access to files

## Datanode - Slave
- Contains data
- Informs namenode

# Map

```
def sqr(n):
    return n * n

list = [1,2,3,4]

map(sqr, list) -> [1,4,9,16]
```

# Reduce

```
def add(i, j):
    return i + j

list = [1,2,3,4]

reduce(add, list) -> 10
```

# MapReduce

```
def MapReduce(data, mapper, reducer):
    return reduce(reducer, map(mapper, data))

MapReduce(list, sqr, add) -> 30
```

```java
public static class Map extends MapReduceBase implements Mapper<LongWritable, Text, Text,
IntWritable> {
  private final static IntWritable one = new IntWritable(1);
  private Text word = new Text();

  public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output,
Reporter reporter) throws IOException {
    String line = value.toString();
    StringTokenizer tokenizer = new StringTokenizer(line);
    while (tokenizer.hasMoreTokens()) {
     word.set(tokenizer.nextToken());
     output.collect(word, one);
    }
   }
  }
```

# MAP

wordcount.java

```java
public static class Reduce extends MapReduceBase implements Reducer<Text, IntWritable, Text,
IntWritable> {
  public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text,
IntWritable> output, Reporter reporter) throws IOException {
    int sum = 0;
    while (values.hasNext()) {
     sum += values.next().get();
    }
    output.collect(key, new IntWritable(sum));
  }
}
```

# REDUCE

wordcount.java

```
$ mkdir wordcount_classes
$ javac -classpath ${HADOOP_HOME}/hadoop-${HADOOP_VERSION}-core.jar -d wordcount_classes
WordCount.java
$ jar -cvf /usr/joe/wordcount.jar -C wordcount_classes/ .
```

Assuming that:

- /usr/joe/wordcount/input - input directory in HDFS
- /usr/joe/wordcount/output - output directory in HDFS

Sample text-files as input:

```
$ bin/hadoop dfs -ls /usr/joe/wordcount/input/
/usr/joe/wordcount/input/file01
/usr/joe/wordcount/input/file02

$ bin/hadoop dfs -cat /usr/joe/wordcount/input/file01
Hello World Bye World

$ bin/hadoop dfs -cat /usr/joe/wordcount/input/file02
Hello Hadoop Goodbye Hadoop
```

Run the application:

```
$ bin/hadoop jar /usr/joe/wordcount.jar org.myorg.WordCount /usr/joe/wordcount/input
/usr/joe/wordcount/output
```

Output:

```
$ bin/hadoop dfs -cat /usr/joe/wordcount/output/part-00000
Bye 1
Goodbye 1
Hadoop 2
Hello 2
World 2
```

# RUNNING IT

wordcount.java