



Why Scala Is Taking Over the Big Data World



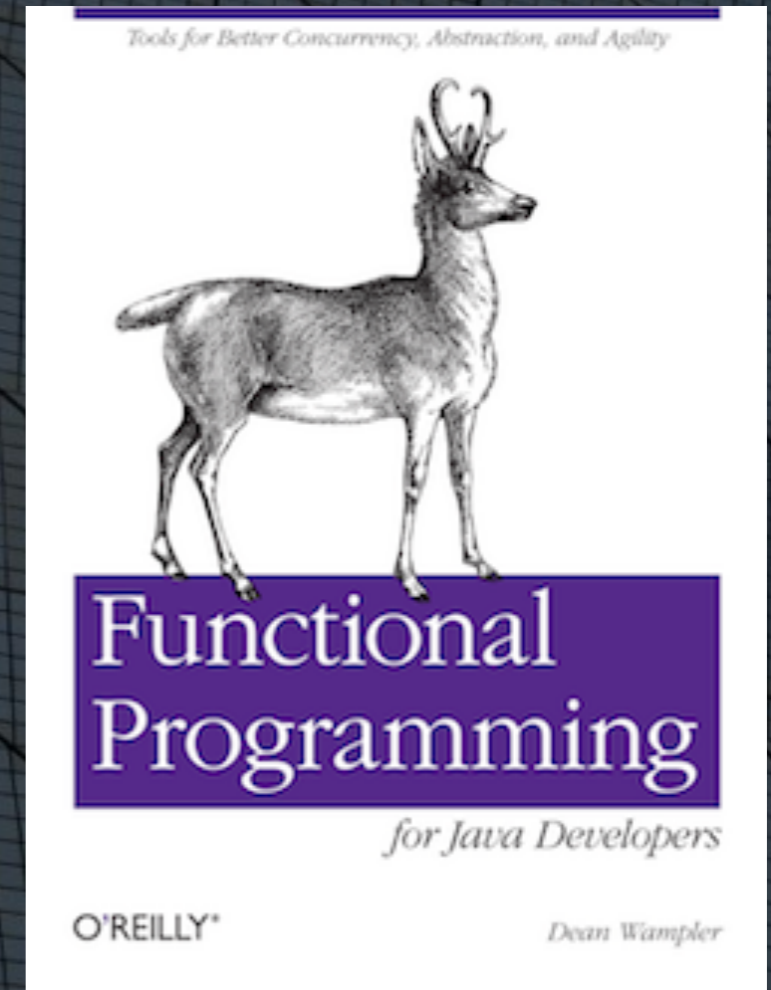
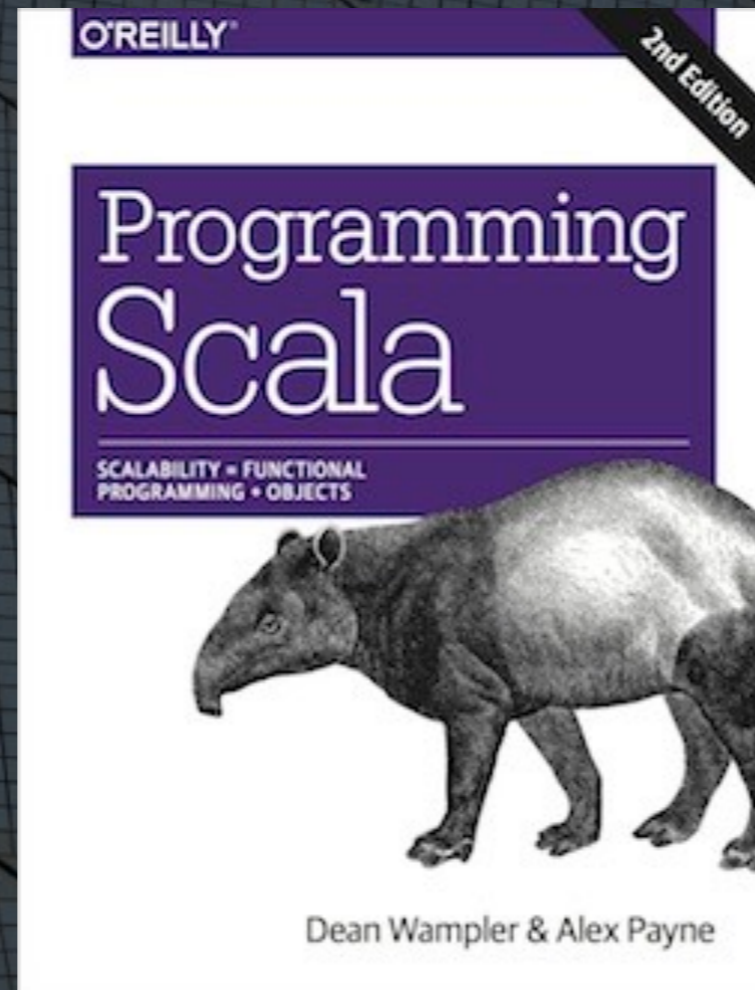
Saturday, January 10, 15

Sydney Opera House photos Copyright © Dean Wampler, 2011-2015, Some Rights Reserved.

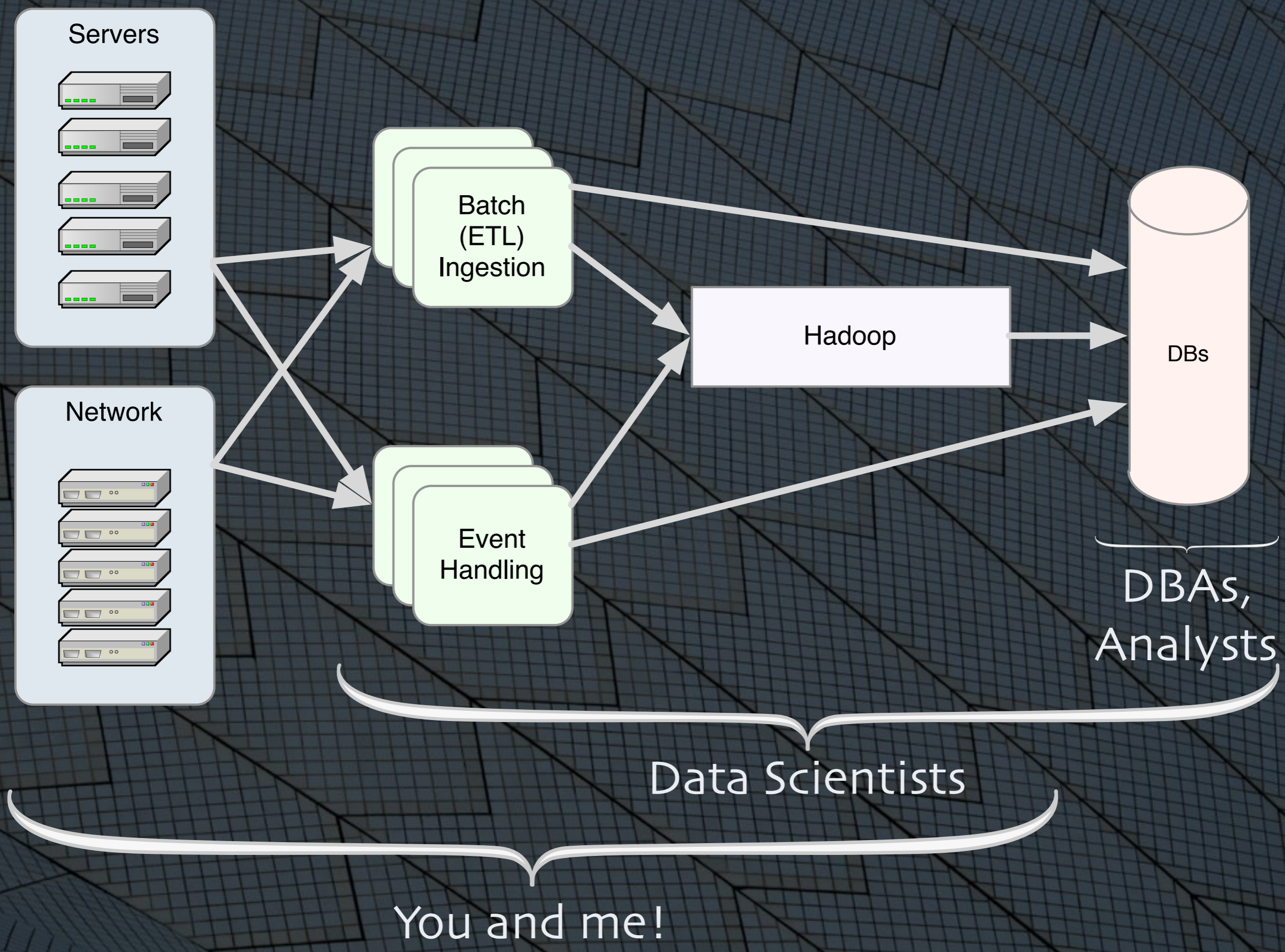
The content is free to reuse, but attribution is requested.

<http://creativecommons.org/licenses/by-nc-sa/2.0/legalcode>

<plug>
<shameless>

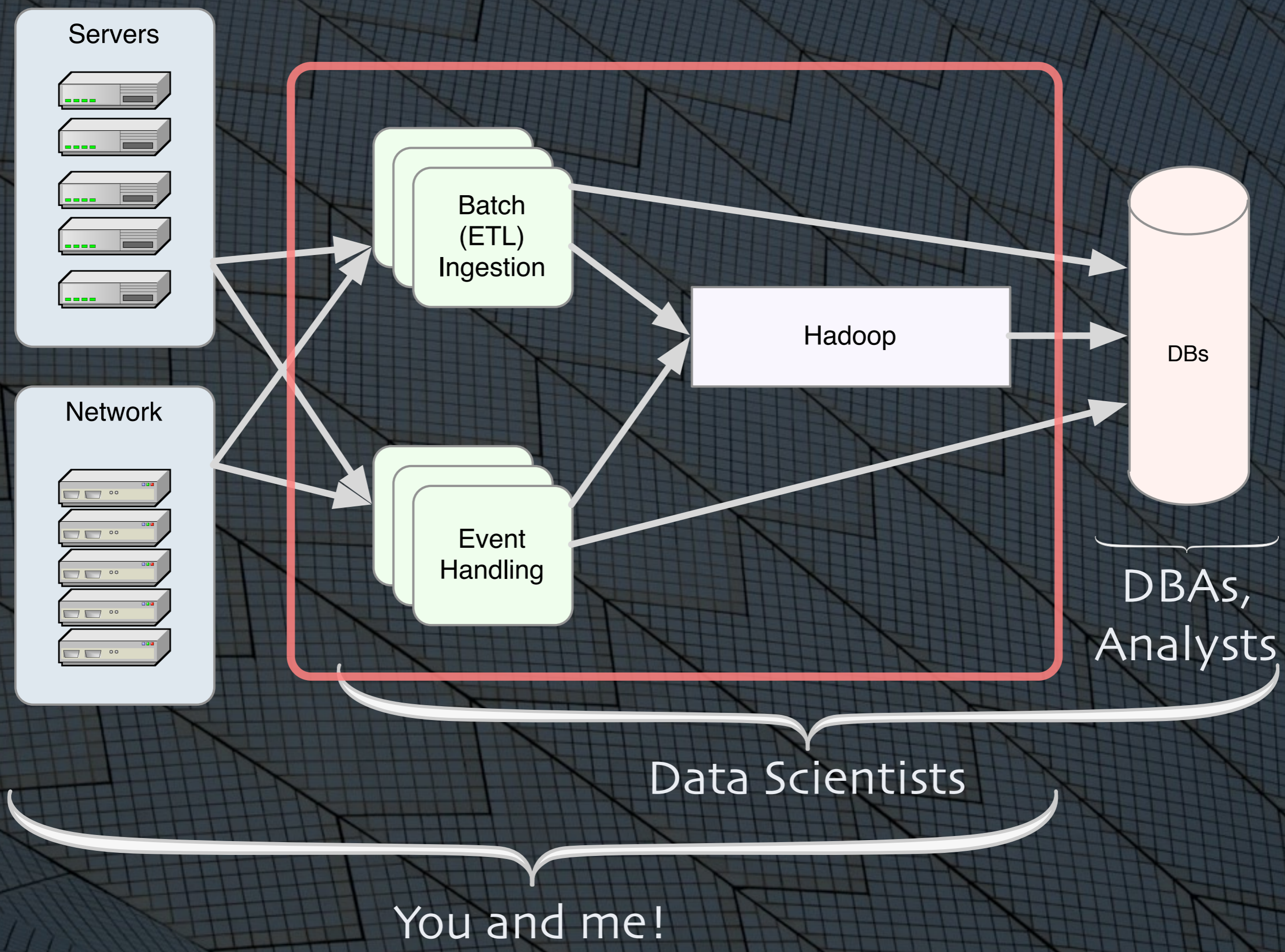


</shameless>
</plug>



Saturday, January 10, 15

DBAs, traditional data analysts: SQL, SAS. "DBs" could also be distributed file systems.
 Data Scientists - Statistics experts. Some programming, especially Python, R, Julia, maybe Matlab, etc.
 Developers like us, who figure out the infrastructure (but don't usually manage it), and write the programs that do batch-oriented ETL (extract, transform, and load), and more real-time event handling.
 Often this data gets pumped into Hadoop or other compute engines and data stores, including various SQL and NoSQL DBs, and file systems.





You are here...

5

Saturday, January 10, 15

Let's put all this into perspective...

http://upload.wikimedia.org/wikipedia/commons/thumb/8/8f/Whole_world_-_land_and_oceans_12000.jpg/1280px-Whole_world_-_land_and_oceans_12000.jpg



... and it's 2008.

Saturday, January 10, 15

Let's put all this into perspective, circa 2008...

http://upload.wikimedia.org/wikipedia/commons/thumb/8/8f/Whole_world_-_land_and_oceans_12000.jpg/1280px-Whole_world_-_land_and_oceans_12000.jpg

Hadoop



Saturday, January 10, 15

Let's drill down to Hadoop, which first gained widespread awareness in 2008-2009, when Yahoo! announced they were running a 10K core cluster with it, Hadoop became a top-level Apache project, etc.

4000

Scaling Hadoop to 4000 nodes at Yahoo!

By aanand – Tue, Sep 30, 2008 10:04 AM EDT

 Recommend

1

 Tweet

0

Tue, Sep 30, 2008

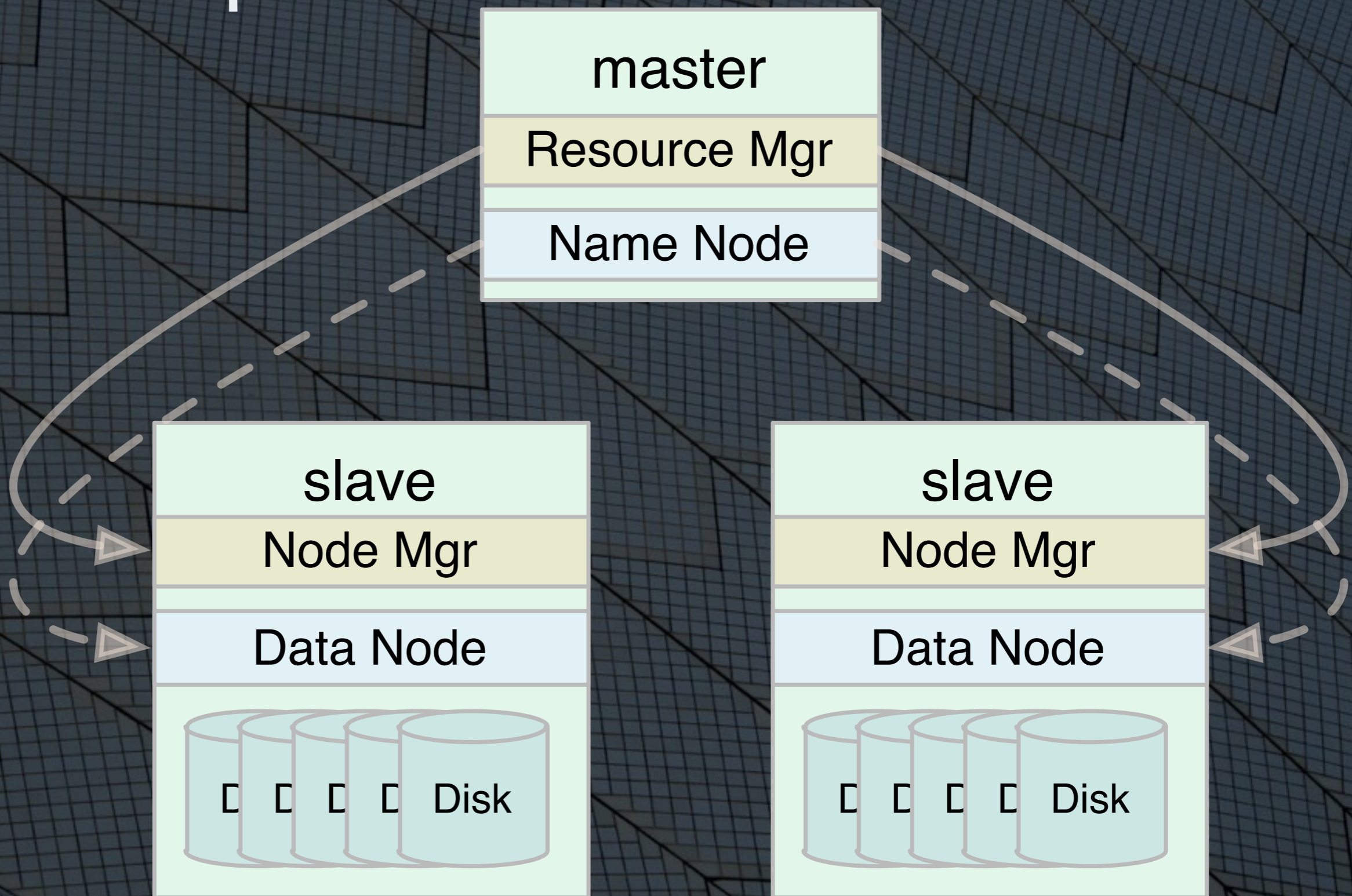
We recently ran Hadoop on what we believe is the single largest Hadoop installation, ever:

- 4000 nodes
- 2 quad core Xeons @ 2.5ghz per node
- 4x1TB SATA disks per node
- 8G RAM per node
- 1 gigabit ethernet on each node
- 40 nodes per rack
- 4 gigabit ethernet uplinks from each rack to the core (unfortunately a misconfiguration, we usually do 8 uplinks)
- Red Hat Enterprise Linux AS release 4 (Nahant Update 5)
- Sun Java JDK 1.6.0_05-b13
- So that's well over 30,000 cores with nearly 16PB of raw disk!

Quant, by
today's standards

16PB

Hadoop



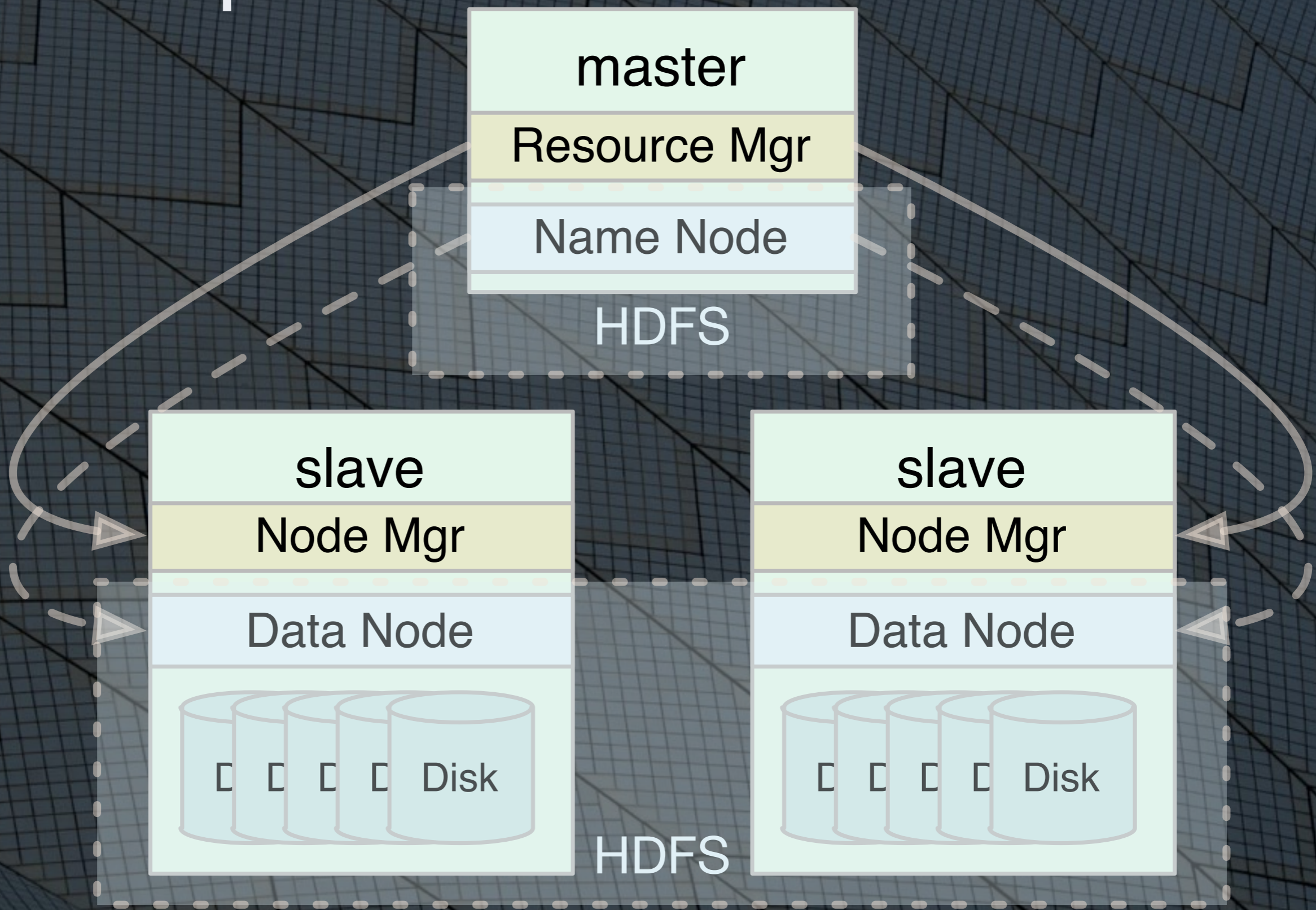
Saturday, January 10, 15

The schematic view of a Hadoop v2 cluster, with YARN (Yet Another Resource Negotiator) handling resource allocation and job scheduling. (V2 is actually circa 2013, but this detail is unimportant for this discussion). The master services are federated for failover, normally (not shown) and there would usually be more than two slave nodes. Node Managers manage the tasks

The Name Node is the master for the Hadoop Distributed File System. Blocks are managed on each slave by Data Node services.

The Resource Manager decomposes each job in to tasks, which are distributed to slave nodes and managed by the Node Managers. There are other services I'm omitting for simplicity.

Hadoop

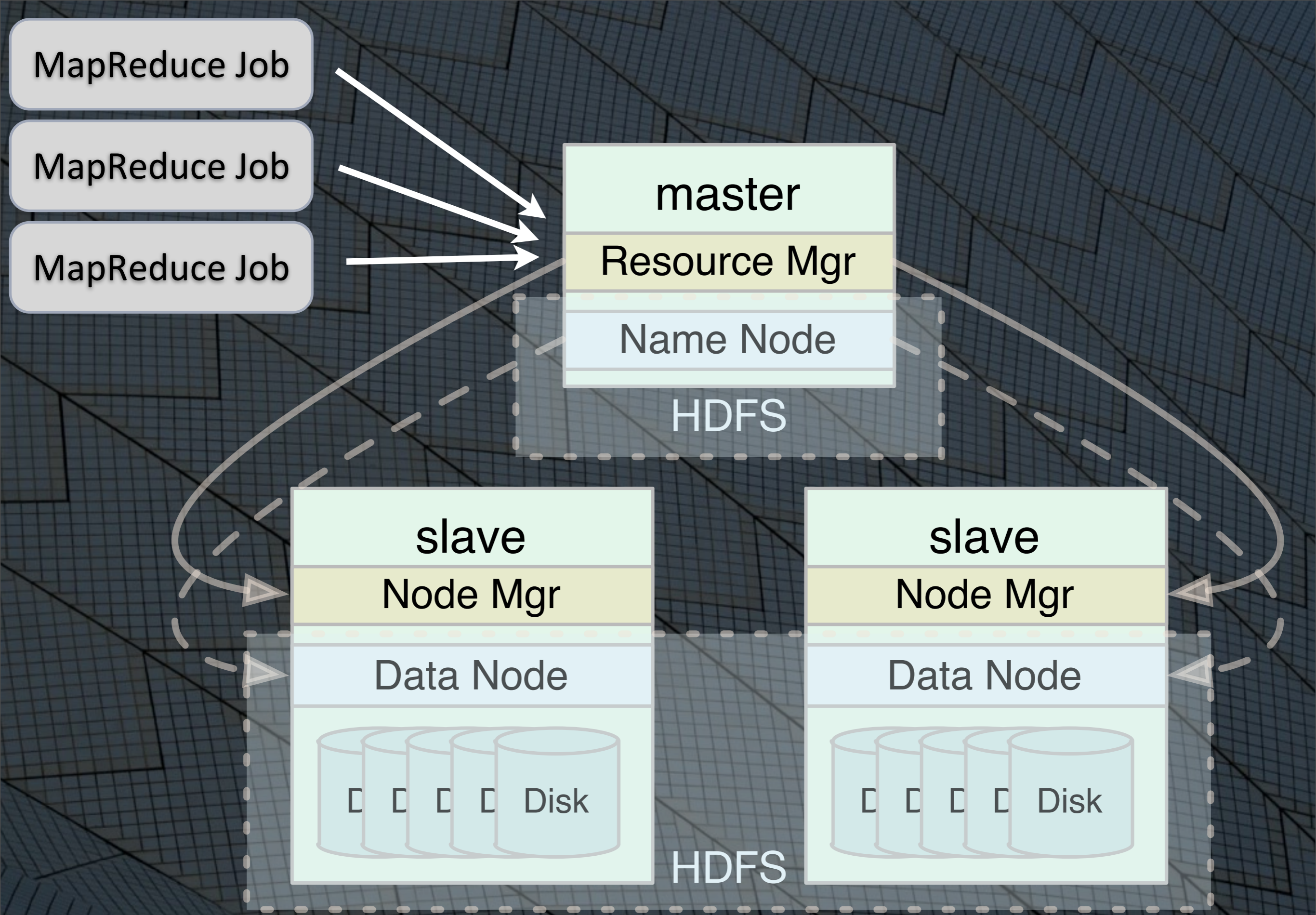


Saturday, January 10, 15

The schematic view of a Hadoop v2 cluster, with YARN (Yet Another Resource Negotiator) handling resource allocation and job scheduling. (V2 is actually circa 2013, but this detail is unimportant for this discussion). The master services are federated for failover, normally (not shown) and there would usually be more than two slave nodes. Node Managers manage the tasks

The Name Node is the master for the Hadoop Distributed File System. Blocks are managed on each slave by Data Node services.

The Resource Manager decomposes each job in to tasks, which are distributed to slave nodes and managed by the Node Managers. There are other services I'm omitting for simplicity.



MapReduce

The classic
compute model
for Hadoop

Example: Inverted Index

wikipedia.org/hadoop
Hadoop provides
MapReduce and HDFS

...

wikipedia.org/hbase
HBase stores data in HDFS

...

wikipedia.org/hive
Hive queries HDFS files and
HBase tables with SQL



inverse index

block

...	...
hadoop	(.../hadoop,1)
hbase	(.../hbase,1),(.../hive,1)
hdfs	(.../hadoop,1),(.../hbase,1),(.../hive,1)
hive	(.../hive,1)
...	...

block

...	...
-----	-----

block

...	...
-----	-----

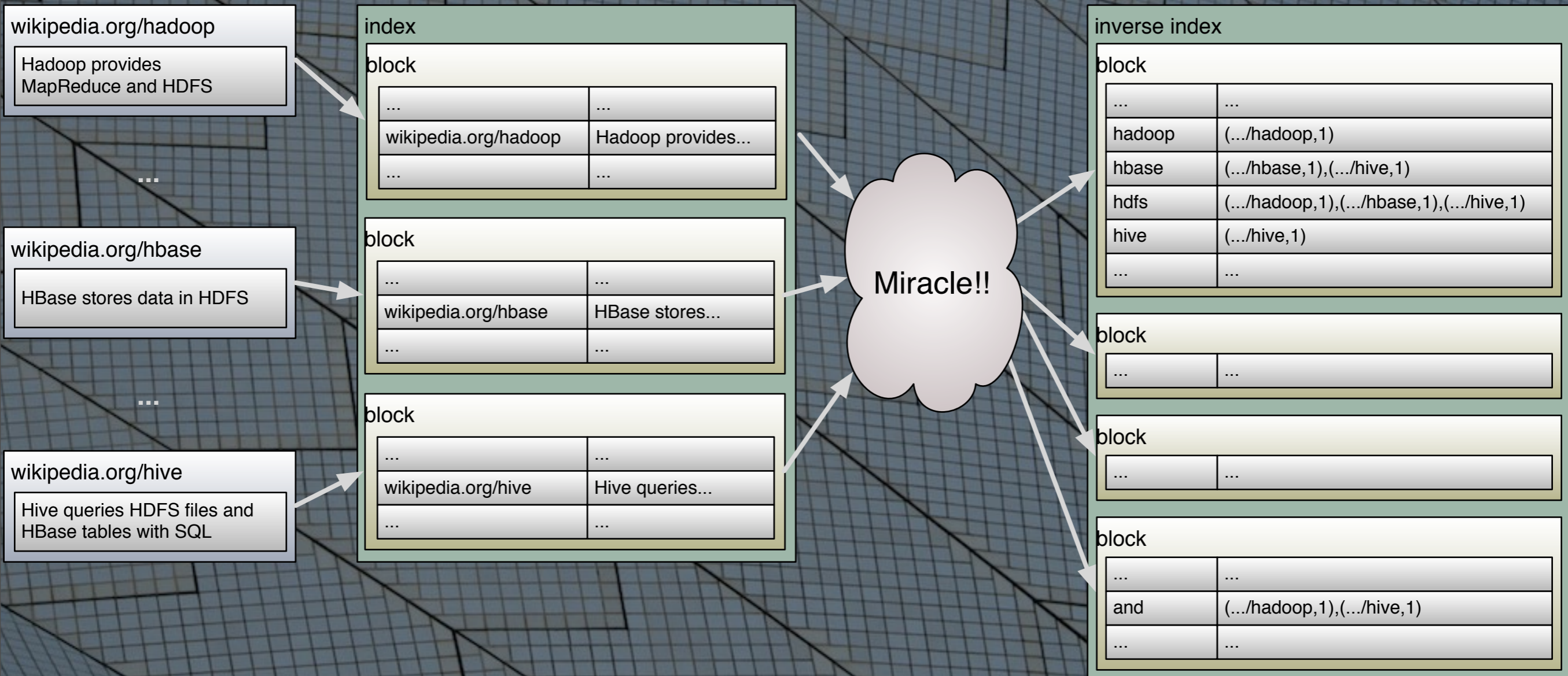
block

...	...
and	(.../hadoop,1),(.../hive,1)

Example: Inverted Index

Web Crawl

Compute Inverted Index



Web Crawl

Compute Inverted Index

wikipedia.org/hadoop
Hadoop provides
MapReduce and HDFS

wikipedia.org/hbase
HBase stores data in HDFS

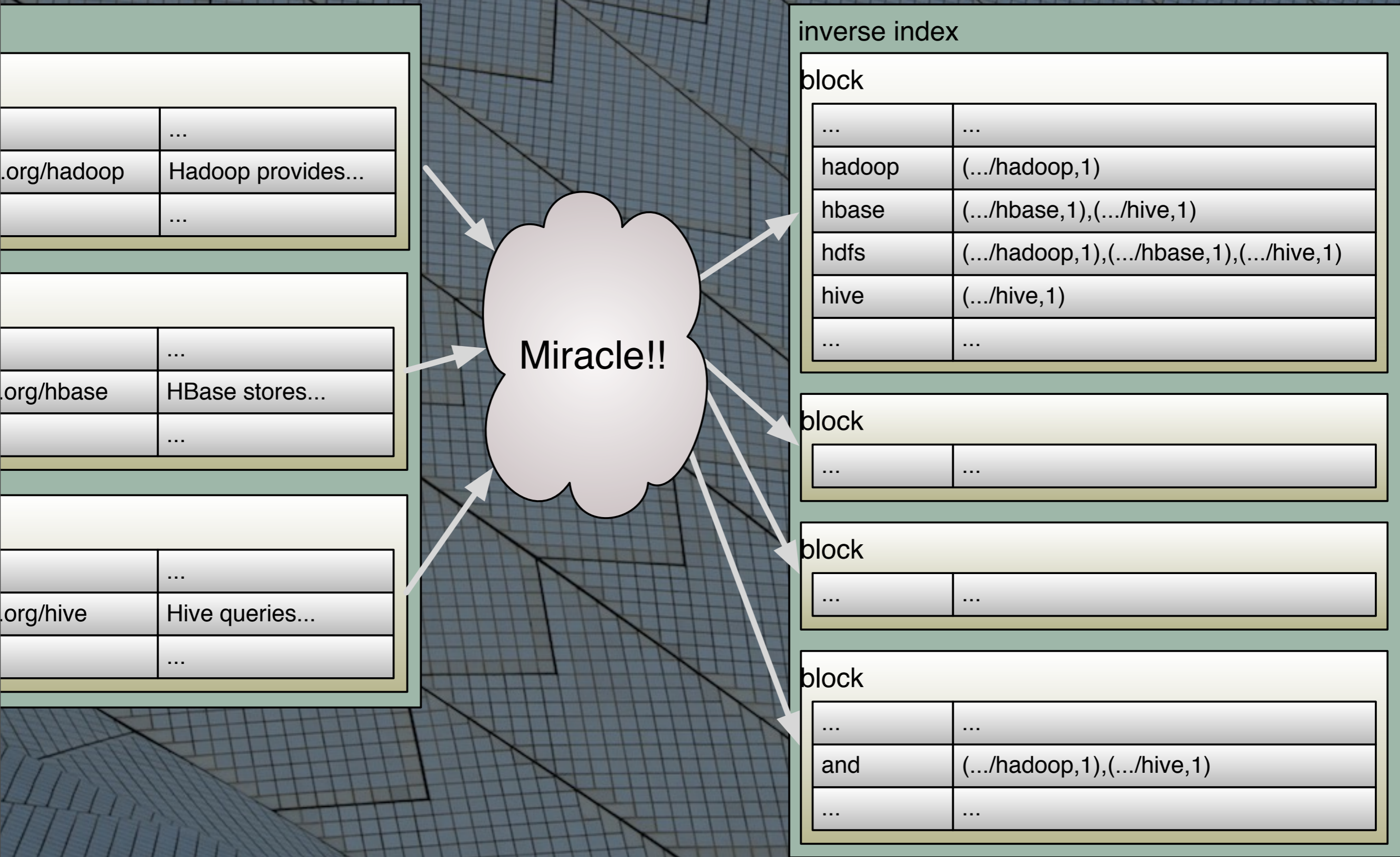
wikipedia.org/hive
Hive queries HDFS files and
HBase tables with SQL

index	
block	
...	...
wikipedia.org/hadoop	Hadoop provides...
...	...
block	
...	...
wikipedia.org/hbase	HBase stores...
...	...
block	
...	...
wikipedia.org/hive	Hive queries...
...	...

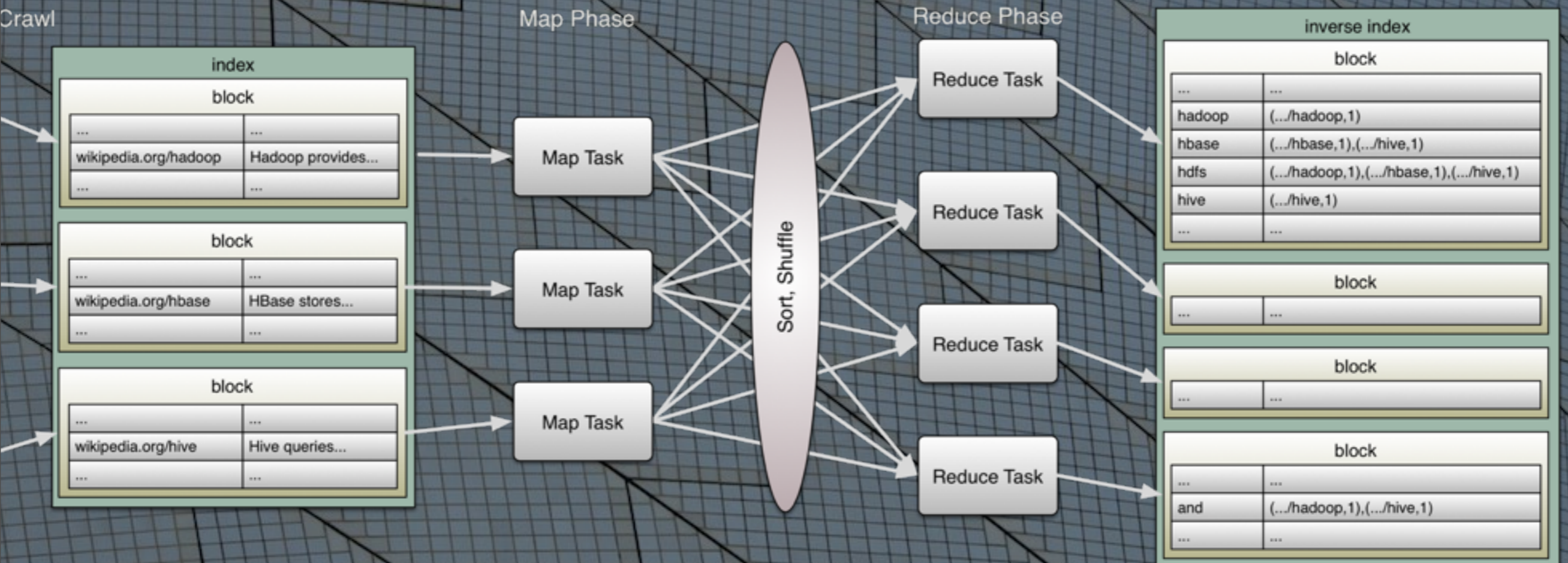


in
bl
bl
bl
bl
bl

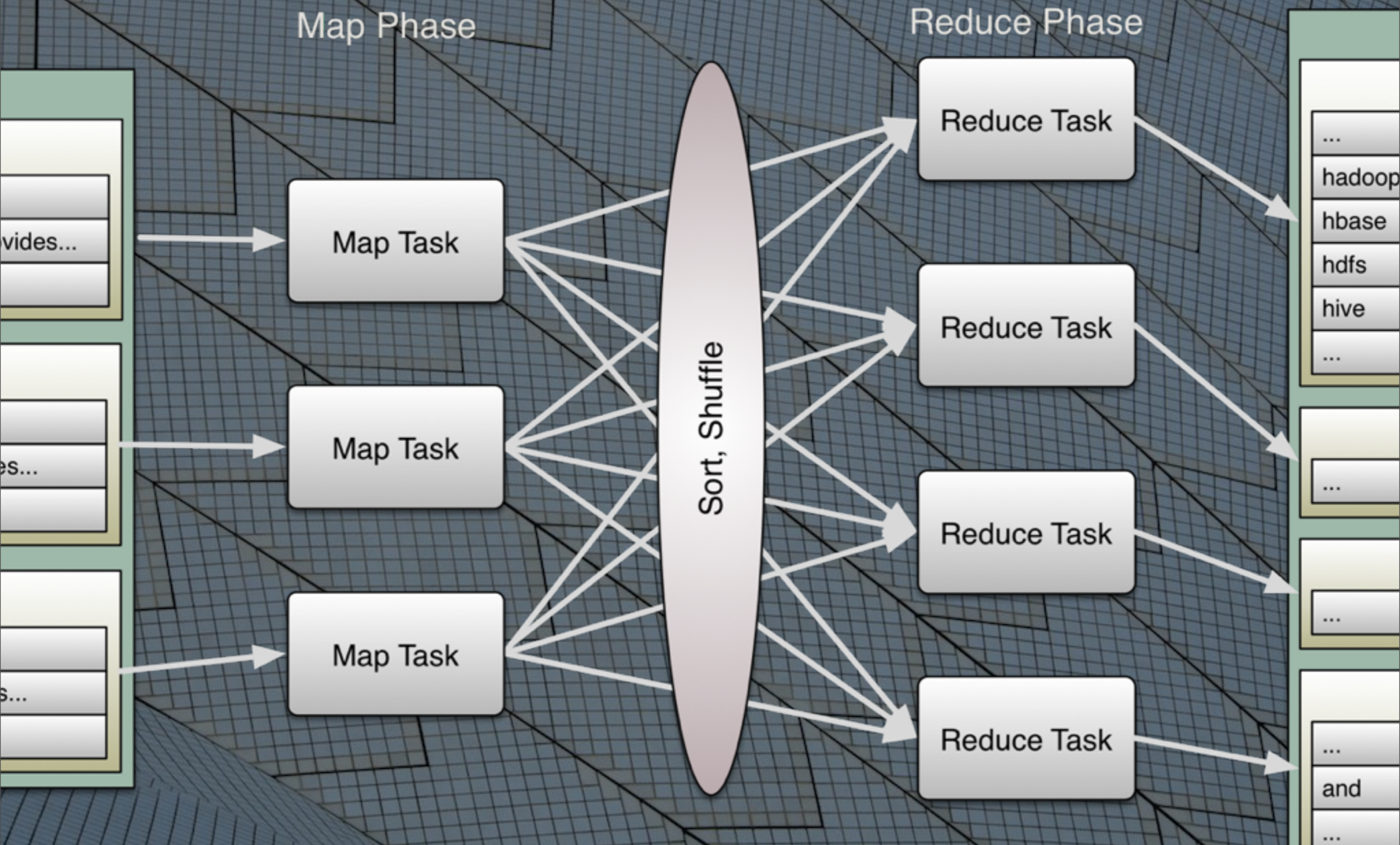
Compute Inverted Index



1 Map step + 1 Reduce step



1 Map step + 1 Reduce step



Problems

Hard to
implement
algorithms...

Problems

... and the
Hadoop API is
horrible:

```
import java.io.IOException;
import java.util.*;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapred.*;

public class LineIndexer {

    public static void main(String[] args) {
        JobClient client = new JobClient();
        JobConf conf =
            new JobConf(LineIndexer.class);

        conf.setJobName("LineIndexer");
        conf.setOutputKeyClass(Text.class);
        conf.setOutputValueClass(Text.class);
```

21

Saturday, January 10, 15

For example, the classic inverted index, used to convert an index of document locations (e.g., URLs) to words into the reverse; an index from words to doc locations. It's the basis of search engines.

I'm not going to explain the details. The point is to notice all the boilerplate that obscures the problem logic.

Everything is in one outer class. We start with a main routine that sets up the job.

I used yellow for method calls, because methods do the real work!! But notice that most of the functions in this code don't really do a whole lot of work for us...

```
JobClient client = new JobClient();
JobConf conf =
    new JobConf(LineIndexer.class);

conf.setJobName("LineIndexer");
conf.setOutputKeyClass(Text.class);
conf.setOutputValueClass(Text.class);
FileInputFormat.addInputPath(conf,
    new Path("input"));
FileOutputFormat.setOutputPath(conf,
    new Path("output"));
conf.setMapperClass(
    LineIndexMapper.class);
conf.setReducerClass(
    LineIndexReducer.class);

client.setConf(conf);
```

```
new Path("output"));
conf.setMapperClass(
    LineIndexMapper.class);
conf.setReducerClass(
    LineIndexReducer.class);

client.setConf(conf);

try {
    JobClient.runJob(conf);
} catch (Exception e) {
    e.printStackTrace();
}
}
```

```
public static class LineIndexMapper
    extends MapReduceBase
```

```
public static class LineIndexMapper
    extends MapReduceBase
    implements Mapper<LongWritable, Text,
                    Text, Text> {
    private final static Text word =
        new Text();
    private final static Text location =
        new Text();

    public void map(
        LongWritable key, Text val,
        OutputCollector<Text, Text> output,
        Reporter reporter) throws IOException {

        FileSplit fileSplit =
            (FileSplit)reporter.getInputSplit();
        String fileName =
```

24

Saturday, January 10, 15

This is the LineIndexMapper class for the mapper. The map method does the real work of tokenization and writing the (word, document-name) tuples.


```
FileSplit fileSplit =
    (FileSplit)reporter.getInputSplit();
String fileName =
    fileSplit.getPath().getName();
location.set(fileName);

String line = val.toString();
StringTokenizer itr = new
    StringTokenizer(line.toLowerCase());
while (itr.hasMoreTokens()) {
    word.set(itr.nextToken());
    output.collect(word, location);
}
}
}
```

`public static class LineIndexProducer`

```

public static class LineIndexReducer
    extends MapReduceBase
    implements Reducer<Text, Text,
                    Text, Text> {
    public void reduce(Text key,
        Iterator<Text> values,
        OutputCollector<Text, Text> output,
        Reporter reporter) throws IOException {
        boolean first = true;
        StringBuilder toReturn =
            new StringBuilder();
        while (values.hasNext()) {
            if (!first)
                toReturn.append(", ");
            first=false;
            toReturn.append(
                values.next().toString());
        }
    }
}

```

26

Saturday, January 10, 15

The reducer class, LineIndexReducer, with the reduce method that is called for each key and a list of values for that key. The reducer is stupid; it just reformats the values collection into a long string and writes the final (word,list-string) output.

```
Reporter reporter) throws IOException {
    boolean first = true;
    StringBuilder toReturn =
        new StringBuilder();
    while (values.hasNext()) {
        if (!first)
            toReturn.append(", ");
        first=false;
        toReturn.append(
            values.next().toString());
    }
    output.collect(key,
        new Text(toReturn.toString()));
}
}
}
```

```

import java.io.IOException;
import java.util.*;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapred.*;

public class LineIndexer {

    public static void main(String[] args) {
        JobClient client = new JobClient();
        JobConf conf =
            new JobConf(LineIndexer.class);

        conf.setJobName("LineIndexer");
        conf.setOutputKeyClass(Text.class);
        conf.setOutputValueClass(Text.class);
        FileInputFormat.addInputPath(conf,
            new Path("input"));
        FileOutputFormat.setOutputPath(conf,
            new Path("output"));
        conf.setMapperClass(
            LineIndexMapper.class);
        conf.setReducerClass(
            LineIndexReducer.class);

        client.setConf(conf);

        try {
            JobClient.runJob(conf);
        } catch (Exception e) {
            e.printStackTrace();
        }
    }

    public static class LineIndexMapper
        extends MapReduceBase
        implements Mapper<LongWritable, Text,
            Text, Text> {
        private final static Text word =
            new Text();
        private final static Text location =
            new Text();

        public void map(
            LongWritable key, Text val,
            OutputCollector<Text, Text> output,
            Reporter reporter) throws IOException {

            FileSplit fileSplit =
                (FileSplit)reporter.getInputSplit();
            String fileName =
                fileSplit.getPath().getName();
            location.set(fileName);

            String line = val.toString();
            StringTokenizer itr = new
                StringTokenizer(line.toLowerCase());
            while (itr.hasMoreTokens()) {
                word.set(itr.nextToken());
                output.collect(word, location);
            }
        }
    }

    public static class LineIndexReducer
        extends MapReduceBase
        implements Reducer<Text, Text,
            Text, Text> {
        public void reduce(Text key,
            Iterator<Text> values,
            OutputCollector<Text, Text> output,
            Reporter reporter) throws IOException {
            boolean first = true;
            StringBuilder toReturn =
                new StringBuilder();
            while (values.hasNext()) {
                if (!first)
                    toReturn.append(", ");
                first=false;
                toReturn.append(
                    values.next().toString());
            }
            output.collect(key,
                new Text(toReturn.toString()));
        }
    }
}

```

Altogether



Saturday, January 10, 15

Let's put all this into perspective, circa 2012...

http://upload.wikimedia.org/wikipedia/commons/thumb/8/8f/Whole_world_-_land_and_oceans_12000.jpg/1280px-Whole_world_-_land_and_oceans_12000.jpg

Dean Wampler



*“Trolling the
Hadoop community
since 2012...”*

NE Scala Symposium
@deanwampler
March 9, 2012

Why Big Data Needs to Be Functional



In which I claimed that:



*Hadoop is the
Enterprise Java Beans
of our time.*

Salvation!

Scalding

Saturday, January 10, 15

Twitter wrote a Scala API, <https://github.com/twitter/scalding>, to hide the mess. Actually, Scalding sits on top of Cascading (<http://cascading.org>) a higher-level Java API that exposes more sensible “combinators” of operations, but is still somewhat verbose due to the pre-Java 8 conventions it must use. Scalding gives us the full benefits of Scala syntax and functional operations, “combinators”.

Scalding (Scala)

Cascading (Java)

MapReduce (Java)

```
import com.twitter.scalding._

class InvertedIndex(args: Args)
  extends Job(args) {

  val texts = Tsv("texts.tsv", ('id, 'text))

  val wordToIds = texts
    .flatMap(('id, 'text) -> ('word, 'id2)) {
      fields: (Long, String) =>
        val (id2, text) =
          text.split("\\s+").map {
            word => (word, id2)
          }
    }
}
```

```
.flatMap((id, text) => (word, id2) => {
  fields: (Long, String) =>
    val (id2, text) =
      text.split("\\s+").map {
        word => (word, id2)
      }
})

val invertedIndex =
  wordToTweets.groupBy('word) {
    _.toList[Long]('id2 -> 'ids)
  }
invertedIndex.write(Tsv("output.tsv"))
}
```

That's it!

Problems


MapReduce is
“batch-mode”
only



Event stream processing.

(actually 2011)

Storm!



Twitter wrote Summingbird for Storm + Scalding

Storm!

Saturday, January 10, 15

Storm is a popular framework for scalable, resilient, event-stream processing.

Twitter wrote a Scalding-like API called Summingbird (<https://github.com/twitter/summingbird>) that abstracts over Storm and Scalding, so you can write one program that can run in batch mode or process events.

(For time's sake, I won't show an example.)

Problems

Flush to disk,
then reread
between jobs

100x perf. hit!

Saturday, January 10, 15

While your algorithm may be implemented using a sequence of MR jobs (which takes specialized skills to write...), the runtime system doesn't understand this, so the output of each job is flushed to disk (HDFS), even if it's TBs of data. Then it is read back into memory as soon as the next job in the sequence starts! This problem plagues Scalding (and Cascading), too, since they run on top of MapReduce (although Cascading is being ported to Spark, which we'll discuss next). However, as of mid-2014, Cascading is being ported to a new, faster runtime called Apache Tez, and it might be ported to Spark, which we'll discuss. Twitter is working on its own optimizations within Scalding. So the perf. issues should go away by the end of 2014.



It's 2013.

Saturday, January 10, 15

Let's put all this into perspective, circa 2013...

http://upload.wikimedia.org/wikipedia/commons/thumb/8/8f/Whole_world_-_land_and_oceans_12000.jpg/1280px-Whole_world_-_land_and_oceans_12000.jpg

Salvation v2.0!

Use Spark

Saturday, January 10, 15

The Hadoop community has realized over the last several years that a replacement for MapReduce is needed. While MR has served the community well, it's a decade old and shows clear limitations and problems, as we've seen. In late 2013, Cloudera, the largest Hadoop vendor officially embraced Spark as the replacement. Most of the other Hadoop vendors followed.

```

import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._

object InvertedIndex {
  def main(args: Array[String]) = {

    val sc = new SparkContext(
      "local", "Inverted Index")

    sc.textFile("data/crawl")
      .map { line =>
        val array = line.split("\t", 2)
        (array(0), array(1))
      }
      .flatMap {
        case (path, text) =>
          text.split(" ") map {

```

This implementation is more sophisticated than the Scalding example. It also computes the count/document of each word. Hence, there are more steps (some of which could be merged).

It starts with imports, then declares a singleton object (a first-class concept in Scala), with a main routine (as in Java).

The methods are colored yellow again. Note this time how dense with meaning they are this time.

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
```

```
object InvertedIndex {
  def main(args: Array[String]) = {
```

```
    val sc = new SparkContext(
      "local", "Inverted Index")
```

```
    sc.textFile("data/crawl")
      .map { line =>
        val array = line.split("\t", 2)
        (array(0), array(1))
      }
```

```
      .flatMap {
        case (path, text) =>
```

```
          text.split(" ") map {
```

This implementation is more sophisticated than the Scalding example. It also computes the count/document of each word. Hence, there are more steps (some of which could be merged).

It starts with imports, then declares a singleton object (a first-class concept in Scala), with a main routine (as in Java).

The methods are colored yellow again. Note this time how dense with meaning they are this time.

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
```

```
object InvertedIndex {
  def main(args: Array[String]) = {
```

```
    val sc = new SparkContext(
      "local", "Inverted Index")
```

```
    sc.textFile("data/crawl")
      .map { line =>
        val array = line.split("\t", 2)
        (array(0), array(1))
      }
```

```
    .flatMap {
      case (path, text) =>
```

```
      text.split("|||w|||") map {
```

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
```

```
object InvertedIndex {
  def main(args: Array[String]) = {
```

```
    val sc = new SparkContext(
      "local", "Inverted Index")
```

```
    sc.textFile("data/crawl")
      .map { line =>
        val array = line.split("\t", 2)
        (array(0), array(1))
      }
```

```
    .flatMap {
      case (path, text) =>
```

```
      text.split("|||w||||") map {
```

```

sc.textFile("data/crawl")
  .map { line =>
    val array = line.split("\t", 2)
    (array(0), array(1))
  }
  .flatMap {
    case (path, text) =>
      text.split("""\W+""") map {
        word => (word, path)
      }
  }
  .map {
    case (w, p) => ((w, p), 1)
  }
  .reduceByKey {
    (n1, n2) => n1 + n2
  }
}

```

Now we begin a sequence of transformations on the input data.

First, we map over each line, a string, to extract the original document id (i.e., file name, UUID), followed by the text in the document, all on one line. We assume tab is the separator. “(array(0), array(1))” returns a two-element “tuple”. Think of the output RDD as having a schema “String fileName, String text”.

flatMap maps over each of these 2-element tuples. We split the text into words on non-alphanumeric characters, then output collections of word (our ultimate, final “key”) and the path. Each line is converted to a collection of (word,path) pairs, so flatMap converts the collection of collections into one long “flat” collection of (word,path) pairs.

```

sc.textFile("data/crawl")
  .map { line =>
    val array = line.split("\t", 2)
    (array(0), array(1))
  }
  .flatMap {
    case (path, text) =>
      text.split("""\W+""") map {
        word => (word, path)
      }
  }
  .map {
    case (w, p) => ((w, p), 1)
  }
  .reduceByKey {
    (n1, n2) => n1 + n2
  }
}

```

Next, flatMap maps over each of these 2-element tuples. We split the text into words on non-alphanumeric characters, then output collections of word (our ultimate, final “key”) and the path. Each line is converted to a collection of (word,path) pairs, so flatMap converts the collection of collections into one long “flat” collection of (word,path) pairs.


```
sc.textFile("data/crawl")
  .map { line =>
    val array = line.split("\t", 2)
    (array(0), array(1))
  }
  .flatMap {
    case (path, text) =>
      text.split("""\W+""") map {
        word => (word, path)
      }
  }
  .map {
    case (w, p) => ((w, p), 1)
  }
  .reduceByKey {
    (n1, n2) => n1 + n2
  }
}
```

```
}  
  .reduceByKey {  
    (n1, n2) => n1 + n2  
  }
```

```
  .groupBy {  
    case ((w, p), n) => w  
  }
```

```
((word1, path1), n1)  
((word2, path2), n2)  
...
```

```
  .map {  
    case (w, seq) =>  
      val seq2 = seq map {  
        case (_, (p, n)) => (p, n)  
      }.sortBy {  
        case (path, n) => (-n, path)  
      }  
      (w, seq2.mkString(", "))  
  }
```

```
  .saveAsTextFile(argz.outpath)
```

```
}  
.reduceByKey {  
  (n1, n2) => n1 + n2  
}
```

```
.groupBy {  
  case (w, p), n => w  
}
```

```
.map {  
  (word, Seq((word, (path1, n1)), (word, (path2, n2)), ...))  
  case ...
```

```
  val seq2 = seq.map {  
    case (_, (p, n)) => (p, n)  
  }.sortBy {  
    case (path, n) => (-n, path)  
  }  
  (w, seq2.mkString(", "))
```

```
}  
.saveAsTextFile(argz.outpath)
```

```

    case ((w, p), n) => w
  }
  .map {
    case (w, seq) =>
      val seq2 = seq map {
        case (_, (p, n)) => (p, n)
      }.sortBy {
        case (path, n) => (-n, path)
      }
      (w, seq2.mkString(", "))
  }
  .saveAsTextFile("word-counts", 1)
  sc.stop()
}
}

```

(word, "(path4, n4), (path3, n3), (path2, n2), ...")

```
    case ((w, p), n) => w
  }
  .map {
    case (w, seq) =>
      val seq2 = seq map {
        case (_, (p, n)) => (p, n)
      }.sortBy {
        case (path, n) => (-n, path)
      }
      (w, seq2.mkString(", "))
  }
  .saveAsTextFile(argz.outpath)
  sc.stop()
}
```

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._

object InvertedIndex {
  def main(args: Array[String]) = {

    val sc = new SparkContext(
      "local", "Inverted Index")

    sc.textFile("data/crawl")
      .map { line =>
        val array = line.split("\t", 2)
        (array(0), array(1))
      }
      .flatMap {
        case (path, text) =>
          text.split("""\W+""") map {
            word => (word, path)
          }
      }
      .map {
        case (w, p) => ((w, p), 1)
      }
      .reduceByKey {
        (n1, n2) => n1 + n2
      }
      .groupByKey {
        case (w, (p, n)) => w
      }
      .map {
        case (w, seq) =>
          val seq2 = seq map {
            case (_, (p, n)) => (p, n)
          }
          (w, seq2.mkString(", "))
      }
      .saveAsTextFile(argz.outpath)

    sc.stop()
  }
}
```

Altogether

```
text.split("\\W+") map {  
  word => (word, path)  
}  
  
}  
  
.map {  
  case (w, p) => ((w, p), 1)  
}  
  
.reduceByKey {  
  (n1, n2) => n1 + n2  
}  
  
.groupBy {  
  case (w, (p, n)) => w  
}  
  
.map {  
  case (w, seq) =>  
    val seq2 = seq map {  
      case (_, (p, n)) => (p, n)  
    }  
}
```

Powerful,
beautiful
combinators

55

Saturday, January 10, 15

Stop for a second and admire the simplicity and elegance of this code, even if you don't understand the details. This is what coding should be, IMHO, very concise, to the point, elegant to read. Hence, a highly-productive way to work!!

$$\nabla \cdot \mathbf{D} = \rho$$

$$\nabla \cdot \mathbf{B} = 0$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}$$

$$\nabla \times \mathbf{H} = \mathbf{J} + \frac{\partial \mathbf{D}}{\partial t}$$

Spark also has a streaming mode for “mini-batch” event handling.

(We’ll come back to it...)



Winning!

Saturday, January 10, 15

Let's recap why is Scala taking over the big data world.

Elegant DSLs

...

```
.map {  
  case (w, p) => ((w, p), 1)  
}
```

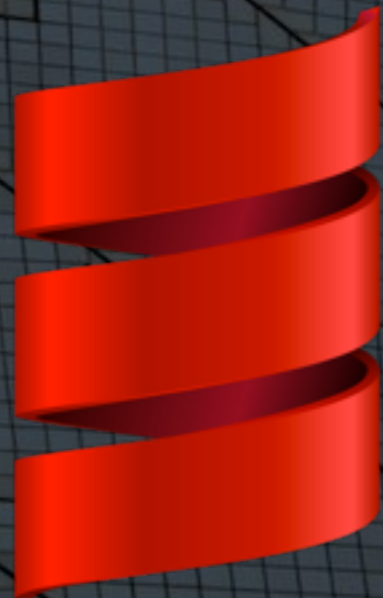
```
.reduceByKey {  
  (n1, n2) => n1 + n2  
}
```

```
.map {  
  case ((w, p), n) => (w, (p, n))  
}
```

```
.groupBy {  
  case (w, (p, n)) => w  
}
```

...

The JVM



eclipse



Functional Combinators

SQL Analog

```
CREATE TABLE inverted_index (  
  word    CHARACTER(64),  
  id1     INTEGER,  
  count1  INTEGER,  
  id2     INTEGER,  
  count2  INTEGER);
```

```
val inverted_index:  
Stream[(String,Int,Int,Int,Int)]
```

Functional Combinators

```
SELECT * FROM inverted_index  
WHERE word LIKE 'sc%';
```

Restrict

```
inverted_index.filter {  
  case (word, _) =>  
    word.startsWith "sc"  
}
```

Functional Combinators

```
SELECT word FROM inverted_index;
```

Projection

```
inverted_index.map {  
  case (word, _, _, _, _) =>  
    word  
}
```

Functional Combinators

```
SELECT count1, COUNT(*) AS size
FROM inverted_index
GROUP BY count1
ORDER BY size DESC;
```

Group By and Order By

```
inverted_index.groupBy {
  case (_, _, count1, _, _) => count1
} map {
  case (count1, words) => (count1, words.size)
} sortBy {
  case (count, size) => -size
}
```


Unification?



Spark Core + Spark SQL + Spark Streaming

```
val sparkContext =  
  new SparkContext("local[*]", "Much Wow!")  
val streamingContext =  
  new StreamingContext(  
    sparkContext, Seconds(60))  
val sqlContext =  
  new SQLContext(sparkContext)  
import sqlContext._
```

```
case class Flight(  
  number: Int,  
  carrier: String,  
  origin: String,  
  destination: String,  
  ...)
```

```
val sparkContext =  
  new SparkContext("local", "connections")  
val streamingContext =  
  new StreamingContext(  
    sparkContext, Seconds(60))  
val sqlContext =  
  new SQLContext(sparkContext)  
import sqlContext._
```

```
case class Flight(  
  number: Int,  
  carrier: String,  
  origin: String,  
  destination: String,  
  ...)
```

```
val sparkContext =  
  new SparkContext("local", "connections")  
val streamingContext =  
  new StreamingContext(  
    sparkContext, Seconds(60))  
val sqlContext =  
  new SQLContext(sparkContext)  
import sqlContext._
```

```
case class Flight(  
  number: Int,  
  carrier: String,  
  origin: String,  
  destination: String,  
  ...)
```

```
import sqlContext._
```

```
case class Flight(  
  number: Int,  
  carrier: String,  
  origin: String,  
  destination: String,  
  ...)
```

```
object Flight {  
  def parse(str: String): Option[Flight]=  
    {...}  
}
```

```
val server = ... // IP address or name  
val port = ... // integer  
val dStream =
```

```
streamingContext.socketTextStream(  
  ...)
```

```
}  
  
val server = ... // IP address or name  
val port = ... // integer  
val dStream =  
    streamingContext.socketTextStream(  
        server, port)
```

```
val flights = for {  
    line <- dStream  
    flight <- Flight.parse(line)  
} yield flight
```

```
flights.foreachRDD { (rdd, time) =>  
    rdd.registerTempTable("flights")  
    sql(s"""  
        SELECT $time, carrier, origin,  
        destination, COUNT(*)
```

```
}  
  
val server = ... // IP address or name  
val port = ... // integer  
val dStream =  
    streamingContext.socketTextStream(  
        server, port)
```

```
val flights = for {  
    line <- dStream  
    flight <- Flight.parse(line)  
} yield flight
```

```
flights.foreachRDD { (rdd, time) =>  
    rdd.registerTempTable("flights")  
    sql(s"""  
        SELECT $time, carrier, origin,  
            destination, COUNT(*)
```



```
flights.foreachRDD { (rdd, time) =>
  rdd.registerTempTable("flights")
  sql(s"""
    SELECT $time, carrier, origin,
      destination, COUNT(*)
    FROM flights
    GROUP BY carrier, origin, destination
    ORDER BY c4 DESC
    LIMIT 20""").foreach(println)
}
```

```
streamingContext.start()
streamingContext.awaitTermination()
streamingContext.stop()
```

```
flights.foreachRDD { (rdd, time) =>
  rdd.registerTempTable("flights")
  sql(s"""
    SELECT $time, carrier, origin,
           destination, COUNT(*)
    FROM flights
    GROUP BY carrier, origin, destination
    ORDER BY c4 DESC
    LIMIT 20""").foreach(println)
}
```

```
streamingContext.start()
streamingContext.awaitTermination()
streamingContext.stop()
```

We Won!



75

Saturday, January 10, 15

The title of this talk is in the present tense (present participle to be precise?), but has Scala already won? Is the game over?

dean.wampler@typesafe.com
polyglotprogramming.com/talks
[@deanwampler](#)



Saturday, January 10, 15

See also the bonus slides that follow.

Bonus Slides



Scala for Mathematics



78

Saturday, January 10, 15

Spire and Algebird. ScalaZ also has some of these data structures and algorithms.

Algebird

Large-scale Analytics

Algebraic types like Monoids,
which generalize addition.

–A set of elements.

–An associative binary operation.

–An identity element.

Efficient approximation algorithms.

– “Add All the Things”,
[infoq.com/presentations/
abstract-algebra-analytics](http://infoq.com/presentations/abstract-algebra-analytics)

Hash, don't Sample!

-- Twitter



Spire

Fast Numerics

- Types: Complex, Quaternion, Rational, Real, Interval, ...
- Algebraic types: Semigroups, Monoids, Groups, Rings, Fields, Vector Spaces, ...
- Trigonometric Functions.
- ...



What to Fix?

85

Saturday, January 10, 15

What could be improved?

Schema Management

Tuples limited to 22 fields.

Schema Management

Instantiate an object for each record?

Saturday, January 10, 15

Do we want to create an instance of an object for each record? We already have support for data formats like Parquet that implement columnar storage. Can our Scala APIs transparently use arrays of primitives for columns, for better performance? Spark has a support for Parquet which does this. Can we do it and should we do it for all data sets?

Schema Management

Remaining
boxing/
specialization
limitations

iPython Notebooks

Need an
equivalent for
Scala

Saturday, January 10, 15

iPython Notebooks are very popular with data scientists, because they integrate data visualization, etc. with code.
github.com/Bridgewater/scala-notebook is a start.