

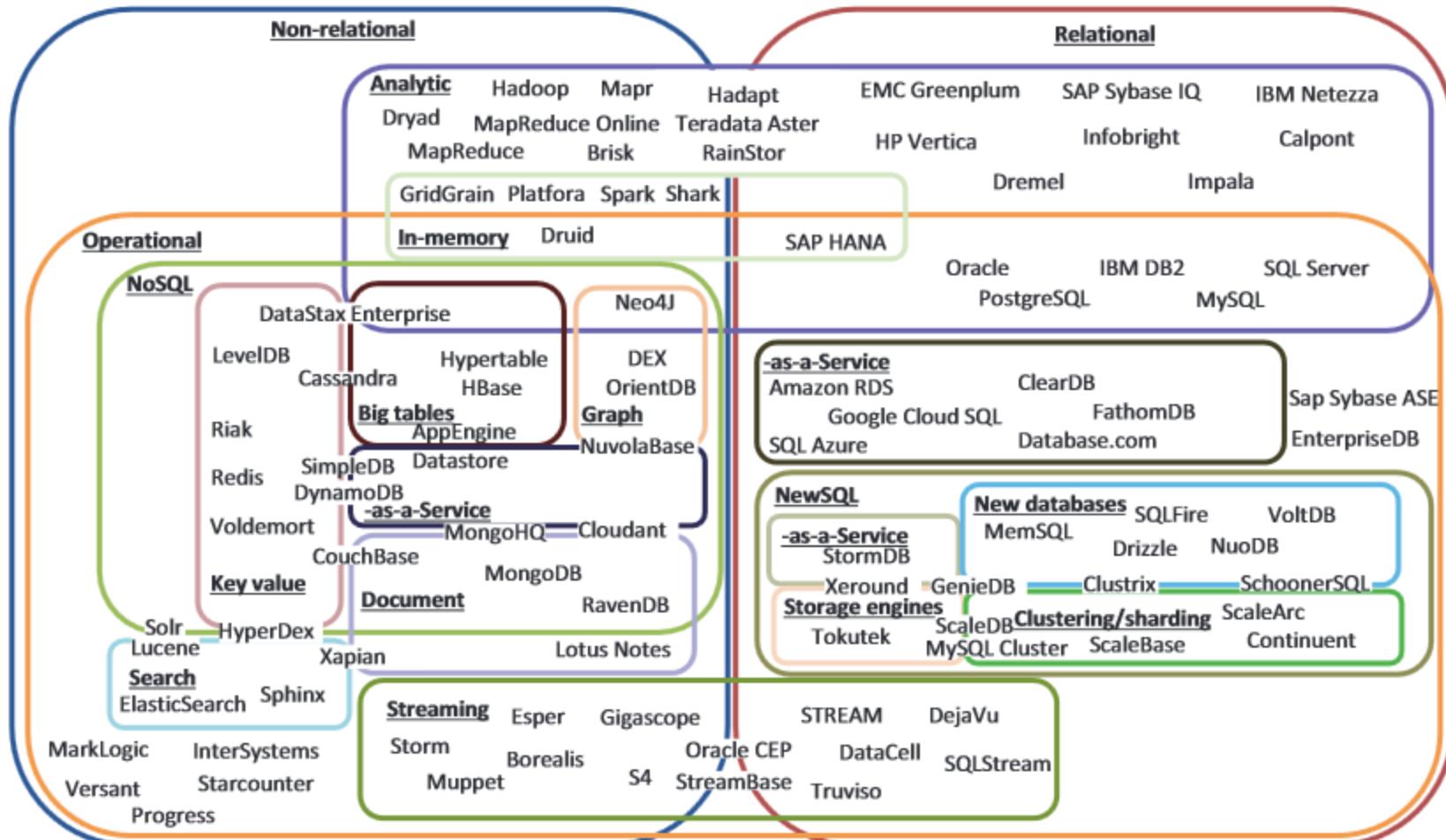
Lecture 20: NoSQL II

Monday, April 13, 2015

Announcements

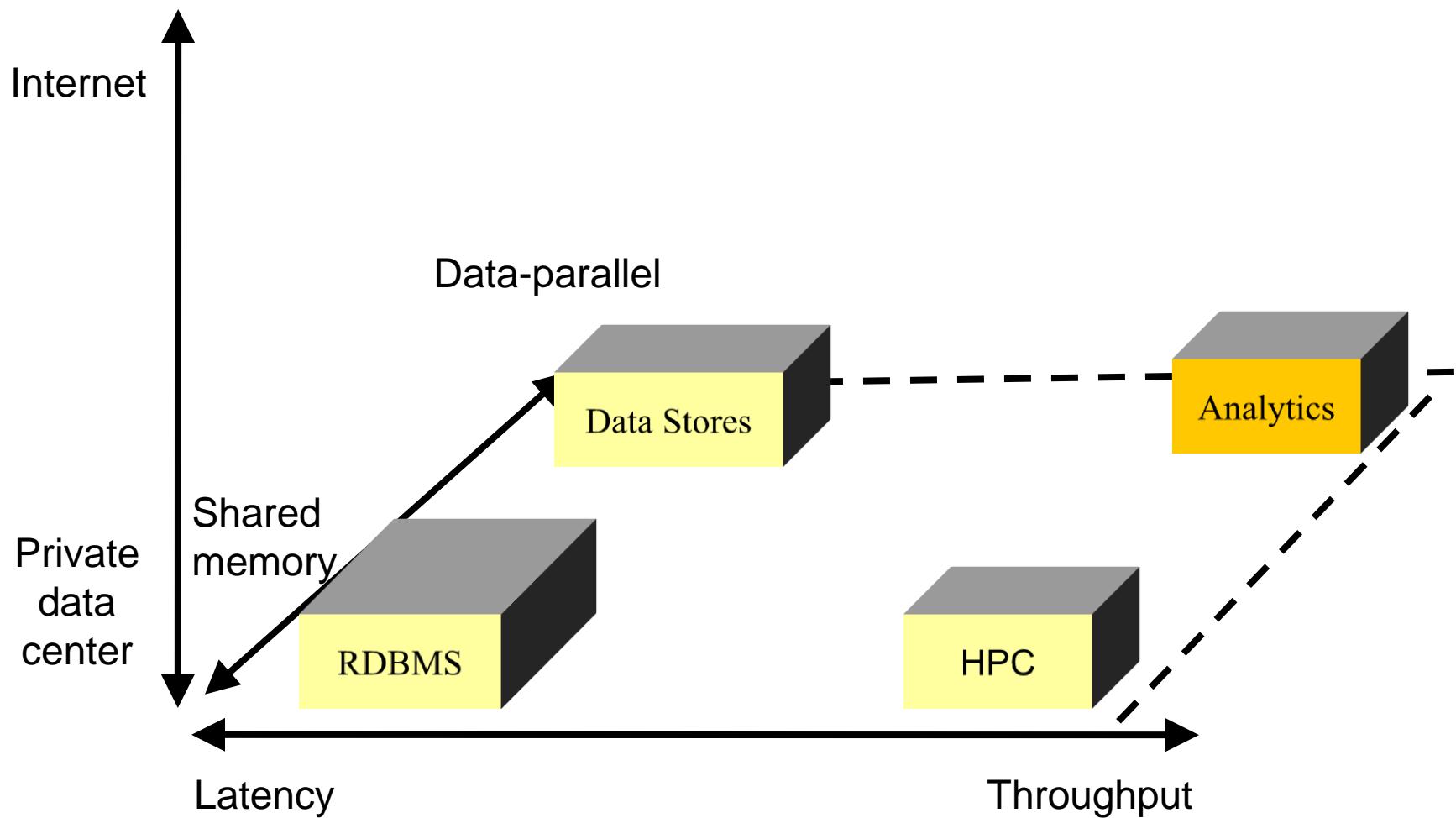
- Today: MapReduce & flavor of Pig
- Next class: Cloud platforms and Quiz #6
- HW #4 is out and will be due 04/27
- Grading questions:
 - Class participation
 - Homeworks
 - Quizzes
 - Class project

“Data Systems” Landscape



Source: Lim et al., "How to Fit when No One Size Fits", CIDR 2013.

Data Systems Design Space

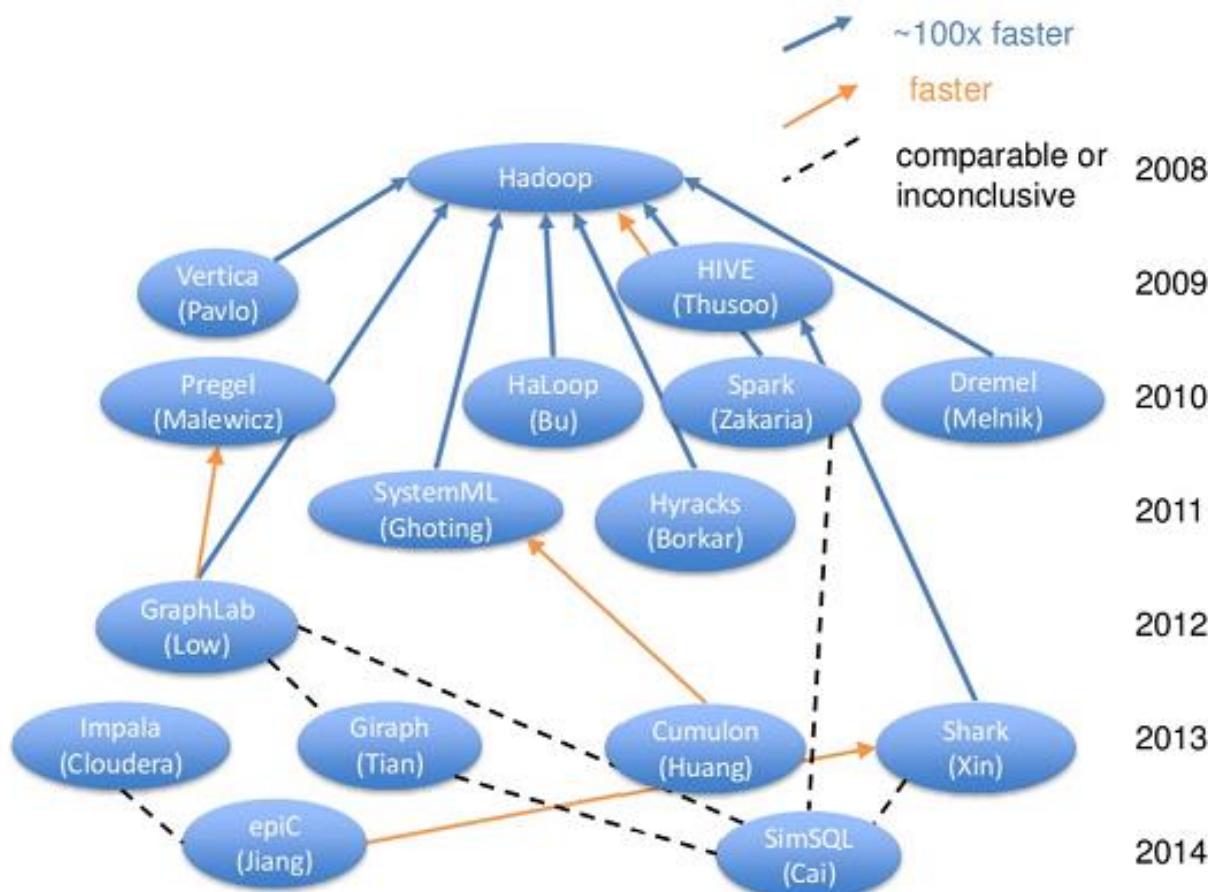


Source: Adapted from Michael Isard, Microsoft Research.

MapReduce

- MapReduce = high-level programming model and implementation for large-scale parallel data processing
- Inspired by primitives from Lisp and other functional programming languages
- History:
 - 2003: built at Google
 - 2004: published in OSDI (Dean & Ghemawat)
 - 2005: open-source version Hadoop
 - 2005 - 2014: very influential in DB community

MapReduce Literature



Source: David Maier and Bill Howe, "Big Data Middleware", CIDR 2015.

Data Model

MapReduce knows files!

A file = a bag of (key, value) pairs

A MapReduce program:

- Input: a bag of (inputkey, value) pairs
- Output: a bag of (outputkey, values) pairs

Step 1: Map Phase

- User provides the map function:
 - Input: one (input key, value) pair
 - Output: bag of (intermediate key, value) pairs
- MapReduce system applies the map function in parallel to all (input key, value) pairs in the input file
- Results from the Map phase are stored to disk and redistributed by the intermediate key during the Shuffle phase

Step 2: Reduce Phase

- MapReduce system groups all pairs with the same intermediate key, and passes the bag of values to the Reduce function
- User provides the Reduce function:
 - Input: (intermediate key, bag of values)
 - Output: bag of output values
- Results from Reduce phase stored to disk

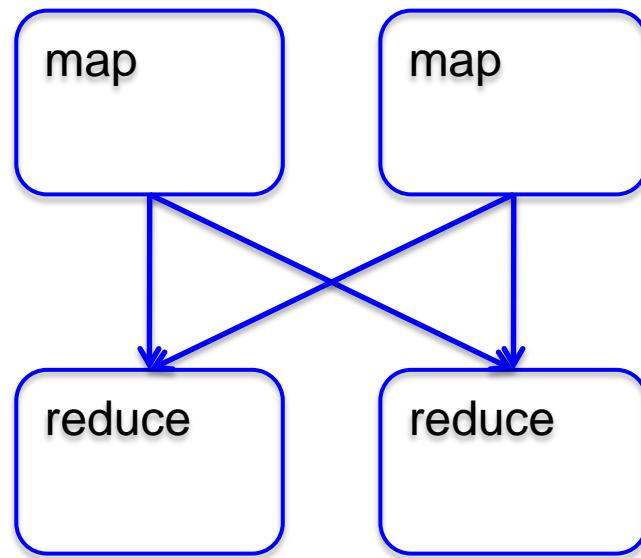
Canonical Example

Pseudocode for counting the number of occurrences of each word in a large collection of documents

```
map(String key, String input_value):
    // key: document name
    // input_value: document contents
    for each word in input_value:
        EmitIntermediate(word, "1");
```

```
reduce(String inter_key, Iterator inter_values):
    // inter_key: a word
    // inter_values: a list of counts
    int sum = 0;
    for each value in inter_values:
        sum += ParseInt(value);
    EmitFinal(inter_key, sum);
```

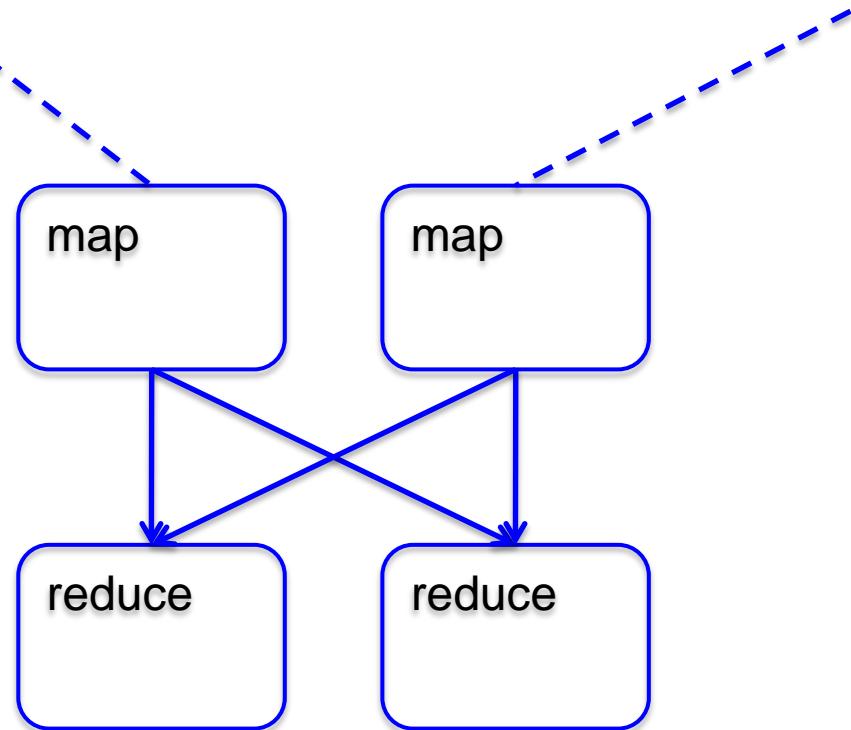
MapReduce Illustrated



MapReduce Illustrated

Romeo, Romeo, wherefore art thou Romeo?

What, art thou hurt?



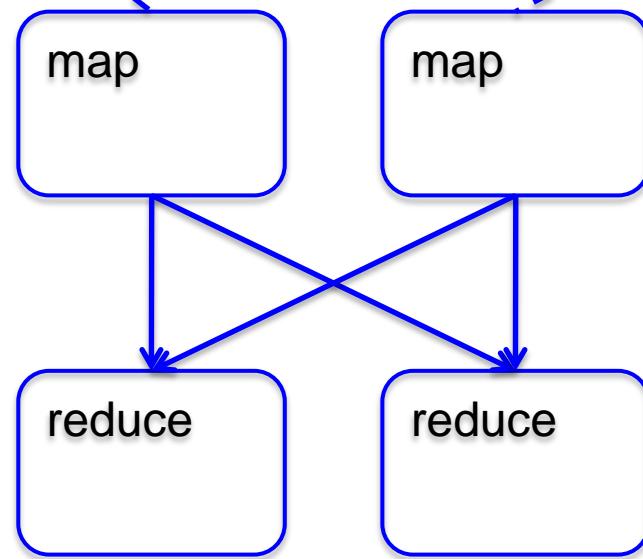
MapReduce Illustrated

Romeo, Romeo, wherefore art thou Romeo?

Romeo, 1
Romeo, 1
wherefore, 1
art, 1
thou, 1
Romeo, 1

What, art thou hurt?

What, 1
art, 1
thou, 1
hurt, 1



MapReduce Illustrated

Romeo, Romeo, wherefore art thou Romeo?

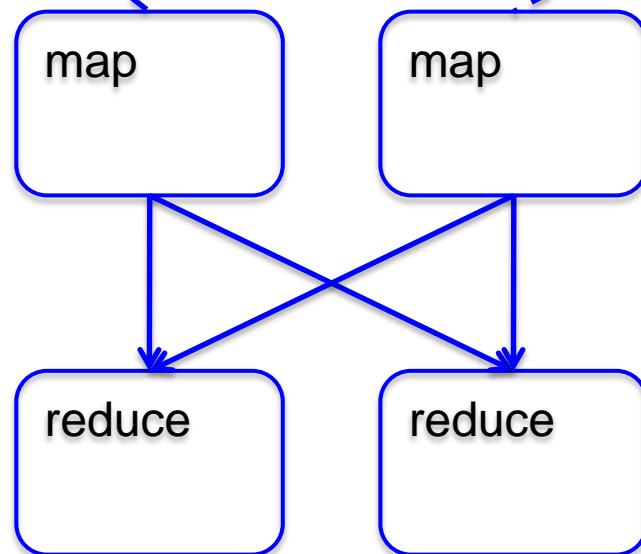
Romeo, 1
Romeo, 1
wherefore, 1
art, 1
thou, 1
Romeo, 1

art, (1, 1)
hurt (1),
thou (1, 1)

What, art thou hurt?

What, 1
art, 1
thou, 1
hurt, 1

Romeo, (1, 1, 1)
wherefore, (1)
what, (1)



MapReduce Illustrated

Romeo, Romeo, wherefore art thou Romeo?

Romeo, 1
Romeo, 1
wherefore, 1
art, 1
thou, 1
Romeo, 1

art, (1, 1)
hurt (1),
thou (1, 1)

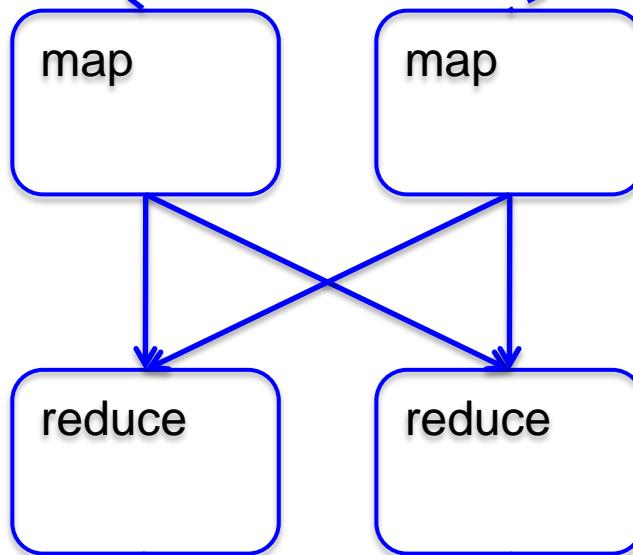
art, 2
hurt, 1
thou, 2

What, art thou hurt?

What, 1
art, 1
thou, 1
hurt, 1

Romeo, (1, 1, 1)
wherefore, (1)
what, (1)

Romeo, 3
wherefore, 1
what, 1



Rewritten as SQL

Documents(document_id, word)

```
SELECT word, COUNT(*)  
FROM Documents  
GROUP BY word
```

Observe:

Map + Shuffle Phases = Group By
Reduce Phase = Aggregate

More generally, each of the SQL operators that we have studied can be implemented in MapReduce

Relational Join

Employees(emp_id, last_name, first_name, dept_id)

Departments(dept_id, dept_name)

```
SELECT *
FROM Employees e, Departments d
WHERE e.dept_id = d.dept_id
```

Relational Join

Employees(emp_id, emp_name, dept_id)

<u>emp_id</u>	emp_name	dept_id
20	Alice	100
21	Bob	100
25	Carol	150

Departments(dept_id, dept_name)

<u>dept_id</u>	dept_name
100	Product
150	Support
200	Sales

```
SELECT e.emp_id, e.emp_name, d.dept_id, d.dept_name  
FROM Employees e, Deparments d  
WHERE e.dept_id = d.dept_id
```

<u>emp_id</u>	emp_name	dept_id	dept_name
20	Alice	100	Product
21	Bob	100	Product
25	Carol	150	Support

Relational Join

Employees(emp_id, emp_name, dept_id)

<u>emp_id</u>	emp_name	dept_id
20	Alice	100
21	Bob	100
25	Carol	150

Departments(dept_id, dept_name)

<u>dept_id</u>	dept_name
100	Product
150	Support
200	Sales

Input:

Employee, 20, Alice, 100
Employee, 21, Bob, 100
Employee, 25, Carol, 150
Departments, 100, Product
Departments, 150, Support
Departments, 200, Sales



Output:

k=100,v=(Employee, 20, Alice, 100)
k=100,v=(Employee, 21, Bob, 100)
k=150, v=(Employee, 25, Carol, 150)
k=100, v=(Departments, 100, Product)
k=150, v=(Departments, 150, Support)
k=200, v=(Departments, 200, Sales)

Relational Join

Employees(emp_id, emp_name, dept_id)

<u>emp_id</u>	emp_name	dept_id
20	Alice	100
21	Bob	100
25	Carol	150

Departments(dept_id, dept_name)

<u>dept_id</u>	dept_name
100	Product
150	Support
200	Sales

Input:

k=100,v=[(Employee, 20, Alice, 100),
(Employee, 21, Bob, 100),
(Departments, 100, Product)]

k=150, v=[(Employee, 25, Carol, 150),
(Departments, 150, Support)]

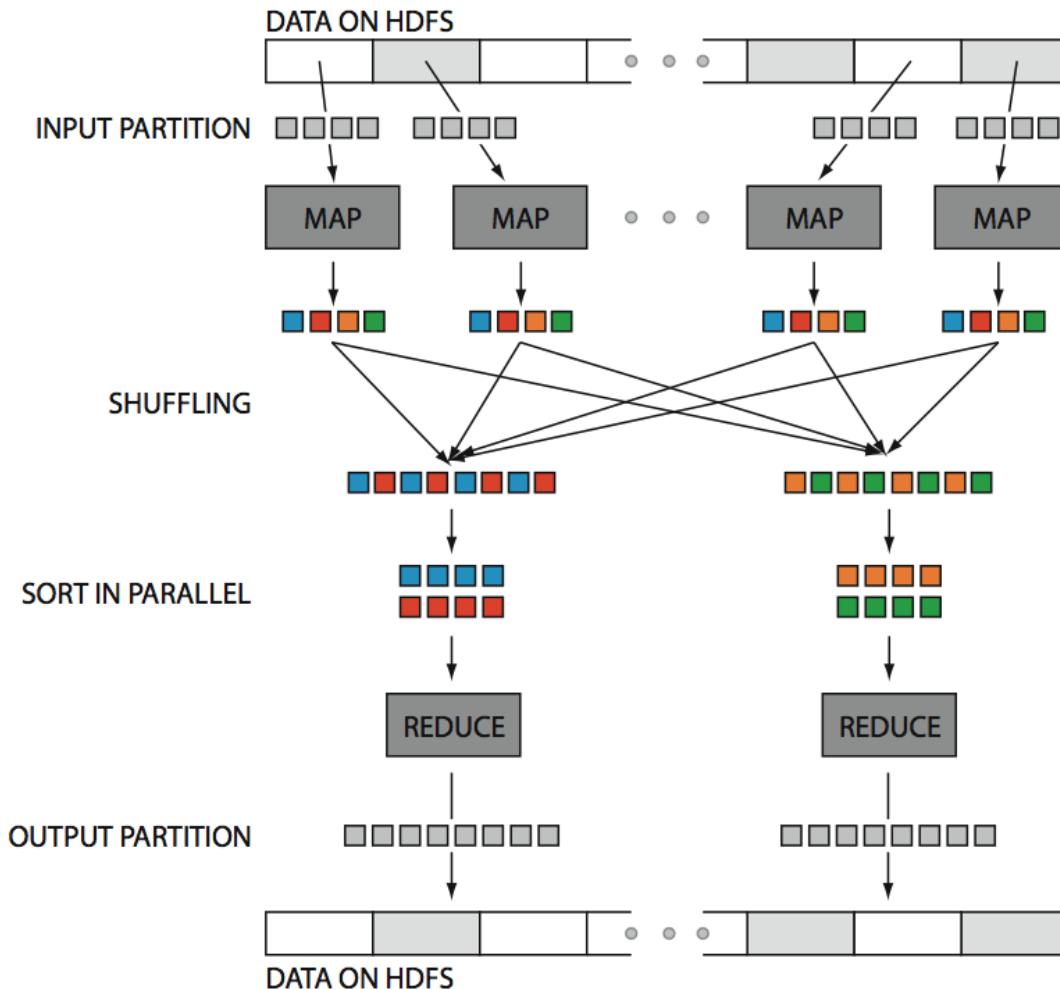
k=200, v=[(Departments, 200, Sales)]

Reduce

Output:

20, Alice, 100, Product
21, Bob, 100, Product
25, Carol, 150, Support

Hadoop on One Slide



Source: Huy Vo, NYU Poly

MapReduce Internals

- Single master node
- Master partitions input file by key into **M splits** (> servers)
- Master assigns **workers** (=servers) to the **M map tasks**, keeping track of their progress
- Workers write their output to local disk, partition into **R regions** (> servers)
- Master assigns workers to the **R reduce tasks**
- Reduce workers read regions from the map workers' local disks

Key Implementation Details

- Worker failures:
 - Master pings workers periodically, looking for stragglers
 - When straggle is found, master reassigns splits to **other** workers
 - Stragglers are a main reason for slowdown
 - Solution: pre-emptive backup execution of last few remaining in-progress tasks
- Choice of M and R:
 - Larger than servers is better for load balancing

MapReduce Summary

- Hides scheduling and parallelization details
- Not most efficient implementation, but has great fault tolerance
- However, limited queries:
 - Difficult to write more complex tasks
 - Need multiple MapReduce operations
- Solution:
 - Use high-level language (e.g. Pig, Hive, Sawzall, Dremel, Tenzing) to express complex queries
 - Need optimizer to compile queries into MR tasks

MapReduce Summary

- Hides scheduling and parallelization details
- Not most efficient implementation, but has great fault tolerance
- However, limited queries:
 - Difficult to write more complex tasks
 - Need multiple MapReduce operations
- Solution:
 - Use high-level language (e.g. Pig, Hive, Sawzall, Dremel, Tenzing) to express complex queries
 - Need optimizer to compile queries into MR tasks

Pig & Pig Latin

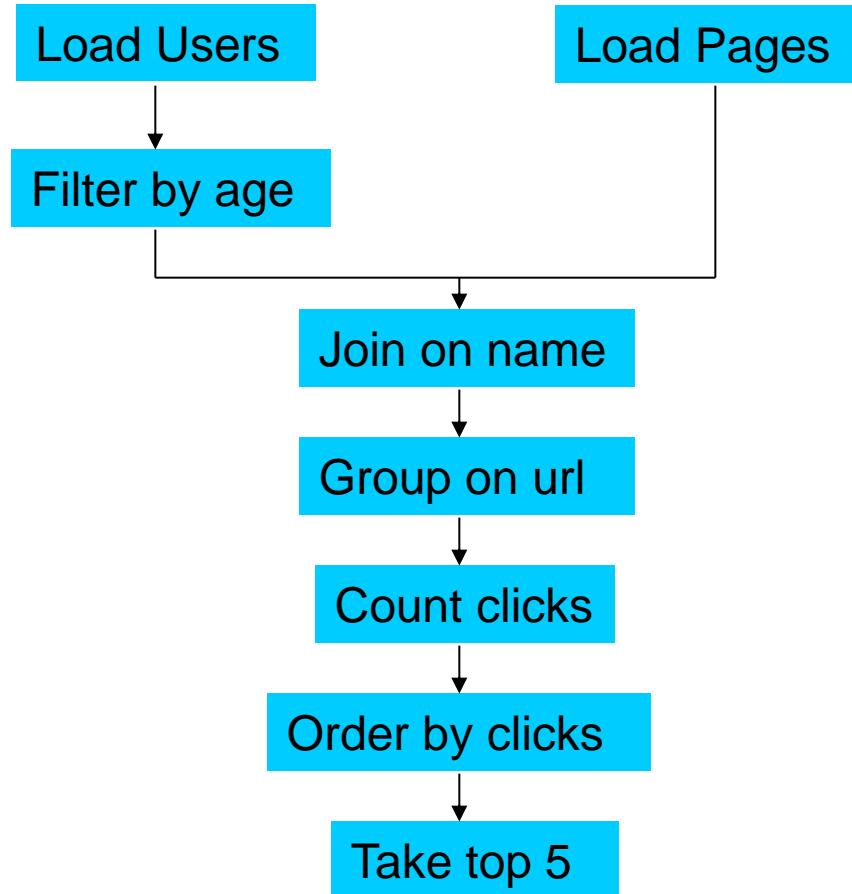
- An engine and language for executing programs on top of Hadoop
- Logical plan → sequence of MapReduce ops
- Free and open-sourced (unlike some others)
<http://hadoop.apache.org/pig/>
- ~70% of Hadoop jobs are Pig jobs at Yahoo!
- Being used at Twitter, LinkedIn, and other companies
- Available as part of Amazon, Hortonworks and Cloudera Hadoop distributions



Why use Pig?

Find the top 5 most visited sites by users aged 18 - 25.

Assume: user data stored in one file and website data in another file.



In MapReduce

```
import java.util.List;
import java.util.Map;
import org.apache.hadoop.io.IOException;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.FileInputFormat;
import org.apache.hadoop.mapred.FileOutputFormat;
import org.apache.hadoop.mapred.JobConf;
import org.apache.hadoop.mapred.KeyValueTextInputFormat;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.MapReduceJobControl;
import org.apache.hadoop.mapred.ReduceCollector;
import org.apache.hadoop.mapred.Reducer;
import org.apache.hadoop.mapred.SequenceFileInputFormat;
import org.apache.hadoop.mapred.TextInputFormat;
import org.apache.hadoop.mapred.jobcontrol.Job;
import org.apache.hadoop.mapred.jobcontrol.IdentityMapper;
public class MRExample extends MapReduceBase {
    public static class LoadPages extends Mapper<Text, Text, Text, Text> {
        public void map(Text key, Text val, Reporter reporter) throws IOException {
            // Pull out the key
            String line = val.toString();
            int firstComma = line.indexOf(',');
            String value = line.substring(firstComma + 1);
            // Prepend an index to the value so we know which
            // it came from.
            Text outVal = new Text("2" + value);
            oc.collect(outKey, outVal);
        }
        public static class ReduceUrls extends Reducer<Text, LongWritable, Text, Text> {
            public void reduce(Text key, Iterator<LongWritable> iter, Reporter reporter) throws IOException {
                long sum = 0;
                while (iter.hasNext()) {
                    sum += iter.next().get();
                }
                oc.collect(key, new LongWritable(sum));
            }
        }
        public static class Join extends MapReduceBase {
            implements Reducer<Text, Text, Text, Text> {
                public void reduce(Text key, Iterator<Text> iter, Reporter reporter) throws IOException {
                    // For each value, figure out which file it's from
                    // according to its first character.
                    List<String> first = new ArrayList<String>();
                    List<String> second = new ArrayList<String>();
                    while (iter.hasNext()) {
                        Text item = iter.next();
                        if (item.charAt(0) == '1') {
                            first.add(item.substring(1));
                        } else {
                            second.add(item.substring(1));
                        }
                    }
                    reporter.setStatus("OK");
                }
            }
        }
        public static class LoadJoined extends MapReduceBase {
            implements Mapper<Text, Text, Text, LongWritable> {
                public void map(Text key, Text val, Reporter reporter) throws IOException {
                    // Do the cross product and collect the values
                    for (String s1 : first) {
                        for (String s2 : second) {
                            String outval = key + "," + s1 + "," + s2;
                            reporter.collect(null, new Text(outval));
                        }
                    }
                    reporter.setStatus("OK");
                }
            }
        }
        public static class LoadJoinedMapper extends Mapper<Text, Text, Text, LongWritable> {
            implements Mapper<Text, Text, LongWritable> {
                public void map(Text key, Text val, Reporter reporter) throws IOException {
                    // just pass a 1 for the combiner/reducer
                    Text outkey = new Text(key);
                    outkey.setLength(1);
                    reporter.collect(outkey, new LongWritable(1));
                }
            }
        }
        public static class LoadAndFilterUsers extends Mapper<Text, Text, Text, Text> {
            public void map(Text key, Text val, Reporter reporter) throws IOException {
                String line = val.toString();
                int firstComma = line.indexOf(',');
                String value = line.substring(firstComma + 1);
                if (value.equals("1")) {
                    reporter.collect(null, new Text("1"));
                }
            }
        }
        public static class LoadAndFilterPages extends Mapper<Text, Text, Text, Text> {
            public void map(Text key, Text val, Reporter reporter) throws IOException {
                String line = val.toString();
                int firstComma = line.indexOf(',');
                String value = line.substring(firstComma + 1);
                if (value.equals("2")) {
                    reporter.collect(null, new Text("2"));
                }
            }
        }
        public static class LoadAndFilterURLs extends Mapper<Text, Text, Text, Text> {
            public void map(Text key, Text val, Reporter reporter) throws IOException {
                String line = val.toString();
                int firstComma = line.indexOf(',');
                String value = line.substring(firstComma + 1);
                if (value.equals("3")) {
                    reporter.collect(null, new Text("3"));
                }
            }
        }
        public static class LoadAndFilterClicks extends Mapper<Text, Text, Text, Text> {
            public void map(Text key, Text val, Reporter reporter) throws IOException {
                String line = val.toString();
                int firstComma = line.indexOf(',');
                String value = line.substring(firstComma + 1);
                if (value.equals("4")) {
                    reporter.collect(null, new Text("4"));
                }
            }
        }
        public static class LimitClicks extends MapReduceBase {
            implements Reducer<Text, LongWritable, Text, LongWritable> {
                public void reduce(Text key, Iterator<LongWritable> iter, Reporter reporter) throws IOException {
                    while (iter.hasNext()) {
                        LongWritable val = iter.next();
                        if (val.get() > 100) {
                            break;
                        }
                        reporter.collect((LongWritable) val, (Text) key);
                    }
                }
            }
        }
        public static class Top100 extends MapReduceBase {
            implements Reducer<Text, LongWritable, Text, LongWritable> {
                public void reduce(Text key, Iterator<LongWritable> iter, Reporter reporter) throws IOException {
                    int count = 0;
                    for (LongWritable val : iter) {
                        if (val.get() > 100) {
                            break;
                        }
                        count++;
                    }
                    if (count > 100) {
                        break;
                    }
                    reporter.collect(key, new LongWritable(count));
                }
            }
        }
        public static void main(String[] args) throws IOException {
            JobConf ip = new JobConf(MRExample.class);
            ip.setInputFormat(TextInputFormat.class);
            ip.setOutputFormat(TextOutputFormat.class);
            ip.setMapperClass(LoadPages.class);
            ip.setMapperClass(LoadAndFilterPages.class);
            ip.setMapperClass(LoadJoinedMapper.class);
            ip.setMapperClass(LoadAndFilterURLs.class);
            ip.setMapperClass(LoadAndFilterClicks.class);
            ip.setMapperClass(LimitClicks.class);
            ip.setMapperClass(Top100.class);
            JobControl jc = new JobControl();
            jc.addJob(ip);
            jc.run();
        }
    }
}
```

170 lines of code, 4 hours to write

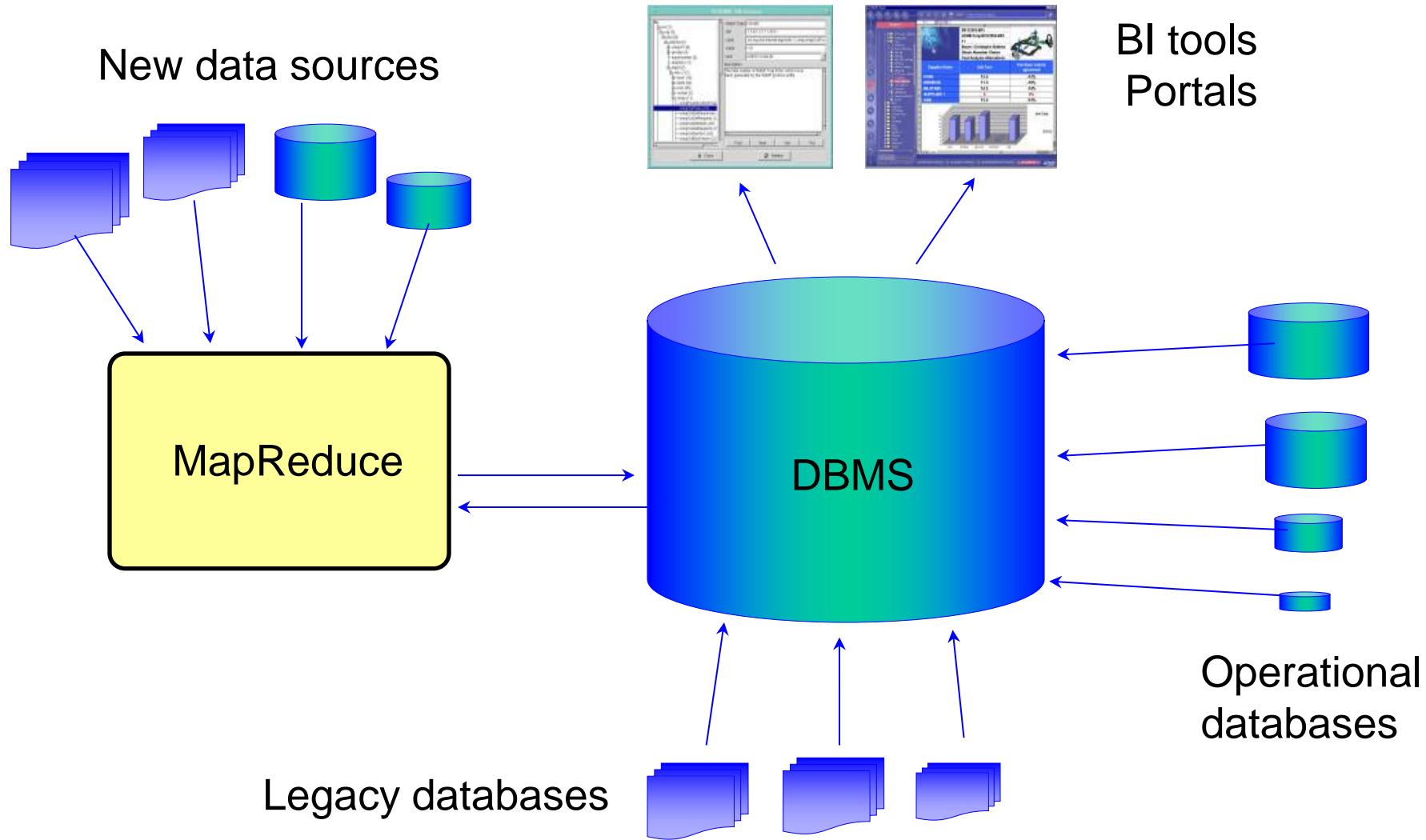
Source: Yahoo! Pig Team

In Pig Latin

```
Users = load 'users' as (name, age);  
Fltrd = filter Users by  
    age >= 18 and age <= 25;  
Pages = load 'pages' as (user, url);  
Jnd = join Fltrd by name, Pages by user;  
Grpd = group Jnd by url;  
Smmd = foreach Grpd generate group,  
    COUNT(Jnd) as clicks;  
Srted = order Smmd by clicks desc;  
Top5 = limit Srted 5;  
store Top5 into 'top5sites';
```

9 lines of code, 15 minutes to write

Emerging Analytics Pipeline



Optional References

MapReduce: Simplified Data Processing on Large Clusters [Dean & Ghemawat OSDI '04]

Pig Latin: A Not-So-Foreign Language for Data Processing [Olston et. al. SIGMOD '08]

Hive – A Petabyte Scale Data Warehouse Using Hadoop [Thusoo VLDB '09]

Designs, Lessons and Advice from Building Large Distributed Systems [Dean LADIS '09]

Tenzing: A SQL Implementation On The MapReduce Framework [Chattpadhyay et. al. VLDB '11]

Next Class

- Cloud platforms (guest speaker Jacob Walcik)
- Quiz #6