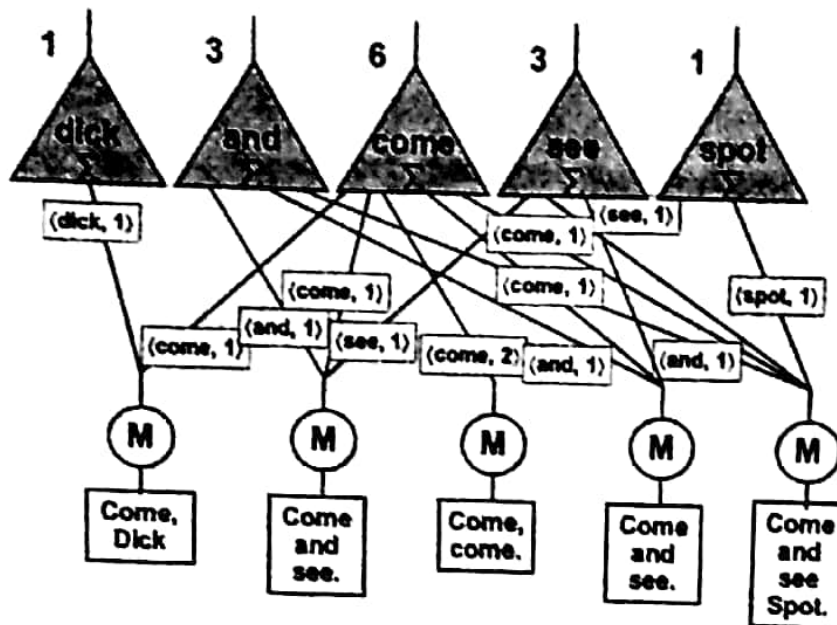


Map/Reduce Example



Word-Count Pairs

Extract

- Create an word index of set of documents
- Map: generate $\langle \text{word}, \text{count} \rangle$ pairs for all words in document
- Reduce: sum word counts across documents

- 3 -

Hadoop API

Requirements

- Programmer must supply Mapper & Reducer classes

Mapper

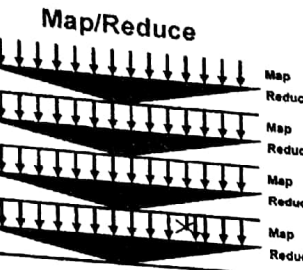
- Steps through file one line at a time
- Code generates sequence of $\langle \text{key}, \text{value} \rangle$
 - Call `output.collect(key, value)`
- Default types for keys & values are strings
 - Lots of low-level machinery to convert to & from other data types
 - But can use anything "writable"

Reducer

- Given key + iterator that generates sequence of values
- Generate one or more $\langle \text{key}, \text{value} \rangle$ pairs
 - Call `output.collect(key, value)`

- 5 -

Fault Tolerance



Map/Reduce

Map
Reduce
Map
Reduce
Map
Reduce
Map
Reduce

Data Integrity

- Store multiple copies of each file
- Including intermediate results of each Map / Reduce
 - Continuous checkpointing

Recovering from Failure

- Simply recompute lost result
 - Localized effect
- Dynamic scheduler keeps all processors busy

- 9 -

Example: Sparse Matrices with Map/Reduce

$$\begin{bmatrix} 10 & & 20 \\ & 30 & 40 \\ 50 & 60 & 70 \end{bmatrix} \times \begin{bmatrix} -1 & \\ -2 & -3 \\ & -4 \end{bmatrix} = \begin{bmatrix} -10 & -80 \\ -60 & -250 \\ -170 & -460 \end{bmatrix}$$

- Task: Compute product $C = A \cdot B$
- Assume most matrix entries are 0

Motivation

- Core problem in scientific computing
- Challenging for parallel execution
- Demonstrate expressiveness of Map/Reduce

- 10 -

Computing Sparse Matrix Product

A

10	20
	30 40
50	60 70

Row Col

1 $\frac{10}{A} \rightarrow 1$

1 $\frac{20}{A} \rightarrow 3$

2 $\frac{30}{A} \rightarrow 2$

2 $\frac{40}{A} \rightarrow 3$

3 $\frac{50}{A} \rightarrow 1$

3 $\frac{60}{A} \rightarrow 2$

3 $\frac{70}{A} \rightarrow 3$

B

-1	
-2	-3
	-4

Row Col

1 $\frac{-1}{B} \rightarrow 1$

2 $\frac{-2}{B} \rightarrow 1$

2 $\frac{-3}{B} \rightarrow 2$

3 $\frac{-4}{B} \rightarrow 2$

- Represent matrix as list of nonzero entries (row, col, value, matrixID)
- Strategy
 - Phase 1: Compute all products $a_{i,k} \cdot b_{k,j}$
 - Phase 2: Sum products for each entry i,j
 - Each phase involves a Map/Reduce

- 11 -

Phase 1 Map of Matrix Multiply

Row Col

1 $\frac{10}{A} \rightarrow 1$

1 $\frac{20}{A} \rightarrow 3$

2 $\frac{30}{A} \rightarrow 2$

2 $\frac{40}{A} \rightarrow 3$

3 $\frac{50}{A} \rightarrow 1$

3 $\frac{60}{A} \rightarrow 2$

3 $\frac{70}{A} \rightarrow 3$

Key = col

Key = 1

1 $\frac{10}{A} \rightarrow 1$

3 $\frac{50}{A} \rightarrow 1$

Key = 2

2 $\frac{30}{A} \rightarrow 2$

3 $\frac{60}{A} \rightarrow 2$

Key = 3

1 $\frac{20}{A} \rightarrow 3$

2 $\frac{40}{A} \rightarrow 3$

3 $\frac{70}{A} \rightarrow 3$

- Group values $a_{i,k}$ and $b_{k,j}$ according to key k

- 12 -

Phase 1 "Reduce" of Matrix Multiply

Key = 1

1 $\frac{10}{A} \rightarrow 1$

3 $\frac{50}{A} \rightarrow 1$

Key = 2

2 $\frac{30}{A} \rightarrow 2$

3 $\frac{60}{A} \rightarrow 2$

Key = 3

1 $\frac{20}{A} \rightarrow 3$

2 $\frac{40}{A} \rightarrow 3$

3 $\frac{70}{A} \rightarrow 3$

Key = 1

1 $\frac{-10}{C} \rightarrow 1$

3 $\frac{-50}{C} \rightarrow 1$

Key = 2

2 $\frac{-60}{C} \rightarrow 1$

2 $\frac{-90}{C} \rightarrow 2$

3 $\frac{-120}{C} \rightarrow 1$

3 $\frac{-180}{C} \rightarrow 2$

Key = 3

1 $\frac{-80}{C} \rightarrow 2$

2 $\frac{-160}{C} \rightarrow 2$

3 $\frac{-280}{C} \rightarrow 2$

- Generate all products $a_{i,k} \cdot b_{k,j}$

- 13 -

Phase 2 Map of Matrix Multiply

Row Col

1 $\frac{-10}{C} \rightarrow 1$

3 $\frac{-50}{C} \rightarrow 1$

2 $\frac{-60}{C} \rightarrow 1$

2 $\frac{-90}{C} \rightarrow 2$

3 $\frac{-120}{C} \rightarrow 1$

3 $\frac{-180}{C} \rightarrow 2$

1 $\frac{-80}{C} \rightarrow 2$

2 $\frac{-160}{C} \rightarrow 2$

3 $\frac{-280}{C} \rightarrow 2$

Key = row,col

Key = 1,1

1 $\frac{-10}{C} \rightarrow 1$

Key = 1,2

1 $\frac{-80}{C} \rightarrow 2$

Key = 2,1

2 $\frac{-60}{C} \rightarrow 1$

Key = 2,2

2 $\frac{-90}{C} \rightarrow 2$

2 $\frac{-160}{C} \rightarrow 2$

Key = 3,1

3 $\frac{-120}{C} \rightarrow 1$

3 $\frac{-50}{C} \rightarrow 1$

Key = 3,2

3 $\frac{-280}{C} \rightarrow 2$

3 $\frac{-180}{C} \rightarrow 2$

- Group products $a_{i,k} \cdot b_{k,j}$ with matching values of i and j

- 14 -

Phase 2 Reduce of Matrix Multiply

Key = 1,1	1 $\frac{-10}{C} \rightarrow 1$	1 $\frac{-10}{C} \rightarrow 1$
Key = 1,2	1 $\frac{-80}{C} \rightarrow 2$	1 $\frac{-80}{C} \rightarrow 2$
Key = 2,1	2 $\frac{-60}{C} \rightarrow 1$	2 $\frac{-60}{C} \rightarrow 1$
Key = 2,2	2 $\frac{-90}{C} \rightarrow 2$	2 $\frac{-250}{C} \rightarrow 2$
	2 $\frac{-160}{C} \rightarrow 2$	
Key = 3,1	3 $\frac{-120}{C} \rightarrow 1$	3 $\frac{-170}{C} \rightarrow 1$
	3 $\frac{-50}{C} \rightarrow 1$	
Key = 3,2	3 $\frac{-280}{C} \rightarrow 2$	3 $\frac{-460}{C} \rightarrow 2$
	3 $\frac{-180}{C} \rightarrow 2$	

C

$$\begin{bmatrix} -10 & -80 \\ -60 & -250 \\ -170 & -460 \end{bmatrix}$$

■ Sum products to get final entries

Matrix Multiply Phase 1 Mapper

```
public class P1Mapper extends MapReduceBase implements Mapper {
    public void map(WritableComparable key, Writable values,
        OutputCollector output, Reporter reporter) throws
        IOException {
        try {
            GraphEdge e = new GraphEdge(values.toString());
            IntWritable k;
            if (e.tag.equals("A"))
                k = new IntWritable(e.toNode);
            else
                k = new IntWritable(e.fromNode);
            output.collect(k, new Text(e.toString()));
        } catch (BadGraphException e) {}
    }
}
```

Matrix Multiply Phase 1 Reducer

```
public class P1Reducer extends MapReduceBase implements Reducer {

    public void reduce(WritableComparable key, Iterator values,
                      OutputCollector output, Reporter reporter)
                      throws IOException

    {
        Text outv = new Text(""); // Don't really need output values
        /* First split edges into A and B categories */
        LinkedList<GraphEdge> alist = new LinkedList<GraphEdge>();
        LinkedList<GraphEdge> blist = new LinkedList<GraphEdge>();
        while(values.hasNext()) {
            try {
                GraphEdge e =
                    new GraphEdge(values.next().toString());
                if (e.tag.equals("A")) {
                    alist.add(e);
                } else {
                    blist.add(e);
                }
            } catch (BadGraphException e) {}
        }
        // Continued
    }
}
```

Matrix Multiply Phase 2 Mapper

```
public class P2Mapper extends MapReduceBase implements Mapper {  
    public void map(WritableComparable key, Writable values,  
                    OutputCollector output, Reporter reporter)  
                    throws IOException {  
        String es = values.toString();  
        try {  
            GraphEdge e = new GraphEdge(es);  
            // Key based on head & tail nodes  
            String ks = e.fromNode + " " + e.toNode;  
            output.collect(new Text(ks), new Text(e.toString()));  
        } catch (BadGraphException e) {}  
    }  
}
```

MapReduce Implementation

Built on Top of Parallel File System

- **Google: GFS, Hadoop: HDFS**
- **Provides global naming**
- **Reliability via replication (typically 3 copies)**

Breaks work into tasks

- **Master schedules tasks on workers dynamically**
- **Typically #tasks >> #processors**

Net Effect

- **Input: Set of files in reliable file system**
- **Output: Set of files in reliable file system**
- **Can write program as series of MapReduce steps**

Reducing

Shuffle

- Each reducer fetches its share of key, value pairs from each mapper using RPC
- Sort data according to keys
 - Use disk-based ("external") sort if too much data for memory

Reduce Operation

- Step through key-value pairs in sorted order
- For each unique key, call reduce function for all values
- Append result to output file

Result

- R output files
- Typically supply to next round of MapReduce

Example Parameters

Sort Benchmark

- 10^{10} 100-byte records
- Partition into $M = 15,000$ 64MB pieces
 - Key = value
 - Partition according to most significant bytes
- Sort locally with $R = 4,000$ reducers

Machine

- 1800 2Ghz Xeons
- Each with 2 160GB IDE disks
- Gigabit ethernet
- 891 seconds total

Interesting Features

Fault Tolerance

- Assume reliable file system
- Detect failed worker
 - Heartbeat mechanism
- Rescheduled failed task

Stragglers

- Tasks that take long time to execute
- Might be bug, flaky hardware, or poor partitioning
- When done with most tasks, reschedule any remaining executing tasks
 - Keep track of redundant executions
 - Significantly reduces overall run time