

# Introduction to Hadoop

Owen O'Malley

Yahoo Inc!

[omalley@apache.org](mailto:omalley@apache.org)



# Hadoop: Why?

- Need to process 100TB datasets with multi-day jobs
- On 1 node:
  - scanning @ 50MB/s = 23 days
  - MTBF = 3 years
- On 1000 node cluster:
  - scanning @ 50MB/s = 33 min
  - MTBF = 1 day
- Need framework for distribution
  - Efficient, reliable, easy to use

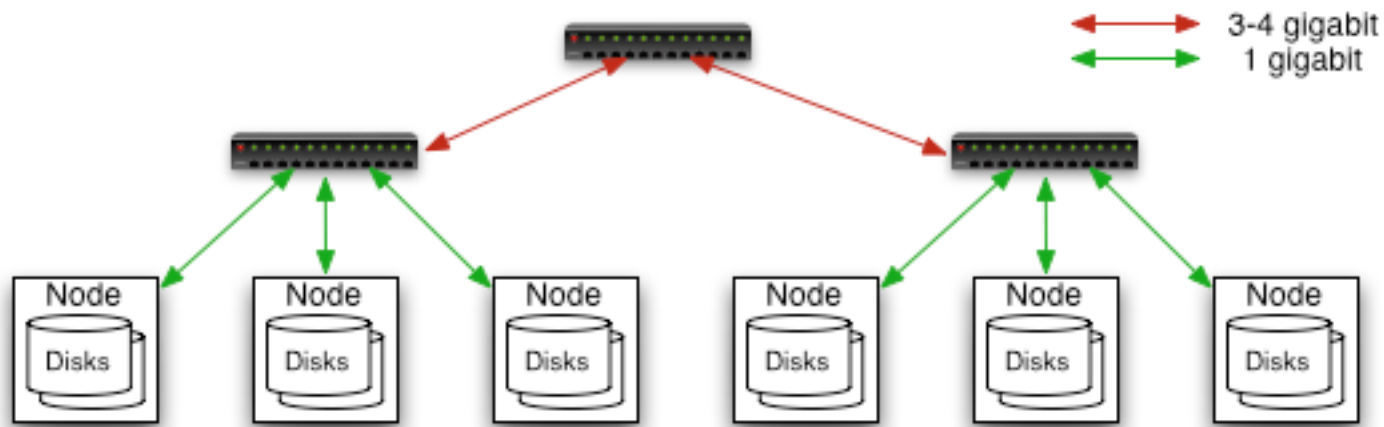


# Hadoop: How?

- Commodity Hardware Cluster
- Distributed File System
  - Modeled on GFS
- Distributed Processing Framework
  - Using Map/Reduce metaphor
- Open Source, Java
  - Apache Lucene subproject



# Commodity Hardware Cluster



- Typically in 2 level architecture
  - Nodes are commodity PCs
  - 30-40 nodes/rack
  - Uplink from rack is 3-4 gigabit
  - Rack-internal is 1 gigabit



# Distributed File System

- Single namespace for entire cluster
  - Managed by a single *namenode*.
  - Hierarchical directories
  - Optimized for streaming reads of large files.
- Files are broken in to large blocks.
  - Typically 64 or 128 MB
  - Replicated to several *datanodes*, for reliability
  - Clients can find location of blocks
- Client talks to both namenode and datanodes
  - Data is not sent through the namenode.



# Distributed Processing

- User submits Map/Reduce *job* to *JobTracker*
- System:
  - Splits job into lots of *tasks*
  - Schedules tasks on nodes close to data
  - Monitors tasks
  - Kills and restarts if they fail/hang/disappear
- Pluggable file systems for input/output
  - Local file system for testing, debugging, etc...



# Map/Reduce Metaphor

- Abstracts a very common pattern (munge, regroup, munge)
- Natural for
  - Building or updating offline databases (eg. indexes)
  - Computing statistics (eg. query log analysis)
- Software framework
  - Frozen part: distributed sort, and reliability via re-execution
  - Hot parts: input, map, partition, compare, reduce, and output



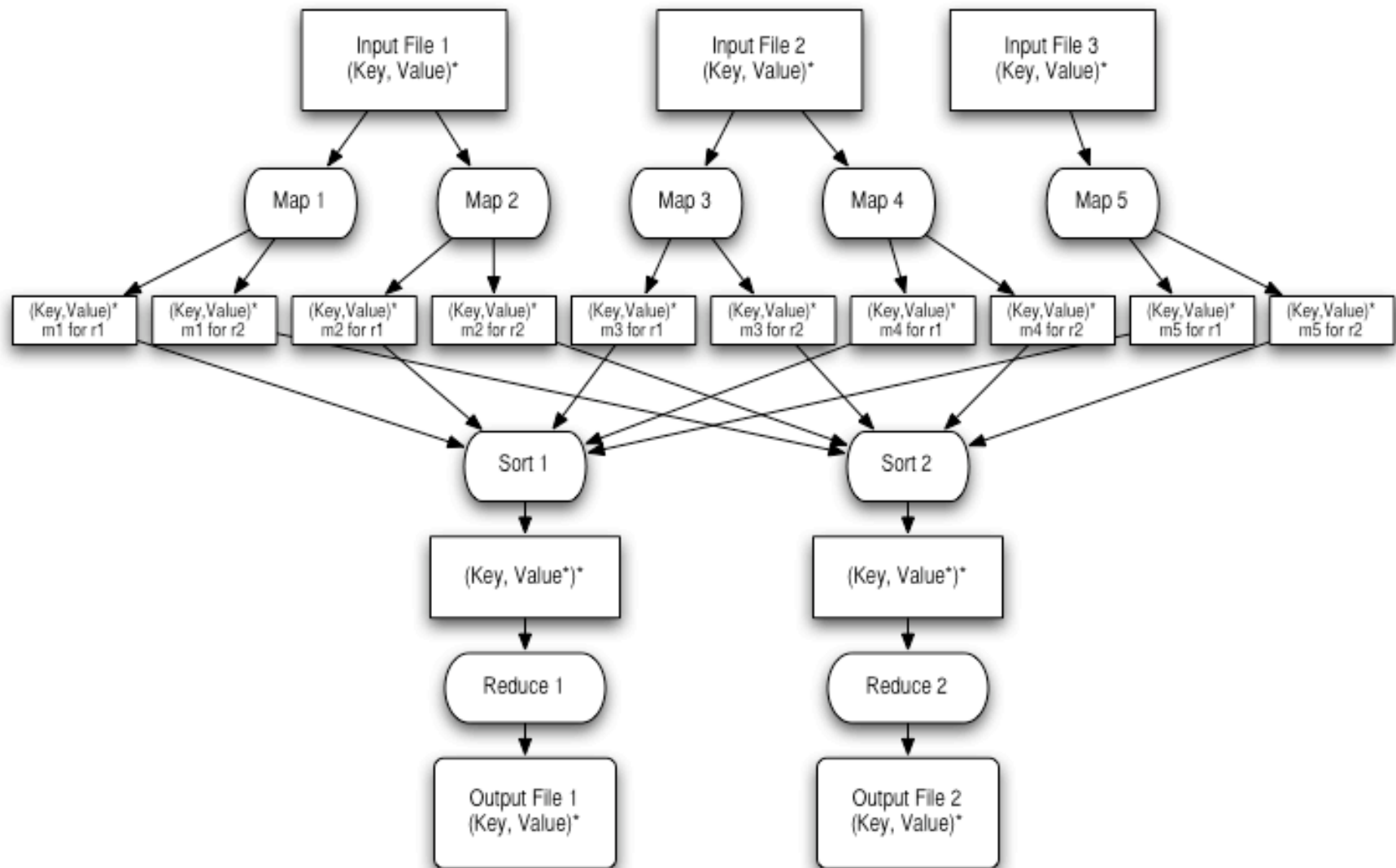
# Map/Reduce Metaphor

- Data is a stream of *keys* and *values*
- Mapper
  - Input: *key1, value1* pair
  - Output: *key2, value2* pairs
- Reducer
  - Called once per a key, in sorted order
  - Input: *key2, **stream** of value2*
  - Output: *key3, value3* pairs
- Launching Program
  - Creates a JobConf to define a job.
  - Submits JobConf and waits for completion.





# Map/Reduce Dataflow



# Map/Reduce Optimizations

- Overlap of maps, shuffle, and sort
- Mapper locality
  - Schedule mappers close to the data.
- Combiner
  - Mappers may generate duplicate keys
  - Side-effect free reducer run on mapper node
  - Minimize data size before transfer
  - Reducer is still run
- Speculative execution
  - Some nodes may be slower
  - Run duplicate task on another node



# HOWTO: Setting up Cluster

- Modify **hadoop-site.xml** to set directories and master hostnames.
- Create a **slaves** file that lists the worker machines one per a line.
- Run **bin/start-dfs** on the namenode.
- Run **bin/start-mapred** on the jobtracker.



# HOWTO: Write Application

- To write a distributed word count program:
  - Mapper: Given a line of text, break it into words and output the word and the count of 1:
    - “hi Apache bye Apache” ->
    - (“hi”, 1), (“Apache”, 1), (“bye”, 1), (“Apache”, 1)
  - Combiner/Reducer: Given a word and a set of counts, output the word and the sum
    - (“Apache”, [1, 1]) -> (“Apache”, 2)
  - Launcher: Builds the configuration and submits job



# Word Count Mapper

```
public class WCMapper extends MapReduceBase implements Mapper {  
  
    private static final IntWritable ONE = new IntWritable(1);  
  
    public void map(WritableComparable key, Writable value,  
                   OutputCollector output,  
                   Reporter reporter) throws IOException {  
        StringTokenizer itr = new StringTokenizer(value.toString());  
        while (itr.hasMoreTokens()) {  
            output.collect(new Text(itr.nextToken()), ONE);  
        }  
    }  
}
```



# Word Count Reduce

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



# Word Count Launcher

```
public static void main(String[] args) throws IOException {
    JobConf conf = new JobConf(WordCount.class);
    conf.setJobName("wordcount");

    conf.setOutputKeyClass(Text.class);
    conf.setOutputValueClass(IntWritable.class);

    conf.setMapperClass(WCMap.class);
    conf.setCombinerClass(WCReduce.class);
    conf.setReducerClass(WCReduce.class);

    conf.setInputPath(new Path(args[0]));
    conf.setOutputPath(new Path(args[1]));

    JobClient.runJob(conf);
}
```



```
Shell No. 2 - Konsole
Session Edit View Bookmarks Settings Help

oom@kry2000> bin/hadoop dfs -ls demo
Found 1 items
/user/oom/demo/books    <dir>
oom@kry2000> bin/hadoop dfs -ls demo/books | head
Found 1000 items
/user/oom/demo/books/part-0000 <r 3> 1120139520
/user/oom/demo/books/part-0001 <r 3> 1120139520
/user/oom/demo/books/part-0002 <r 3> 1120139520
/user/oom/demo/books/part-0003 <r 3> 1120139520
/user/oom/demo/books/part-0004 <r 3> 1120139520
/user/oom/demo/books/part-0005 <r 3> 1120139520
/user/oom/demo/books/part-0006 <r 3> 1120139520
/user/oom/demo/books/part-0007 <r 3> 1120139520
/user/oom/demo/books/part-0008 <r 3> 1120139520
oom@kry2000> bin/hadoop jar hadoop-0.12.4-dev-examples.jar wordcount demo/books demo/words
07/05/03 07:16:17 INFO mapred.FileInputFormat: Total input paths to process : 1000
07/05/03 07:16:27 INFO mapred.JobClient: Running job: job_0001
07/05/03 07:16:28 INFO mapred.JobClient:  map 0% reduce 0%
07/05/03 07:16:39 INFO mapred.JobClient:  map 1% reduce 0%
```





# kry2001 Hadoop Map/Reduce Administration

**Started:** Thu May 03 07:12:07 UTC 2007

**Version:** 0.12.4-dev, r534234

**Compiled:** Wed May 2 11:46:34 UTC 2007 by oom

## Cluster Summary

Maps	Reduces	Tasks/Node	Nodes
941	23	2	<a href="#">913</a>

## Running Jobs

Running Jobs								
Jobid	User	Name	Map % complete	Map total	Maps completed	Reduce % complete	Reduce total	Reduces
<a href="#">job_0001</a>	oom	wordcount	28.48%	17000	4830	0.00%	1000	0

## Completed Jobs



Hadoop reduce task list for job_0001 on kry2001 - Mozilla Firefox				
File Edit View History Bookmarks Tools Help				
<a href="#">tip_0001_r_000024</a>	33.29%	reduce > copy (16979 of 17000 at 0.00 MB/s) >	3-May-2007 07:16:51	
<a href="#">tip_0001_r_000025</a>	33.29%	reduce > copy (16982 of 17000 at 0.00 MB/s) >	3-May-2007 07:16:51	
<a href="#">tip_0001_r_000026</a>	33.29%	reduce > copy (16980 of 17000 at 0.00 MB/s) >	3-May-2007 07:16:51	
<a href="#">tip_0001_r_000027</a>	33.32%	reduce > copy (16994 of 17000 at 0.00 MB/s) >	3-May-2007 07:16:51	
<a href="#">tip_0001_r_000028</a>	33.30%	reduce > copy (16983 of 17000 at 0.00 MB/s) >	3-May-2007 07:16:51	
<a href="#">tip_0001_r_000029</a>	33.25%	reduce > copy (16959 of 17000 at 0.00 MB/s) >	3-May-2007 07:16:51	
<a href="#">tip_0001_r_000032</a>	33.33%	reduce > copy (16999 of 17000 at 0.00 MB/s) >	3-May-2007 07:16:51	
<a href="#">tip_0001_r_000033</a>	33.26%	reduce > copy (16967 of 17000 at 0.00 MB/s) >	3-May-2007 07:16:51	
<a href="#">tip_0001_r_000034</a>	33.31%	reduce > copy (16989 of 17000 at 0.00 MB/s) >	3-May-2007 07:16:51	



Kind	% Complete	Num Tasks	Pending	Running	Complete	Killed	<a href="#">Failed/Killed Task Attempts</a>
<a href="#">map</a>	100.00%	17000	0	0	17000	0	0 / 0
<a href="#">reduce</a>	100.00%	1000	0	0	1000	0	0 / 0

	Counter	Map	Reduce	Total
org.apache.hadoop.examples.WordCount\$Counter	WORDS	11,903,802,000	0	11,903,802,000
	VALUES	11,903,802,000	129,009,000	12,032,811,000
Map-Reduce Framework	Map input records	1,339,416,000	0	1,339,416,000
	Map output records	11,903,802,000	0	11,903,802,000
	Map input bytes	68,157,030,000	0	68,157,030,000
	Map output bytes	113,957,261,000	0	113,957,261,000
	Combine input records	11,903,802,000	0	11,903,802,000
	Combine output records	129,009,000	0	129,009,000



# Running on Amazon EC2/S3

- Amazon sells cluster services
  - EC2: \$0.10/cpu hour
  - S3: \$0.20/gigabyte month
- Hadoop supports:
  - EC2: cluster management scripts included
  - S3: file system implementation included
- Tested on 400 node cluster
- Combination used by several startups



# Hadoop On Demand

- Traditionally Hadoop runs with dedicated servers
- Hadoop On Demand works with a batch system to allocate and provision nodes dynamically
  - Bindings for Condor and Torque/Maui
- Allows more dynamic allocation of resources



# Scalability

- Runs on 1000 nodes
- 5TB sort on 500 nodes takes 2.5 hours
- Distributed File System:
  - 150 TB
  - 3M files



# Thank You

- Questions?
- For more information:
  - <http://lucene.apache.org/hadoop/>

