

CS 425 / ECE 428
Distributed Systems
Fall 2016

Indranil Gupta (Indy)

Aug 30-Sep 1, 2016

Lecture 4: Mapreduce and Hadoop

All slides © IG

WHAT IS MAPREDUCE?

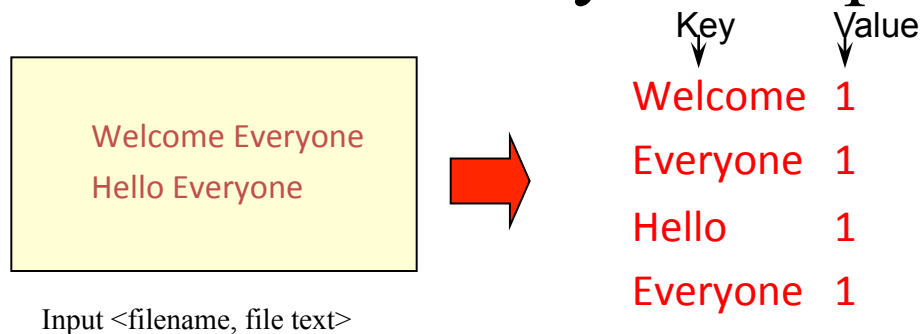
- Terms are borrowed from Functional Language (e.g., Lisp)

Sum of squares:

- (map square '(1 2 3 4))
 - Output: (1 4 9 16)
 - [processes each record sequentially and independently]
- (reduce + '(1 4 9 16))
 - (+ 16 (+ 9 (+ 4 1)))
 - Output: 30
 - [processes set of all records in batches]
- Let's consider a sample application: **Wordcount**
 - You are given a huge dataset (e.g., Wikipedia dump or all of Shakespeare's works) and asked to list the count for each of the words in each of the documents therein

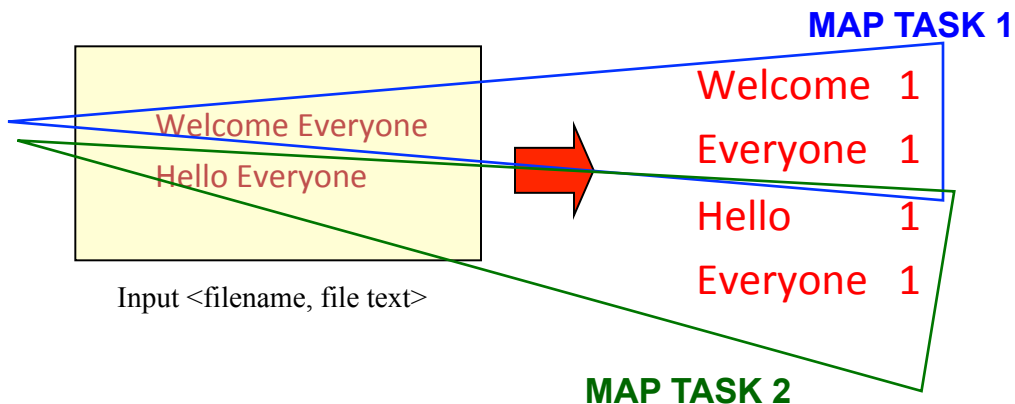
MAP

- Process individual records to generate intermediate key/value pairs.



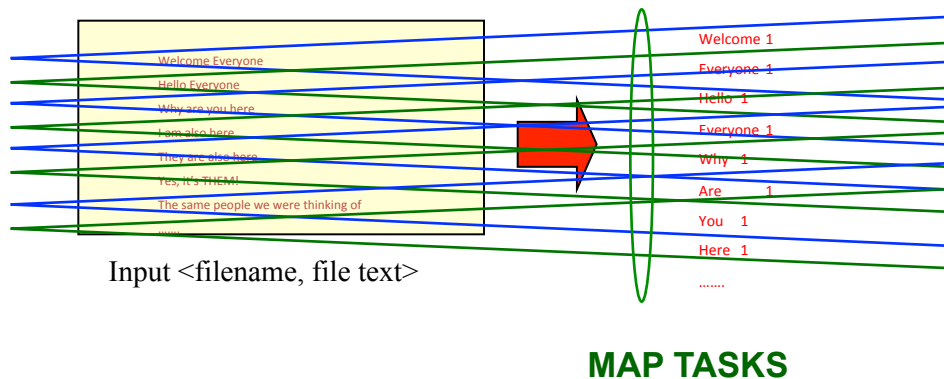
MAP

- **Parallely** Process individual records to generate intermediate key/value pairs.



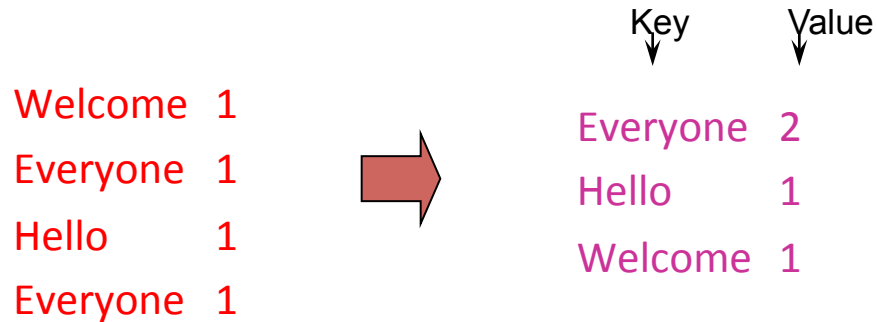
MAP

- **Parallely** Process a large number of individual records to generate intermediate key/value pairs.



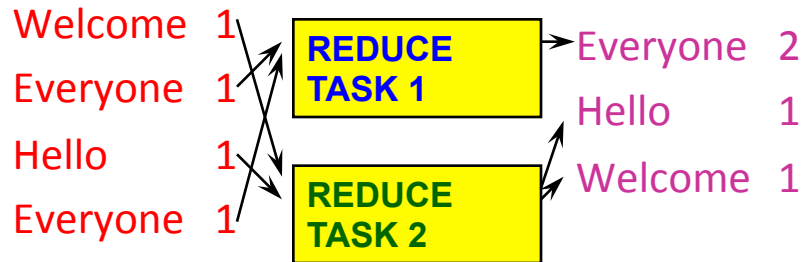
REDUCE

- Reduce processes and merges all intermediate values associated per key



REDUCE

- Each key assigned to one Reduce
- Parallely Processes and merges all intermediate values by partitioning keys



- Popular: *Hash partitioning, i.e.,* key is assigned to
 - $\text{reduce \#} = \text{hash}(\text{key}) \% \text{number of reduce tasks}$

HADOOP CODE - MAP

```
public static class MapClass 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)
        // key is empty, value is the line
        throws IOException {
        String line = value.toString();
        StringTokenizer itr = new StringTokenizer(line);
        while (itr.hasMoreTokens()) {
            word.set(itr.nextToken());
            output.collect(word, one);
        }
    }
}
// Source: http://developer.yahoo.com/hadoop/tutorial/module4.html#wordcount
```


HADOOP CODE - REDUCE

```
public static class ReduceClass extends MapReduceBase implements
Reducer<Text, IntWritable, Text, IntWritable> {
    public void reduce(
        Text key,
        Iterator<IntWritable> values,
        OutputCollector<Text, IntWritable> output,
        Reporter reporter)
    throws IOException {
        // key is word, values is a list of 1's
        int sum = 0;
        while (values.hasNext()) {
            sum += values.next().get();
        }
        output.collect(key, new IntWritable(sum));
    }
} // Source: http://developer.yahoo.com/hadoop/tutorial/module4.html#wordcount
```

HADOOP CODE - DRIVER

```
// Tells Hadoop how to run your Map-Reduce job
public void run (String inputPath, String outputPath)
    throws Exception {
    // The job. WordCount contains MapClass and Reduce.
    JobConf conf = new JobConf(WordCount.class);
    conf.setJobName("mywordcount");
    // The keys are words
    (strings) conf.setOutputKeyClass(Text.class);
    // The values are counts (ints)
    conf.setOutputValueClass(IntWritable.class);
    conf.setMapperClass(MapClass.class);
    conf.setReducerClass(ReduceClass.class);
    FileInputFormat.addInputPath(
        conf, newPath(inputPath));
    FileOutputFormat.setOutputPath(
        conf, new Path(outputPath));
    JobClient.runJob(conf);
}
```

// Source: <http://developer.yahoo.com/hadoop/tutorial/module4.html#wordcount>

SOME APPLICATIONS OF MAPREDUCE

Distributed Grep:

- Input: large set of files
- Output: lines that match pattern

- Map – *Emits a line if it matches the supplied pattern*
- Reduce – *Copies the intermediate data to output*

SOME APPLICATIONS OF MAPREDUCE (2)

Reverse Web-Link Graph

- Input: Web graph: tuples (a, b) where (page a \rightarrow page b)
- Output: For each page, list of pages that link *to* it

- Map – *process web log and for each input $\langle source, target \rangle$, it outputs $\langle target, source \rangle$*
- Reduce - *emits $\langle target, list(source) \rangle$*

SOME APPLICATIONS OF MAPREDUCE (3)

Count of URL access frequency

- Input: Log of accessed URLs, e.g., from proxy server
- Output: For each URL, % of total accesses for that URL

- Map – *Process web log and outputs <URL, 1>*
- Multiple Reducers - *Emits <URL, URL_count>*
(So far, like Wordcount. But still need %)
- Chain another MapReduce job after above one
- Map – *Processes <URL, URL_count> and outputs <1, (<URL, URL_count>)>*
- 1 Reducer – Does two passes. In first pass, sums up all *URL_count's* to calculate overall_count. In second pass calculates %'s
Emits multiple <URL, URL_count/overall_count>

SOME APPLICATIONS OF MAPREDUCE (4)

Map task's output is sorted (e.g., quicksort)

Reduce task's input is sorted (e.g., mergesort)

Sort

- Input: Series of (key, value) pairs
- Output: Sorted <value>s

- Map – $\langle \text{key}, \text{value} \rangle \rightarrow \langle \text{value}, _ \rangle$ (*identity*)
- Reducer – $\langle \text{key}, \text{value} \rangle \rightarrow \langle \text{key}, \text{value} \rangle$ (*identity*)
- Partitioning function – partition keys across reducers based on **ranges** (can't use hashing!)
 - Take data distribution into account to balance reducer tasks

PROGRAMMING MAPREDUCE

Externally: For **user**

1. Write a Map program (short), write a Reduce program (short)
2. Specify number of Maps and Reduces (parallelism level)
3. Submit job; wait for result
4. Need to know very little about parallel/distributed programming!

Internally: For the Paradigm and Scheduler

1. Parallelize Map
2. Transfer data from Map to Reduce (**shuffle data**)
3. Parallelize Reduce
4. Implement Storage for Map input, Map output, Reduce input, and Reduce output

(Ensure that no Reduce starts before all Maps are finished. That is, ensure the ***barrier*** between the Map phase and Reduce phase)

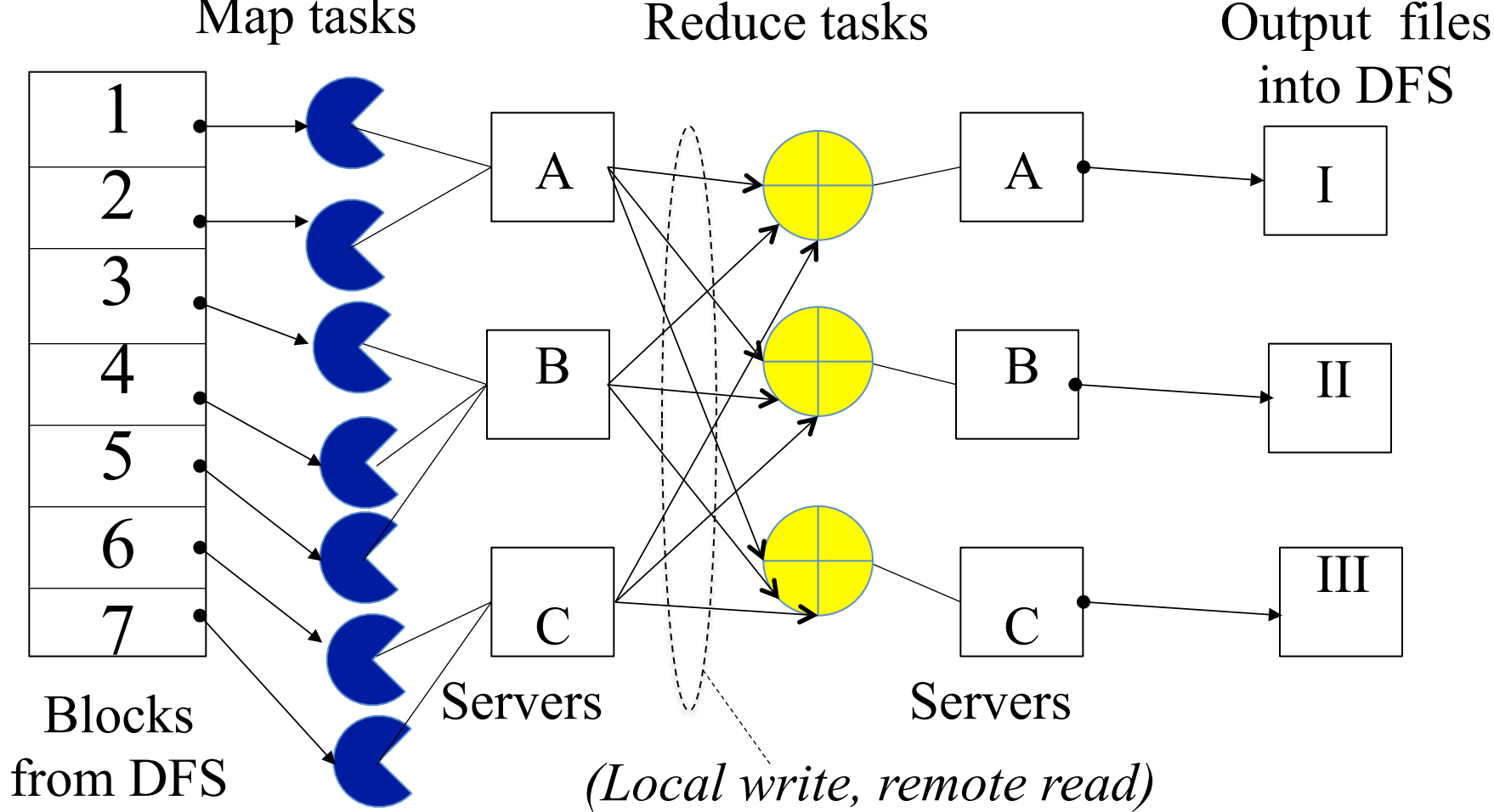
INSIDE MAPREDUCE

For the cloud:

1. Parallelize Map: **easy!** each map task is independent of the other!
 - All Map output records with same key assigned to same Reduce
2. Transfer data from Map to Reduce:
 - Called Shuffle data
 - All Map output records with same key assigned to same Reduce task
 - use **partitioning function, e.g., $\text{hash}(\text{key})\% \text{number of reducers}$**
3. Parallelize Reduce: **easy!** each reduce task is independent of the other!
4. Implement Storage for Map input, Map output, Reduce input, and Reduce output
 - Map input: from **distributed file system**
 - Map output: to local disk (at Map node); uses **local file system**
 - Reduce input: from (multiple) remote disks; uses local file systems
 - Reduce output: to distributed file system

local file system = Linux FS, etc.

distributed file system = GFS (Google File System), HDFS (Hadoop Distributed File System)



Resource Manager (assigns maps and reduces to servers)

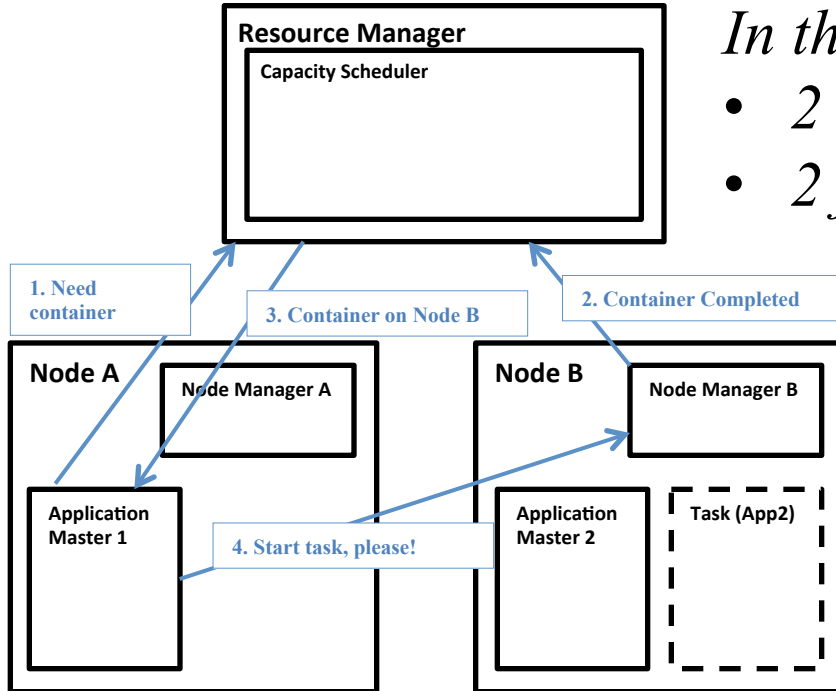
THE YARN SCHEDULER

- Used underneath Hadoop 2.x +
- YARN = Yet Another Resource Negotiator
- Treats each server as a collection of *containers*
 - Container = fixed CPU + fixed memory (think of Linux cgroups, but even more lightweight)
- Has 3 main components
 - Global *Resource Manager (RM)*
 - Scheduling
 - Per-server *Node Manager (NM)*
 - Daemon and server-specific functions
 - Per-application (job) *Application Master (AM)*
 - Container negotiation with RM and NMs
 - Detecting task failures of that job

YARN: HOW A JOB GETS A CONTAINER

In this figure

- 2 servers (A, B)
- 2 jobs (1, 2)



FAULT TOLERANCE

- Server Failure
 - NM heartbeats to RM
 - If server fails: RM times out waiting for next heartbeat, RM lets all affected AMs know, and AMs take appropriate action
 - NM keeps track of each task running at its server
 - If task fails while in-progress, mark the task as idle and restart it
 - AM heartbeats to RM
 - On failure, RM restarts AM, which then syncs it up with its running tasks
- RM Failure
 - Use old checkpoints and bring up secondary RM
- Heartbeats also used to piggyback container requests
 - Avoids extra messages

SLOW SERVERS

Slow tasks are called **Stragglers**

- The slowest task slows the entire job down (why?)
 - Due to Bad Disk, Network Bandwidth, CPU, or Memory
 - Keep track of “progress” of each task (% done)
 - Perform proactive backup (replicated) execution of some straggler tasks
 - A task considered done when its first replica complete (other replicas can then be killed)
 - Approach called **Speculative Execution**.
- Barrier at the end of Map phase!

LOCALITY

- Locality
 - Since cloud has hierarchical topology (e.g., racks)
 - For server-fault-tolerance, GFS/HDFS stores 3 replicas of each of chunks (e.g., 64 MB in size)
 - For rack-fault-tolerance, on different racks, e.g., 2 on a rack, 1 on a different rack
 - Mapreduce attempts to schedule a map task on
 1. a machine that contains a replica of corresponding input data, or failing that,
 2. on the same rack as a machine containing the input, or failing that,
 3. Anywhere
 - Note: The 2-1 split of replicas is intended to reduce bandwidth when writing file.
 - Using more racks does not affect overall Mapreduce scheduling performance

MAPREDUCE: SUMMARY

- Mapreduce uses parallelization + aggregation to schedule applications across clusters
- Need to deal with failure
- Plenty of ongoing research work in scheduling and fault-tolerance for Mapreduce and Hadoop

ANNOUNCEMENTS

- HW1 due Sep 20th
- MP1 due Sep 11th
 - VMs distributed (check Piazza)
 - Demos will be Monday Sep 12th (schedule and details will be posted next week on Piazza)
- Please fill out Student Survey by today (course webpage).
- Check Piazza often! It's where all the announcements are at!
- If your MP partner dropped, please form new group and let us know (cs-425-staff)
 - TODAY is the last day
 - DO NOT change MP groups unless your partner has dropped