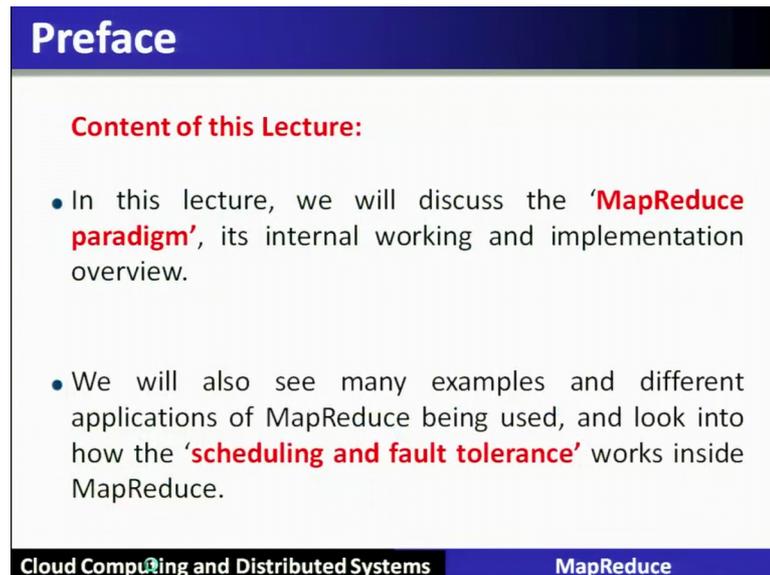


**Cloud Computing and Distributed Systems**  
**Dr. Rajiv Misra**  
**Department of Computer Science and Engineering**  
**Indian Institute of Technology, Patna**

**Lecture – 19**  
**Map Reduce**

(Refer Slide Time: 00:15)



**Preface**

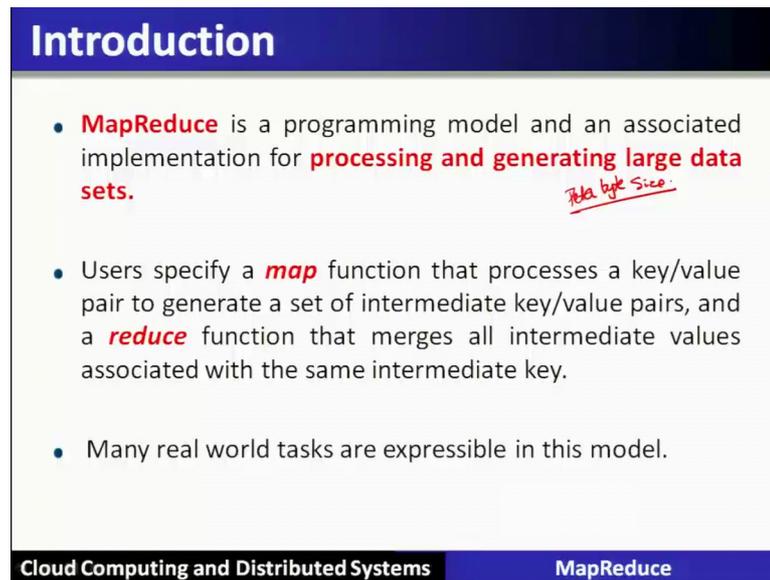
**Content of this Lecture:**

- In this lecture, we will discuss the '**MapReduce paradigm**', its internal working and implementation overview.
- We will also see many examples and different applications of MapReduce being used, and look into how the '**scheduling and fault tolerance**' works inside MapReduce.

Cloud Computing and Distributed Systems      MapReduce

MapReduce, preface content of this lecture we will discuss 'MapReduce paradigm' it is internal working and it is implementation overview. We will also see many examples, using MapReduce and different applications of MapReduce. And also we will look into the aspects how the scheduling and fault tolerance is done inside the MapReduce.

(Refer Slide Time: 00:42)



## Introduction

- **MapReduce** is a programming model and an associated implementation for **processing and generating large data sets**. *700 byte Size*
- Users specify a **map** function that processes a key/value pair to generate a set of intermediate key/value pairs, and a **reduce** function that merges all intermediate values associated with the same intermediate key.
- Many real world tasks are expressible in this model.

Cloud Computing and Distributed Systems      MapReduce

Introduction: MapReduce is a programming model and is associated and implementation for processing and generating large data sets. So, the large data sets we mean that if the data cannot be stored and processed with within one computer system that is in the range of petabytes size.

Therefore, we require a new paradigm that is called MapReduce, which we will enable to process such a large data sets, which are useful in today's different applications. So, here in this simple programming model users have to specify a function which is called a map that processes key value pair to generate a set of intermediate key value pairs, which in turn will be taken up by another function which is called a reduce that merges all the intermediate values associated with the same intermediate keys. And therefore, together map and reduce will work together to basically solve all the many applications in the real world. So, many real world tasks are expressible in this particular model.

(Refer Slide Time: 02:18)

**Contd...**

- Programs written in this functional style **are automatically parallelized and executed** on a large cluster of commodity machines.
- The **run-time system** takes care of the details of partitioning the input data, scheduling the program's execution across a set of machines, handling machine failures, and managing the required inter-machine communication.
- This allows programmers without any experience with parallel and distributed systems to easily utilize the resources of a large distributed system.
- A **typical MapReduce computation processes** many terabytes of data on thousands of machines. Hundreds of MapReduce programs have been implemented and upwards of one thousand MapReduce jobs are executed on **Google's clusters** every day.

Cloud Computing and Distributed Systems      MapReduce

Now, programs written in this functional style are automatically parallelized and executed on a large cluster of commodity machines. The run time system takes care of the details of partitioning the input data, scheduling the programs execution across the set of machines, handling machine failures, managing the required inter machine communication that is all will be made abstraction from the programmers. So, this allows the programmers without having any experience with the parallel and distributed system to easily utilize these resources for solving the problem which comprises of a large data set.

A typical MapReduce computation processes many terabytes of data on thousands of machines. So, that is the scale on which this MapReduce computation can work. Hundreds of MapReduce programs have been implemented and upwards of one thousand of MapReduce jobs are executed on Google's cluster every day.

(Refer Slide Time: 03:29)

**Distributed File System**

- Chunk Servers**
  - File is split into contiguous chunks
  - Typically each chunk is 16-64MB
  - Each chunk replicated (usually 2x or 3x)
  - Try to keep replicas in different racks
- Master node**
  - Also known as Name Nodes in HDFS
  - Stores metadata
  - Might be replicated
- Client library for file access**
  - Talks to master to find chunk servers
  - Connects directly to chunkservers to access data

Cloud Computing and Distributed Systems      MapReduce

This MapReduce internally uses a distributed file system, which is known as HDFS and there are certain components which will be useful in the design. So, we will discuss this HDFS in a separate lecture, but for this point of discussion at this point of time we will understand that the distributed file system has the master node which also known as the name node in HDFS. And this stores the metadata and these may be replicated as per different file system and here we are using the HDFS file system.

There are some client libraries for the file accesses. So, they will talk to the master to find out the available chunk servers, which connects directly to the chunk server to access the data. So, chunk servers are the servers who manages the splits into the contiguous chunks. And each chunk is of the size 16-64 MB and these chunks are replicated normally by default it is 3 times replicated in HDFS file system. This replication is also ensures that these replication also goes into different racks. So, that it will become a rack tolerant, rack fault, tolerant system it can build up in this process of replication.

(Refer Slide Time: 05:16)

## Motivation for Map Reduce (Why ?)

- **Large-Scale Data Processing**
  - Want to use 1000s of CPUs
  - But don't want hassle of managing things ✓
- **MapReduce Architecture provides**
  - Automatic parallelization & distribution ✓
  - Fault tolerance
  - I/O scheduling
  - Monitoring & status updates

Cloud Computing and Distributed Systems      MapReduce

So, let us see the motivation of using MapReduce for large scale data processing. So, to handle this scale of data processing we required 1000s of CPUs and also we do not want the hassle of managing such a large system. Therefore, the MapReduce is a programming model and its architecture will provide automatic parallelization. And the distribution, fault tolerance, and I/O scheduling, monitoring, and status updates we will see their more details in further slides.

(Refer Slide Time: 06:04)

## 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 ✓

Cloud Computing and Distributed Systems      MapReduce

So, MapReduce paradigm let us understand from the scratch about this paradigm. These terms map and reduce their borrow from the functional language which is called a lisp. So, for example, the map function which is there in the list, if it is having the operation or a function called square and a Lisp is given. Then this particular function will be applied on each and every element of this particular list.

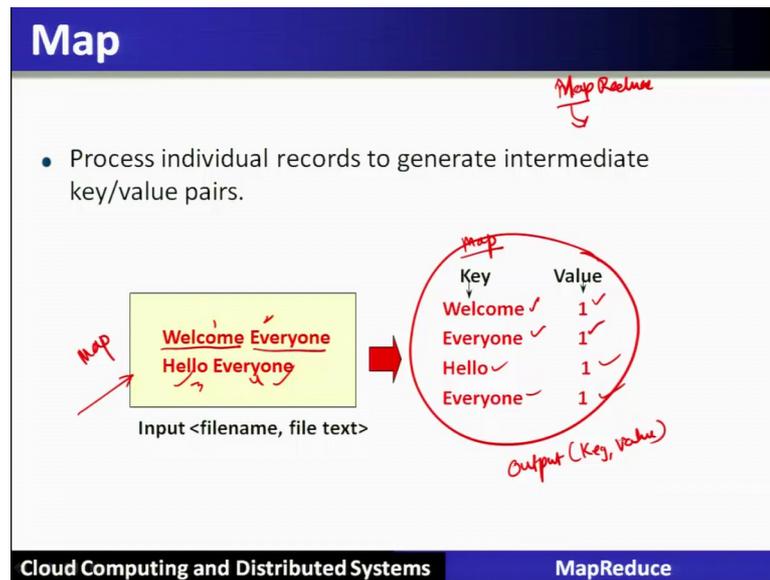
Now this particular function can be executed in parallel why because, this particular function if it is applied on these input values they can be able to generate the output that is the square of 1 is 1 square of 2 is 4 3 is 9 4 is 16. So, if four different machines are involved or four different threads are executing this is square at four different elements. So, in one go it will be able to generate this.

So, that is how the inherent parallelism is being exploited here in the functional languages and that is how the MapReduce will also do the same way without botheration of internally how this parallel programming is to be done by the programmer end. Similarly there will be a reduce function so, in the reduce function there is an operation that is the function which is defined.

So, function is a plus which involves all the elements, which are to be added and given the output. So, internally it will be represented in this particular equation form and two operators at a time because, plus is a binary operation two at a time this will form a complete expression and this will be evaluated to give the output 30. So, these map and reduce functions which are there in functional language like lisp, it is being inspired and it is not that map and reduce of lisp function, but it is a separate construct which is called a MapReduce which is now available to solve this kind of problem.

So, the problems of application for example, a large document large data set for example, it is a dump of a Wikipedia. And here you want to find out the word Shakespeare or some other keyword and you want to list the count for each of these words in a given document there in if that is the application. Then a MapReduce construct will be very easy and its program is called a word count. So, let us see how the word count program will work, where in we can submit a huge data set and using the word count program of MapReduce it can able to solve.

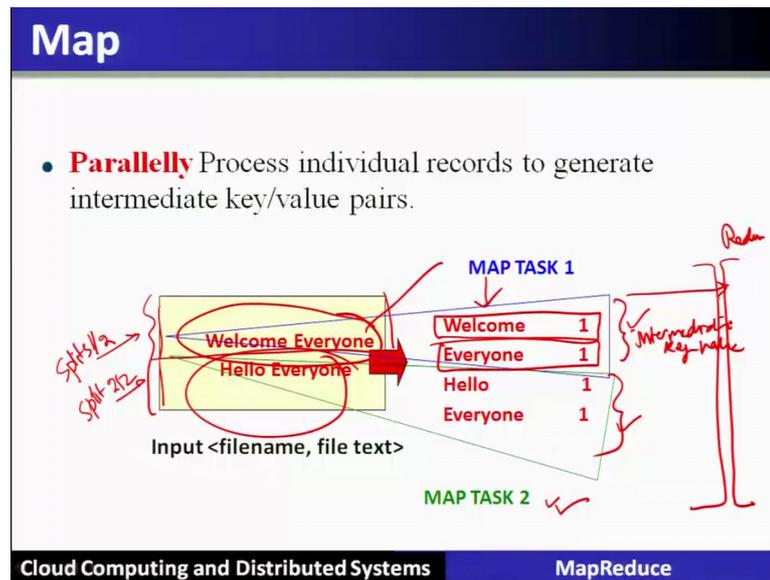
(Refer Slide Time: 09:35)



So, now in the MapReduce there are two functions map and reduce. So, out of them the map function let us discuss first. So, the process individual record to generate the intermediate key value pairs so, take this snapshot of the entire big document, let us say it has four different words in this particular sentence that is welcome everyone hello everyone. Now this particular document when it is given as the input. Then this map function should be able to generate the key value pairs. So, key means every keyword it will generate and it is count it will emit in the map that value of 1. So, for every word it will emit 1 count.

Similarly, hello also will emit 1 count and everyone is appearing again so, it will be generating that particular count so that is the function of map. So, map will produce this as an output. So, when map is given this particular filename, it will generate this particular output that is in the form of a key value pair.

(Refer Slide Time: 11:23)



Now, this particular map function will processes these individual records to generate the intermediate key value pairs that is what we have shown. Now, this particular task can be executed by two different map task or we can divide this file into two different, two different chunks or we can call them as a splits. So, there are two splits and we intend to execute in parallel. So, that this particular document can be processed big size document can be processed.

So, in that process these particular splits are given to different servers, they are called as a map task 1 and the other one is called map task 2. So, this particular input is given to the map task 1 and this map will generate the intermediate key value pairs they are called intermediate key value pairs. So, also the map task 2 also will do the same thing. These intermediate key value pairs will now be given as the input to the next phase which is called reduced phase. So, for reduced phase the output of the map phase will become the input.

(Refer Slide Time: 13:08)

## Map

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

The diagram illustrates the Map phase. On the left, a box labeled 'Input <filename, file text>' contains several lines of text. A red arrow labeled 'MAP TASKS' points from this input to a series of vertical lines representing intermediate key-value pairs. The pairs are: 'Welcome', 'Everyone', 'Hello', 'Welcome', 'Hello', 'Everyone', 'Hello', 'Everyone', 'Hello', 'Everyone'. The bottom of the slide contains the text 'Cloud Computing and Distributed Systems' and 'MapReduce'.

So, the parallelly processing a large number of individual records to generate the intermediate key value pair is basically, inherent in this paradigm that is called map using the map tasks.

(Refer Slide Time: 13:26)

## Reduce

- Reduce processes and merges all intermediate values associated per **key**

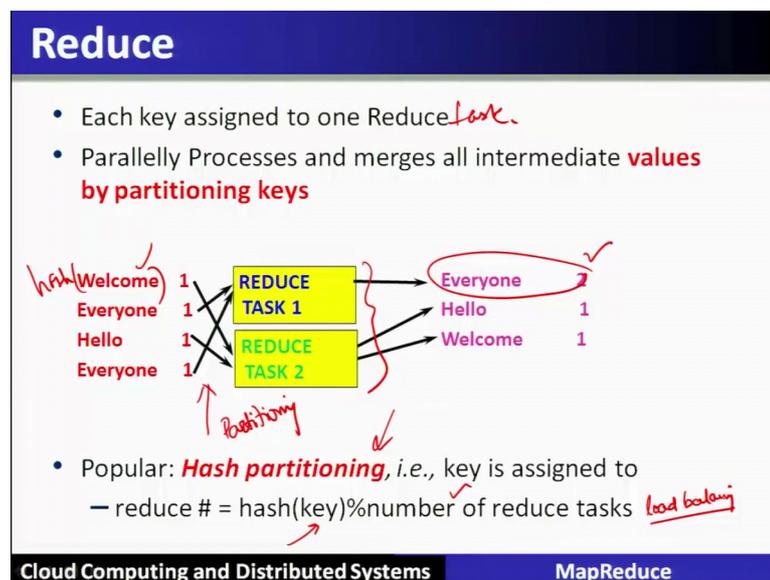
The diagram illustrates the Reduce phase. On the left, a box labeled 'Intermediate Key Value' contains several lines of text: 'Welcome 1', 'Everyone 1', 'Hello 1', 'Everyone 1'. A red arrow labeled 'Reduce' points from this input to a series of vertical lines representing the final key-value pairs. The pairs are: 'Everyone 2', 'Hello 1', 'Welcome 1'. The bottom of the slide contains the text 'Cloud Computing and Distributed Systems' and 'MapReduce'.

Then this particular output of map phase is given to the reducer. That is the key value pairs are given this is called intermediate key value pairs, key value pairs are given to the reduce function. So, reduce processes and merges all the intermediate values associated per key together in this phase. So, here you can see that everyone is appearing twice. So,

they should be merged in the reduce phase and therefore, it will give the final output associated with every key.

So, let us see when this happens using the reduce function what it does? It combines or it merges, these two keywords and it will generate the combined values. Similarly, hello is appearing only once and welcome is also appearing once. Now to merge these different keys which are having the same keys they have to be merged together in the reducer phase. So, they have to give be given to the same reducer. And this hello keyword may be given to the separate reducer. Now there are requirements of three different reducers, but if you have only two machines available, which can be assigned to these reducer. So, a particular machine can handle two different or more than one different reducers also.

(Refer Slide Time: 15:40)



So, let us see the details of this reduce function. So, each key will be assigned to one reduce to one reducer or to one reduce task. So, parallely they will be processed and it will merge all the intermediate values associated of the same key to one reducer by partitioning this particular keys. So, here in this case you can see the partitioning or these keys are basically able to map to reducers which you have here shown in this example.

So, for example, map 1 will go to reduce task 2, everyone will go to reduce task 1 hello will go to reduce task 2 and everyone will go to the reduce task 1. And reduced task 1 will then combine these 2 everyone into that combined function that is it will aggregate.

Now this kind of partitioning of the keys to the reduce task can be done using hash partitioning.

So, using hash partitioning we can hash a particular key. And this particular key when it is hashed then it will be taken the modulo of the number of reduce tasks. In this manner this particular keys will be partitioned equally across all the reduce task. So, there will be a proper load balancing of these keywords which are assigned to this reduce function.

(Refer Slide Time: 17:40)

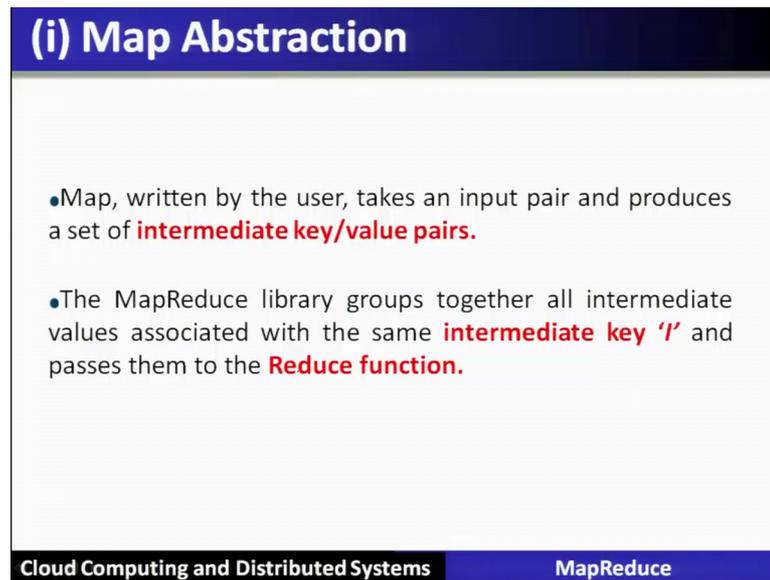
**Programming Model**

- The computation takes a set of **input key/value pairs**, and produces a set of **output key/value pairs**.
- The user of the MapReduce library expresses the computation as two functions:
  - (i) The Map ✓
  - (ii) The Reduce ✓

Cloud Computing and Distributed Systems      MapReduce

So, let us see the programming model in more detail of this MapReduce. So, the computation takes a set of input key and values pairs and produces a set of output key value pairs. So, the users they have to use the MapReduce library to express their computations using these two functions.

(Refer Slide Time: 18:06)



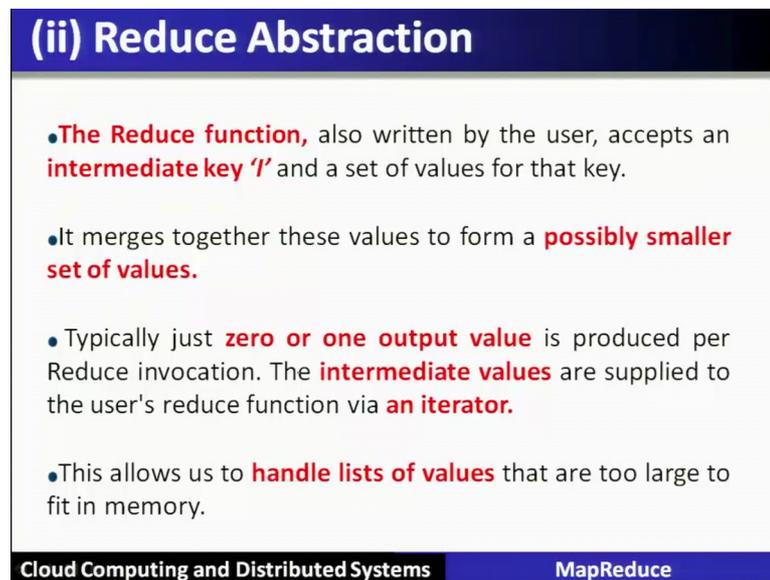
**(i) Map Abstraction**

- Map, written by the user, takes an input pair and produces a set of **intermediate key/value pairs**.
- The MapReduce library groups together all intermediate values associated with the same **intermediate key 'k'** and passes them to the **Reduce function**.

Cloud Computing and Distributed Systems      MapReduce

The map abstraction the map written by the user takes the input pairs. Pair and produces a set of intermediate key value pairs. MapReduce library a groups together all intermediate key values associated with that intermediate key and passes them to the reduce function.

(Refer Slide Time: 18:23)



**(ii) Reduce Abstraction**

- **The Reduce function**, also written by the user, accepts an **intermediate key 'k'** and a set of values for that key.
- It merges together these values to form a **possibly smaller set of values**.
- Typically just **zero or one output value** is produced per Reduce invocation. The **intermediate values** are supplied to the user's reduce function via **an iterator**.
- This allows us to **handle lists of values** that are too large to fit in memory.

Cloud Computing and Distributed Systems      MapReduce

Reduce abstraction reduce function is also written by the user accepts the intermediate keys, which are generated from the previous phase that is from the map phase. And the set of values for that it merges together these values to form a possibly a smaller set of

values. Typically zero or one output value is produced per reduce invocation. The intermediate values are supported to the users reduce functions where it later. This allows us to handle the list of values that are too large to be fit in the memory.

(Refer Slide Time: 19:00)

The slide is titled "Map-Reduce Functions for Word Count" in a blue header. It contains two code blocks. The first block is the map function: `map(key, value):` followed by a comment `// key: document name; value: text of document`, then `for each word w in value:` and `emit(w, 1)`. The second block is the reduce function: `reduce(key, values):` followed by a comment `// key: a word; values: an iterator over counts`, then `result = 0`, `for each count v in values:`, `result += v`, and `emit(key, result)`. Handwritten red annotations include arrows pointing from "key" to "document name" and "value" to "text of document" in the map function; and from "key" to "a word" and "values" to "an iterator over counts" in the reduce function. There are also checkmarks and a "word" label with an arrow pointing to the `emit(w, 1)` line.

```
map(key, value):  
// key: document name; value: text of document  
for each word w in value:  
    emit(w, 1)  
  
reduce(key, values):  
// key: a word; values: an iterator over counts  
result = 0  
for each count v in values:  
    result += v  
emit(key, result)
```

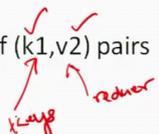
Cloud Computing and Distributed Systems      MapReduce

Let us see through an example of a word count. Here the map function will take the key values and key is nothing, but the name of the document and value is the text of that particular document. So, when given this particular document the name of a document and this is the text of a document or the entire document is given as a value then for each word which are separated by either the blanks are other separator. So, each word will be identified in the document or in the value of that is called the document. And it will emit that word  $w$  with a value 1 in the word count program.

Similarly, this output will be now taken up by the reducer, now here the key is a word and the values are the values which are being generated over here. So, values and iterator will count these values. So, let us begin this process with initialization of result is equal to 0. So, for each count  $v$  in these particular values, the result will be incremented by that value and it will emit for key this corresponding result and that will be output from this MapReduce function.

(Refer Slide Time: 20:35)

## Map-Reduce Functions

- **Input:** a set of key/value pairs ✓
- User supplies two functions:  
map(k,v) → list(k1,v1) ✓  
reduce(k1, list(v1)) → v2 ✓
- (k1,v1) is an intermediate key/value pair
- **Output** is the set of (k1,v2) pairs  


Cloud Computing and Distributed Systems MapReduce

So, in general input will be a set of key value pairs, that we have already seen and the user has to specify its own application through the map function; which will emit a list of intermediate key value pairs are now given as the input to the reduce function which will generate the output that is the values of all keys. So, output will be that keys and the final values which will be generated by the reduce function. And so, let us see some of the applications.

(Refer Slide Time: 21:24)

## Applications

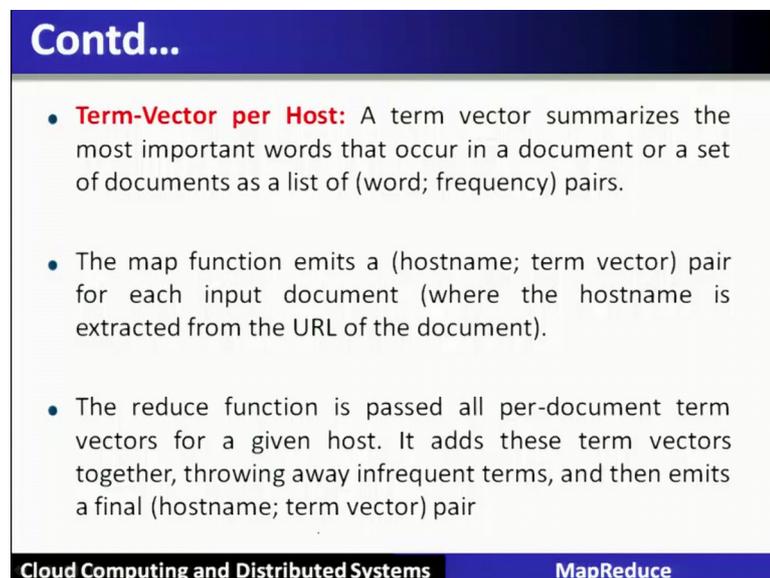
- Here are a few simple applications of interesting programs that can be easily expressed as **MapReduce computations**.
- **Distributed Grep:** The map function emits a line if it matches a supplied pattern. The reduce function is an identity function that just copies the supplied intermediate data to the output.
- **Count of URL Access Frequency:** The map function processes logs of web page requests and outputs (URL; 1). The reduce function adds together all values for the same URL and emits a (URL; total count) pair.
- **ReverseWeb-Link Graph:** The map function outputs (target; source) pairs for each link to a target URL found in a page named source. The reduce function concatenates the list of all source URLs associated with a given target URL and emits the pair: (target; list(source)) ✓

Cloud Computing and Distributed Systems MapReduce

The applications of this MapReduce paradigms are many some of them we will here list them to understand the details of these applications. So, distributed graph application; that means, distributed graph can be applied on the log files or an error files wherein we can check or we can find out that pattern and find out the what kind of error it is incurring. So, since the log files are too big therefore, this distributed graph is an ideal application. Similarly, the count of URL access frequency is also one of the important application. And the list of URL's are being generated by any proxy server or in the cloud who is running the websites.

So, how many such accesses are being made across different websites accesses by the user. So, they are total count also becomes a MapReduce application. Similarly a reverse web link graph that means, a map that means, we want to find out the graph of those websites and who are accessing those websites. So, that is called link graph. So, we want to find out the reverse web link graph means to that website how many people they are accessing it. So, if we want to know that that is that application is called reverse web link graph.

(Refer Slide Time: 23:06)



**Contd...**

- **Term-Vector per Host:** A term vector summarizes the most important words that occur in a document or a set of documents as a list of (word; frequency) pairs.
- The map function emits a (hostname; term vector) pair for each input document (where the hostname is extracted from the URL of the document).
- The reduce function is passed all per-document term vectors for a given host. It adds these term vectors together, throwing away infrequent terms, and then emits a final (hostname; term vector) pair

Cloud Computing and Distributed Systems      MapReduce

Similarly, term vector per host, a term vector summarizes the important words that occur in a document or a set of documents as the list of words and their frequency which is called a term vector per host.

(Refer Slide Time: 23:27)

**Contd...**

- **Inverted Index:** The map function parses each document, and emits a sequence of (word; document ID) pairs. The reduce function accepts all pairs for a given word, sorts the corresponding document IDs and emits a (word; list(document ID)) pair. The set of all output pairs forms a simple inverted index. It is easy to augment this computation to keep track of word positions.
- **Distributed Sort:** The map function extracts the key from each record, and emits a (key; record) pair. The reduce function emits all pairs unchanged.

Cloud Computing and Distributed Systems      MapReduce

So, this particular application can also be programmed using MapReduce inverted index that is the map function parses each document and emit a sequence of words. So, the reduce function accepts all pair of the given words and sorts the corresponding document ID and emits the world. So, the set of all output pair forms a inverted index and it is easy to augment this computation to keep track of what position. Similarly, distributed sort that is the sorting algorithm in a large scale is called distributed sort.

(Refer Slide Time: 23:59)

**Applications of MapReduce**

**(1) 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*

Cloud Computing and Distributed Systems      MapReduce

So, let us see some of the use of MapReduce in these applications. So, distributed graph here the input is a large set of files and the output is the lines into the file that matches that particular given pattern. So, map will emit a line if that pattern is matches the supplied pattern and reduce will copy the intermediate data to the output.

(Refer Slide Time: 24:28)

## Applications of MapReduce

**(2) Reverse Web-Link Graph:**

- **Input:** Web graph: tuples (a, b) where (page a → page b)
- **Output:** For each page, list of pages that link to it
- Map – *process web log and for each input <source, target>, it outputs <target, source>*
- Reduce - *emits <target, list(source)>*

*Example*

map key: a, b, c, a, b, c, a, b, c, a, b, c

map value: a → b, b → c, c → a, b → a, a → c, c → b, b → a, a → c, c → b, b → a

emit: a → c, b → a, c → b

reduce: a → c, b → a, c → b

output: a → (c, b), c → b, b → a

target source

Page a → page b

Cloud Computing and Distributed Systems
MapReduce

So, this can be written using a simple program like word count. Another application is called as reverse web link graph here the input is the web graph in the form of a tuples a, b where the page a will basically accessing or having a link to the page b. Now, given this particular input the output for each page is the list of pages that link to it. To do this, the map function what it will do? It will process this particular web log for each input and it outputs the reverse of that that is the target and the source.

Then reducer will collect the target and form the list of source. Let us take an example; let us see a graph, a b and c. So, here in this particular graph, there are three nodes a b c or let us say there are three pages a b c, a is now referring to b or accessing to b and b is accessing to a and a is accessing to c. So, the list of pages here it will be in the form of a b then, then b c, then c a and then b a. Now the task is to find out for example, in a how many links are being accessed? So, here it is showing there are two links. Let us see how using MapReduce we can do this.

So, this will be the input this value will be the input, this is the key and this is the value this will be input to the map function. So, the map function will be given this input. Now

the map will emit value, comma key so, that means it will emit this particular output. Let us say here it is given as source and target, this is called source and this is called target. So, this will be given as the values to this particular thing. Now this will emit the target, comma source. So, MapReduce output will emit b, comma a it will output c, comma b it will output a, comma c it will output a, comma b. Then this particular output will be now fed to the reducer, reducer will now generate based on the key as the target it will find out the list of source.

So, this particular key for example, a will be combined and the list will be generated for a c comma b this will be list for c there will be list called b and for b the list will be a this will be generated out of the reducer. So, here we can see the node a is having two links pointed one by b and the other by c. So, using MapReduce program we can easily write down the map and reduce function which is very small program. But here we have to know the tricks to solve these particular problems. And when a big programmer or a big file is given it will be able to generate these values.

(Refer Slide Time: 29:01)

## 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>*

Cloud Computing and Distributed Systems
MapReduce

Now, let us see another MapReduce program to count URL access frequency. Now here the input is a log of accessed URL's from the proxy server and output will be the each URL and the percentage of total accesses for that URL so, this will be a bit difficult. Let us see the details how we can write down a program for this particular application of URL access frequency. So, here first what we will do is that map function will process

the web log and output URL with a value 1. So that means, all the URL's will be given to the map function, URL's will be give raise the input and what it will emit here is URL and 1, that means, for every URL it will emit 1.

Now, the next phase comprises of multiple reducers. So, the first phase what it will do? It will take it will take this particular output. Now reducer number 1, it will take this output URL 1 and then it will it will count the total URL's, that is it will generate a reducer 1, comma URL and the count. So, with 1 and all the output whatever is generated, so that means, it will generate the total number of URL's present. So, it will generate the total values of URL present.

Now, then another reducer will be basically, done in the second pass which will count the percentage of the URL in the form that this particular reducer will take the URL's, which are omitted by the map function. And for every URL count whatever is output by the reducer in this particular phase, it will divide by the overall count whatever is calculated in the previous pass of the reducer. Therefore, with one map function and 2 times application of a reducer function, we will be able to calculate the count of URL access frequency. So, this particular application has shown that it is possible to cascade multiple MapReduce job to solve an application.

(Refer Slide Time: 32:40)

**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

Cloud Computing and Distributed Systems      MapReduce

Now let us see another application, here in this application it will do the sorting. Now the sorting using MapReduce; that means, if a very big file is given: how are you going to

sort using MapReduce is the task. Now here the map tasks outputs they result that is intermediate results in a sorted form, that is using the quick sort already is implemented inside the MapReduce, inside the map function.

Similarly, the reduce task also input in the form of sorted manner that is also done using some sorting algorithm for example, using merge sort. So, using this information let us see how the sorting can be done using MapReduce in a very simple manner. So, let us see that if the input is given to the map. So, it will generate that particular value which is output from the map function. And the second one that is the value of will be any value that is not that important.

So, you know that these particular values are already sorted so, the map will generate these particular keys in a sorted order. Then as far as the reducer is concerned this as a input of this map function and it will generate the same value again. The only thing we have to do is that in the partitioning, we have to provide the range instead of hash function because, hashing will randomize them. Since the output of the map function is already in the form of sorted and we do not want to disturb it so, that we have to provide the range values the range of keys. So, range of keys means the reduce task will be assigned according to the ranges.

So, that this particular range will not be disturbed and the output of this map function will be preserved which is in the sorted list and the output will be sorted. So, this particular example or application of MapReduce could make this sorting algorithm easy why because, it has exploited the way map function is giving the output it is in the sorted order. So, we have to just preserve that sorted order without being disturbed by the partitioning function so, that the entire file will be in a sorted order output.

(Refer Slide Time: 36:06)

## Programming MapReduce

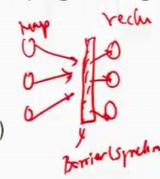
**Externally:** For user

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)



Cloud Computing and Distributed Systems
MapReduce

Now, let us see the MapReduce scheduling how it is done? Now, as far as programming of a MapReduce is concerned externally for the user we have to only write down the map program, these are very short programs. And they have to specify the map and reduce function and submit the job and wait for the results. The users need not have to know the details of parallel and distributed concepts and the programming.

As far as everything is done internally, so, for internally for the paradigm and the scheduler has to build everything and give the output. So, it has to parallelize the map function execution; that means, it will transfer the data from map to the reduce, that is done by the shuffling that we have seen earlier that a partitioning and shuffling using hash function was being carried out.

Then parallelize the reduce function and implement the storage for map input map output reduce input and reduce output. One thing is there that we have to ensure that, no reducer starts before all the maps are finished, that is being done inside MapReduce. So, this ensures that a barrier between the map phase and the reduce phase is there. So, barrier in the sense this is called a barrier is nothing, but a synchronization. So, when all the map function will finish, some are finishing early some are finishing quite late. So, earlier who ever finish they have to wait till all the map function completes their execution then only the next phase called reducer phase will begin afterwards so, this is called a barrier.

In some of the implementations this barriers can also be in some cases being a synchronized; that means, they are not synchronizing may be in the next phase without completing all the reduce phase the next phase can begun.

(Refer Slide Time: 38:34)

### 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

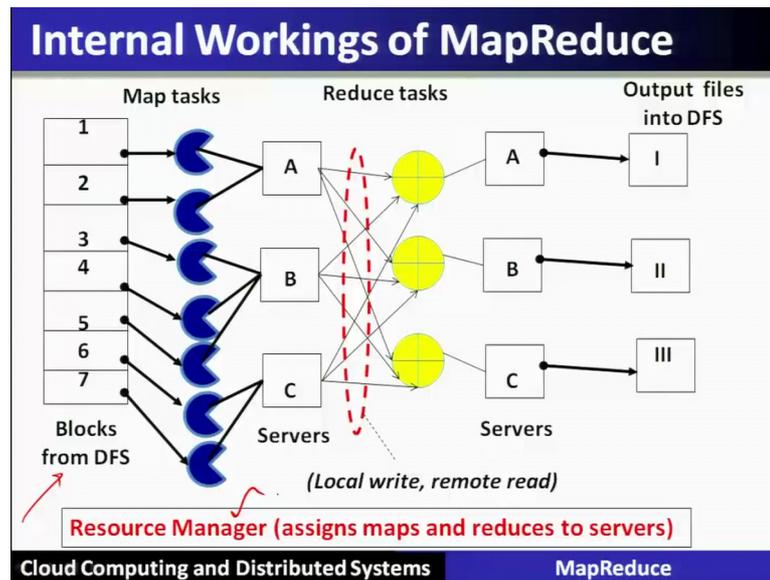
**local file system** = Linux FS, etc.  
**distributed file system** = GFS (Google File System), HDFS (Hadoop Distributed File System)

Cloud Computing and Distributed Systems      MapReduce

So, inside the MapReduce we have seen that parallelize map is quite easy now why because, each task in the map they are all independent. So, they can parallelly execute inside a cloud. To transfer the data from map to reduce a shuffle function is to be called and all map output with the same key is assigned to the same reduce job and this requires to use the partitioning function which is hash based portioning. So, that it will be uniformly distribute across all the servers which are running the reduce task.

Third step is to parallelize reduce is quiet easy why because, reduce is independent of each other, implement the storage for the map input, map output, reduce input, reduce output. So, here we can see that the map input is done from the distributed file system. Map output also uses the local file system and reduce input also from multiple remote disks using the local file system and reduce output node is done through the distributed file system. Distributed file systems which are used here in the MapReduce are either the Google File System or Hadoop Distribute File System.

(Refer Slide Time: 40:11)



Let us see the internal working of the MapReduce. So, first thing is whenever the input file is given. So, this particular input file is split into different splits seven different splits are shown over here. Each and every split is given to different map tasks now these map tasks are to be assigned to servers or the machines. So, for example, here it is shown that two different map tasks are assigned to one server. A three different map tasks are assigned to the server B and two map tasks are assigned to the server C. So, it depends upon how many servers are there these map tasks will be assigned.

Similarly, the output of these servers will be shuffled; that means, they will be grouped by the keys and will be given to the reducer functions. So, here this example shows that there are three different reduce tasks being defined over here. And these reduce tasks will be assigned to different servers. So, there are three different servers assigned to the reduce task which will be executing these reducers. At the output of these particular servers are stored in a file in a distributed file system, which are given different names or there may be only one file will also depend upon.

Now, who allocates who assigns these map tasks to the servers? And the reduce task also to the different servers may be the same servers. So, it is the resource manager, it is the resource manager of a distributed file system or a yarn server that we will see scheduler which will basically internally do this.

(Refer Slide Time: 48:24)

**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
- 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

Cloud Computing and Distributed Systems      MapReduce

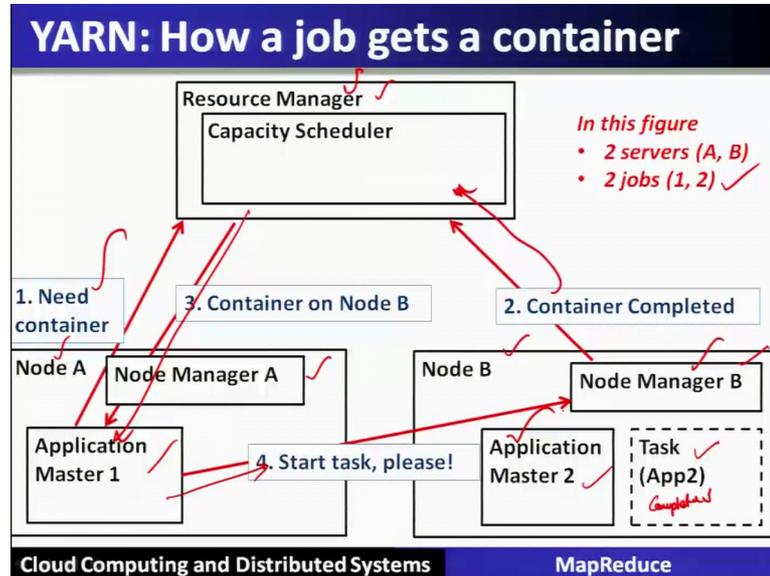
Let us see the yarn scheduler which is combined with the Hadoop 2.0 versions onwards. So, yarn full form is yet another resource negotiator. So, this particular yarn will treat each server as a collection of container. So, it means that, if the server is having some number of course, let us say 3 and it is having 3 GB of internal memory. So that means, one core and 1 GB will form one container so, it has now three containers. So, three different task can be assigned in this manner. So, this yarn scheduler has 3 different components, the use of one component we have seen in the previous slide that is called resource manager. Which assigns these containers to the different tasks that is to the map task and to the reduce task so that is called scheduling.

So, this there is a global resource manager which performs the scheduling of this particular task map or reduce task to this particular machines or the servers or to the container here in this terminology. Then another component of yarn is called the node manager. So, for every server there is a daemon process called node manager which basically, will manage all the server related functions so, that means, a server if it is giving three different containers so the node manager will manage it.

Then per application there is a application manager, application master. So, for every application there is a application master and this particular application master will negotiate with the resource manager and the node manager for allocating their tasks to

the machines. Also the application master is responsible to deal with the task failures at the time of execution of the job.

(Refer Slide Time: 44:45)



Let us see how the yarn will get a job of a container? So, here there is a resource manager and this is the example which is shown that there are two nodes node A and node B. So, the node A is having it is manager node A manager and node manager B will manage this server node B. So, there is resource manager which will manage all the nodes which are available in this scenario that is through the yarn.

Now, for there are two applications 1 jobs job one and job 2. So, there are two different application master they will manage these two tasks. Now let us see the example of these interactions with the resource manager. Now let us assume that the application master to it is task that is the application 2 task is completed. Therefore, this will return back through the node manager that it is free.

Similarly the application one master negotiates with the resource manager asking for a container, now since this container is available why because node B has completed it is task. So, that particular container is now free so it will be assigned back to the master 1 and it will start the task execution.

So, this particular interaction of different resource manager, the node manager and the application master, they will ensure that the resources which are required to execute the applications are done. Let us see the fault tolerance, which is handled in map reduce.

(Refer Slide Time: 46:52)

## Fault Tolerance

- **Server Failure**
  - **NM heartbeats to RM**
    - If server fails, 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 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

Cloud Computing and Distributed Systems      MapReduce

So, MapReduce handles different kind of faults at different level. For example, whenever there is a server failure so, this will affect the name node and resource manager, name manager, resource manager and application manager, will be affected in this particular scenario whenever there is a server failures. Similarly, whenever there is a resource manager failure then basically, it uses the secondary resource manager to failover and start working at.

(Refer Slide Time: 47:38)

## 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 backup (**replicated**) execution of straggler tasks
  - A task considered done when its first replica complete called **Speculative Execution**.

Cloud Computing and Distributed Systems      MapReduce

These are detected using the heart beats which are contained in the request and also avoid. There is a possibility of slow servers called stragglers. So, the slowest task will slow down the entire job that we have already seen because, it is using the barriers concept. So, slow servers are due to the bad disk, network bandwidth, CPU, memory, there can be many different issues. So, it has to keep track of the progress of the task and performs the backup that is replicated execution of staggered task which is called as a speculative execution.

(Refer Slide Time: 48:18)

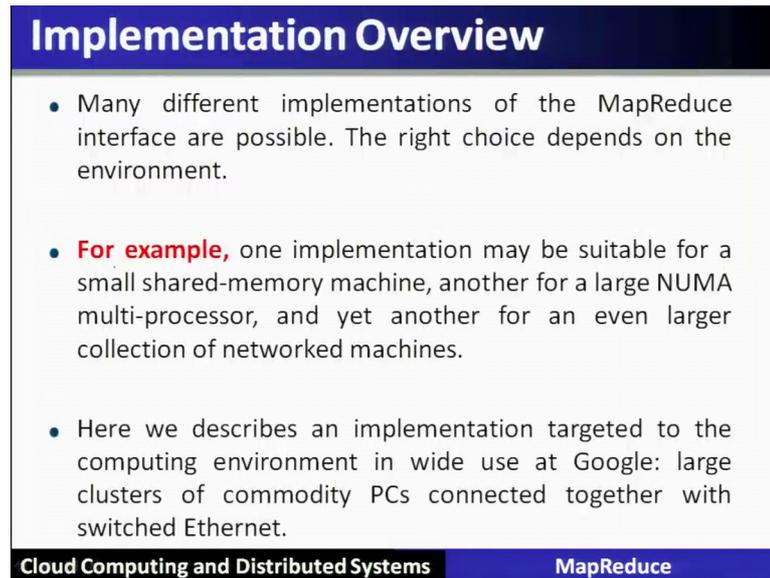
## Locality

- **Locality**
  - Since cloud has hierarchical topology (**e.g., racks**)
  - GFS/HDFS stores 3 replicas of each of chunks (e.g., 64 MB in size)
    - Maybe 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

Cloud Computing and Distributed Systems      MapReduce

Locality; that means, that the cloud has the hierarchical topology of a data centre that is the servers racks and different racks are connected through the top of racks switch. So, the locality has to be exploited, that means, MapReduce will try to attempt to assign the server for a particular task which is very close to it so, that the network delay can be avoided here in this case.

(Refer Slide Time: 48:58)



**Implementation Overview**

- Many different implementations of the MapReduce interface are possible. The right choice depends on the environment.
- **For example**, one implementation may be suitable for a small shared-memory machine, another for a large NUMA multi-processor, and yet another for an even larger collection of networked machines.
- Here we describes an implementation targeted to the computing environment in wide use at Google: large clusters of commodity PCs connected together with switched Ethernet.

Cloud Computing and Distributed Systems      MapReduce

So, let us see some of the implementation of MapReduce, there are many implementation of MapReduce an interface an interfaces are available. So, let us see one such example.

(Refer Slide Time: 49:10)

### Contd...

- (1) Machines are typically dual-processor x86 processor running Linux, with 2-4 GB of memory per machine.
- (2) Commodity networking hardware is used . Typically either 100 megabits/second or 1 gigabit/second at the machine level, but averaging considerably less in overall bisection bandwidth.
- (3) A cluster consists of hundreds or thousands of machines, and therefore machine failures are common.
- (4) Storage is provided by inexpensive IDE disks attached directly to individual machines.
- (5) Users submit jobs to a scheduling system. Each job consists of a set of tasks, and is mapped by the scheduler to a set of available machines within a cluster.

Cloud Computing and Distributed Systems      MapReduce

Here where we assume that the machines are dual processor commodity networks hardware is available and the clusters consist of hundred and thousands of machine is there.

(Refer Slide Time: 49:25)

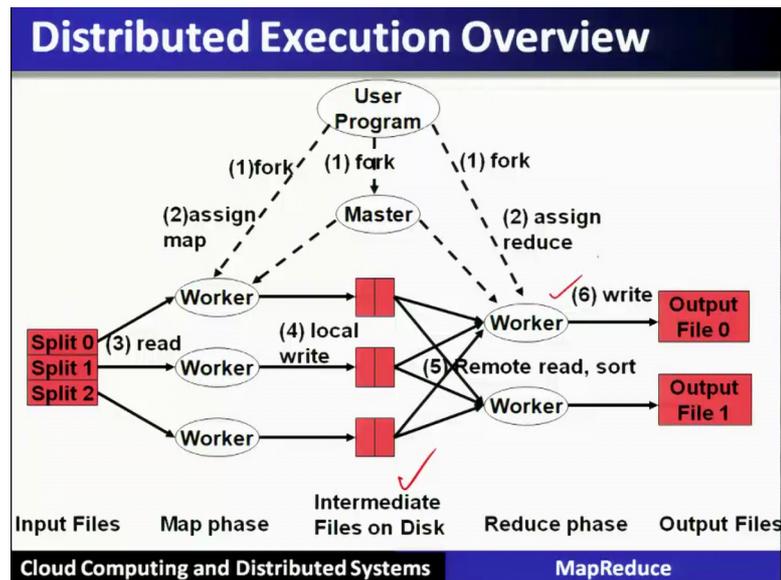
### Distributed Execution Overview

- The **Map invocations** are distributed across multiple machines by automatically partitioning the input data into a set of M splits.
- The input splits can be processed in parallel by different machines.
- **Reduce invocations** are distributed by partitioning the intermediate key space into R pieces using a partitioning function (e.g.,  $\text{hash}(\text{key}) \bmod R$ ).
- The number of partitions (R) and the partitioning function are specified by the user.

Cloud Computing and Distributed Systems      MapReduce

And storage is also provided and user can submit the job to this kind of system. Let us see the distributed execution of this MapReduce invocation.

(Refer Slide Time: 49:32)



So, here we can see that the user program, user can submit their job in this particular manner and this particular user also gives the splits that is the input file is divided into different splits. And these particular user program, that is the map phase will generate different workers and these workers are assigned to the different machines. And when these splits argument to the worker they will output and they will write down on the intermediate files on the disks using some distributed file system.

Now, the next task is that these particular intermediate output is given to the input to the next phase that is to the reducer phase. And this reducer phase that is called worker will be assigned by the user program that is to the master. Now, then these particular workers are assigned by the user program that is through the master. Now, then these particular workers are assigned the servers in the reduce phase and they will output through the to the output file in the distributed file system. So, this is the typical scenario of a distributed execution now here we can see the sequence of actions in this distributed execution.

(Refer Slide Time: 50:51)

## Sequence of Actions

When the user **program calls the MapReduce function**, the following **sequence of actions** occurs:

1. The MapReduce library in the user program first splits the input files into M pieces of typically 16 megabytes to 64 megabytes (MB) per piece. It then starts up many copies of the program on a cluster of machines.
2. One of the copies of the program is special- the master. The rest are workers that are assigned work by the master. There are M map tasks and R reduce tasks to assign. The master picks idle workers and assigns each one a map task or a reduce task.
3. A worker who is assigned a map task reads the contents of the corresponding input split. It parses key/value pairs out of the input data and passes each pair to the user-defined Map function. The intermediate key/value pairs produced by the Map function are buffered in memory.

Cloud Computing and Distributed Systems      MapReduce

So, the program calls the MapReduce function and the sequence of actions, which we have told is being listed over here. For example, MapReduce library will split the user file into different pieces of 16 MB to 64 MB and then these one of these copies is called one of the program copies is called a master. And the rest are it will generate the different worker process and these workers are assigned the map task.

(Refer Slide Time: 51:31)

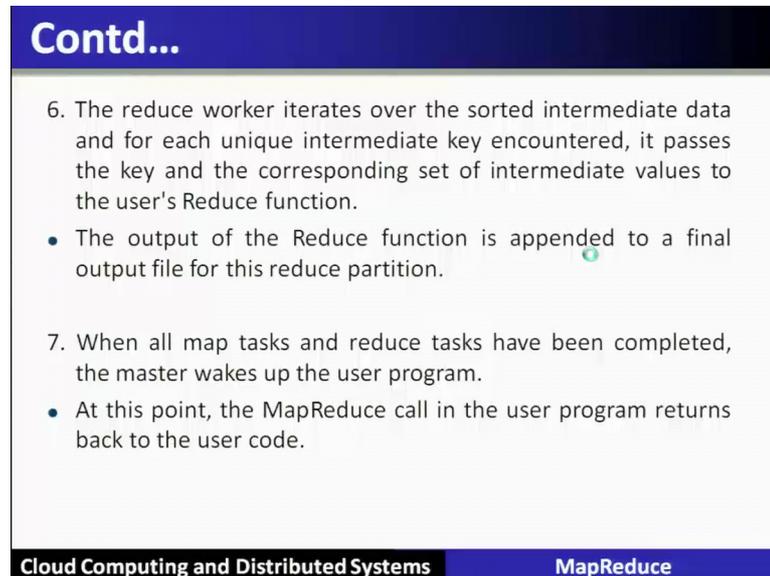
## Contd...

4. Periodically, the buffered pairs are written to local disk, partitioned into R regions by the partitioning function.
  - The locations of these buffered pairs on the local disk are passed back to the master, who is responsible for forwarding these locations to the reduce workers.
5. When a reduce worker is notified by the master about these locations, it uses remote procedure calls to read the buffered data from the local disks of the map workers. When a reduce worker has read all intermediate data, it sorts it by the intermediate keys so that all occurrences of the same key are grouped together.
  - The sorting is needed because typically many different keys map to the same reduce task. If the amount of intermediate data is too large to fit in memory, an external sort is used.

Cloud Computing and Distributed Systems      MapReduce

And into the end each split will be assigned to these workers and periodically they are being buffered and partitioned. These reduce and these particular intermediate results will be stored in the file system.

(Refer Slide Time: 51:51)



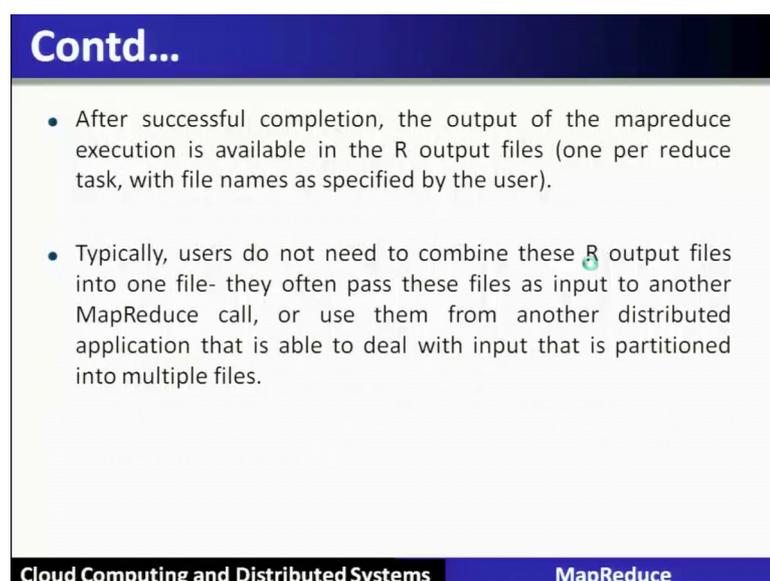
**Contd...**

6. The reduce worker iterates over the sorted intermediate data and for each unique intermediate key encountered, it passes the key and the corresponding set of intermediate values to the user's Reduce function.
  - The output of the Reduce function is appended to a final output file for this reduce partition.
7. When all map tasks and reduce tasks have been completed, the master wakes up the user program.
  - At this point, the MapReduce call in the user program returns back to the user code.

Cloud Computing and Distributed Systems      MapReduce

Now, the reducer worker is notified by the master and this particular output, intermediate output of the previous stage will be given as the input and it will be reducer phase will give the output.

(Refer Slide Time: 52:03)



**Contd...**

- After successful completion, the output of the mapreduce execution is available in the R output files (one per reduce task, with file names as specified by the user).
- Typically, users do not need to combine these R output files into one file- they often pass these files as input to another MapReduce call, or use them from another distributed application that is able to deal with input that is partitioned into multiple files.

Cloud Computing and Distributed Systems      MapReduce

(Refer Slide Time: 52:04)

## Master Data Structures

- The master keeps several data structures. For each map task and reduce task, it stores the **state (idle, in-progress, or completed)**, and the identity of the worker machine (**for non-idle tasks**).
- The master is the conduit through which the location of intermediate file regions is propagated from map tasks to reduce tasks. Therefore, for each completed map task, the master stores the locations and sizes of the R intermediate file regions produced by the map task.
- Updates to this location and size information are received as map tasks are completed. The information is pushed incrementally to workers that have in-progress reduce tasks.

Cloud Computing and Distributed Systems

MapReduce

So here the master data structure is like this, that it basically has different states.

(Refer Slide Time: 52:18)

## Fault Tolerance

- Since the MapReduce library is designed to help process very large amounts of data using hundreds or thousands of machines, the library must tolerate machine failures gracefully.
- **Map worker failure**
  - Map tasks completed or in-progress at worker are reset to idle
  - Reduce workers are notified when task is rescheduled on another worker
- **Reduce worker failure**
  - Only in-progress tasks are reset to idle
- **Master failure**
  - MapReduce task is aborted and client is notified

Cloud Computing and Distributed Systems

MapReduce

And it will trace through different states internally in the execution.

(Refer Slide Time: 52:27)

## Task Granularity

- The **Map phase is subdivided into M pieces and the reduce phase into R pieces.**
- Ideally, M and R should be much larger than the number of worker machines.
- Having each worker perform many different tasks improves dynamic load balancing, and also speeds up recovery when a worker fails: the many map tasks it has completed can be spread out across all the other worker machines.
- There are practical bounds on how large M and R can be, since the master must make  **$O(M + R)$  scheduling decisions** and keeps  **$O(M * R)$  state in memory.**
- Furthermore, R is often constrained by users because the output of each reduce task ends up in a separate output file.

Cloud Computing and Distributed Systems      MapReduce

(Refer Slide Time: 52:29)

## Partition Function

- Inputs to map tasks are created by contiguous splits of input file
- For reduce, we need to ensure that records with the same intermediate key end up at the same worker
- System uses a default partition function e.g.,  **$\text{hash}(\text{key}) \bmod R$**
- Sometimes useful to override
  - E.g.,  **$\text{hash}(\text{hostname}(\text{URL})) \bmod R$**  ensures URLs from a host end up in the same output file

Cloud Computing and Distributed Systems      MapReduce

Fault tolerance we have already covered, locality is being exploited task granularity we have assumed partition function we have covered.

(Refer Slide Time: 52:31)

## Ordering Guarantees

- It is guaranteed that within a given partition, the intermediate key/value pairs are processed in increasing key order.
- This ordering guarantee makes it easy to generate a sorted output file per partition, which is useful when the output file format needs to support efficient random access lookups by key, or users of the output and it convenient to have the data sorted.

Cloud Computing and Distributed Systems      MapReduce

(Refer Slide Time: 52:38)

## Combiners Function (1)

- In some cases, there is significant repetition in the intermediate keys produced by each map task, and the user specified Reduce function is commutative and associative.
- A good example of this is the word counting example. Since word frequencies tend to follow a Zipf distribution, each map task will produce hundreds or thousands of records of the form <the, 1>.
- All of these counts will be sent over the network to a single reduce task and then added together by the Reduce function to produce one number. We allow the user to specify an optional Combiner function that does partial merging of this data before it is sent over the network.

Cloud Computing and Distributed Systems      MapReduce

And ordering guarantees is also there that we have used we have shown in the sorting application and this is the combined function.

(Refer Slide Time: 52:44)

## Combiners Function (2)

- **The Combiner function is executed on each machine that performs a map task.**
- Typically the same code is used to implement both the combiner and the reduce functions.
- The only difference between a reduce function and a combiner function is how the MapReduce library handles the output of the function.
- The output of a reduce function is written to the final output file. The output of a combiner function is written to an intermediate file that will be sent to a reduce task.
- Partial combining significantly speeds up certain classes of MapReduce operations.

Cloud Computing and Distributed Systems      MapReduce

(Refer Slide Time: 52:47)

## Example: 1 Word Count using MapReduce

```
map(key, value):
// key: document name; value: text of document
for each word w in value:
    emit(w, 1)

reduce(key, values):
// key: a word; values: an iterator over counts
result = 0
for each count v in values:
    result += v
emit(key, result)
```

Cloud Computing and Distributed Systems      MapReduce

(Refer Slide Time: 52:51)

### Count Illustrated

```

map(key=url, val=contents):
  For each word w in contents, emit (w, "1")
reduce(key=word, values=uniq_counts):
  Sum all "1"s in values list
  Emit result "(word, sum)"

```

see bob run  
see spot throw

→

see 1  
bob 1  
run 1  
see 1  
spot 1  
throw 1

→

bob 1  
run 1  
see 2  
spot 1  
throw 1

Cloud Computing and Distributed Systems
MapReduce

So, MapReduce more examples are available.

(Refer Slide Time: 52:52)

### Example 2: Counting words of different lengths

- The map function takes a value and outputs key:value pairs.
- For instance, if we define a map function that takes a string and outputs the length of the word as the key and the word itself as the value then
  - map(steve) would return 5:steve and
  - map(savannah) would return 8:savannah.

emit (length, word)

This allows us to run the map function against values in parallel and provides a huge advantage.

Cloud Computing and Distributed Systems
MapReduce

This word count we have already done. Now the counting words of different length, if we want to find out. Then what we will do is? Given a particular input that is then we will written its length and that particular word also. So, it we will written the length and the word both will be emit in the map function.

(Refer Slide Time: 53:27)

### Example 2: Counting words of different lengths

Before we get to the reduce function, the mapreduce framework groups all of the values together by key, so if the map functions output the following **key:value pairs**:

3 : the ✓	✓ They get grouped as:
3 : and ✓	
3 : you ✓	✓ 3 : [the, and, you]
4 : then ✓	✓ 4 : [then, what, when]
4 : what ✓	5 : [steve, where]
4 : when ✓	8 : [savannah, research]
5 : steve ✓	
5 : where ✓	
8 : savannah ✓	
8 : research ✓	

Cloud Computing and Distributed Systems MapReduce

So, that is the case then for example, all these particular words, first their length and their words will be written. And then reduce function will combine them using that particular key that is the length of the words. And in this case for every length of word every length different keywords will be combined in the form of a list.

(Refer Slide Time: 53:58)

### Example 2: Counting words of different lengths

- Each of these lines would then be passed as an argument to the reduce function, which accepts a key and a list of values.
- In this instance, we might be trying to figure out how many words of certain lengths exist, so our reduce function will just count the number of items in the list and output the key with the size of the list, like:

3 : 3 ✓
4 : 3 ✓
5 : 2 ✓
8 : 2 ✓

Cloud Computing and Distributed Systems MapReduce

And this shows that the word that is having the length 3 there are 3 different words and length 4 there are 3 different words, length 5 there are 2 different and length 8 there are 2 different words. So, here we have counted the words of different length.

(Refer Slide Time: 54:10)

### Example 2: Counting words of different lengths

- The reductions can also be done in parallel, again providing a huge advantage. We can then look at these final results and see that there were only two words of length 5 in the corpus, etc...
- **The most common example of mapreduce is for counting the number of times words occur in a corpus.**

Cloud Computing and Distributed Systems MapReduce

(Refer Slide Time: 54:14)

### Example 3: Finding Friends

- Facebook has a list of friends (note that friends are a bi-directional thing on Facebook. If I'm your friend, you're mine).
- They also have lots of disk space and they serve hundreds of millions of requests everyday. They've decided to pre-compute calculations when they can to reduce the processing time of requests. **One common processing request is the "You and Joe have 230 friends in common" feature.**
- When you visit someone's profile, you see a list of friends that you have in common. This list doesn't change frequently so it'd be wasteful to recalculate it every time you visited the profile (sure you could use a decent caching strategy, but then we wouldn't be able to continue writing about mapreduce for this problem).
- We're going to use mapreduce so that we can calculate everyone's common friends once a day and store those results. Later on it's just a quick lookup. We've got lots of disk, it's cheap.

Cloud Computing and Distributed Systems MapReduce

(Refer Slide Time: 54:16)

### Example 3: Finding Friends

- Assume the friends are stored as **Person->[List of Friends]**, our friends list is then:
  - A -> B C D
  - B -> A C D E
  - C -> A B D E
  - D -> A B C E
  - E -> B C D

Another example is to find a common friends, list of common friends. So, here every input every person will provide a list of their friends. So, A B C D E they have provided their own list of friends, if this is the input given to the map function.

(Refer Slide Time: 54:29)

### Example 3: Finding Friends

For map(A -> B C D):

(A B) -> B C D

(A C) -> B C D

(A D) -> B C D

For map(B -> A C D E): (Note that A comes before B in the key)

(A B) -> A C D E

(B C) -> A C D E

(B D) -> A C D E

(B E) -> A C D E

Map function will output, for every such list it will give a pair that is A B will list out this particular list of friends. And A C will list out friends and similarly A D. Similarly the person number B will give it is list and this also will be emitted in this particular format.

(Refer Slide Time: 54:56)

### Example 3: Finding Friends

**For map(C -> A B D E) :**

(A C) -> A B D E  
(B C) -> A B D E  
(C D) -> A B D E  
(C E) -> A B D E

**For map(D -> A B C E) :**

(A D) -> A B C E  
(B D) -> A B C E  
(C D) -> A B C E  
(D E) -> A B C E

**And finally for map(E -> B C D):**

(B E) -> B C D  
(C E) -> B C D  
(D E) -> B C D

Cloud Computing and Distributed Systems      MapReduce

So, when all these list are output.

(Refer Slide Time: 55:00)

### Example 3: Finding Friends

- Before we send these key-value pairs to the reducers, we group them by their keys and get:

(A B) -> (A C D E) (B C D)  
(A C) -> (A B D E) (B C D)  
(A D) -> (A B C E) (B C D)  
(B C) -> (A B D E) (A C D E)  
(B D) -> (A B C E) (A C D E)  
(B E) -> (A C D E) (B C D)  
(C D) -> (A B C E) (A B D E)  
(C E) -> (A B D E) (B C D)  
(D E) -> (A B C E) (B C D)

*Intersection*

Cloud Computing and Distributed Systems      MapReduce

Then as far as the reducer is concerned, what reducer will do? Reducer will now for every keyword which is there. And the list which is basically available in the intermediate they will take the intersection. Intersection of these particular list of common friends.

(Refer Slide Time: 55:25)

### Example 3: Finding Friends

- Each line will be passed as an argument to a reducer.
- The **reduce function will simply intersect the lists of values** and output the same key with the result of the intersection.
- For example, **reduce((A B) -> (A C D E) (B C D))**  
will **output (A B) : (C D)**
- **and means that friends A and B have C and D as common friends.**

Cloud Computing and Distributed Systems      MapReduce

So, intersection will find out the list of common friends and that will be output.

(Refer Slide Time: 55:27)

### Example 3: Finding Friends

- The result after reduction is:

(A B) -> (C D)

(A C) -> (B D)

(A D) -> (B C)

(B C) -> (A D E)

(B D) -> (A C E)

(B E) -> (C D)

(C D) -> (A B E)

(C E) -> (B D)

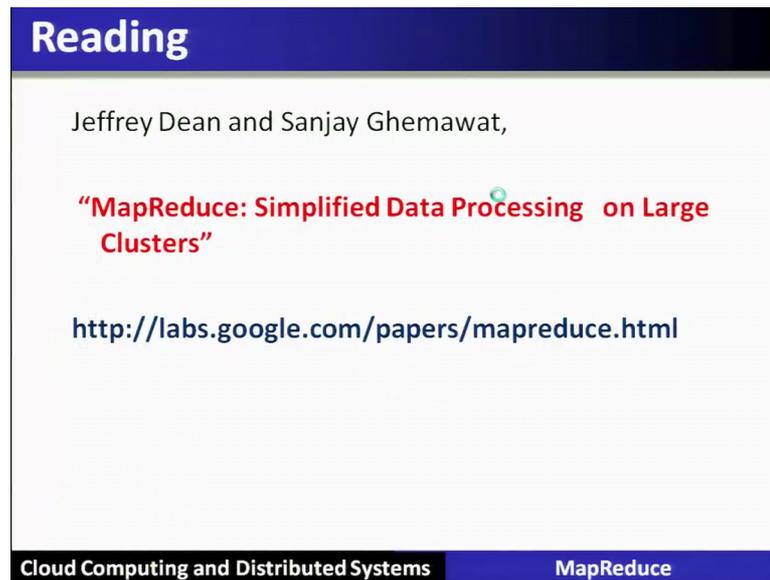
(D E) -> (B C)

Now when D visits B's profile, we can quickly look up (B D) and see that they have three friends in common, (A C E).

Cloud Computing and Distributed Systems      MapReduce

So, this is the result which is shown over here. So, using MapReduce all these particular program become simple, but writing or thinking about this program is not so, trivial.

(Refer Slide Time: 55:39)



**Reading**

Jeffrey Dean and Sanjay Ghemawat,

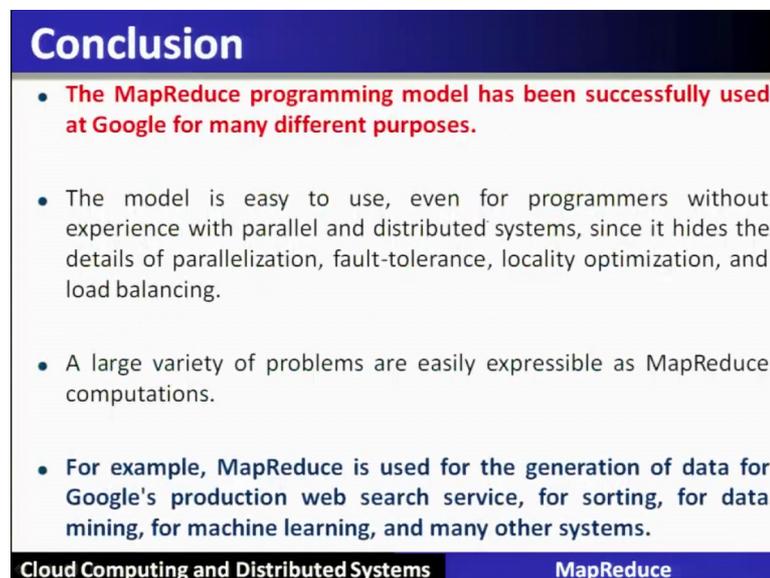
**“MapReduce: Simplified Data Processing on Large Clusters”**

<http://labs.google.com/papers/mapreduce.html>

Cloud Computing and Distributed Systems      MapReduce

More details of this MapReduce you can refer to this particular paper by Jeffrey Dean and Sanjay Ghemawat that is the MapReduce simplified data processing on large clusters.

(Refer Slide Time: 55:57)



**Conclusion**

- **The MapReduce programming model has been successfully used at Google for many different purposes.**
- The model is easy to use, even for programmers without experience with parallel and distributed systems, since it hides the details of parallelization, fault-tolerance, locality optimization, and load balancing.
- A large variety of problems are easily expressible as MapReduce computations.
- For example, MapReduce is used for the generation of data for Google's production web search service, for sorting, for data mining, for machine learning, and many other systems.

Cloud Computing and Distributed Systems      MapReduce

So, conclusion MapReduce programming model has been successfully used in the Google and many different purposes.

(Refer Slide Time: 56:08)

**Conclusion**

- 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.**

Cloud Computing and Distributed Systems      MapReduce

So, we have seen that for a large size data set, this particular simple programming paradigm we have also seen the internal details. We have also shown how it is automatically dealing with the parallelization and dealing with the failures. There are many ongoing research work happening in the in the scheduling fault tolerance in the scenarios of MapReduce and Hadoop.

Thank you.