

## What is MapReduce

What is MapReduce?

- MapReduce in simple words is a programming paradigm to help solve Big Data problems.
- Hadoop framework works on this paradigm.
- This is great for tasks that are sorting or disc read intensive.

Ideally, you would need two functions or pieces of logic:

- Mapper - which converts every record from the input into "key-value" pairs.
- Reducer - which aggregates value for each key as defined by Mapper or the Map phase.

It is also supported by many systems such as MongoDB / CouchDB / Cassandra, Apache Spark.

MapReduce in Hadoop can be written in Java, Shell, Python or any binaries.

Let's take a look at how MapReduce gets executed. In this diagram, we have three machines containing data on which Map functions are getting executed.

Mapper is the logic that you have defined. This logic takes a record as input and converts it into "key-value" pairs. Please note that Map logic is provided by you. This logic can be very complex or very simple based on your need.

These key-value pairs are sorted and then grouped together by Hadoop, based on the key. All of the values for each key are aggregated by your Reducer logic. So if you want to group data based on some criteria, that criteria would be expressed in the mapper logic and how to combine all these values for each key is governed by your logic of Reducer. The result of reducer is saved in the HDFS.

Let's imagine for a moment that we would like to prepare a veg. Burger on a very large scale. As you can see in the diagram, the function `CutIntoPieces()` will be executed on each vegetable, chopping vegetable into pieces and the result will be reduced to form a Burger.

# Thinking in MR - Unix Pipeline

Third approach is to use Unix command in Pipeline or in Chain. Let us first try to understand what does it mean by Pipeline.

As we discussed earlier, when we run a program, it may take input from you. In other words, you may provide input to a program by typing. A program or command may also print some output on the screen. In UNIX, you can provide output of one program as the input to another. This is known as piping. A pipe is denoted by "|".

command1 | command2 -> means the output of command 1 will be the input to command 2.

For eg. echo is a UNIX command that prints on the screen. Whatever argument is passed to it, for eg. echo "Hi" prints Hi on the screen.

"wc" is a command that prints the number of characters words and lines out of whatever you type on the standard input. It would print "no lines | words | characters"

eg. echo "Hello, World" | wc

Output: 1 2 13 - lines, words, characters

Let us try to understand this pipeline of command for word counting in parts.

The first command cat myfile prints the content of the file myfile.

The second command in the chain is "sed", which stands for **streaming editor**. It is used to replace the text with something else in the input. It is very similar to the search and replace option of text editors. You can use regular expression with sed by providing an option "-E"

```
sed -E 's/[t ]+/\n/g'
```

this replaces spaces and tabs with new line, essentially it converts the text into one word per line. This 1 word per line text can be further sent to a command called "sort" which can order lines in the input. The sort command takes various options. The option '-S' makes it use only limited memory. In our case, we are using "-S 1G" option to sort data using only 1 GB of memory.

The last command is "uniq". It finds unique line in the input. It expects data to be ordered already. In case the input to uniq is not sorted, the result is not correct. Uniq command has "-C" which prints the count of each unique word. So "uniq -C" would print counts of each unique word in the sorted input. So the entire pipeline consisting of cat, sed, sort followed by uniq prints the word count of unique words in the text file.

### **Problems in approach 2 SQL and approach 3 Unix**

The CPU, Disk Speed and Disk Space become bottlenecks

### **Approach 4 - external sort**

Split files into many machines

Use approach in 2&3 approach to find freq

Then merge and sum of frequencies

Problems with Approach 4

Time consumed in transport of data

Would require lot of engineering

For each requirement we would need to special purpose network oriented program.

### **What is Map/Reduce**

We don't have to transfer data to other computers

Mapper - convert input record into key value pairs

Reducers - Aggregates all the values for a key