MapReduce is framework using which we can write application to process huge amount of data, in parallel, on large cluster of commodity hardware in a reliable manner.

## What is MapReduce?

- MapReduce is processing technique and program model for distributed computing based on java.

- MapReduce paradigm is based on setting the computer to write where the data resides.

There are 2 stages in MapReduce

Stage 1 :- Map
Stage 2 :- Reduce.
Both Map and Reduce only works on (key, value) pair

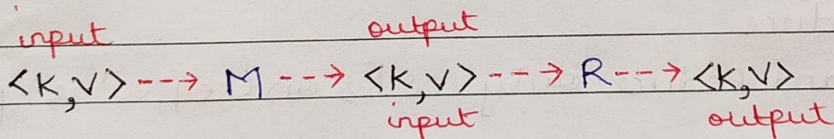Map stage :- The map or mapper job is to process the input data.
- Generally the input data is in the form of file or directory and is stored in the HDFS.

- The input file is passed to the mapper function line by li

- The mapper process the data and creates several small chunks of data.

**Reduce Stage:-** This stage is the combination of the shuffle stage and the Reduce stage.

- The Reducer's job is to process that data comes from the mapper.

- After processing, it produces a new set of output, which will be stored in the HDFS.

## What is (key, value)?

| Key | Value |
|-----|-------|
| Id | 101 |
| Name | Ram |
| Designation | Developer |

input     output

$\langle K,V\rangle \dashrightarrow M \dashrightarrow \langle K,V\rangle \dashrightarrow R \dashrightarrow \langle K,V\rangle$

      input     output

## Record Reader

- The role of Record Reader is to convert each input line into (key, value) pair suitable for reading by Mapper.
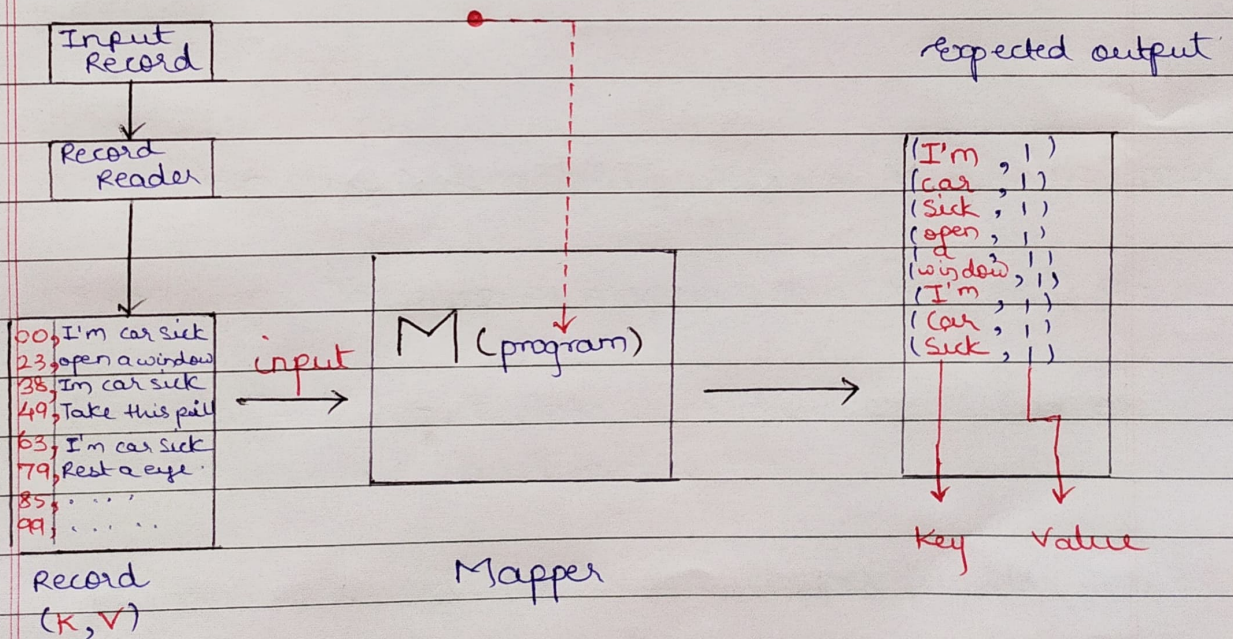
**Input Record**

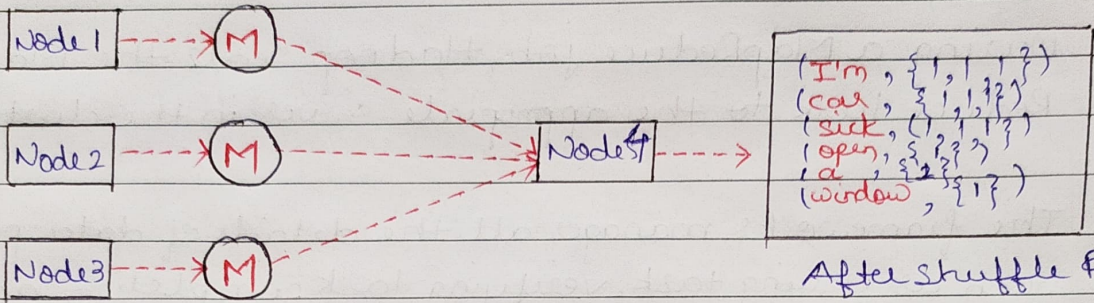| Hello how are You |
| Hello World |
| ... |

Record Reader

**Output Record**

| 00, | Hello how are You |
| 23, | Hello world |
| 58, | |
| 99, | ↑ |
| ↑ | value |
| key | |

- During a MapReduce job, Hadoop sends the Map and Reduce task to the appropriate server in the cluster.

- The frameworks manages all the details of data-passing such as issuing task, verifying task completion, and copy data around the cluster between the nodes.

- Most of the computing takes places on nodes with the data on local disks that reduces the network traffic.

- After completion of the given tasks, the cluster collects and reduces the data to from an appropriate result, and sends it back to the Hadoop server.

How Mapper Works.
                    Reducers



| Input Record |                          | Expected output |

| Record Reader |

| 00, I'm car sick | input | M (program) | (I'm , 1) |
| 23, open a window |       |             | (car , 1) |
| 38, Im car sick   |       |             | (Sick , 1) |
| 49, Take this pill |      |             | (open , 1) |
| 63, I'm car sick  |       |             | (a , 1) |
| 79, Rest a eye    |       |             | (window, 1) |
| 85, . . . .       |       |             | (I'm , 1) |
| 99, . . . . .     |       |             | (Car , 1) |
|                   |       |             | (Sick , 1) |

Record          Mapper                          Key    Value
(K, V)

Map.

Node1 --→ M

Node2 --→ M  ⟶ Node4 ---→

Node3 --→ M

(I'm, {1,1,1,1?})
(car, {1,1,1?})
(sick, {1,1,1,1?})
(open, {1,1,1})
(a, {1,1,1?})
(window, {1?})

After shuffle & sort

Final output

Node1 --→ M

Node2 --→ M  ⟶ Node4 ---→ R ---→

Node3 --→ M

Reduce

(I'm, {3?})
(car, {2?})
(sick, {3?})
(open, {1?})
(a, {2?})
(window, {1?})