

# Big Data Analytics

Academic Year 2022-23

## Index

1. MapReduce
2. MapReduce Workflow
3. MapReduce Word Count



---

Subject Teacher: **Prof. Prakash Parmar, Assistant Professor, CMPN**

- MapReduce is a framework used for writing application that processes large data sets in a distributed manner with parallel algorithms.
- It is core component of the Apache Hadoop ecosystem
- MapReduce have main two function
  1. Map
  2. Reduce
- Both these Map & Reduce works only on key-value pairs
  - Ex. (Roll\_No, Name)
  - (1, Sohan)
  - (2, Mohan)
- In both function, input is Key-value pair and output also Key, Value pair
  - (k, v) -> Map -> (k, v) -> Reduce -> (k, v)

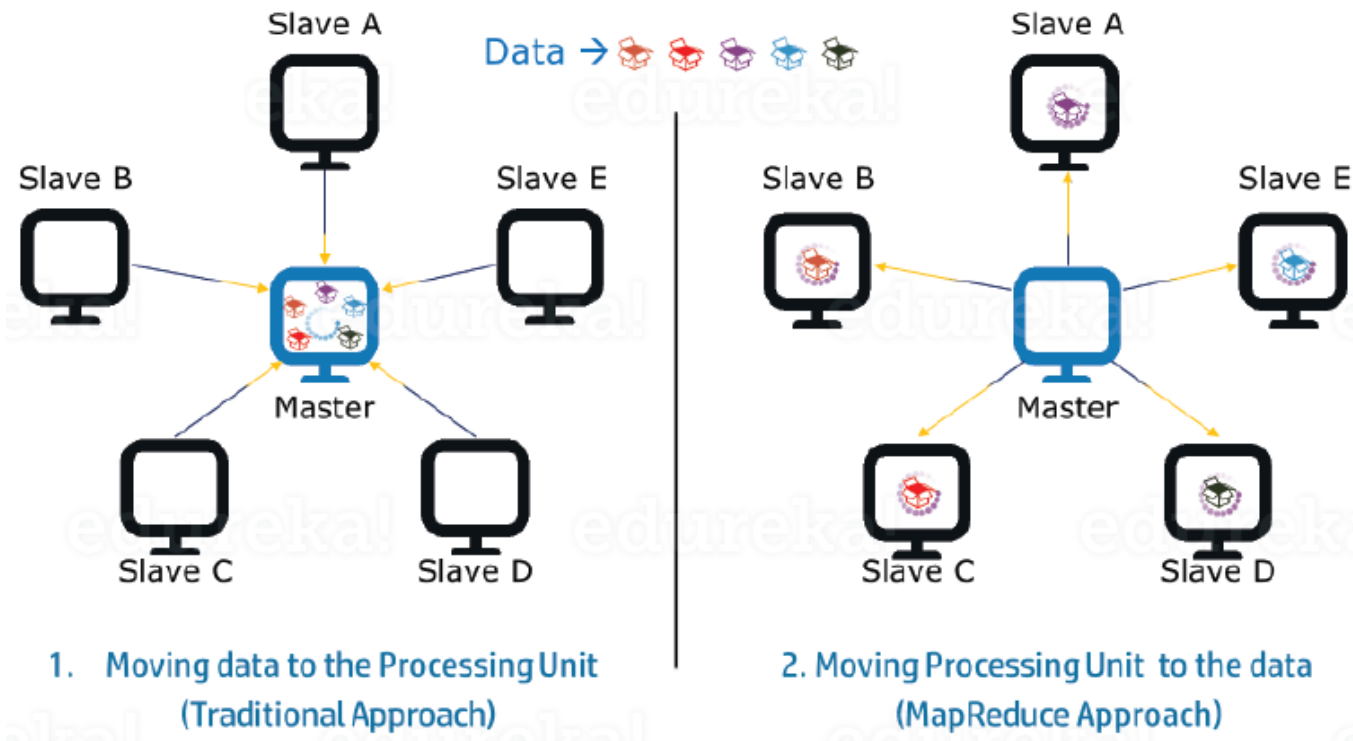
# MapReduce, Why?

- MapReduce is a programming paradigm
- Traditional programming models work only when data is kept on a single machine and if data kept on multiple machine in a distributed manner than we require a new programming model to solve the problems.
- ***MapReduce is a computing paradigm for processing data that resides on many machines.***

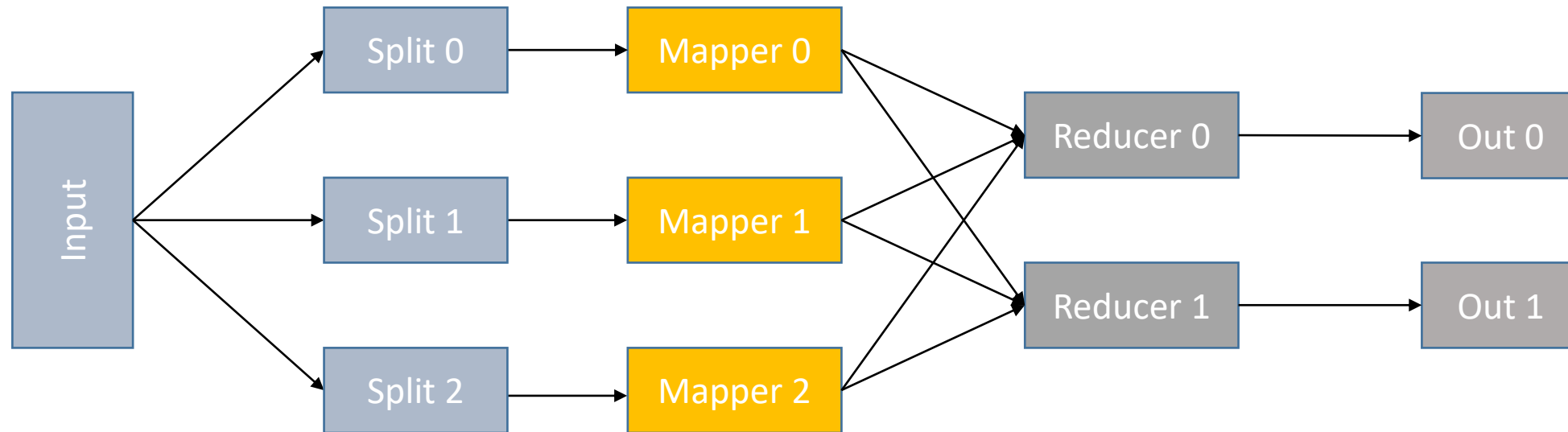


# MapReduce: Working Mechanism

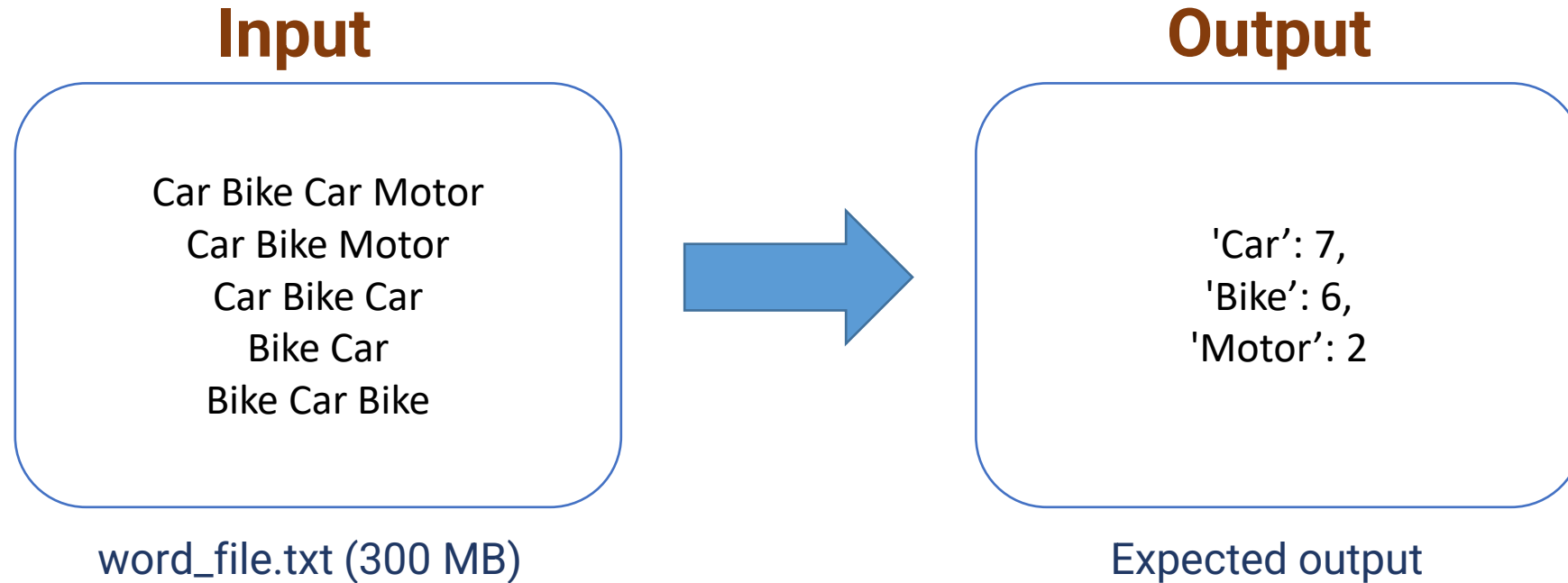
- Hadoop works on principle of data Locality, i.e., the data is processed where it is kept.
- In MapReduce code is move to data, data is not coming towards code, it makes fast processing.
- Data Block – 128 MB but code around 20 KB.



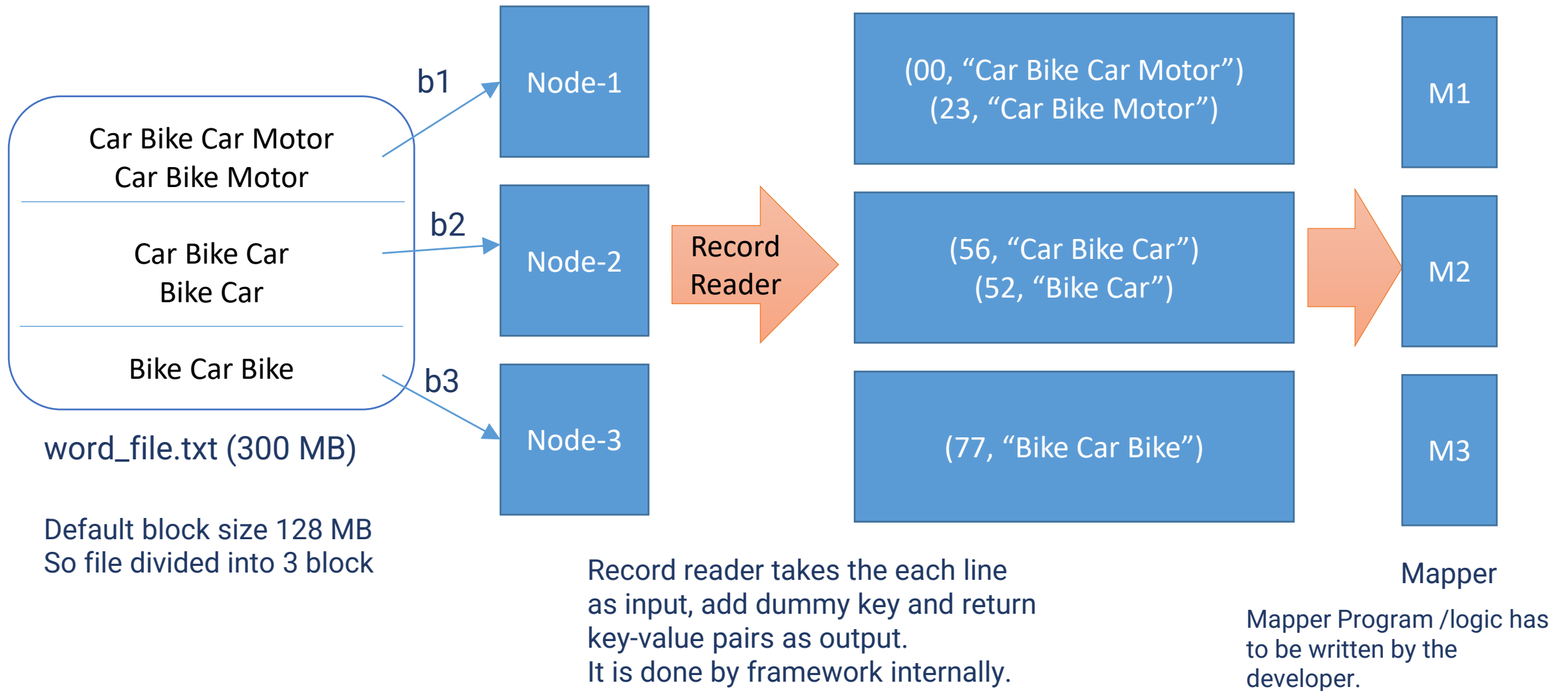
# MapReduce: Work Flow



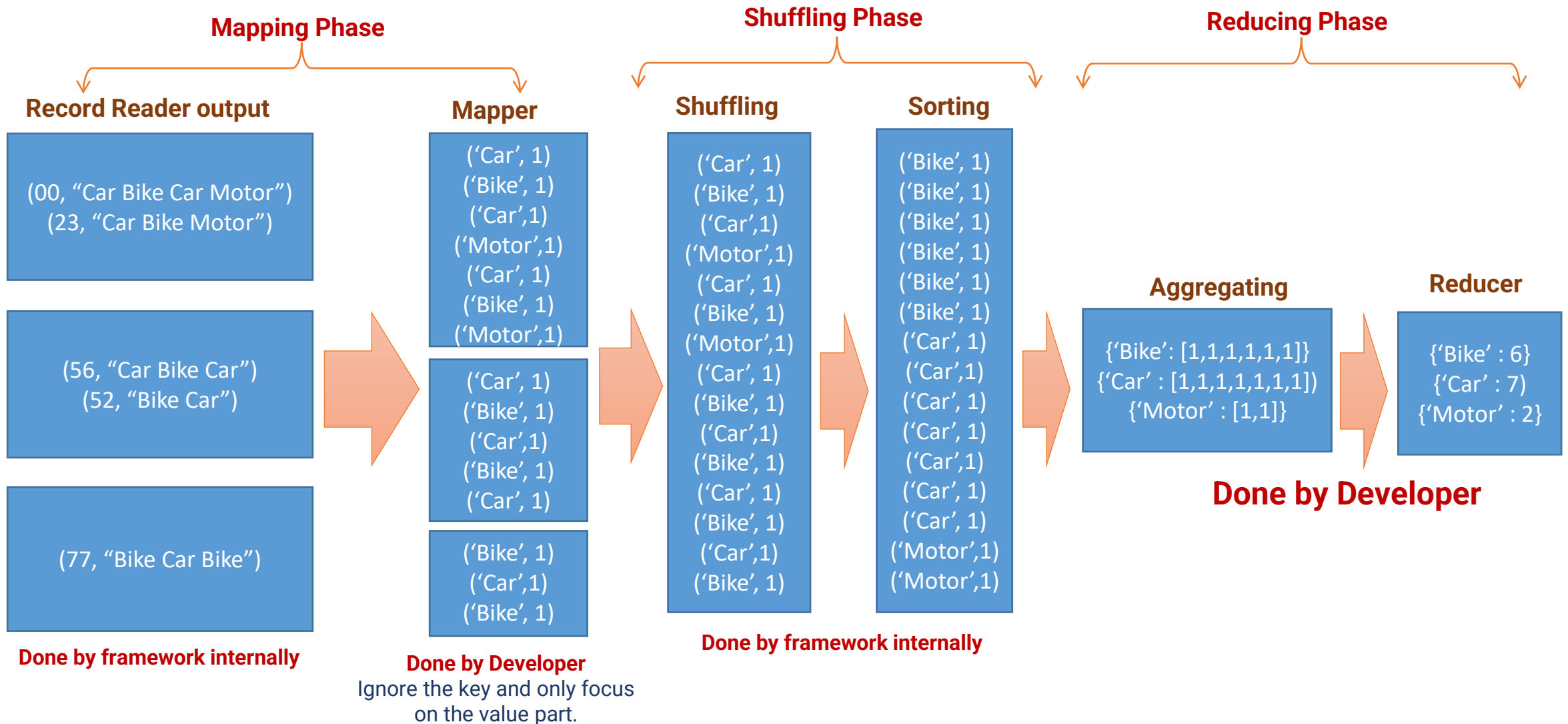
# MapReduce: Word Count



# MapReduce: Word Count



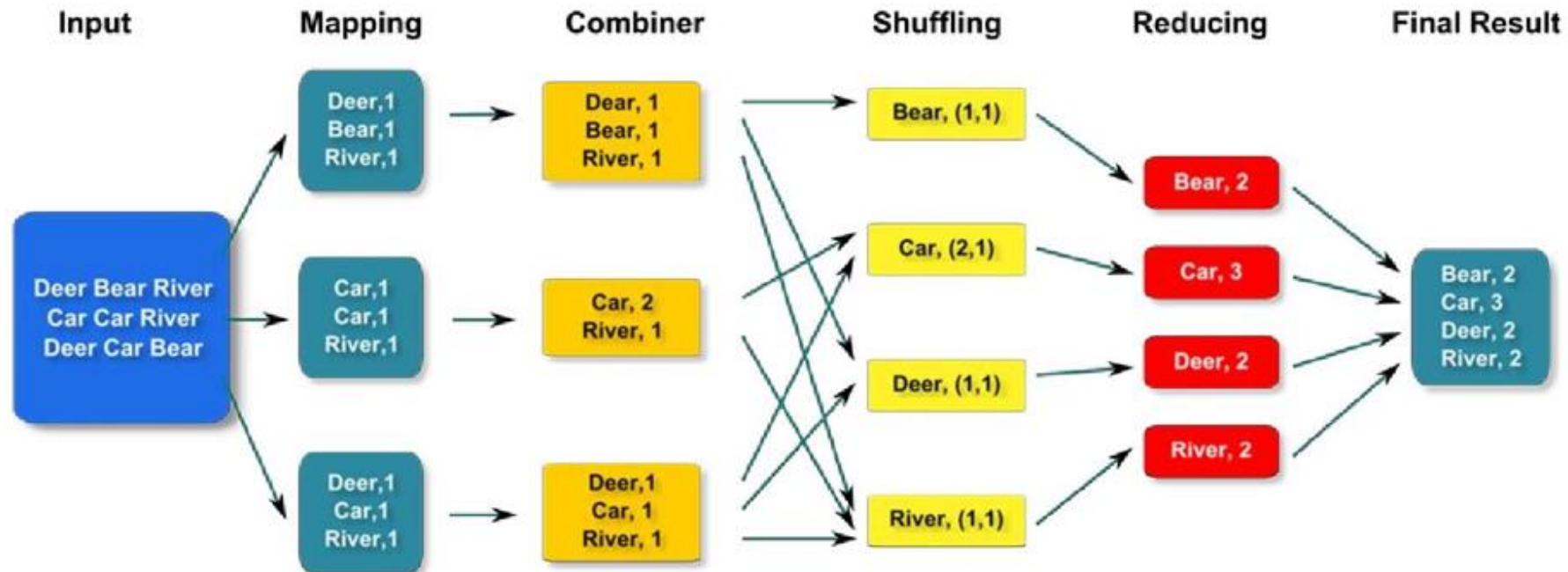
# MapReduce: Word Count





- Combiner is also known as “**Mini-Reducer**” that synopsisizes the Mapper output record with the same Key before passing to the Reducer.
- On a huge dataset when we run a MapReduce job. So Mapper creates large chunks of intermediate data. Then the framework passes this intermediate data on the Reducer for further handling. This leads to huge network congestion. The Hadoop framework offers a function known as Combiner that plays a key role in reducing network congestion.
- The main job of Combiner a “Mini-Reducer is to handle the output data from the Mapper, before passing it to Reducer. It works after the mapper and before the Reducer. Its usage is optional.

## Combiner - Local Reduce



## Advantages

- Use of combiner decreases the time taken for data transfer between mapper and reducer.
- Combiner increases the overall performance of the reducer.
- It reduces the amount of data that the reducer has to process.

## Disadvantages

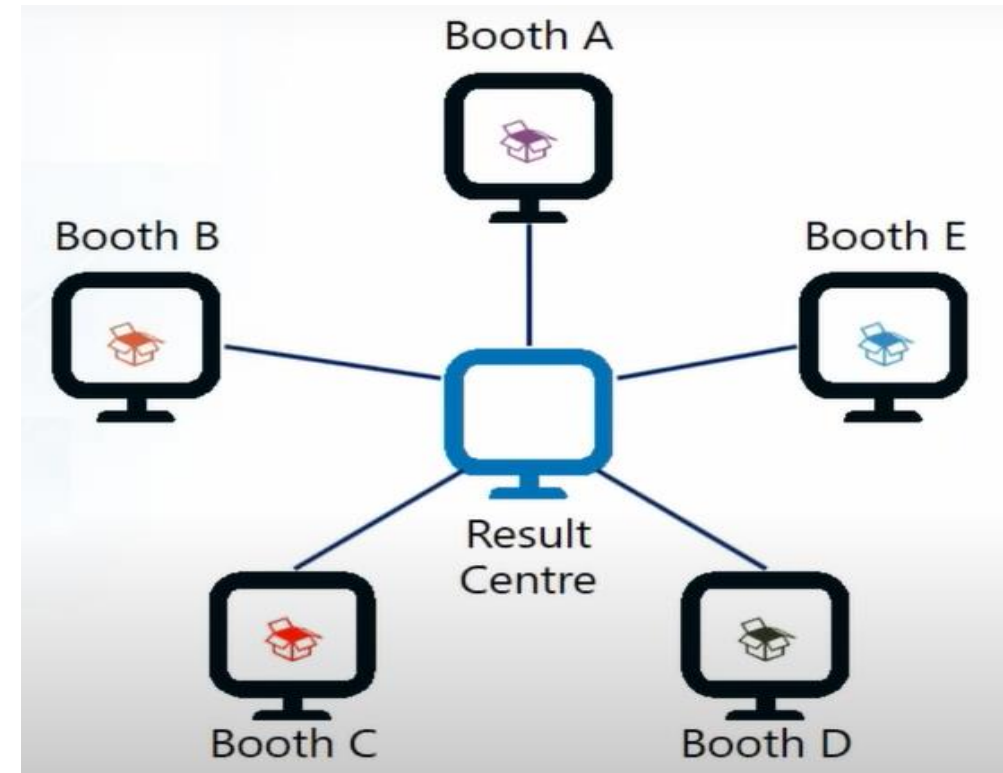
- In the native filesystem, when Hadoop stores the key-value pairs and runs the combiner later this will result in expensive disk IO.
- MapReduce jobs can't rely on the combiner execution as there is no guarantee in its execution.

# MapReduce Use Case: Election Vote Counting

Votes is stored at different Booths and Result centre has the details of all the booths

## Counting: Traditional Approach

- Votes are moved to result centre for counting.
- Moving all the votes to centre is costly.
- Results centre is over burdened.
- Counting takes time.

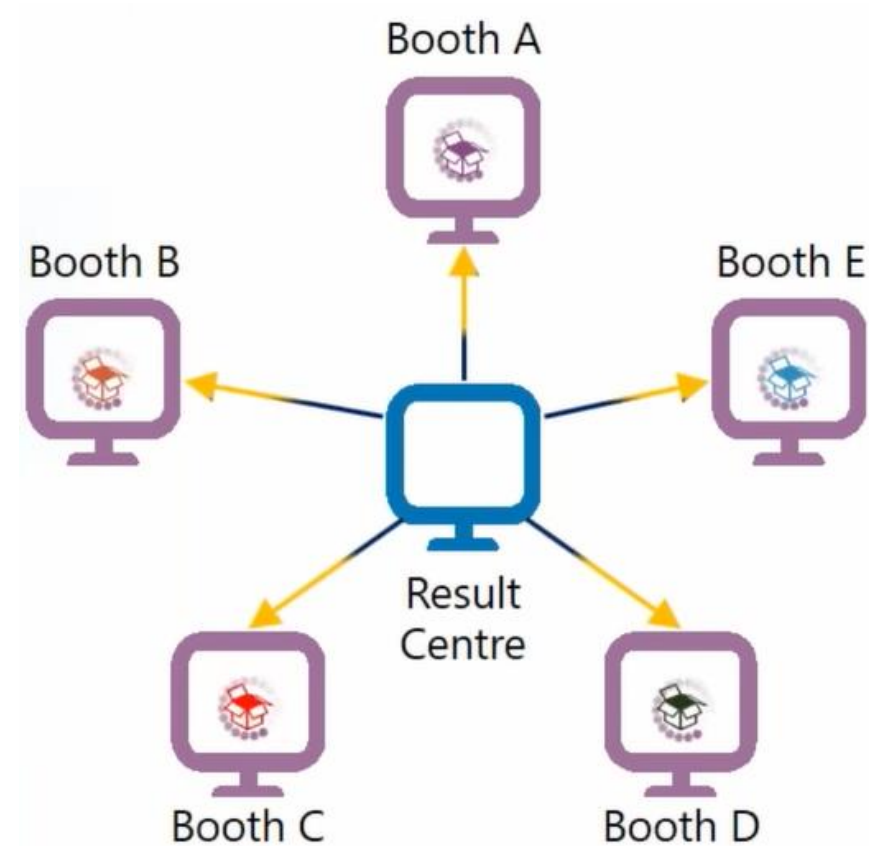


# MapReduce Use Case: Election Vote Counting

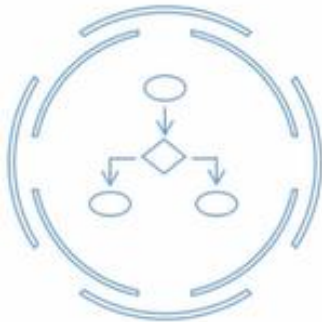
Votes is stored at different Booths and Result centre has the details of all the booths

## Counting: MapReduce Approach

- Votes are counted at individual booths.
- Booth-wise results are sent to the result centre.
- Final result is declared easily and quickly using this way.



## Real-Time Uses of MapReduce



Algorithms



Sorting



Data mining



Search engine operations



Enterprise analytics



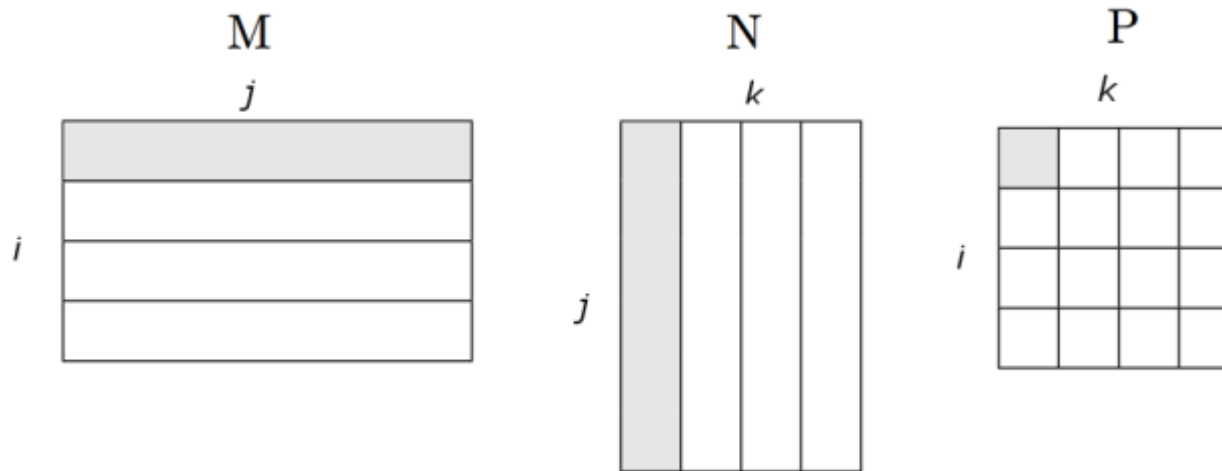
Gaussian analysis



Semantic web 3.0

# Algorithm uses MapReduce : Matrix Vector Multiplication

Matrix-vector and matrix-matrix calculations fit nicely into the MapReduce style of computing.



## Matrix Data Model for MapReduce



---

## Algorithm 1: The Map Function

---

```
1 for each element  $m_{ij}$  of  $M$  do
2   produce (key, value) pairs as  $((i, k), (M, j, m_{ij}))$ , for  $k = 1, 2, 3, \dots$  up
   to the number of columns of  $N$ 
3 for each element  $n_{jk}$  of  $N$  do
4   produce (key, value) pairs as  $((i, k), (N, j, n_{jk}))$ , for  $i = 1, 2, 3, \dots$  up
   to the number of rows of  $M$ 
5 return Set of (key, value) pairs that each key,  $(i, k)$ , has a list with
   values  $(M, j, m_{ij})$  and  $(N, j, n_{jk})$  for all possible values of  $j$ 
```

---

---

## Algorithm 2: The Reduce Function

---

```
1 for each key  $(i, k)$  do
2   sort values begin with  $M$  by  $j$  in  $list_M$ 
3   sort values begin with  $N$  by  $j$  in  $list_N$ 
4   multiply  $m_{ij}$  and  $n_{jk}$  for  $j_{th}$  value of each list
5   sum up  $m_{ij} * n_{jk}$ 
6 return  $(i, k), \sum_{j=1} m_{ij} * n_{jk}$ 
```

---



## Matrix Data Model for MapReduce

$$M_{ij} = \begin{bmatrix} 10 & 20 \\ 30 & 40 \end{bmatrix} \quad N_{jk} = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$$

As input data files, we store matrix M and N on HDFS in following format

$(M, i, j, m_{ij})$  and  $(N, i, j, m_{ij})$

$(M, 1, 1, 10)$  ,  $(M, 1, 2, 20)$ ,  $(M, 2, 1, 30)$  ... ..

$(N, 1, 1, 5)$  ,  $(N, 1, 2, 6)$ ,  $(N, 2, 1, 7)$  ... ..

Note:

Most matrices are sparse so large amount of cells have value zero. When we represent matrices in this form, we do not need to keep entries for the cells that have values of zero to save large amount of disk space.

# Matrix Multiplication using MapReduce

for each element  $m_{ij}$  of M do  
 produce (key, value) pairs as  $((i, k), (M, j, m_{ij}))$   
 for  $k=1, 2, \dots$  Up to column of matrix N

## Mapper code for Matrix M

k=#columns of N	i = row number of M	j = column number of M	formula (key, value) $((i, k), (M, j, m_{ij}))$
k=1	i=1	j=1	$((1, 1), (M, 1, 10))$
		j=2	$((1, 1), (M, 2, 20))$
	i=2	j=1	$((2, 1), (M, 1, 30))$
		j=2	$((2, 1), (M, 2, 40))$
k=2	i=1	j=1	$((1, 2), (M, 1, 10))$
		j=2	$((1, 2), (M, 2, 20))$
	i=2	j=1	$((2, 2), (M, 1, 30))$
		j=2	$((2, 2), (M, 2, 40))$

for each element  $n_{jk}$  of N do  
 produce (key, value) pairs as  $((i, k), (N, j, n_{jk}))$   
 for  $i=1, 2, \dots$  Up to row of matrix M

## Mapper code for Matrix N

i=#columns of M	j = row number of N	k = column number of N	formula (key, value) $((i, k), (N, j, n_{ij}))$
i=1	j=1	k=1	$((1, 1), (N, 1, 5))$
		k=2	$((1, 2), (N, 1, 6))$
	j=2	k=1	$((1, 1), (N, 2, 7))$
		k=2	$((1, 2), (N, 2, 8))$
i=2	j=1	k=1	$((2, 1), (N, 1, 5))$
		k=2	$((2, 2), (N, 1, 6))$
	j=2	k=1	$((2, 1), (N, 2, 7))$
		k=2	$((2, 2), (N, 2, 8))$

Return Set of (key, value) pairs that same key has a list with values M & N

$((1,1), [(M, 1, 10)), (M, 2, 20)]),$

$((2,1), [(M, 1, 30)), (M, 2, 40)]),$

$((1,2), [(M, 1, 10)), (M, 2, 20)]),$

$((2,2), [(M, 1, 30)), (M, 2, 40)])$

$((1,1), [(N, 1, 5)), (N, 2, 7)]),$

$((1,2), [(N, 1, 6)), (N, 2, 8)]),$

$((2,1), [(N, 1, 5)), (N, 2, 7)]),$

$((2,2), [(N, 1, 6)), (N, 2, 8)]),$

# Matrix Multiplication using MapReduce

---

## Algorithm 2: The Reduce Function

---

```
1 for each key (i,k) do
2   sort values begin with M by j in listM
3   sort values begin with N by j in listN
4   multiply  $m_{ij}$  and  $n_{jk}$  for  $j_{th}$  value of each list
5   sum up  $m_{ij} * n_{jk}$ 
6 return (i,k),  $\sum_{j=1} m_{ij} * n_{jk}$ 
```

---

## Reduce Code:

Sort values in M and N list according to J and perform multiplication.

$((1,1), 10*5+20*7) \Rightarrow ((1,1), 190)$

$((1,2), 10*6+20*8) \Rightarrow ((1,2), 220)$

$((2,1), 30*5+40*7) \Rightarrow ((2,1), 430)$

$((2,2), 30*6+40*8) \Rightarrow ((2,2), 500)$

## Result :

$$M_{ij} = \begin{bmatrix} 10 & 20 \\ 30 & 40 \end{bmatrix} \quad N_{jk} = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \quad M \times N = \begin{bmatrix} 190 & 220 \\ 430 & 500 \end{bmatrix}$$

- **Failures** are **norm** in commodity hardware
- **Worker** failure
  - Detect failure via periodic **heartbeats**
  - **Re-execute** in-progress map/reduce tasks
- **Master** failure
  - Single point of failure; Resume from Execution Log
- **Robust**
  - Google's experience: **lost 1600 of 1800 machines once!**, but **finished fine**

# *Learn Fundamentals & Enjoy Engineering*

*Thank You*



Prof. Prakash Parmar  
Assistant Professor  
Computer Engineering Department  
Vidyalankar Classes CSE GATE Faculty