
2017



Map-Reduce Working in Hadoop

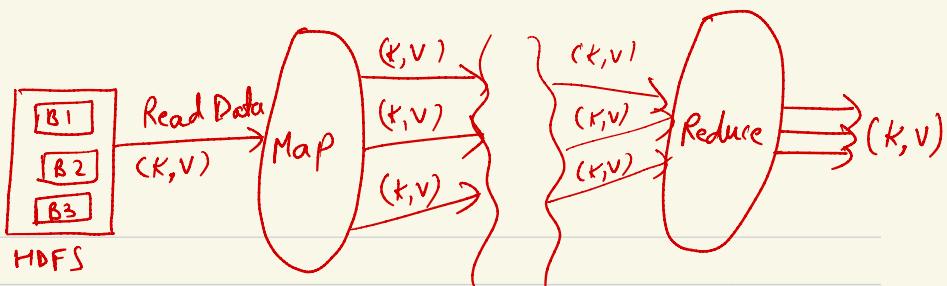
① HDFS ↗

② YARN ↗

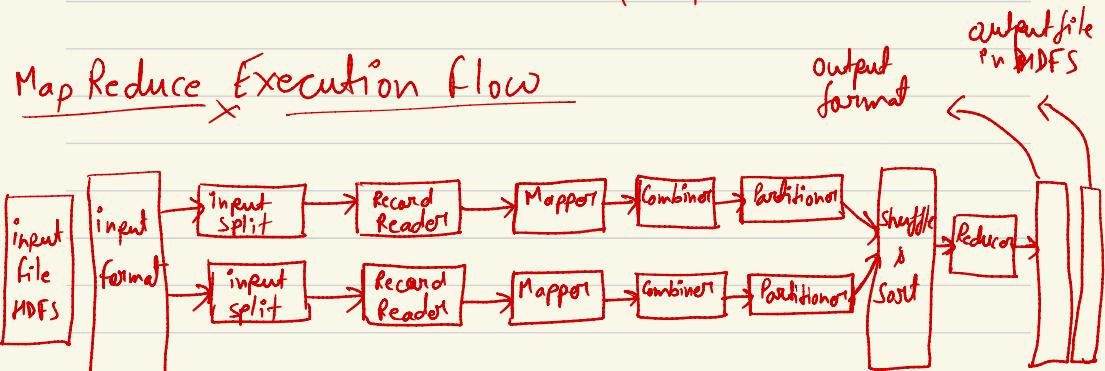
③ MAP-REDUCE

→ MAP → Reads Blocks/data and generate (Key, Value) output.

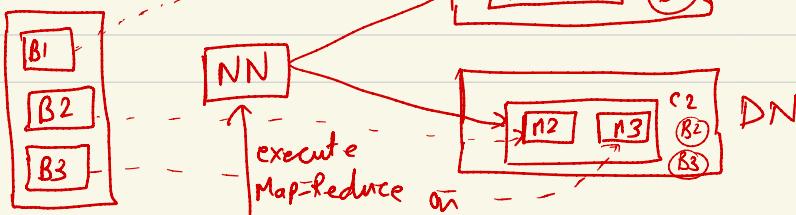
→ REDUCE → Reads (Key, Value) data generated by map phase and aggregates the result. Final result will be also in (Key, Value)



Map Reduce Execution Flow



file.txt



HDFS

Client

file.txt

word-count.py

def mapper ({k, v}):

return {k1, v1}

Coming from HDFS
Input file from
HDFS

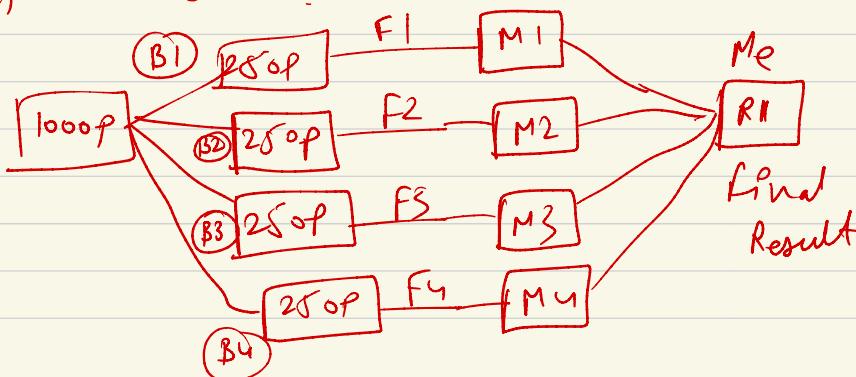
def Reducer ({k1, v1}):

print ({k2, v2})

Output file to
HDFS

book of 1000 pages.

Q) Time for a human to calculate frequency of each word?



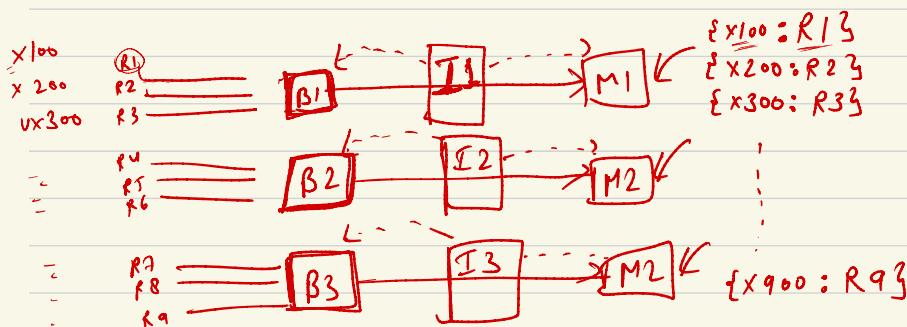
① Input Files: Simply the files or references for the data stored in HDFS

② Input Format: defines how those input files are read from source or how to serialize & deserialize them properly.

③ Input Splits: Logical representation of the data which will be processed by an individual mapper task. One map task is created for each input split.

* Block: Actual unit whole raw data will be stored physically.

* Input Split: It is like a reference to the original Block of Data.

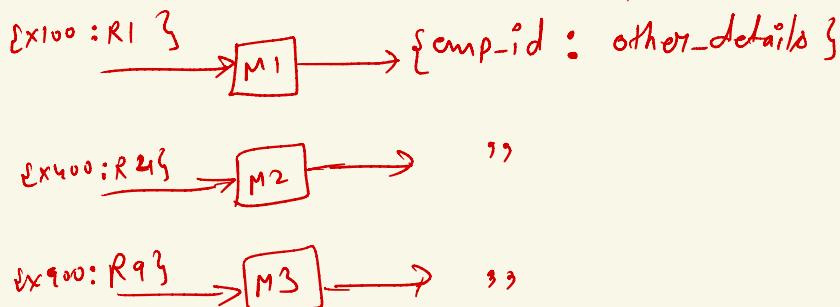


* Record Reader

↳ It communicates with the Input Splits in Hadoop Map Reduce and converts the data into key-value pairs.

Mapper:

It processes each input record from record Reader and generates new Key-Value pair which is different from input Key-Value pair.

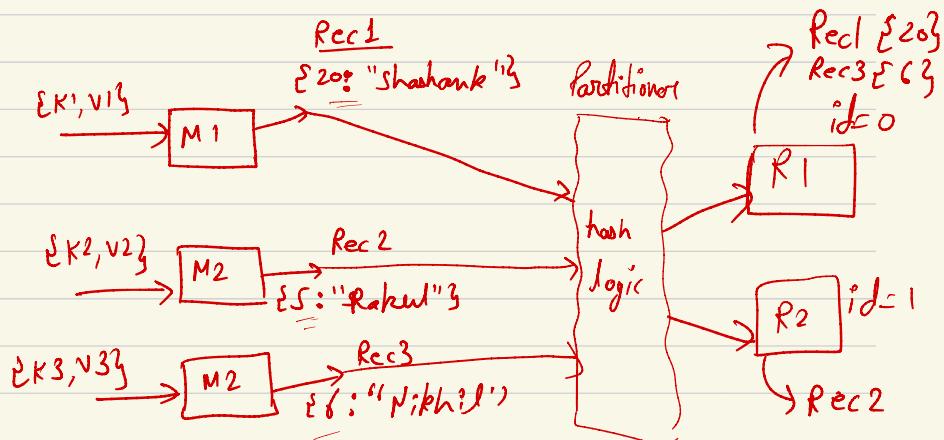


Combiner

It is like a mini-reducer. It will aggregate results locally generated by each mapper task.

Partitioner:

It will decide in which reducer records of mapper task will go to.



Partitioner logic or hash function.

$$\text{hash}(\text{Key}) = \text{Key \% total_num_of_Reducers}$$

R12
 $\{20: \text{"Shashank"}\} \Rightarrow \text{Key} = 20$
 $= \text{hash}(20)$
 $= 20 \% 2$
 $= 0$

Rec2

$$\{5: \text{"Rahul"}\} \Rightarrow \text{Key} = 5$$

 $= \text{hash}(5)$
 $= 5 \% 2$
 $= 1$

x) Shuffling & Sorting:

It will arrange the data for each key on individual reducer.

x) Reducer: It will receive mapper task's data as input and will generate resultant output.

x) Output Format: In which format we want to write data in HDFS

* Output Files: Generated Reference for output data.

Word Count Example in Map-Reduce.

file.txt



Combiner

Partitioner

Shuffle & Sort

Reduce

Cat, 2 mat, 1
Sat, 1 bat, 1

hash()

R1
Cat, 2
mat, 1
Owl, 2
Cat, 1
Cat, 1
Mat, 2
Cat, 2

Owl, 2
Cat, 1 bat, 1

Sat, 3
Cat, 1

Mat, 2
Cat, 2

Cat, 2
Cat, 1
Cat, 1
Cat, 2
Mat, 1
Mat, 2
Owl, 2

Cat, 6
Mat, 3
Owl, 2

R2
Cat, 1
Bat, 1
Bat, 1
Sat, 3

Bat, 1
Bat, 1
Sat, 1
Mat, 3

Bat, 2
Sat, 4

