

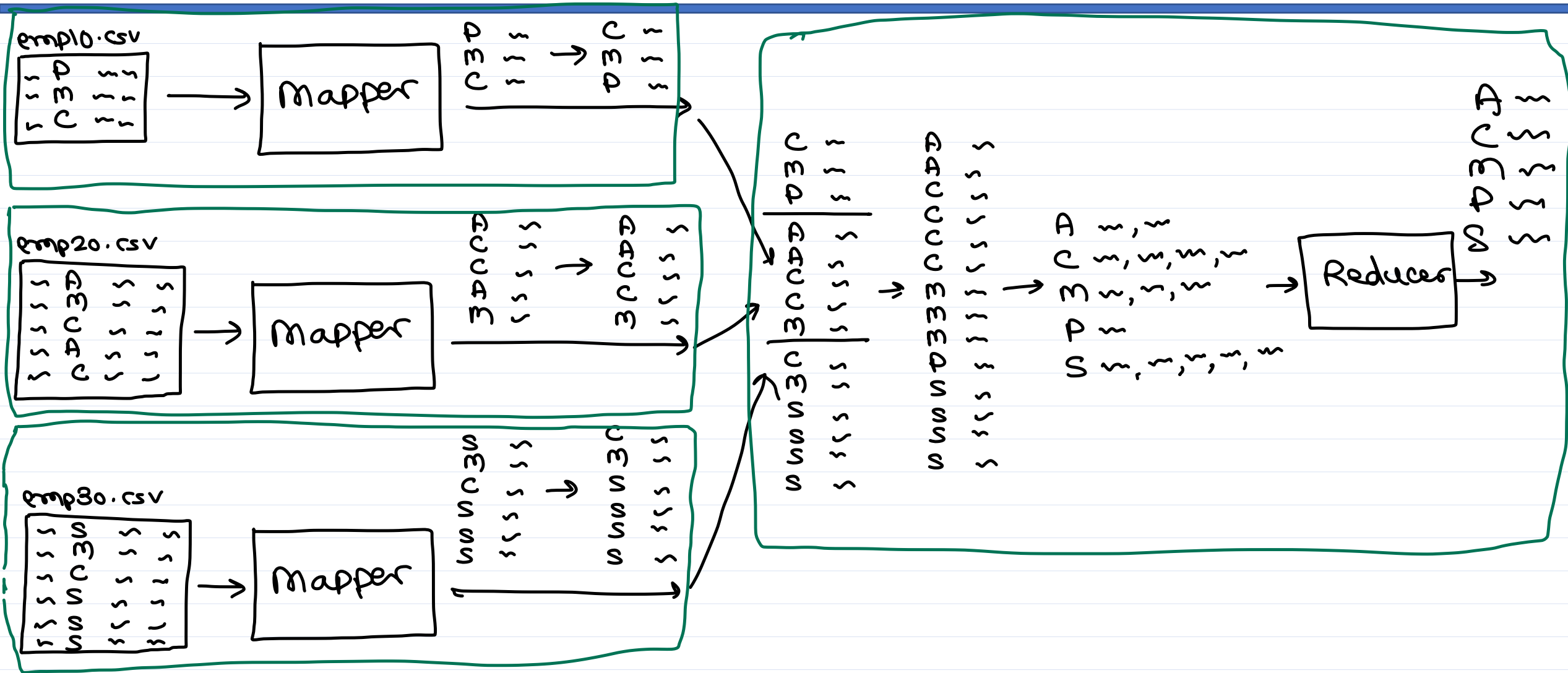


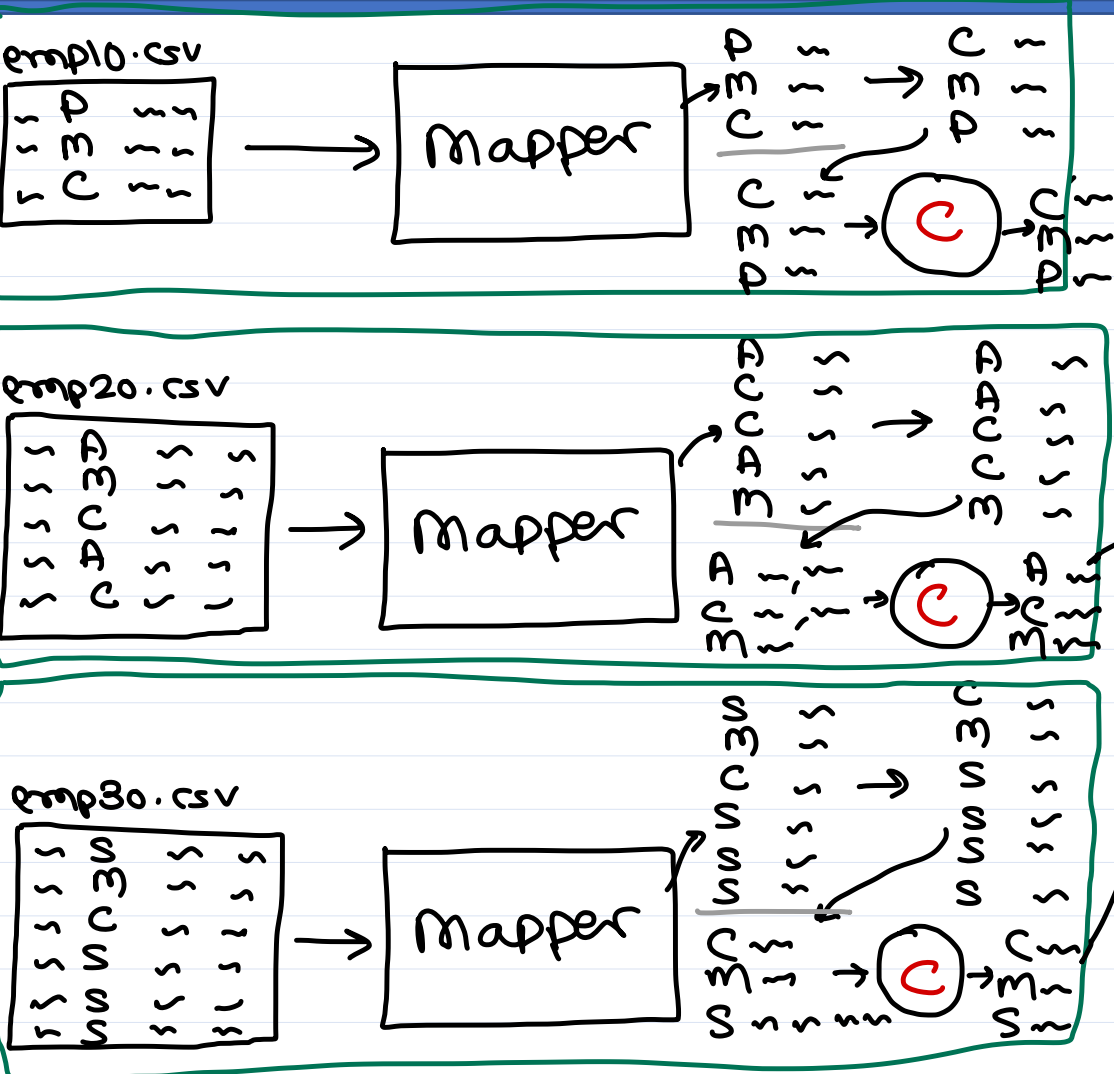
Big Data Technologies

Trainer: Mr. Nilesh Ghule.

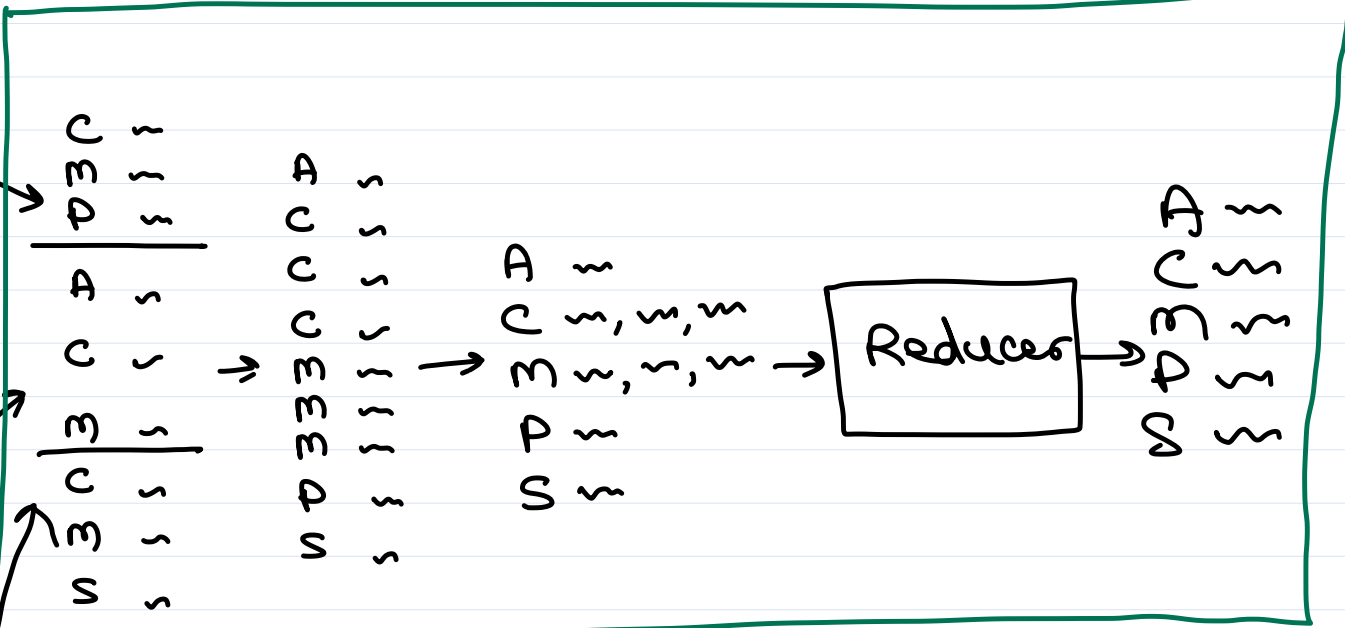


MR execution flow Multiple Mappers and Single Reducer





job.setCombinerClass(EmpReducer.class);

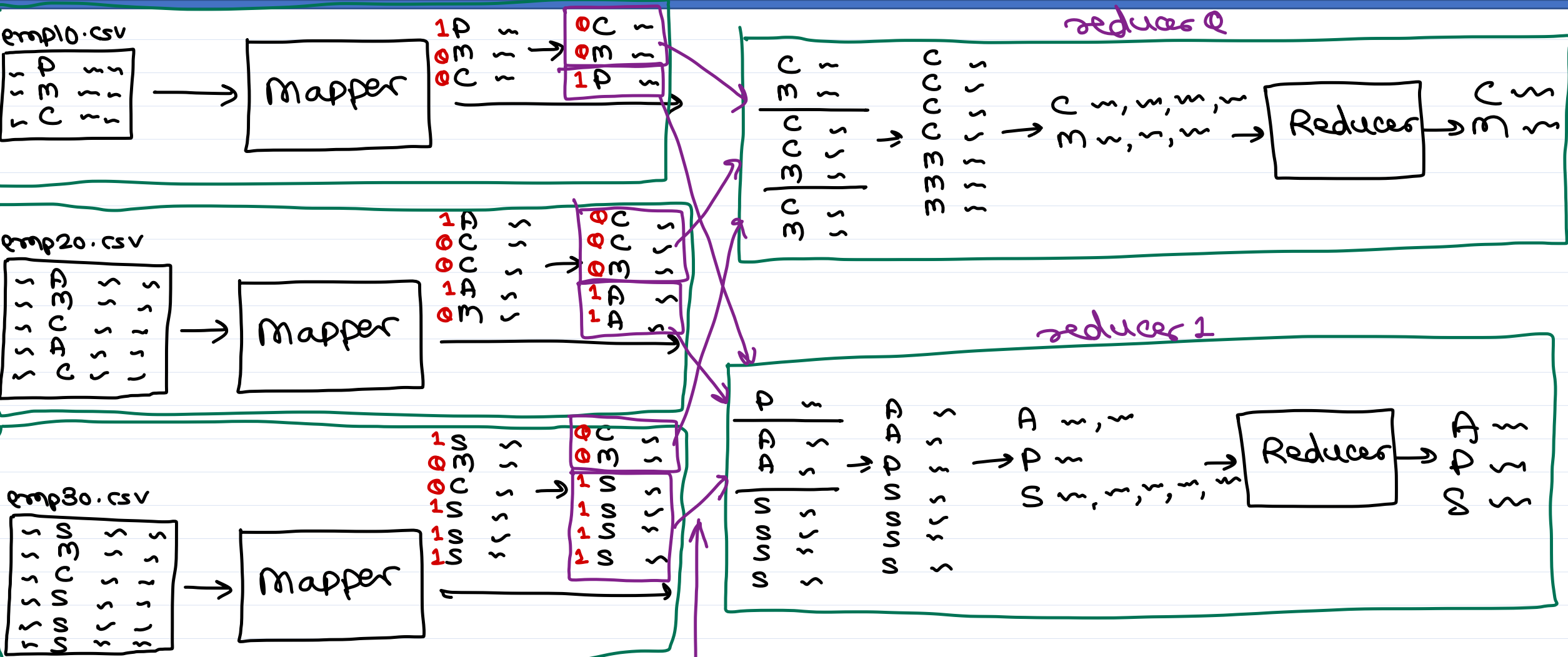


Advantages: ① Reduce network traffic betn mapper & reducer.
② Reduce work load/resources on reducer.

We can use Combiner if:

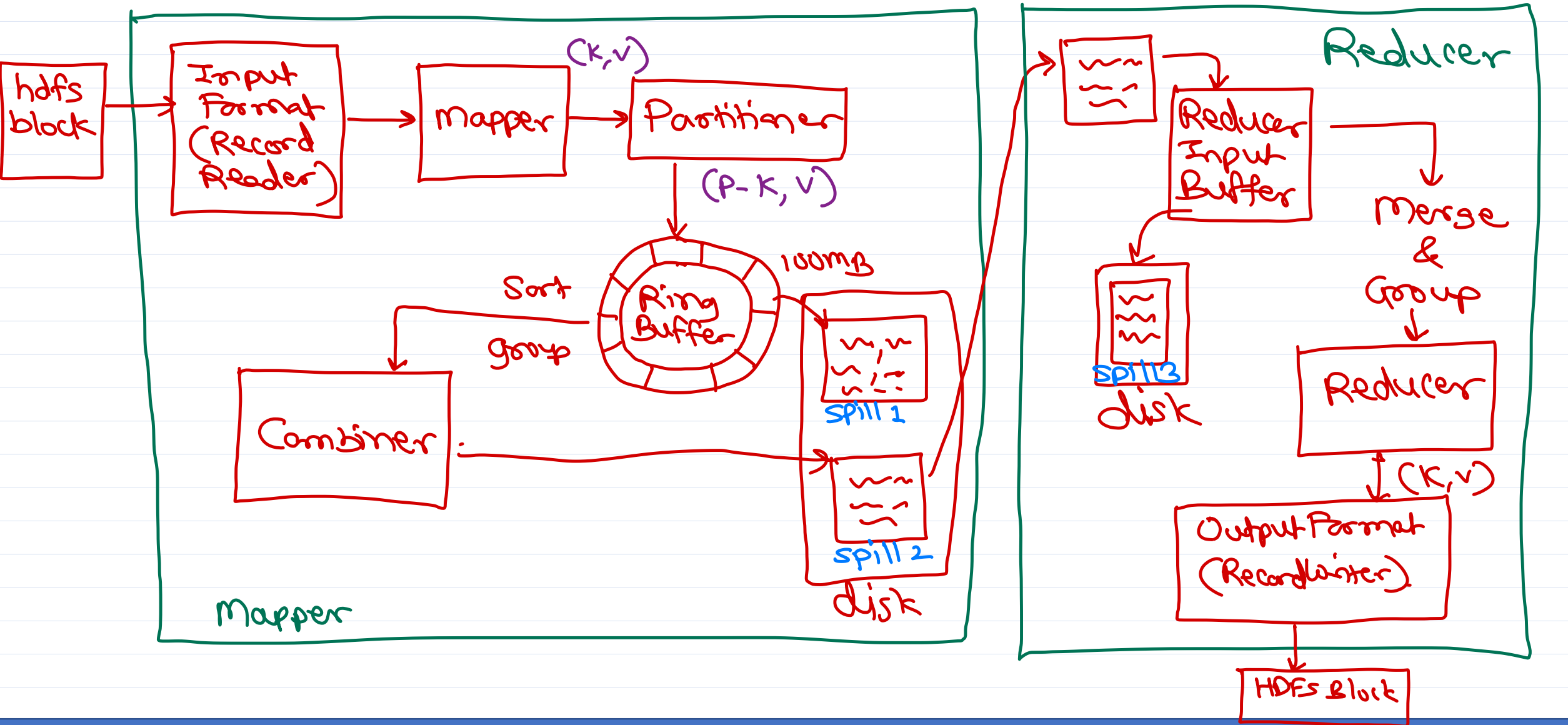
- ① agg op is associative & commutative
- ② combiner/reducer class in key value type should be same as out key value type.

MR execution flow Multiple Mappers and Multiple Reducer

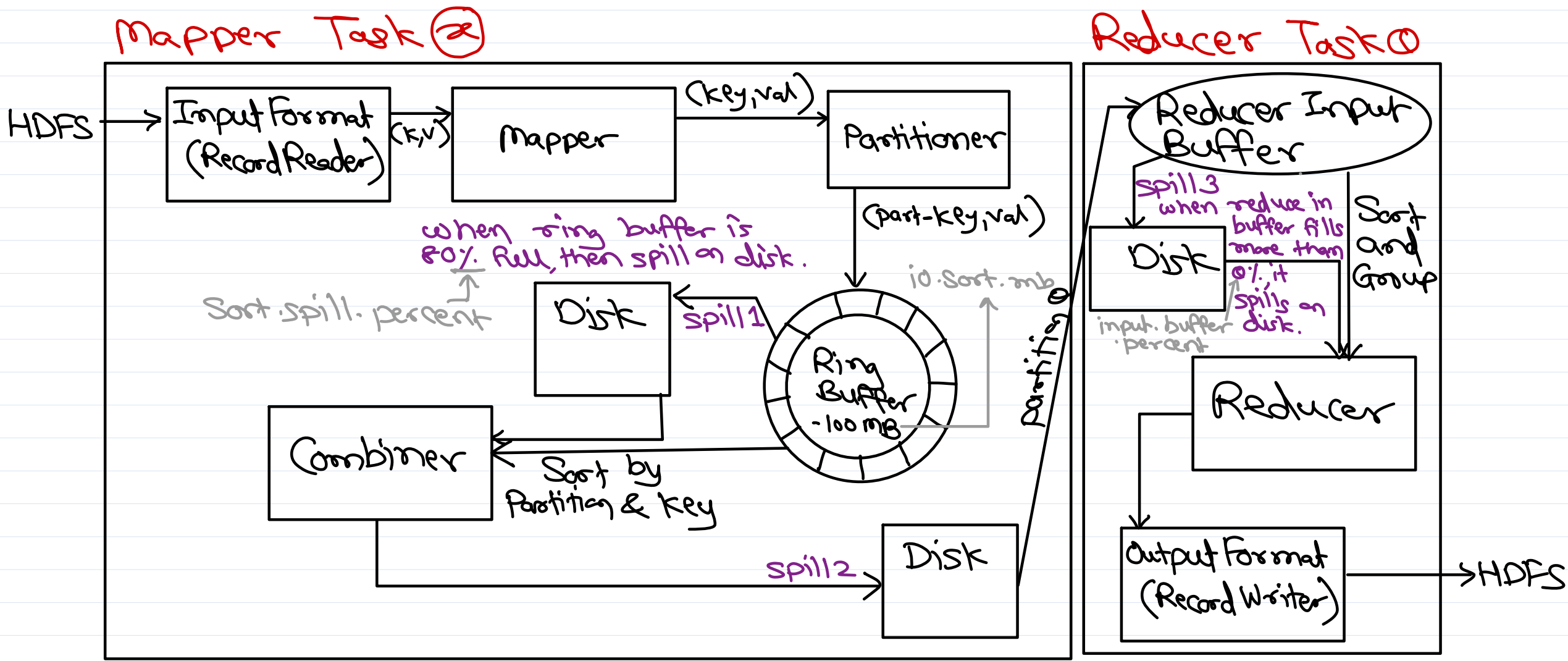


✓ Num of partitions = Num of Reducers. Shuffle-Sort-Merge
✓ HashPartitioner: part no = key.hashCode() % Num of Reducers.

Map Reduce Flow



Map Reduce Data Flow



Input Splits

Number of mappers = Number of input splits

Number of Input splits \approx Number of HDFS blocks

Ex1: 20 files each of 1MB approx

→ Num of HDFS blocks = 20

→ Num of input splits = 20

→ Num of mappers = 20.

Ex2: 2 files - $f1 = 200\text{MB}$, $f2 = 300\text{MB}$

→ Num of HDFS blocks = $2 + 3 = 5$

→ Num of input splits = 5

→ Num of mappers = 5

Input format is responsible for reading records one by one. It has an internal RecordReader.

Text Input Format → 1 record = 1 line → (offset, line).

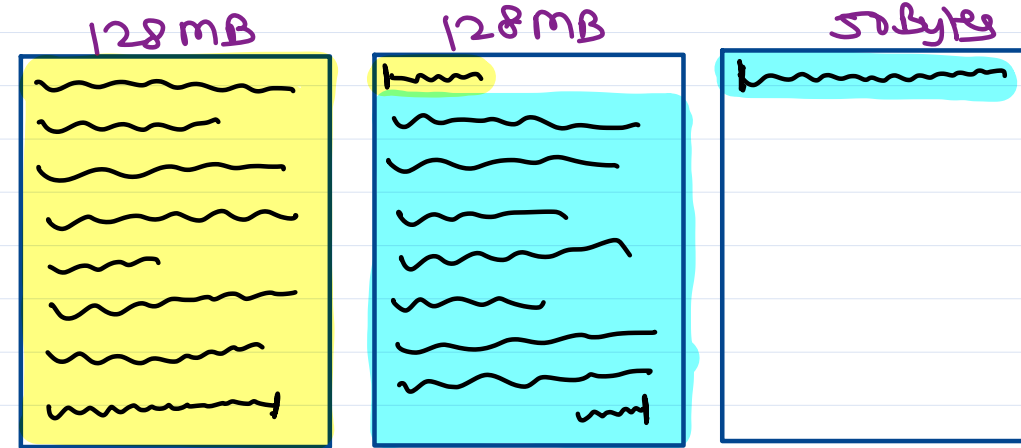
KeyValueText Input Format → 1 record = 1 line → (key, value)
key | value

DB Input Format → 1 record = 1 row (RDBMS)

Combined Text Input Format → Combine multiple text files into single input split.

$f1 = 256\text{MB} + 50\text{bytes}$.

HDFS blocks = 3.



Input split = 2

Input split is logical part of data that is to be processed by a single mapper.

Mostly 1 split = 1 hdfs block.





Thank you!

Nilesh Ghule <nilesh@sunbeaminfo.com>

