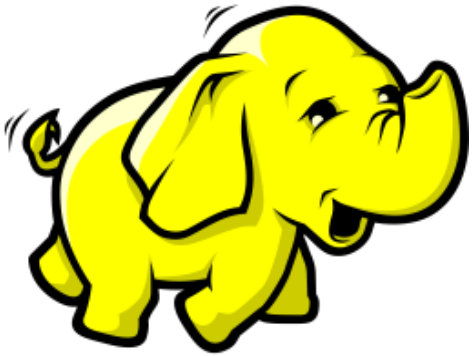# Big Data – Hadoop

Trainer: Mr. Nilesh Ghule.

# Data flow of MR job (Single reducer) Job wise total salary

**input HDFS**

emp10.csv

split 0 → map → sort → [ ]

copy

emp20.csv

split 1 → map → [ ]

emp30.csv

split 2 → map → [ ]

C M P (sort box)

+A

+C

+Sum

A C C C M M P P (column)

A C M P S (column)

A C M P S (red column)

merge → [ ] → reduce → part 0 → HDFS replication

**output HDFS**

A C M P S (right red column)

local reducer
on mapper node → aggregation

Combiner can be used if aggregate operation

$$A + (B + C) = (A + B) + C \rightarrow associative$$

$$A + B = B + A \rightarrow commutative$$

e.g. sum, max, min.

job.set Combiner Class ( EmpReducer.class );
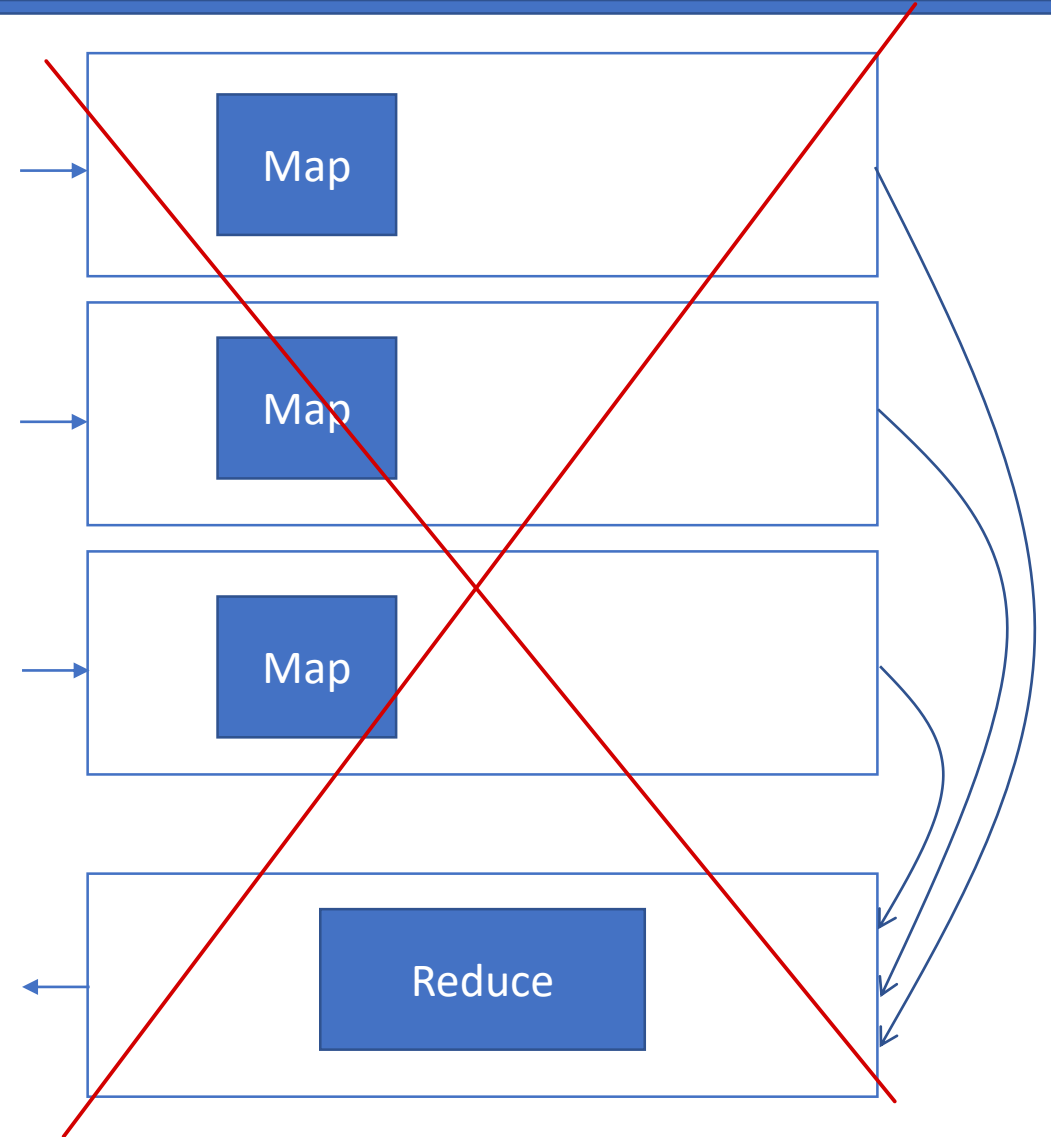
K → [ Combiner ] → K
V → V

# Combiner

- Combiner is a local reducer i.e. runs reducer (aggregation logic) within mapper task process.

- Minimize output for mapper task
  - Less merge & shuffle
  - Less network transfer
  - Less aggregation in reducer

- Combiner is optional.

- Works only for commutative & associative aggregate functions only.
  - A + B = B + A
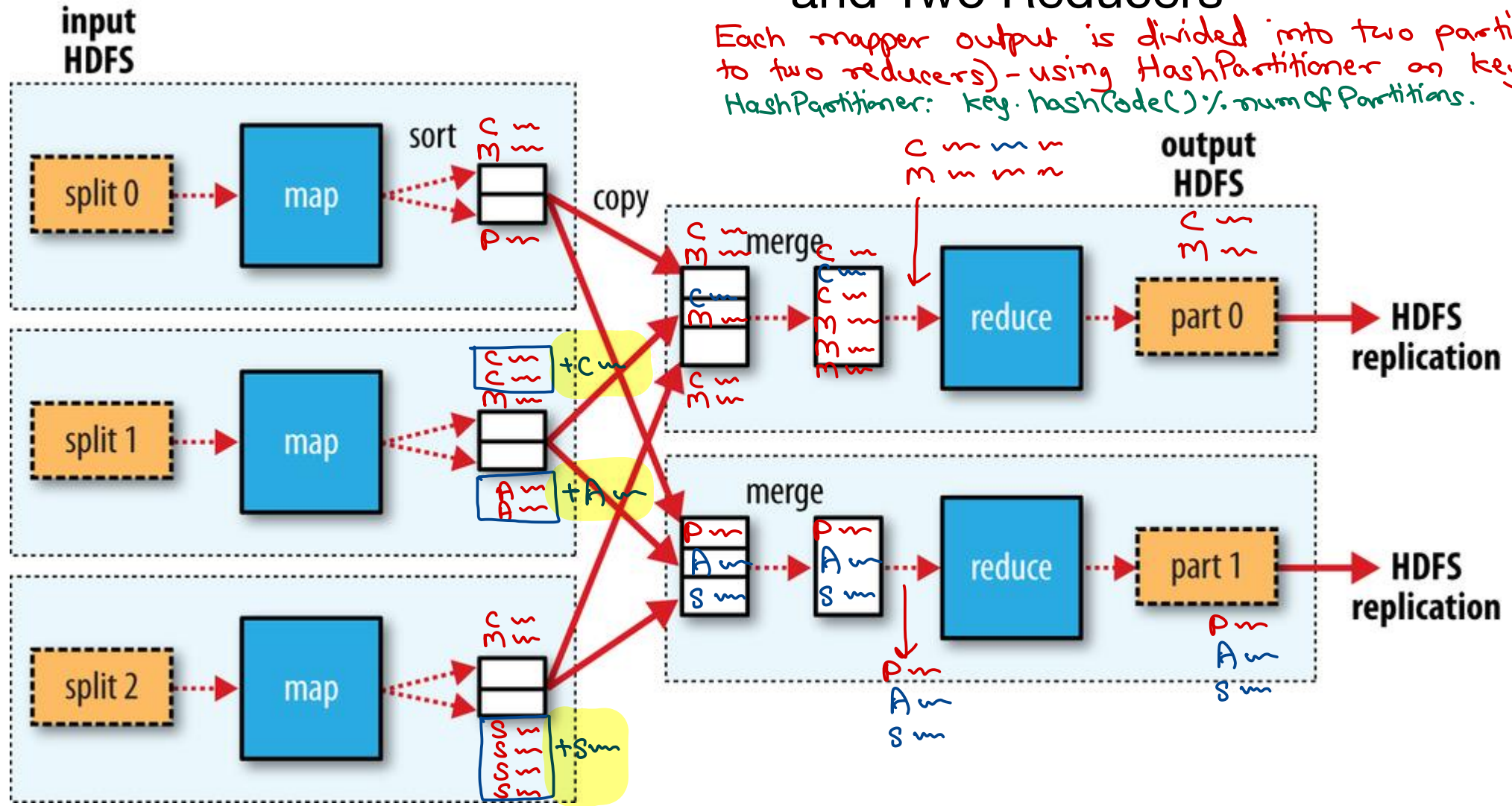  - A + (B + C) = (A + B) + C

# Partitioner

- By default MR job have single reducer.
- Having huge data for aggregation may lead to out of memory error.
- Number of reducers can be configured in job configuration file or in driver code.
  - job.setNumReduceTasks(2);
  - mapreduce.job.reduces = 2
- Number of partitions = Number of reducers
- Output of mapper is divided into multiple partitions based produced key
- By default HashPartitioner is used, that distributes mapper output in number of partitions uniformly.
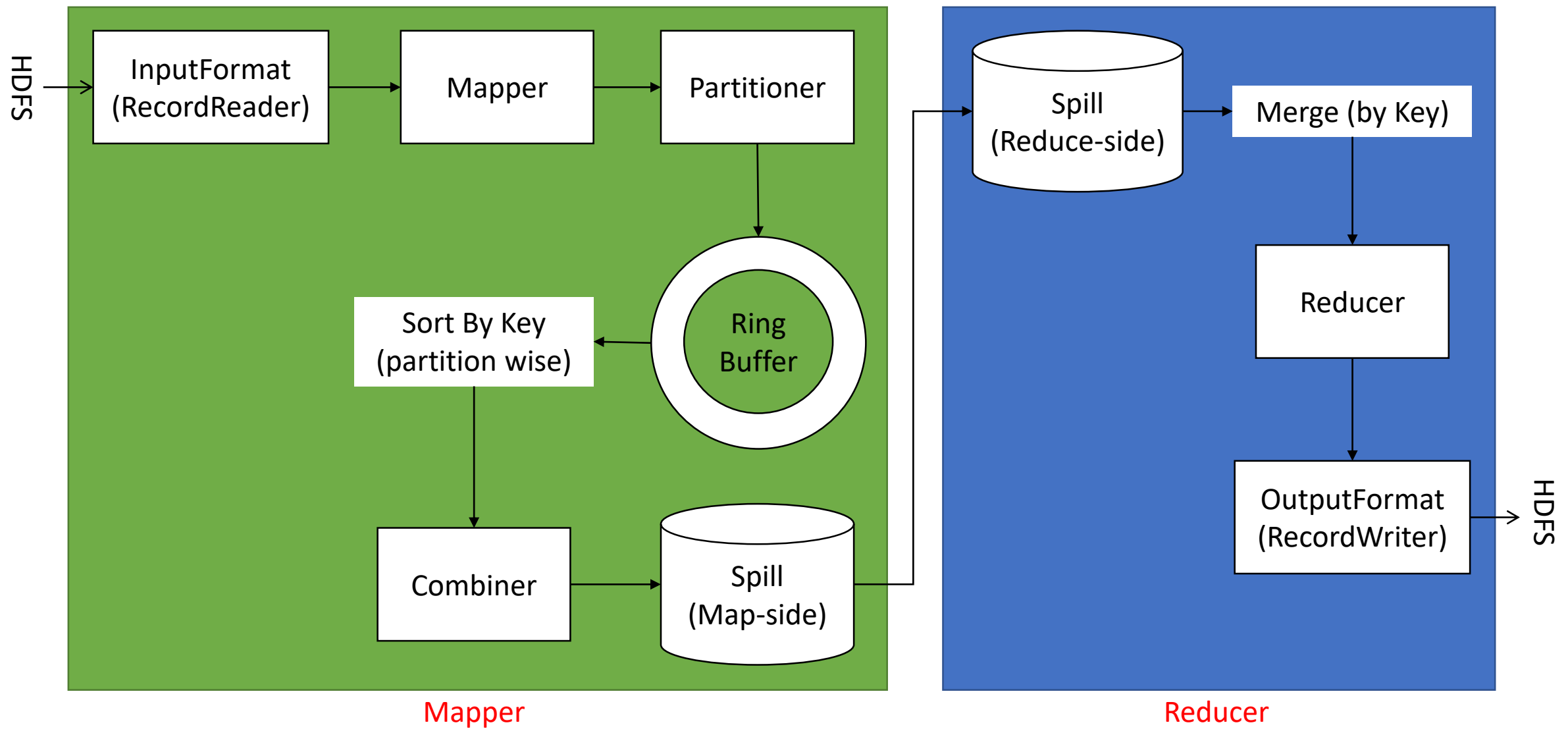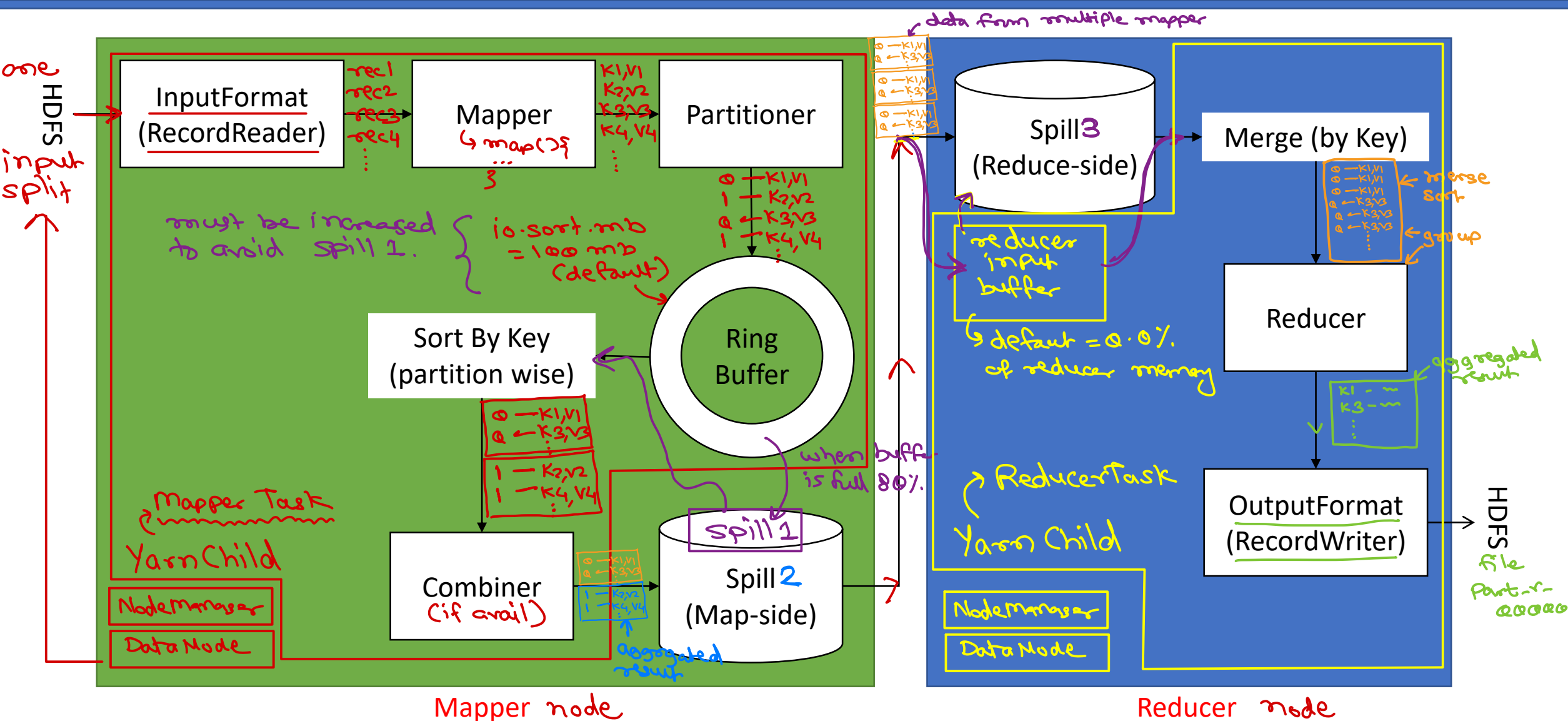
Map

Map

Map

Reduce

## Job wise total salary using Combiner and Two Reducers

Each mapper output is divided into two partitions (due to two reducers) - using HashPartitioner on key.

HashPartitioner: key.hashCode() % num Of Partitions.

# Hadoop MR data flow (detailed)



InputFormat (RecordReader) → Mapper → Partitioner → Ring Buffer → Sort By Key (partition wise) → Combiner → Spill (Map-side)

HDFS → InputFormat (RecordReader)

Spill (Reduce-side) → Merge (by Key) → Reducer → OutputFormat (RecordWriter) → HDFS

Mapper

Reducer

# Hadoop MR data flow (detailed)



Mapper node

Reducer node

Sunbeam Infotech

www.sunbeaminfo.com

# MR on YARN



hadoop jar app.jar

RunJar

Submit() @

Ⓒ Submit()

ⓑ Submit()

app id

jobjar
job.xml

/tmp/yarn/staging/
app-id/

job.jar
job.xml

nudes andadata

# Hadoop Streaming

```python
#!/usr/bin/python3
# mapper.py
import sys
for line in sys.stdin:
    words = line.split()
    for word in words:
        print(f"{word}\t1")
```

Stdout

```python
#!/usr/bin/python3
# reducer.py
di = dict()
for line in sys.stdin:
    (word,cnt) = line.split()
    newcnt = di.get(word, 0) + int(cnt)
    di[word] = newcnt
for word,total in di.items():
    print print(f"{word}\t{total}")
```

Stdout

```
hadoop jar $HADOOP_HOME/share/.../hadoop-streaming-2.7.3.jar \
-files mapper.py,reducer.py \
-input /user/nilesh/wc/input \
-output /user/nilesh/wc/output \
-mapper mapper.py -reducer reducer.py
```

*Contains driver code & other helper code.*



NodeManager

launch

**task JVM**

Child

run

MapTask or ReduceTask

input key/values — std in

output key/values — std out

launch

**Streaming process**

*mapper.py or reducer.py*

*python process*