

## Assignment - Benefit / Limitation Case of S.E

Date: \_\_\_\_\_ Page: 15

- \* Task complete - inform job tracker.
- \* whichever copy finishes first - definitive copy
- \* Kill other task (Slow)
- \* Enabled by default
- \* - mapred.map.task.Speculative.execution [config]
- mapred.reduce.task.Speculative.execution

### Combiner

- Mapper produce large amount of intermediate data
- Must be passed to Reducers
  - ↳ Can result in lot of n/w traffic
- Like a mini reducer
- Run locally on single Mapper's o/p
- O/P from Combiner is sent to reducers.
- Combiner - Reducer are often identical
  - ↳ possible only if operation performed is commutative and associative
  - ↳ I/P and o/p data type for C/R must be identical
- Combiner func is an optimization function

Block1

↓  
Inputformat  
↓

Mapper

↓  
Combiner

↓  
Partitioner

Block2

↓  
Inputformat  
↓

Mapper

↓  
Combiner

↓  
Partitioner

→ shuffle/sort ←

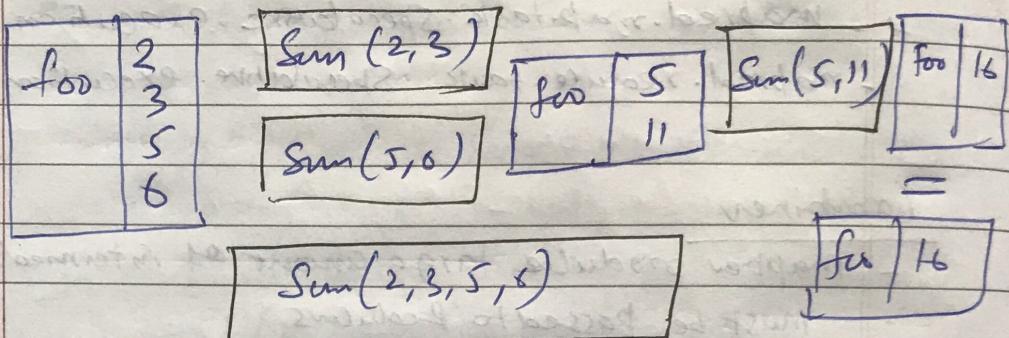
extramarks

Reducer

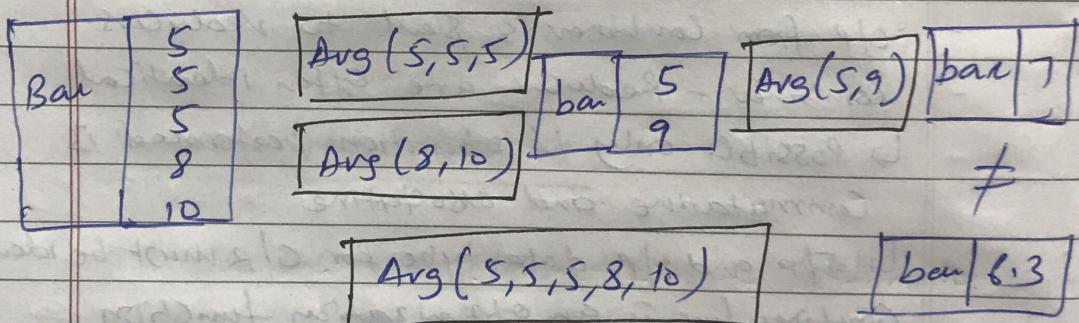
Reducer

eg., in wordcount with Combiner (KIS)

- Some Reducers may be used as Combiner if operation is associative + commutative eg. SumReducer



- Some Reducers Cannot be used as a combiner  
eg. Average Reducer



### Specify

- JobSetCombinerClass (SumReducer.class)
- Combiner may run once, more than once on the o/p from given Mapper
- Don't put code in Combiner that influence result if run more than once.
- It doesn't replace reduce function
- It help cut down the amount of data shuffle b/w mapper and reducer (Always worth)

## Configuration file (mapred-site.xml)

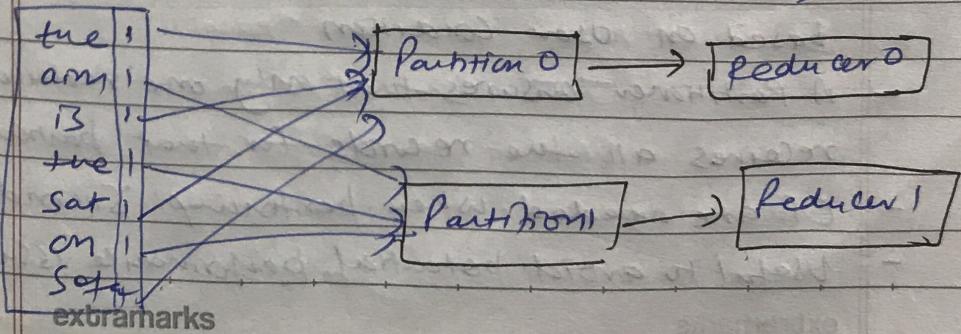
- \* Stored in /etc folder
  - \* Contains configuration info
- 1- mapreduce.framework.name - Yarn
  - 2- mapreduce.map.memory.mb - (1024) Large resource limit for mappers
  - 3- mapreduce.reduce.memory.mb (3072) Large resource limit for reducers
  - 4- mapreduce.task.io.sort.mb (512) Higher memory limit while sorting data for efficiency
  - 5- mapreduce.maps.java.opts - ( $Xmx1024m$ ) Larger heap size for child JVMs of map
  - 6- mapreduce.reduce.java.opts - ( $Xmx2560m$ ) Larger heap size for child JVM of reduce
  - 7- mapreduce.reduce.shuffle.parallelcopies (50) Parallel copies run by reducer.

## Partitioner

- Determine which reducer each intermediate key and its associated value goes to

Get Partition:

(inter-key, inter-value, num-reducers)  $\rightarrow$  Partition



- Partitioning of keys of intermediate o/p is controlled by Partitioner
- Each mapper o/p is partitioned and records having same key value go into same partition
- Then each Partition is sent to Reducer.
- Partition phase takes place after map and before reduce
- Default Partitioner
- The default partitioner is HashPartitioner
  - ↳ Use Java hashCode method
  - Guarantees all pair with same key go to same Reducer

public class HashPartitioner<K,V> extends Partitioner<K,V> {

    public int getPartition(K key, V value, int numReducers) {
 ↳ return (key.hashCode() % Integer.MAX\_VALUE) / numReducers;
 }

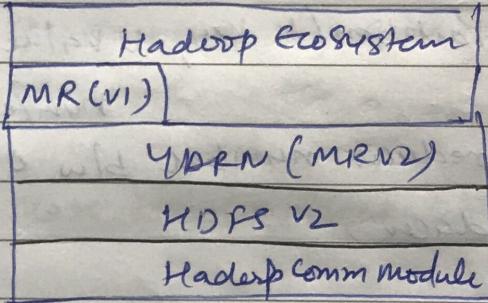
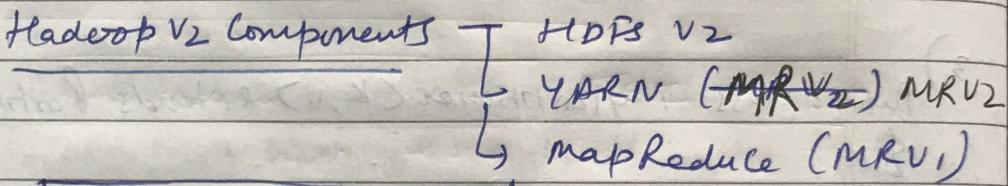
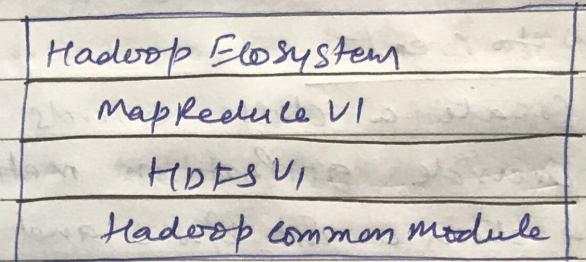
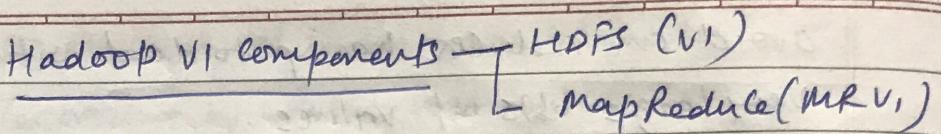
### Custom Partitioner (CP)

- Process allow to store the result in different reducers based on user condition
- A Partitioner ensures that only one reducer receives all the records for that particular key
- CP are needed when performing Secondary Sort
- Useful to avoid potential performance issues - extramarks

避免 one reducer having to deal with many very large lists of values.

### How to Create

- 1) Create a class extends Partitioner
- 2) Override getPartition method  
↳ return an int b/w 0 and 1 less than <sup>no. of</sup> Reducers  
e.g - 10 Reducers (0 - 9)
- 3) Public class MyPartitioner<K,V> extends Partitioner<K,V>  
 {
 public int getPartition(K key, V value, int numReduceTasks)
 {
 // determine reducer number b/w 0 - numReduceTasks
 return reducer; -> 290N
 }
 }
- 4) Specify C.P in your driver code  
job.setPartitionClass(MyPartitioner.class)



### Hadoop V1 limitations

- 1) Batch Processing
- 2) Not suitable for Real time Data Processing
- 3) Support 4000 nodes per cluster
- 4) More burden on JT (Scheduling, Monitoring, Client request, Job history)
- 5) SPOF  $\rightarrow$  NN - Cluster down

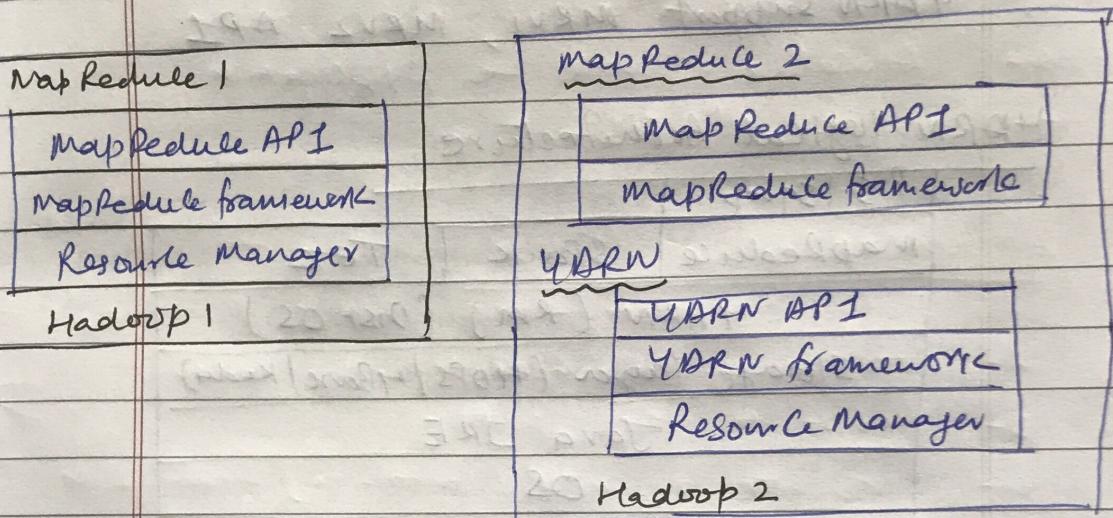
JT - Can't exec the program

SNN - Can't get the latest file image

- (6) Scalability - 4000 (node)(4K)  $\begin{cases} YARN \\ -10\text{ K Node} \\ -10\text{ K jobs} \\ -100\text{ K task} \end{cases}$
- 40,000 Tasks (40K)

extramarks

- 7) fixed slot size - utilization
- 8) Availability - rapidly changing Complex State in JT memory, make it diff to retrofit HA in JT Service
- 9) Doesn't support multi-tenancy



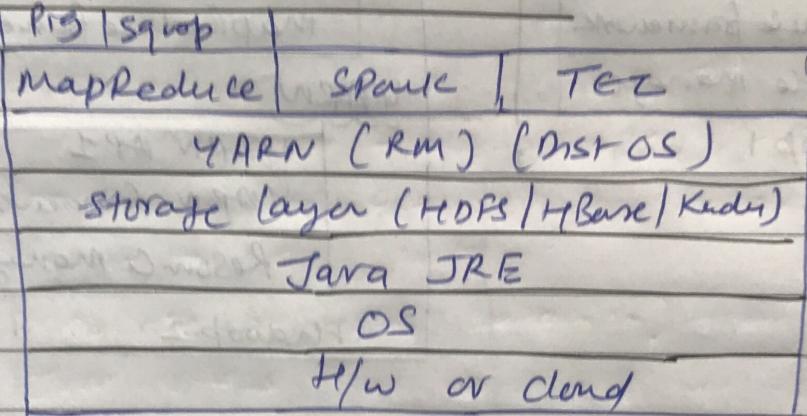
## YARN

- Yet Another Resource Negotiator
- MRV2 - Next generation MR
- Distributed operating sys.
- General framework to support other distributed coupling paradigm
- Yahoo 2010 began to increase performance by smaller memory utilization, enhance scalability, flexibility so to run many dist framework in parallel.
- Hide resource mgt details from user.

## Main Idea

- Shift Job tracker responsibility
  - 1) Resource mgt (Job scheduling)
  - 2) Diff master (Task monitoring / status / job progress)
- YARN supports MRV1, MRV2 API

## YARN logical Architecture



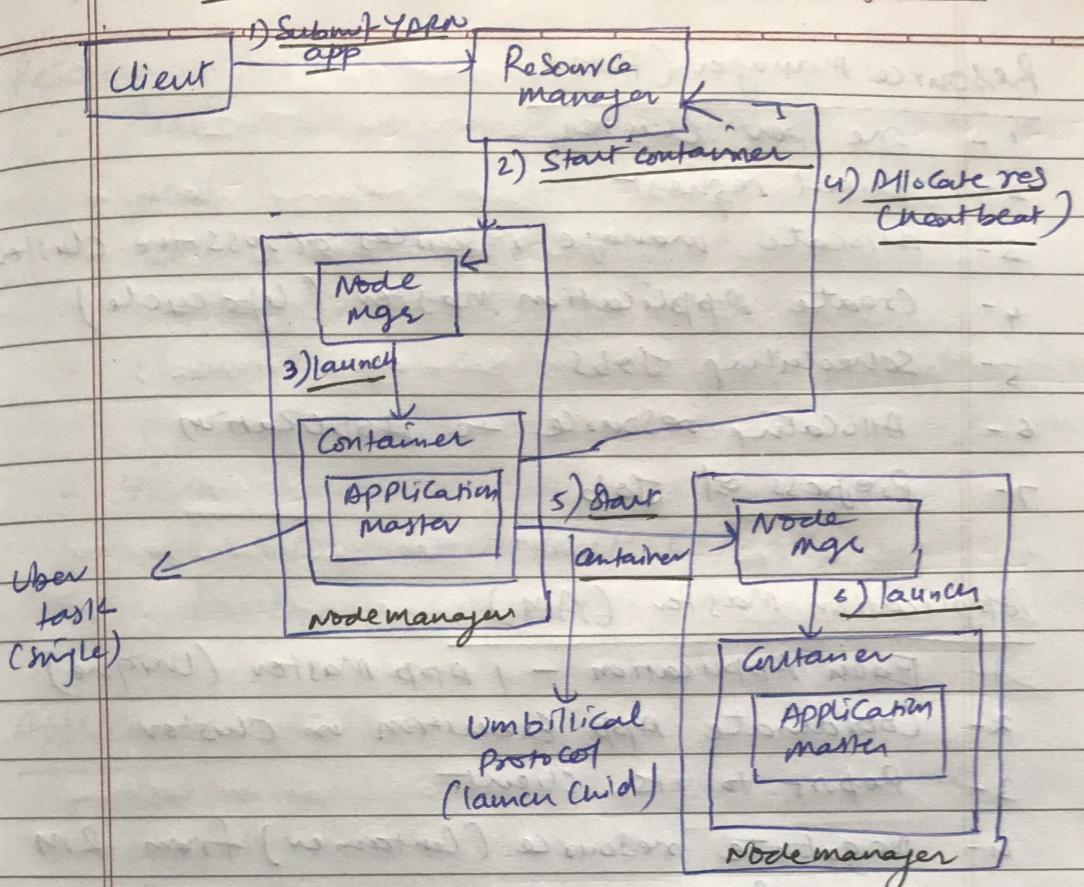
## Benefits

- 1) High Scalability (10K nodes, 100K tasks)
- 2) High Availability (HA both RM, AM)
- 3) Support multiple programming models
- 4) Multi-Tenancy (Beyond MR, even differentiation of MR)
- 5) Multiple Name space
- 6) Cluster utilization (soft slots for map and reduce)

## YARN Architecture

Date: \_\_\_\_\_

Page: 23



### Daemons

- 1) Node Manager
- 2) Resource Manager

### Node Manager (NM)

- 1 - Running on all WORKER nodes
- 2 - launch and monitor Container
- 3 - Keep RM update (Heart Beat)
- 4 - Container life cycle mgmt
- 5 - Monitor resource (CPU, memory)
- 6 - tracking node health
- 7 - Kill container direct by RM

## Resource Manager (RM)

- 1 - one per cluster
- 2 - Client request
- 3 - Allocate manage resources across the cluster
- 4 - Create Application Master (Lifecycle)
- 5 - Scheduling Jobs
- 6 - Allocating resource to Application
- 7 - Progress of Job

## Application Master (AM)

- 1 - Each Application - 1 App Master (unique)
- 2 - Coordinate App Execution in Cluster
- 3 - Report to its Client
- 4 - Negotiate resource (Container) from RM
- 5 - JVM process
- 6 - Run a computation in the container (return to client)
- 7 - <sup>or</sup> Request more container from RM to run a distributed computation

## Container

- 1 - Physical resources (RAM, CPU cores)
- 2 - Container may be Unix process or Linux cgroup
- 3 - Can have multiple containers in single node