Schedulers:-

- In hadoop, we can receive multiple jobs from differed clients to perform.

+ The Map Reduce framework is used to perform multiple tasks in parallel in a typical Hadoop cluster to provide process large size datasets at a fast rate.

→ This Map-Reduce framework is responsible for scheduling and monitoring the tasks given by different clients in a Hadoop cluster. But this method of scheduling jobs is used prior to Hadoop 2.

→ Now in Hadoop 2, we have YARN (Yet Another Resource Negotiator)

→ In Yarn we have separate Daemons for performing job scheduling, Monitoring and Resource Management as Application Master, Node Manager and Resource Manager respectively.

Resource Manager:- Resource Manager is the Master Daemon responsible for tracking or providing the resources required by any application within the cluster, a

Node Manager:- Node Manager is the Slave Daemon which monitors and keep track of the resources used by an application and sends the feedback to Resource Manager.

The Schedular in YARN is totally dedicated to Scheduling the it cannot track the Status of the application.

⇒ On basis of required resources, the Schedular performs or we can say schedule the Jobs.

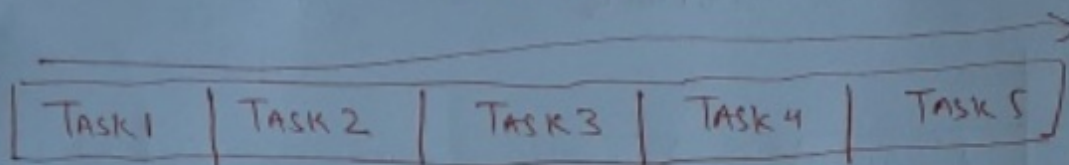## Types of SCHEDULER

FIFO SCHEDULAR.

FAIR SCHEDULER

CAPACITY SCHEDULAR

These Schedulars are actually a kind of algorithms that we use to Schedule tasks in Hadoop Clusters when we receive requests from different-different clients.

A Job queue is nothing but the collection of various tasks that we have received from our various client. The tasks are available in the queue and we need to schedule this task on the basis of our requirements.

2

## JOB QUEUE

| TASK 1 | TASK 2 | TASK 3 | TASK 4 | TASK 5 |

① FIFO SCHEDULAR:

→ As the name suggests fifo i.e. first In first Out, So the tasks or application that comes first will be Served first.

→ This is the default Schedular we use in Hadoop. Two tasks are placed in a queue and the tasks are performed in their

Submission order.

→ In this method, once the job is Scheduled, no intervention is allowed.

→ So Sometimes the high-priority process has to wait for a long time since the priority of the task does not matter in this method.
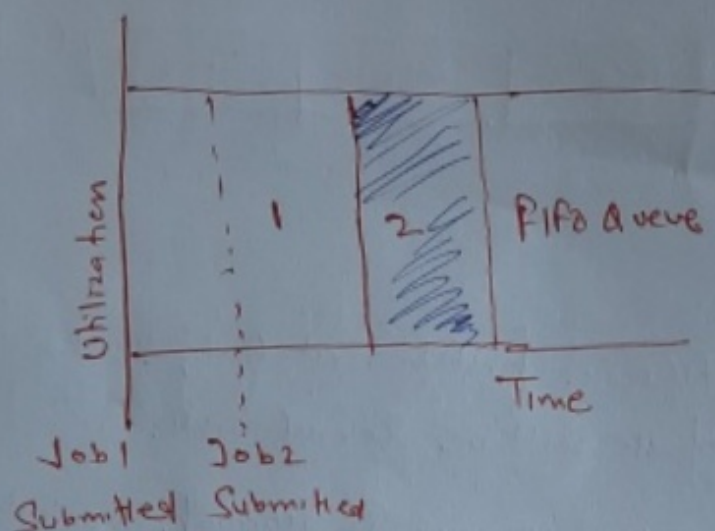
Advantage:
• No need for Configuration.
• First Come First Serve
- Simple to execute.

Disadvantage!
• Priority of task doesn't matter, So high priority jobs need to wait.

• Not suitable for shared cluster.

# fifo SCHEDUER



Time

Utilization

Job1          Job2
Submitted   Submitted

2. Capacity Scheduler (default Scheduler with YARN)

(client      Same cluster
              can be sorted
                     out)

1) In Capacity Schedular we have multiple job queues for Scheduling our tasks.

→ The Capacity Schedular allows multiple occupants to share a large Size Hadoop cluster.

→ In Capacity Scheduler corresponding for each job queue, we provide Some Slots or Cluster resources for performing job operation

→ Each job queue has its own slots to perform its task.

→ In case we have tasks to perform in only one queue then the tasks of that queue can access the Slots of other queues also as they are free to use, and when the new task enters to some other queue then jobs in running in its own slots of the cluster are replaced with its own job.

4

Capacity Scheduler also provids a level of abstrachen to know which occupant is utilizy the more Cluster resource or slots, so that the single user or application doesn't take disappropnate on unnecessary slots in the cluster.

→ The Capacity Scheduler mainly contains 3 types of the queue that are root, parent and leaf which are used to represent Cluster, organisation, or any subgroup, application submission crespectively.
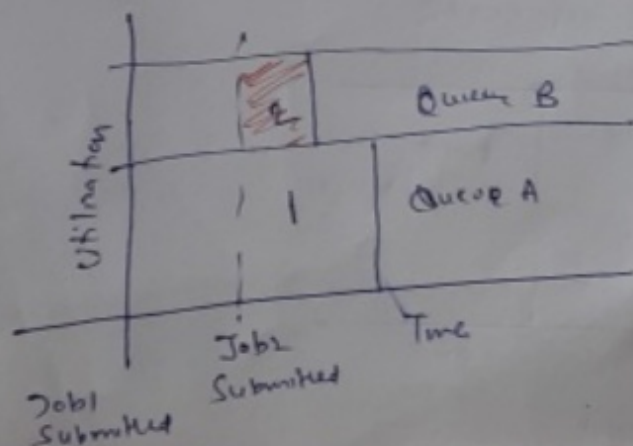
Advantge:

• Best for working with Multiple clients or many jobs in Hadoop Cluster.
• Maximizes throughput in Hadoop Cluster.

Disadvantge:

• More Complex    • Not easy to configure for everyone

CAPACITY SCHEDULER

Queue B

Queue A

Utilization

Jobs Submitted

Time

Jobl Submitted

# 3. Fair Schedular

→ The Fair Schedular is very much Similar to that of the Capacity Schedular.

→ The priority of the job is kept Consideration.

→ With the help of fair Schedular, the YARN applications can Share the resources in the large Hadoop cluster and these resources are maintained dynamically So no need for prior capacity.

→ The resources are distributed in Such a manner that all applications within a Cluster get an equal amount of time.

→ Fair Schedular takes Scheduling decisions on the basis of memory, we can configure it to work with CPU also.

→ Fair Schedular whenever any high priority job arises in the same queue, the task is processed in parallel by replacing Some portion from the already dedicated Slots.
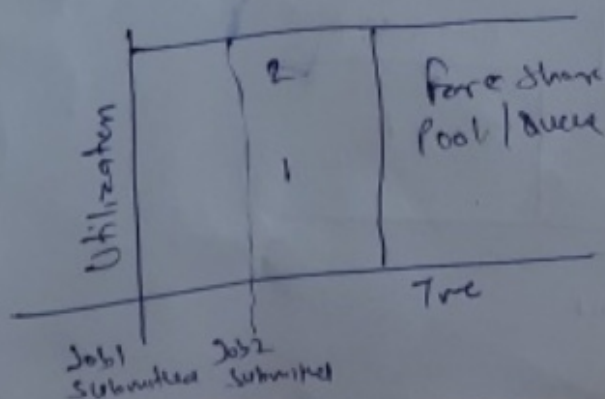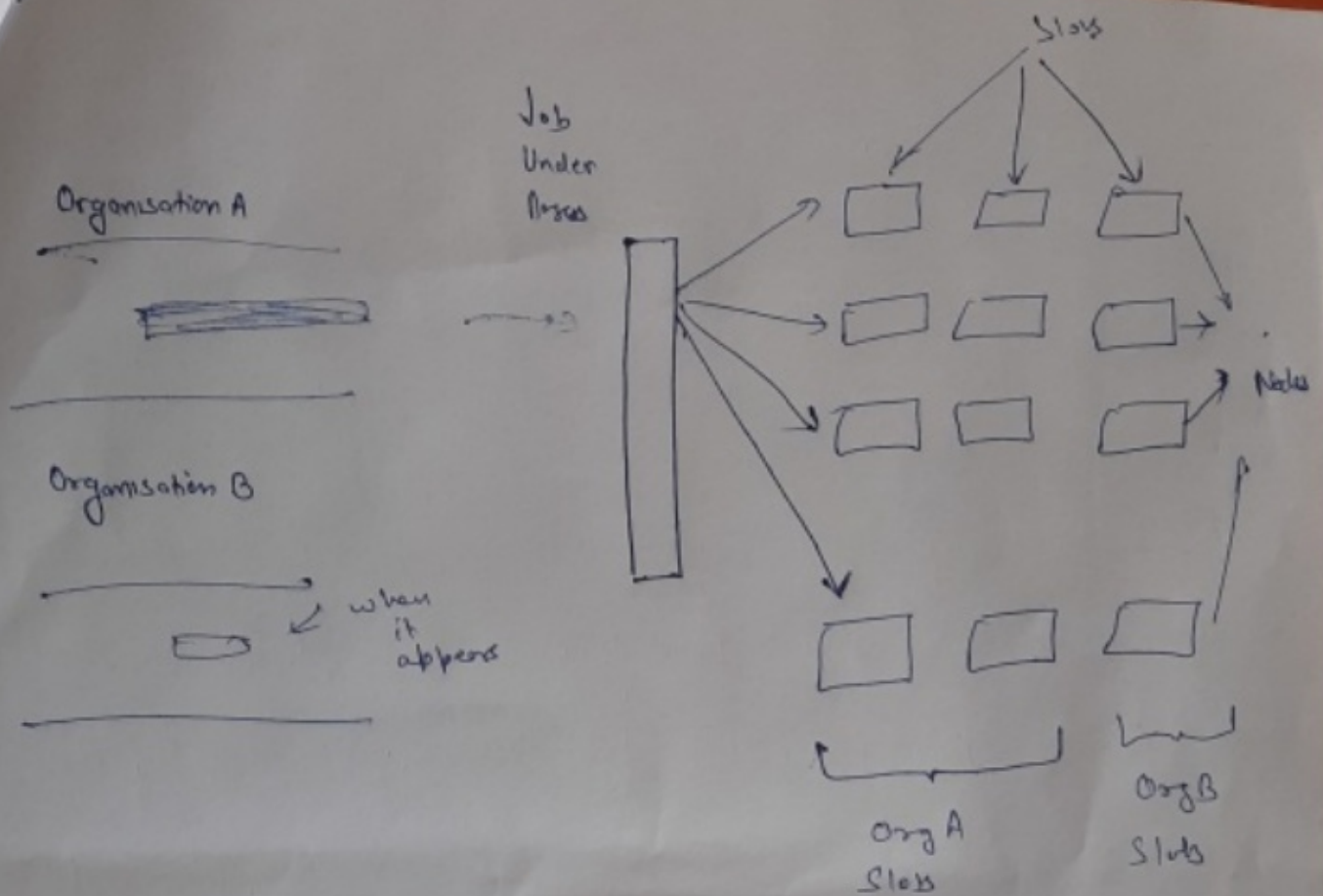
Advantages:
• Resources assigned to each application depend on its priority.
• It can limit the concurrent running task in a particular pool or queue

Disadvantage: The Configuration is required.



FAIR SCHEDULAR

Organisation A

Organisation B

when it appears

Job Under Process

Slots

Node

Org A Slots

Org B Slots

CAPACITY SCHEDULAR

FAIR SCHEDULAR (Jobs are divided into pools)

In Capacity Schedular a Small job having higher priority have to wait but in fair Schedular the job runs parallely and the resources are allocated accordingly.