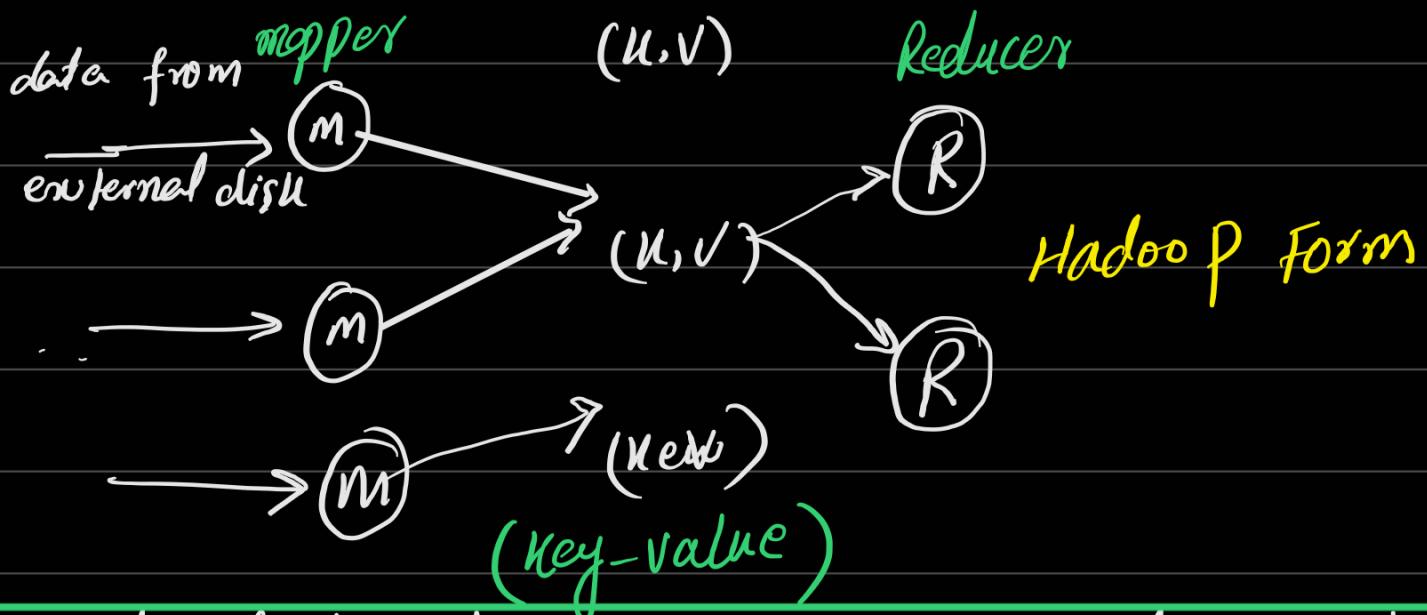


Apache Spark:-

Hadoop vs Spark:

Hadoop Drawback:-

- It was only made for batch data processing.
- slow processing because not in memory computation(Ram).



- Not fault tolerance in term of computation.
∴ fault tolerance = not affected when some node fail.
- NO mechanism for cashing and persist.



Map-Reduce

HDFS → split → map → shuffle → Reduce → output

(1) Apache Spark

It is an distributed computation framework used for large scale data processing. It does not have any separate storage or file system.

Features -

- Speed : spark is ~~100^{times}~~ times faster than hadoop, because it does in **memory processing**. It hold data in memory for processing.
- Powerful - Caching : capability to persist data in memory
- Deployment : spark can work with different **resource managers / cluster manager**.
 - 1 → YARN : resource m. design specially for hadoop.
 - 2 → Mesos

3 → Kubernetes

4 → Spark own cluster. Stand Alone

Resource Manager :- ✓

- Good fit for batch and Real time data processing.

- Polyglot :-

Spark provide high level API or libraries for different language like Python, Java, Scala, R.

Apache Spark Ecosystem ✓

Real time

Spark
SQL
②

Spark
Streaming
③

Spark
R
④

Spark
MLlib
⑤

spark
Graph
X

For ML

History

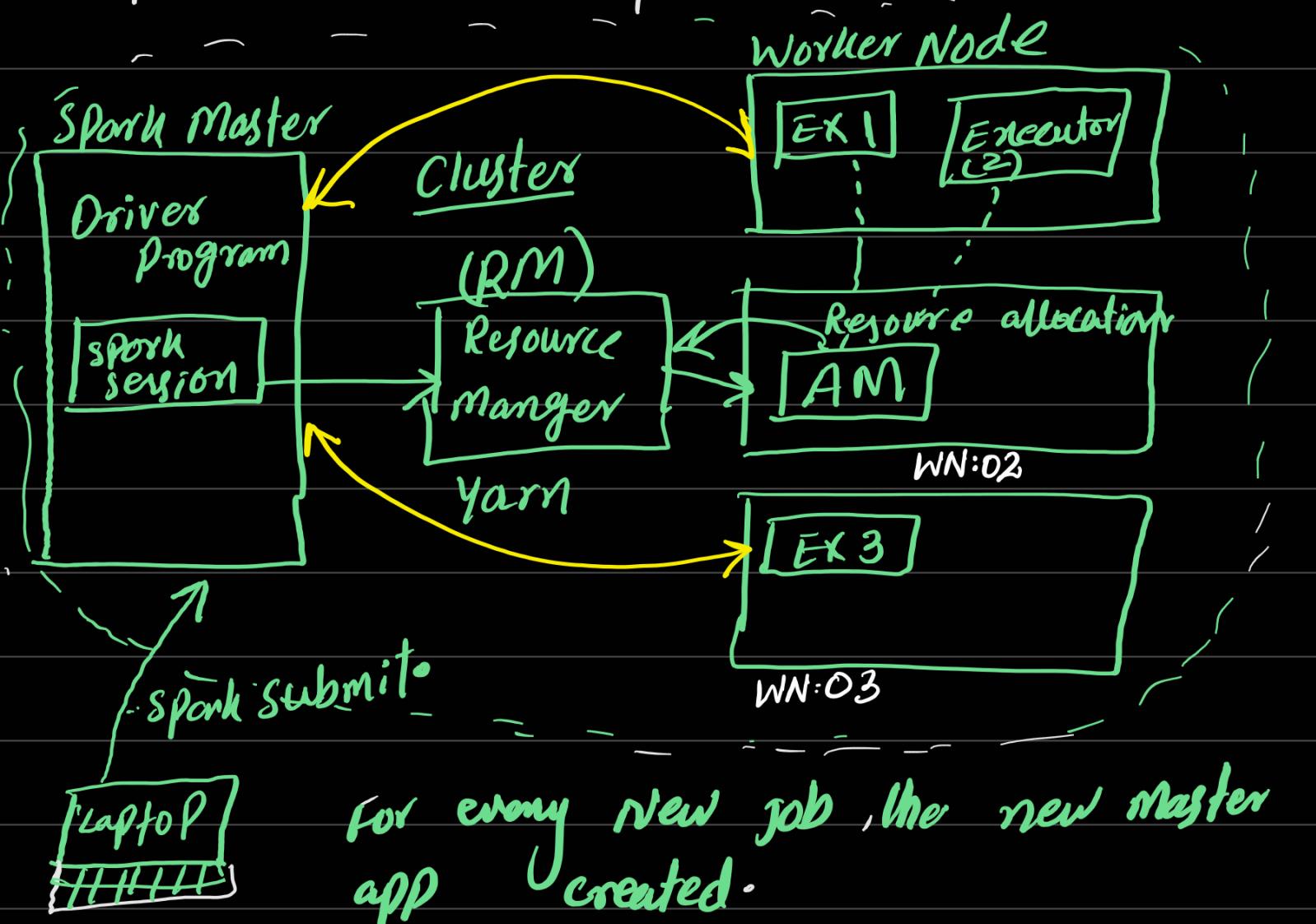
① spark core API
R, Java, Scala, Python

output
data.

Yarn :- work by managing allocation of resources (CPU, memory and disk) across the cluster / scheduling tasks.

(3) Spark Architecture :- ✓

- Spark also follows Master-Slave Architecture. — Do Parallel Processing (Data)
- Spark Master, Spark Worker Node



for every new job, the new master app created.

✓ Driver Program

- It act like a main method which will create Spark context.
The first method which will be called.

If does all management related things for job like lineage Preparation

Scheduling the tasks, data distribution.

✓ Spark Context or Spark Session :- Provide

- entry point in Spark cluster.

spark Context used working for (SQL & DataFrame)

It will connect spark application with execution cluster.

Before Spark 2.0 it was spark context

After 2.0 it is session.

Resource Manager:-

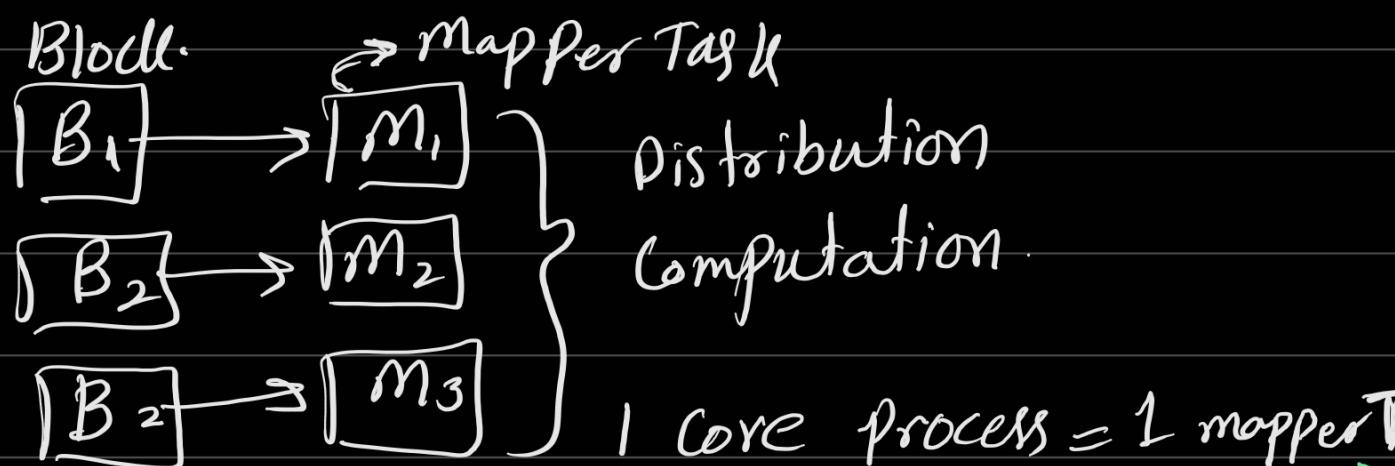
- It will create the very first container and start application master Service in it.
- Application master will request to RM for complete resource allocation or to start required numbers of executor.

Executor:- ✓

If does the Actual computation. They are like virtual containers on worker →

node which has CPU core, RAM, & Network Bandwidth.

(RM) will create required number of executors and offers them. Executor will start interacting with (AM) to send the update for job progress.



✓ (mapper task = Number of Blocks)

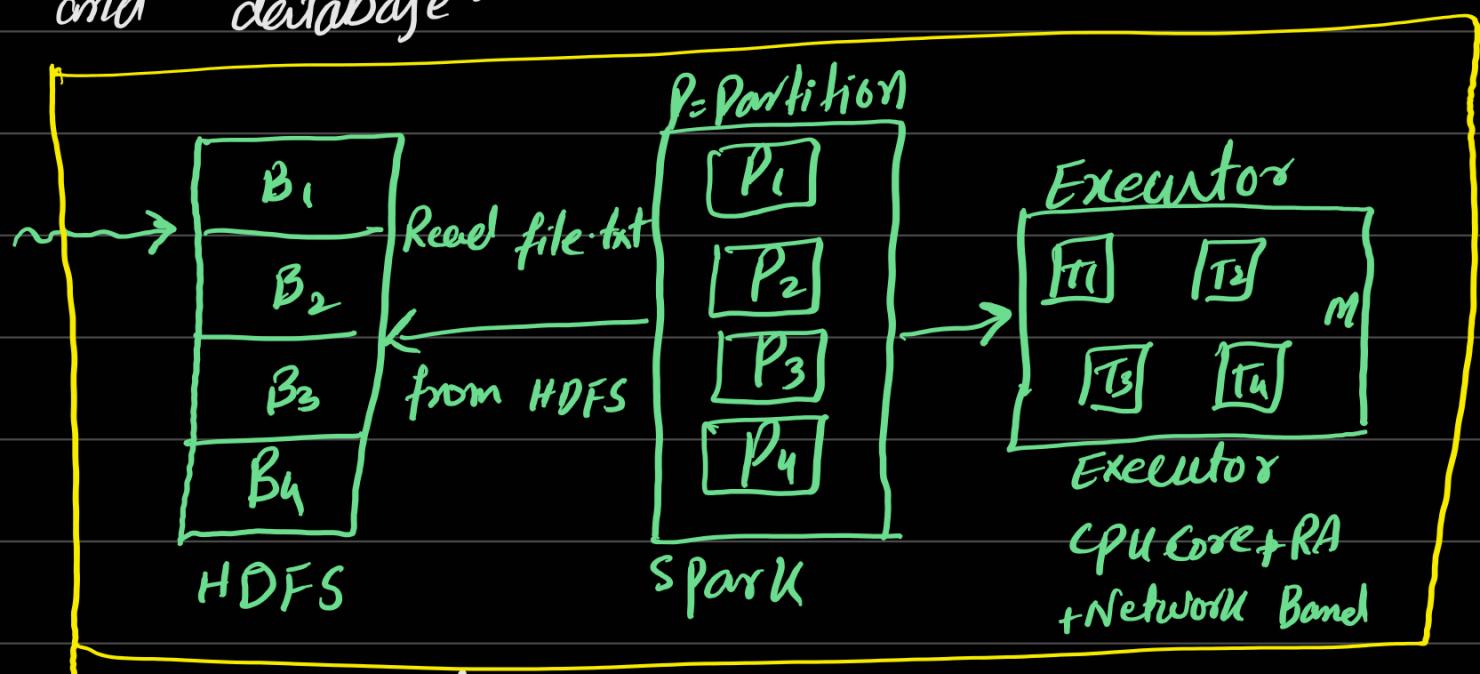
RM :-

RM will create required number of executors.

Executor will start interacting with application master to send the update for job progress

(4) Distributed computation In Spark.

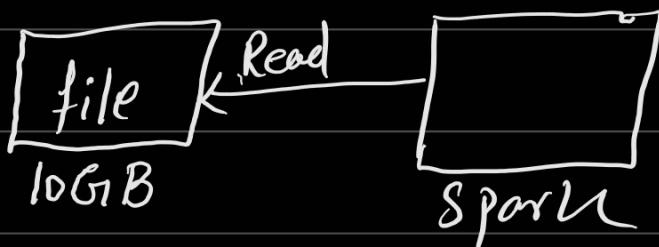
- In spark, data will be divided in number of partition. (NO block)
- How spark does parallel processing?
spark read data from distributed file system and database.



Block act like Partition in spark.

In spark we will have control to decide number of executor, its CPU core.

(5) Partition Explicitly in Spark. ✓



Task :-

- Smallest unit of executor is known as task.
- Each task will process one data partition.
- One task will get executed on 1 CPU core.
- One executor can run many tasks and equal to number of CPU core.
Parameter in Spark Submit Command.
 - Number of executor
 - driver memory (main program to get memory for executor)
 - executor memory
 - Executor core.

(6) Spark with Yarn as RM.

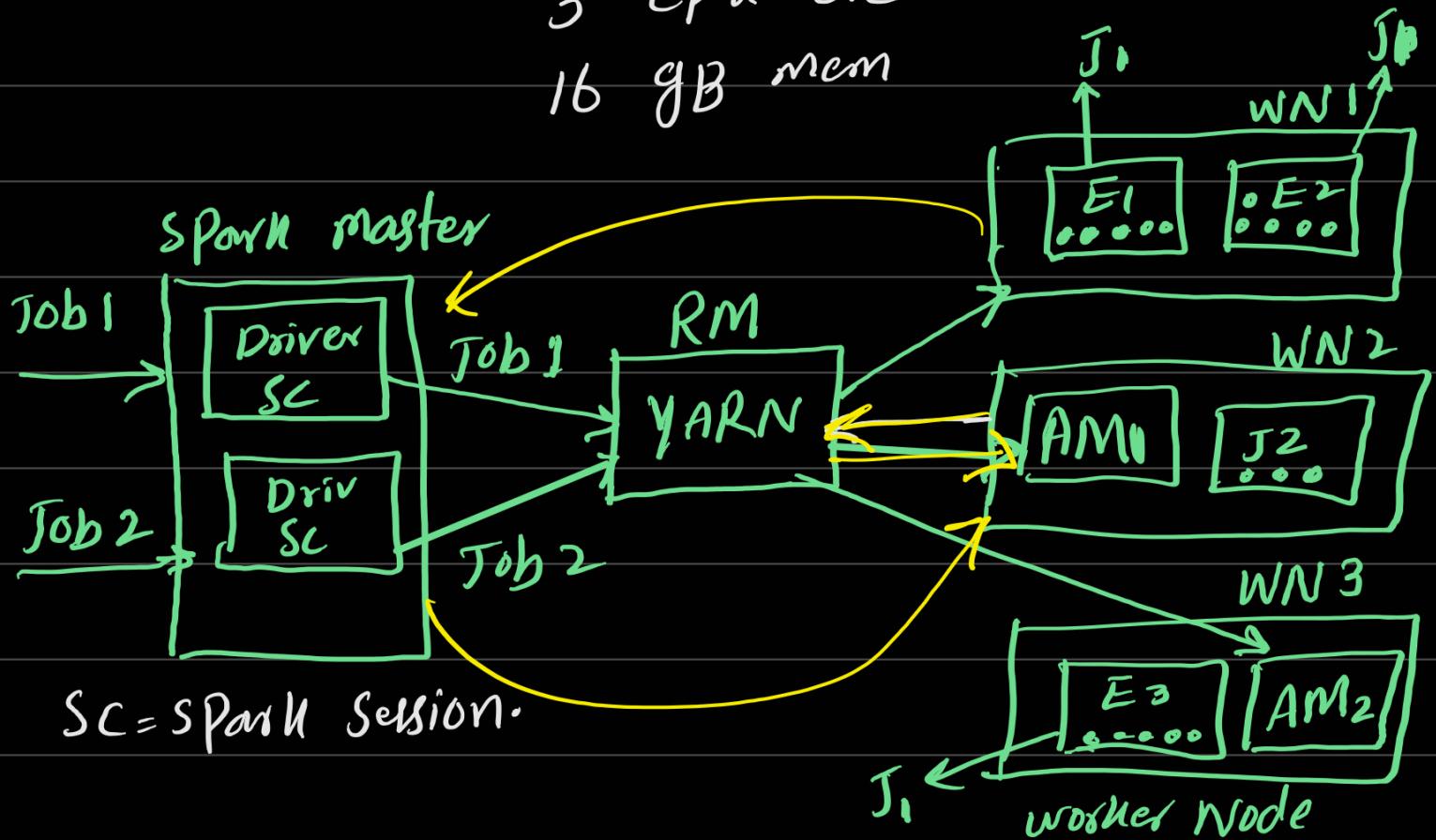
Ex Worker-node \Rightarrow 256GB · 32 CPU core
we want to execute two jobs (Job1, Job2)

Job 1 \Rightarrow 3 executor

5 CPU

Hassan Riaz | Data engineer | MS Data-Science

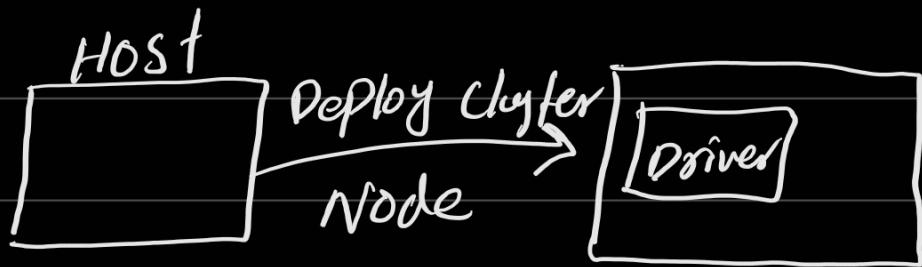
Job 2 => 1 executor
3 CPU core
16 GB mem



- AM request to Yarn, to allocate resources like cores, executor.
- Worker node ^{create} based on cluster capacity, (How much) they need to process in company.

(7) Deployment mode in SPARK ✓

- client mode
- cluster mode



In client when you shutdown the machine task are ended } Job canceled.

In cluster whe you shutdown the machine job running in background and executed.

[Part = 2]

(8)

RDD in spark

RDD → Resilient distributed dataset.

RDD are the main logical data units in the spark. They are distributed collection of objects which stored in memory or disks of different machines on cluster.

Feature of RDD (Data structure)

RDD track data lineage information to recover from faild automatically.

If also known as fault Tolerance.

Ex:-

$\text{RDD} = \text{SparkU.read_CSV("EX.CSV")}$

$\text{RDD2} = \text{RDD.filter("Salary > 1000")}$



if N_2 fail (RDD) automatically check the failed rdd.

• Distributed:-

data present in a RDD resides on multiple nodes.

• Lazy Evaluation.

data does not get loaded in RDD even if you define it. Do Action

• Immutability:

Data stored in an RDD is in the read only mode. we cannot change the value during run time.

$\text{RDD} \xrightarrow{\text{Transform}} \text{RDD1} \xrightarrow{\text{Transform}} \text{RDD2}$

• In memory computation:-

An RDD stores any intermediate data

that is generated in the RAM so that faster access can be provided.

✓ RDD vs Dataframe: ✓

RDD : It can have any type of data (unstructured, semi and structured)

RDD 1 = spark.read("data.txt")

df 1 = spark.read("data.csv")

Dataframe : allow users to create and manipulate data in tabular format, excel, sql

(9) Action and Transformation in SPARK. ✓

Transformation:-

it is a function that produces a new RDD from existing RDD.

It takes RDD as input & generate new RDD as output.

✓ Two types:-

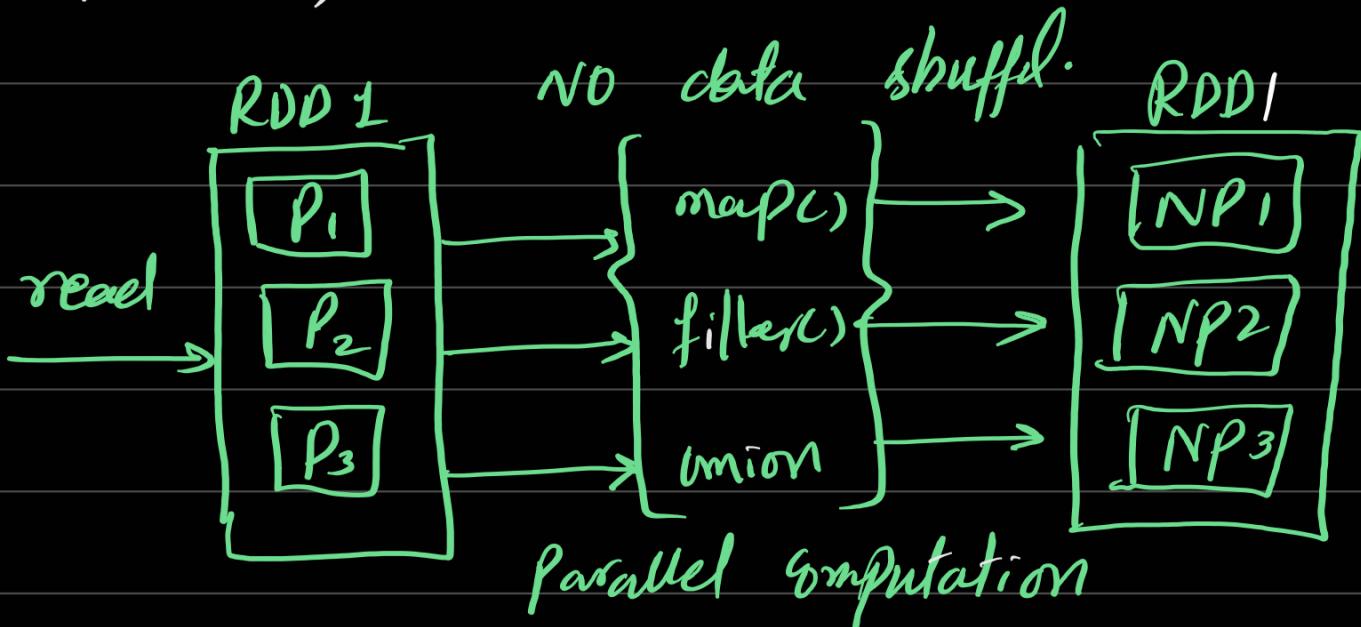
Narrow Transformation (no data shuffle)

Wide Transformation (data shuffling)

Narrow Transformation:-

$$X \xrightarrow{(A)} \left[f(x) \right] \xrightarrow{(A_1)} Y$$

map(), flat map(), union(), filter(), select.



$$a = [1, 2, 3, 4] \quad \text{result} = [1, 4, 9, 16]$$

fun = lambda x: x**2

Wider Transformation: Shuffle:

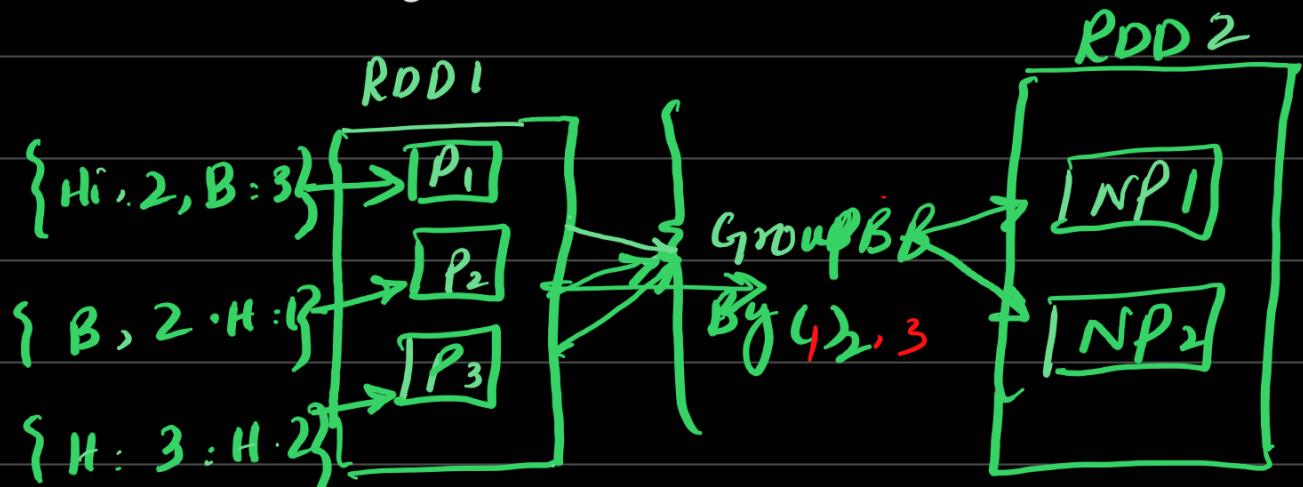
In wide Transformation all the element that are required to compute the records in a single partition, may live in many partition. That's why we need to **shuffle** the data.

Join(), groupby(), repartition()

ReduceByKey(), groupByKey()

get data from many executor/partition.

Example of word count:-



+ Same type of data may be stored in other partition like: Numbers, Int, float, char.

(a)
+ Action In Spark:-

Action is a way of sending data from executor to the driver.

Executor → `collect()`, `take()`, `top()`

Driver is manage the thing.

+ Demo PySpark application (first.py)

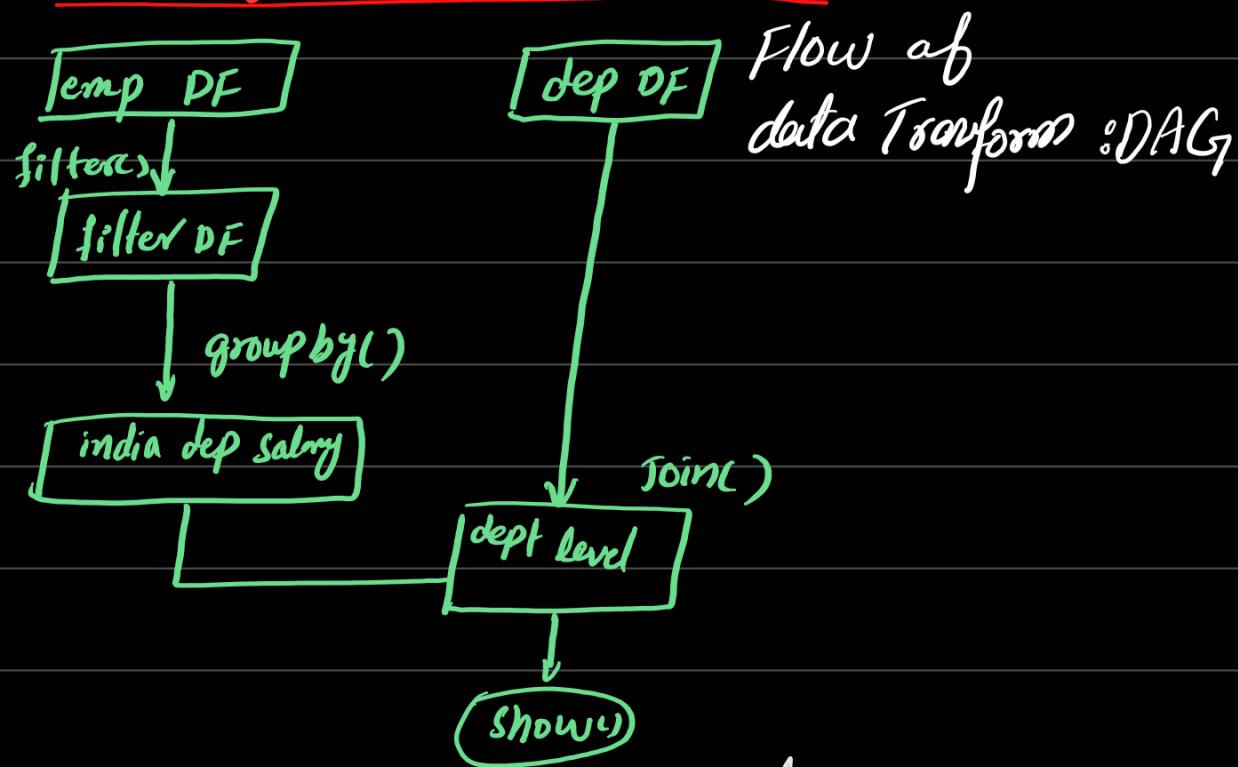
Spark = `SparkSession.getOrCreate()`

```
empDF = spark.read_csv("empby.csv")
```

```
deptDF = spark.read_csv("dep.csv")
```

```
filterDF = empDF.filter("Country = \"INDIA\"")
```

+ Lineage graph(DAG) :-



Two time data shuffling happen
groupby and join.

- DAG created for each job. Job created when action run.

when we do Transformation it optimize
the code in Lazy Transformation.



Phases spark-SQL engine:-

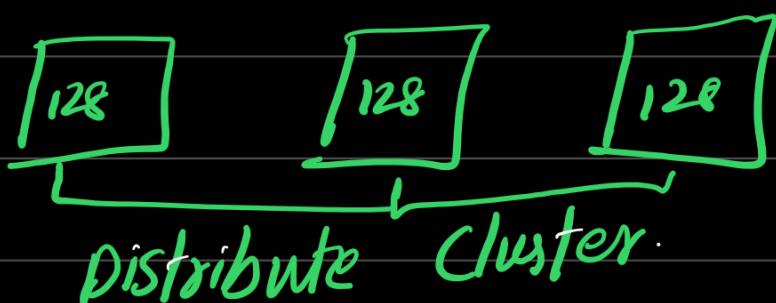
- ① Analysis
- ② Physical
- ③ logical
- ④ code generate.

Catalog:- meta-data of our data:
 { size of file, Name }
 { when upload. }
 all information about original file.

RDD :-

Resilient distributed dataset.

RDD is type of Data structure. It work on our cluster. Suppose = $\frac{500}{128}$ = 4-partition.

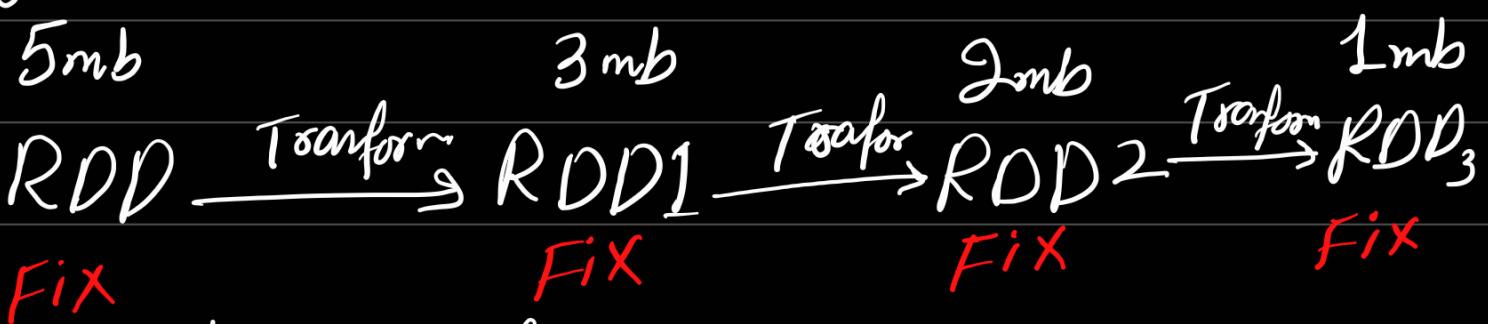


(Resilient) is fault-tolerant.

RDD immutable:-

we cannot change RDD if it create.

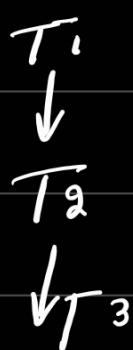
Its have only need operate not right one:-



we create copy form 1RDD to 2nd RDD.

==>

If creat DAG.



• RDD : No optimize

if ~~RDD1~~ Fail then it creat again creat.

• benefit → work with unstructured data.

when to use RDD? when we need full control our data.



Spark Memory Management.

Container [Executor]



Container (Executor)

- CPU , core
- Ram
- Network Bandwidth

Spark - Submit command.

Spark - submit

```
--master < master-url >  
--deploy-mode < cluster or client >  
--config < key : value >  
--driver-memory  
--executor-memory  
--num-executors (N)
```

we have 1 node with 32 GB Ram
and each executor can have 4GB
memory?

$$\text{Total executor} = \frac{32}{4} = 8 \text{ executor}$$

+ Reserved memory:-

This is storage for running executors.
it will store the data of internal
processing which is required to
execute the container.

fixed 300mb

+ User memory:-

It stores the data structure, meta data
and user defined data structure (struct)

+ Storage memory:

It stores cash data, data of broad
cost variable.

+ Execution memory:

It is a storage for the data
required for executions, shuffle-join

aggregation etc.

Example of memory distribution.

Let say we have given 4GB memory to each executor.

$$\text{Executor-memory (Java heap)} = 4\text{GB}$$
$$4\text{GB} = 4096 \text{ MB}$$

Two parameters impact on memory component.

- $\text{spark.memory.fraction} = 75\% (0.75)$
- $\text{spark.memory.storageFraction} = 50\% (0.50)$

Calculation:

+ Reserved memory:-

It is fixed 300MB.

+ User memory:-

$$\text{formula} = (\text{Java heap} - \text{Reserved memory}) \\ \times (1 - \text{spark.memory.fraction})$$

$$= (4096 - 300) \times (1 - 0.75) = 949 \text{ MB}$$

+ Storage memory:-

$$\text{formula} = (\text{Java heap} - \text{Reserved mem}) \times (\text{spark.mem.front} \\ \times \text{spark.mem.Stor.fraction})$$

$$= (4096 - 300) \times (0.75 \times 0.50) = 1423.5 \text{ MB}$$

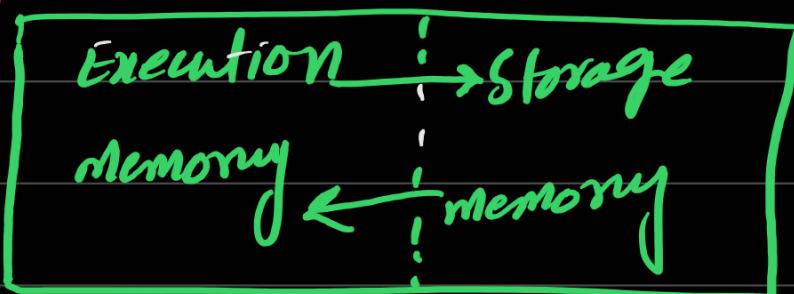
Execution Mem :-

$$\text{formula} = (\text{Java heap} - \text{Reserved mem}) * (\text{spark.mem.frac}) \\ (1 - \text{spark.m.storageFraction}) \\ = (4096 - 300)(0.75)(1 - 0.25) = 1423.5 \text{ MB}$$

Driver memory error
Execution memory error

= 1.4 GB

(ii) Rules for execution and storage memory.



- storage mem can borrow space from execution mem only if memory blocks of execution mem are not in use.
- Execution memory can also borrow space from storage memory if blocks are not in use.
- If block from execution memory is used by storage memory and execution memory need more memory; it can forcefully evict the excess blocks occupied by storage memory.

If block from storage memory is used by execution memory and storage memory needs more memory.

If cannot forcefully evict the excess block occupied by execution memory.

It will end up having less memory

If will wait until spark releases the excess block stored by execution memory and then occupies.

(12) EX Resource Allocation for Spark app:-

Case 1: Hardware Config \rightarrow 6 Nodes in cluster

each node is having 16 cores & 64GB RAM.

* From Recommendation 5 CPU core for executor.

1 GB & 1 core of each node will be required by for operating system. X

15 cores & 63GB

How many executor in our cluster?

$\therefore [5 \text{ cores} = 1 \text{ executor}] ?$

$$15 \text{ cores} = 15/5 = 3$$

$[1 \text{ Node} = 3 \text{ executor}]$

Total executor on cluster = $6 \times 3 = [18 \text{ executor}]$

1 executor will be created by YARN as AM.
so final executors which we should mention in spark submit.

$$\text{num-execut} [= 18 - 1 = 17]$$

Memory for each executor?

$$1 \text{ Node} = 3 \text{ executor}$$

$$1 \text{ Node} = 63 \text{ GB}$$

$$1 \text{ executor} = 63/3 = 21 \text{ GB}$$

spark memory overhead = 10% of exe. memory.

$$= \max(384 \text{ MB}, 0.07 * \text{executor-mem})$$

$$\text{overhead memory} = \max(384 \text{ MB}, 0.07 * 21)$$

$$= \max(384 \text{ MB}, 1.47 \text{ GB})$$

$$[= 1.47 \text{ GB}]$$

Final executor mem = 21 GB - 1.47 GB

memory for each executor is = 19 GB

find values of spark-submit parameter based on (Case - 1) & resource.

5 CPU core per executor
17 executor
19 GB for each executor

(13) How to process 1 TB of data in SPARK.

Here 1TB data store in HDFS

in HDFS the block size is 128 MB.

$$1 \text{ TB} = (1 * 1024 * 1024) / 128 = 8192 \text{ block.}$$

let us take :-

\Rightarrow 20 executor - each executor have 5 cores
 \Rightarrow 5 cores - in each executor
 \Rightarrow 6 GB ram in each executor.

In parallel $20 * 5 = 100$ task will be running.

100 task = 100 data block.

$$100 * 128 = 12800 \text{ MB} = \frac{12800}{1024} = 12.5 \text{ GB}$$

+ 12.5 GB data will get process in first batch.

RAM

$$\Rightarrow 6 \times 20 \text{ exec} = 120 \text{ GB RAM}$$

So at the time 120 GB of Ram will be occupied in cluster.

20 Node = 1.6 GB RAM from each 120 GB executor will be used.

Now Available RAM in each executor ($6 \text{ GB} - 1.6 \text{ GB} = 4.4 \text{ GB}$) will be free to use by other user job.

$$1 \text{ TB} = \frac{1024 \text{ GB}}{120 \text{ GB}} = 85 \text{ Batches.}$$

whole data will process in 85 batches.

Failure Scenarios in Spark:

- ① out of memory error.

- avoid unnecessary data shuffle
- check Broadcast big in size

(2) Network Error

If happen when spark driver & executor are not able to communicate with each other.

Solution = { `spark.network.timeout` }
 { `spark.executor.heartbeatInterval` }

(3) Data Input / Output:

- Make sure data is consistent as per data contract.
- Always put read/write related operation in Try - except block for easy debugging.

(4) Serialization Error:-

Data is not according to the serialized schema.

(5) Disk space Error:

spark job spills intermediate data

into the local machine Disk if no mem
then job fil. Ex (spark.local.dir)

Heap memory :- Part of memory which allocate to (JVM) to store object. JVM used heap memory when spark app owning to store object, data (RDD). Driver/executor memory

(6) Memory configuration.

If we are facing heap space related issue and memory overhead issues then try to configure mem component.

- ⇒ Analyse the data volume
- ⇒ Define driver & executor memory properly
- ⇒ Adjust memory configuration parameter properly.

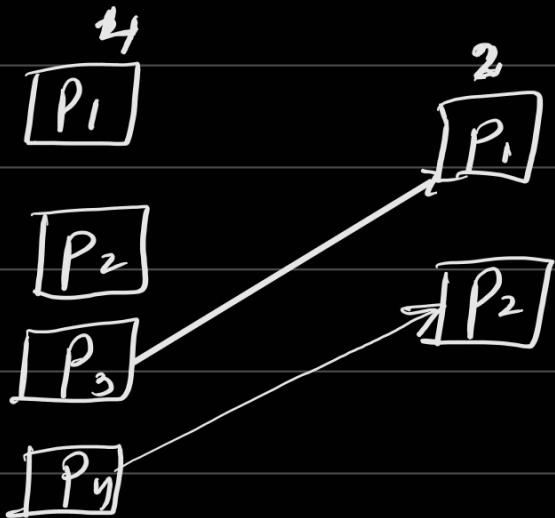
Frequently used optimization in SPARK-

(1) Partition:-

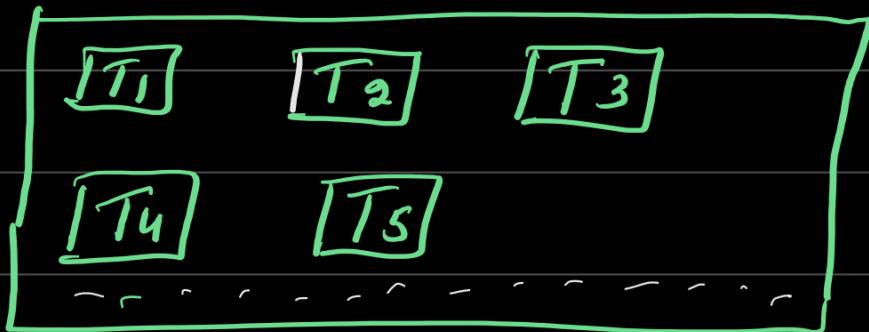
- repartition() increase
- coalesce decrease

emp DF = spark.read.csv("emp.csv")
empDF.partition(10)

JoinDF = empDF.join(deptDF)



spark web UI



(2) optimization in the case of skewed
we can use key-sorting technique.

(3) Broadcast variables:-

If will help to broadcast cost dataset
on multiple executor.

(4) persist() and caching():-

cache() only in memory save

`Persist()` → memory + Disk both.

(5) Try to minimize number of stages in the spark job. That means unoptimized use of wide Transformation

Difference between Job vs Stage is
Task (interview)

Job:-

When Action hit it creates job.
 $N \cdot \text{Action} = N \cdot \text{Job}$.

Stage: wherever data shuffling takes place in the spark app. a new stage gets created.

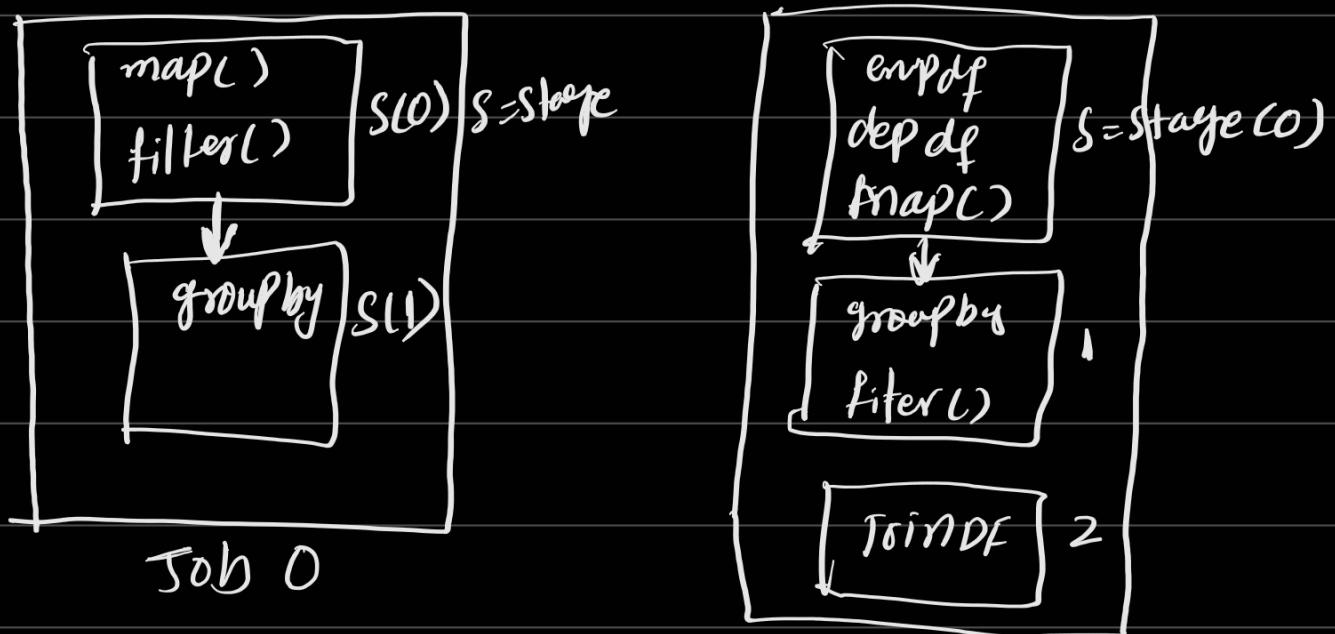
Stage 1 → Stage 2 → Stage 3

Task: small unit of execution. Called Task
emp DF,
dept DF,

Ex: emp DF • filter('INDIA')

Ex: emp DF • groupBy('dept_id') • sum(salary)

Transaction } Join DF = emp DF . JOIN (dept DF)
 { Join DF . filter (dept DF . dept_id not null)
 Join DF . collect (10) } Action



number of action to number ^{Job 1} of job.
and shuffling occurs.

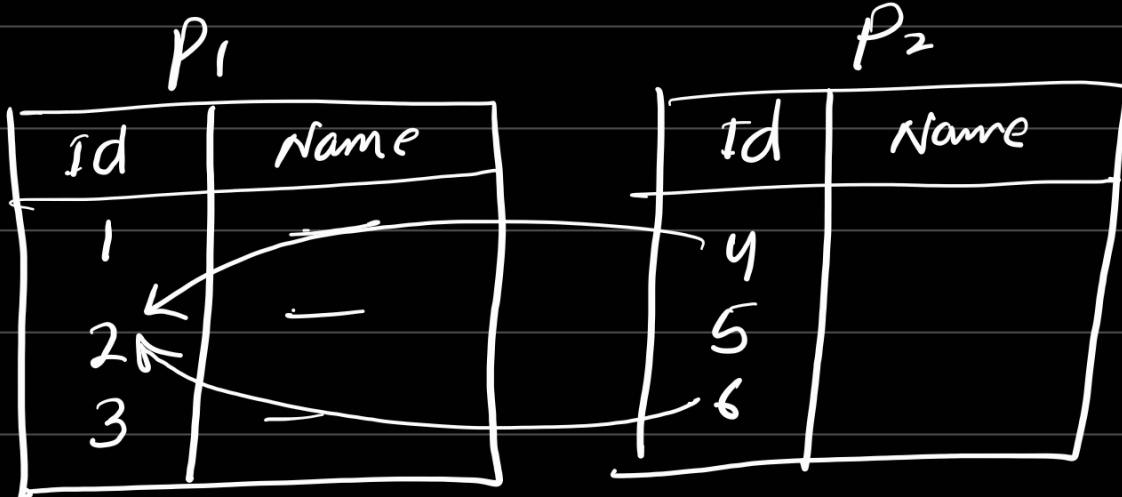
1 Action = 1 Job.



How many job created for 100 GB
of data.

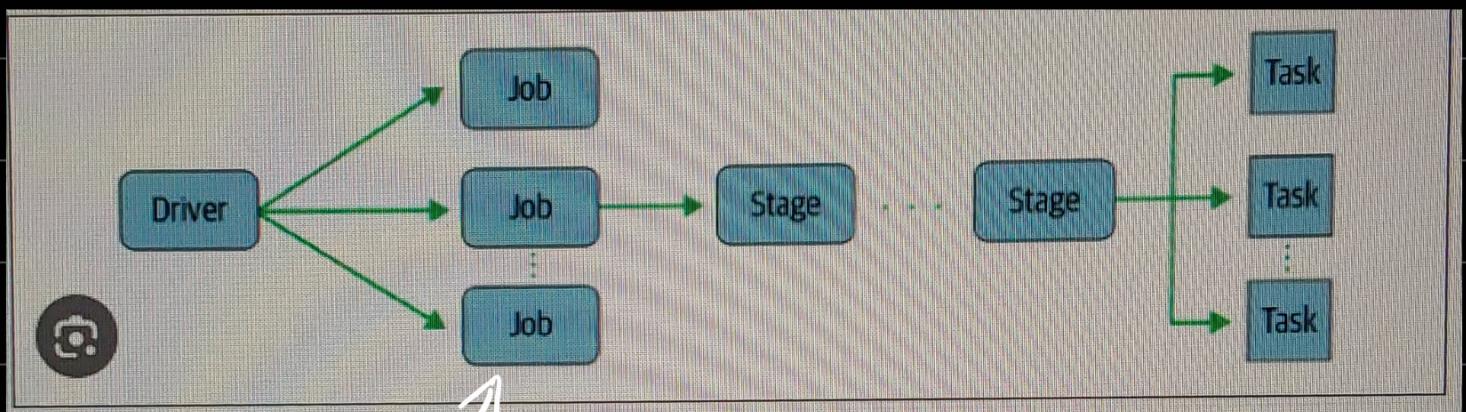
Jobs based on type and transformation
of data.

each job have many stage , each stage have many task



Print where Id is even : (Data Shuffel)
(Point) :-

Exact number of jobs and stages based
on spark app, data and transformation.



Action Hit

Hassan Rafiq || Data engineer || MS Data-Science

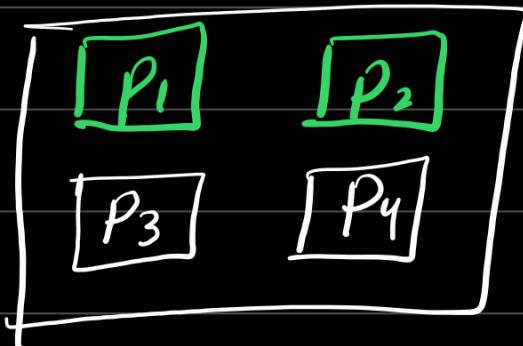
Repartition vs Coalesce :

- = Coalesce used to decrease partition.
- = repartition used to increase or decrease partition.

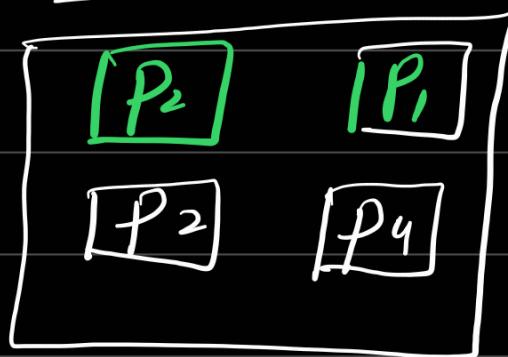
Spark Join :

shuffle data between worker-node

Executor - 1



Executor - 02



Strategies in Spark:

- shuffle Sort Join
- shuffle hash Join
- Broadcast hash Join
- Co-partitioned Join
- Nested Join

Broad Cost Join in Spark :-

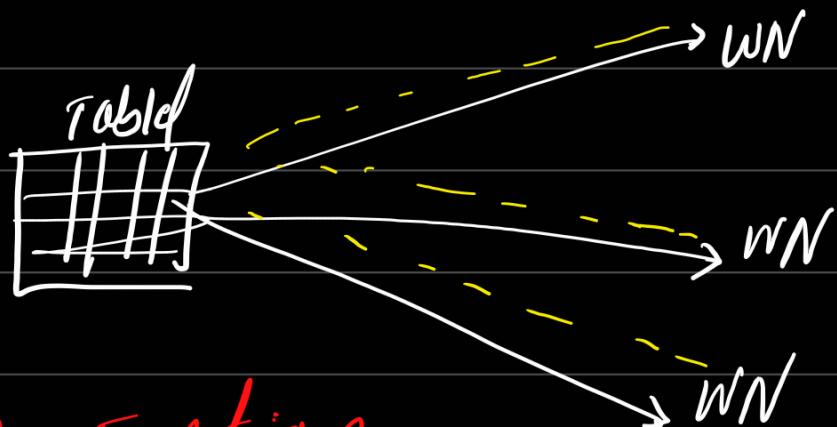
Broad Cost hash Join:-

Broad cost data, during join which **eliminate shuffle**.

If increase the speed of data.

* Broad Cost join sends copy of small table to all worker node.

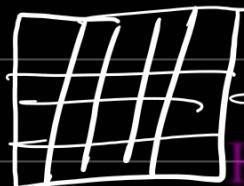
Ensure that all worker node have memory.



Hash Function:-

• meta data of real data.

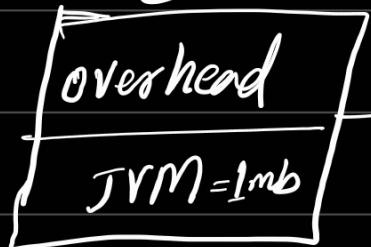
real data store in hash data.



Execution memory error:-

- out-of-memory error:-

Driver



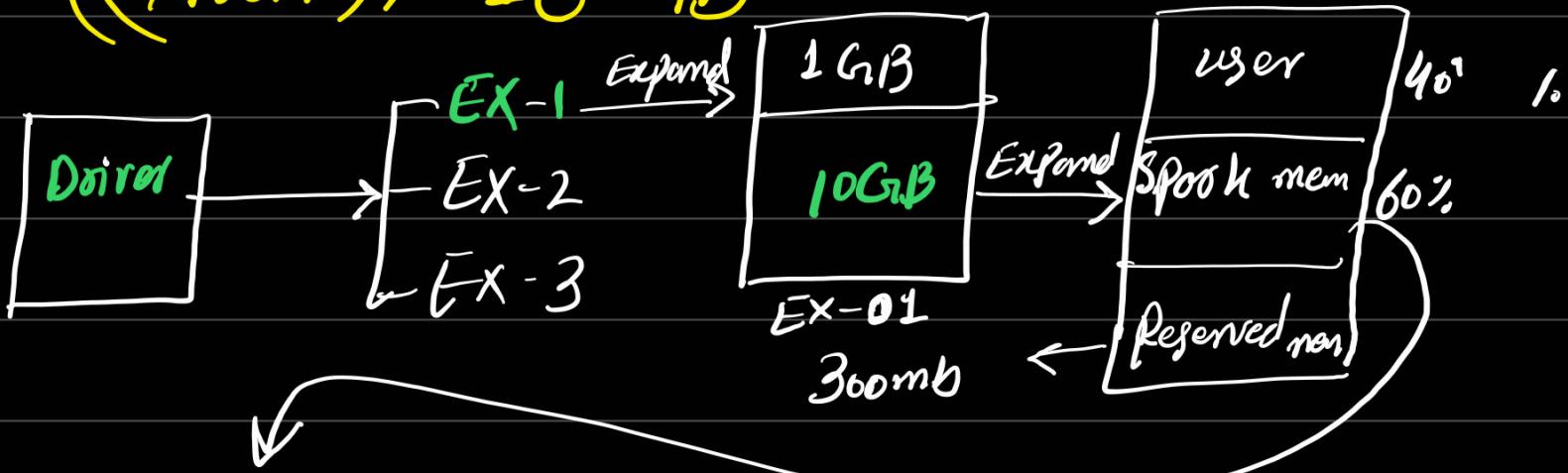
→ Executor-1

→ E2

→ E3

- optimize memory.
 - Caching
 - Partition correctly
 - monitoring
 - data Filter
- } need memory.
} Avoid of error.

((Arch)) 10 GB Data



Spark memory

Storage → 50%

Execution → 50%

Unified memory Management.

4-core , 4-task.

Spark-submit