

## Difference bet'n Apache Spark & Apache Hadoop

### Apache Hadoop

- ① No built-in iterative mode
- ② For data processing Hadoop uses batch engine
- ③
- ④ Doesn't have any abstractions.
- ⑤ Needs Higher storage.

### Apache Spark

- ① Spark is easy to program
- ② has iterative built-in mode
- ③ Uses streaming, machine learning & batch processing.
- ④ Spark completes job processing jobs about 10 to 100 times faster.
- ⑤ Uses Abstraction called RDD, which makes spark feature rich.
- ⑥ Uses low latency.
- ⑦ needs Higher RAM

## Spark Architecture.

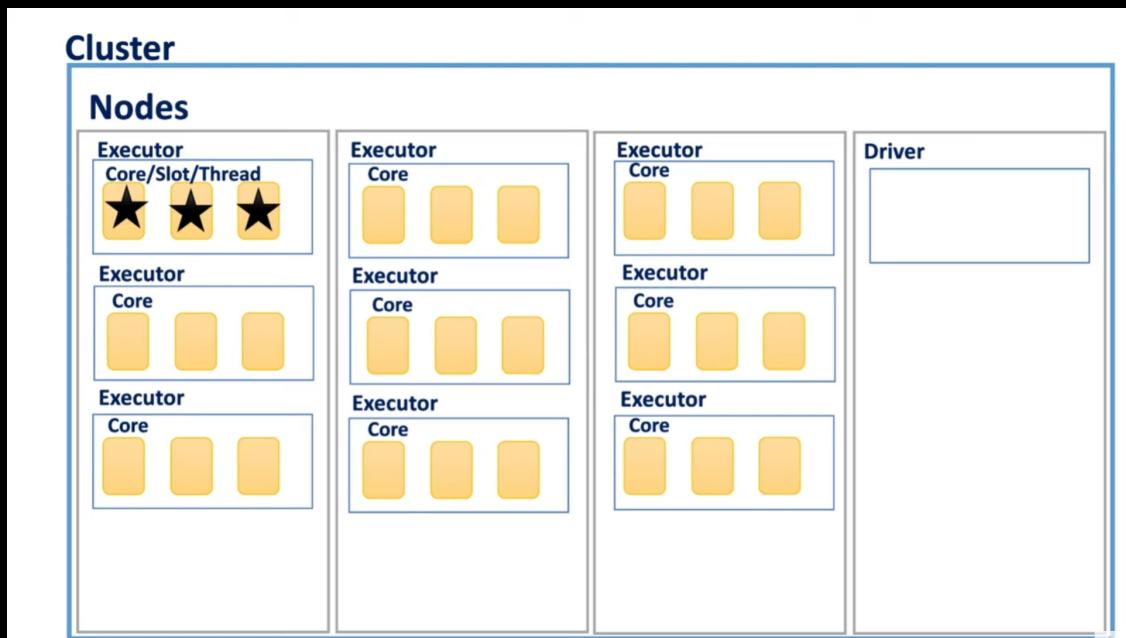
**Driver** - Is responsible to give instructions runs on his own JVM.

**Executors** - Who do execution.

**Dataset** - chunks of data.

**job** - Number of stages

**cluster** - collection of one or more node.



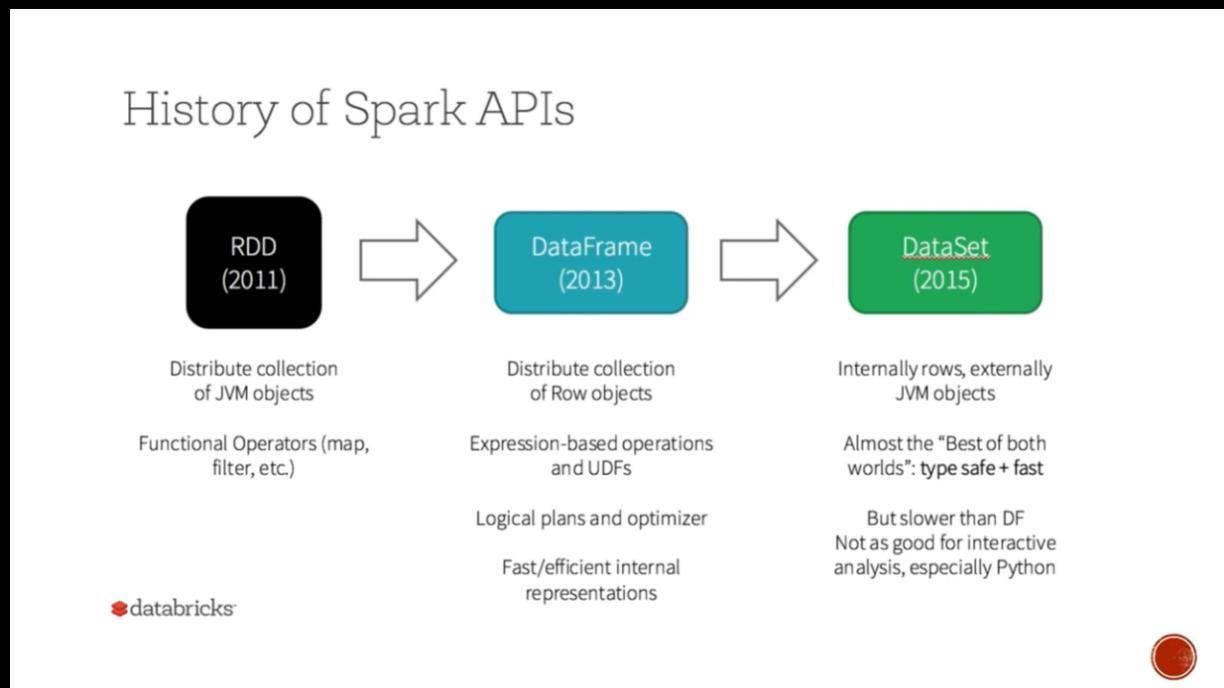
Spark API :-

three types of API

RDD  
(2011)

DataFrame  
(2013)

DataSet  
(2015)



RDD -

Resilient distributed dataset

# RDD

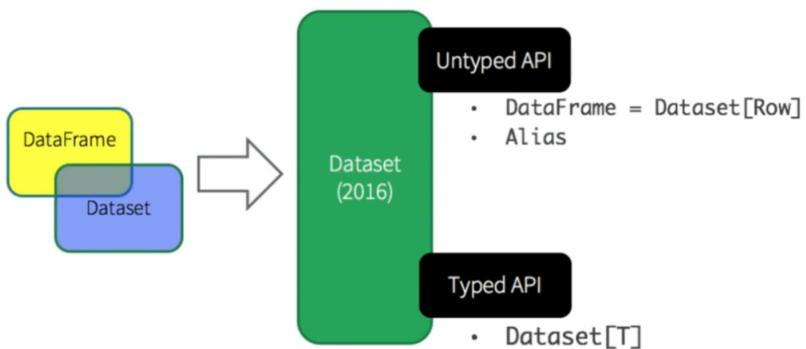
- Resilient distributed dataset
- Fundamental abstraction on which spark was built
- Low level API
- Offers control and flexibility – ‘How’

Resilient  
Immutable  
Compile time type safe  
Lazy Evaluation – materialized in a lazy manner after a action happens  
Doesn't infer schema  
When we do transformation, we get one RDD from another like a lineage

RDD -> T -> RDD-> T->RDD



## Unified Apache Spark 2.0 API



databricks



untypesd - python  
typed - Java.

# DATAFRAME

- High level API
- Immutable distributed collection of data.
- Data is organized into named columns, like a table in a relational database
- DataFrame allows developers to impose a structure onto a distributed collection of data,

Spreadsheet on a single machine



Table or DataFrame partitioned across servers in a data center



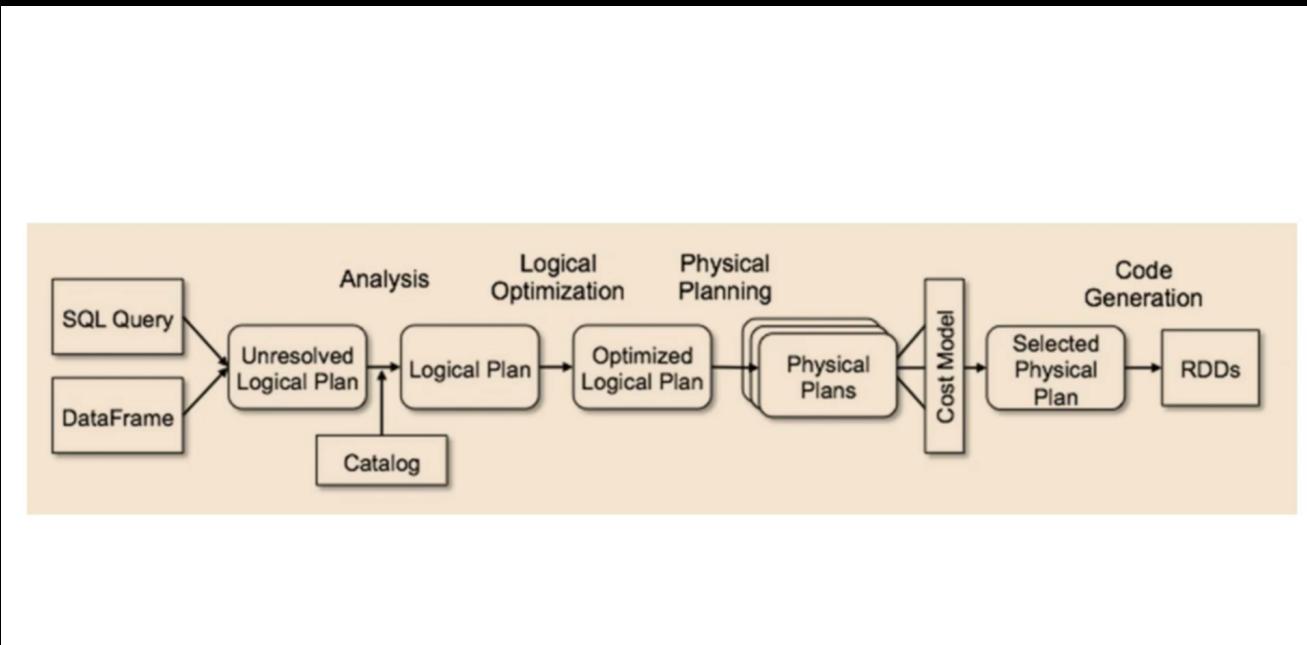
# DATASET

- High level API
- The Apache Spark Dataset API provides a type-safe, object-oriented programming interface.
- Datasets provide compile-time type safety—which means that production applications can be checked for errors before they are run

## WHEN TO USE DATAFRAME AND DATASET

- High-level abstractions, and domain specific APIs
- Processing demands high-level expressions, filters, maps, aggregation, averages, sum, SQL queries, columnar access and use of lambda functions on semi-structured data,
- Higher degree of type-safety at compile time - Dataset
- Unification and simplification of APIs across Spark Libraries
- R user - DataFrames.
- Python user - DataFrames





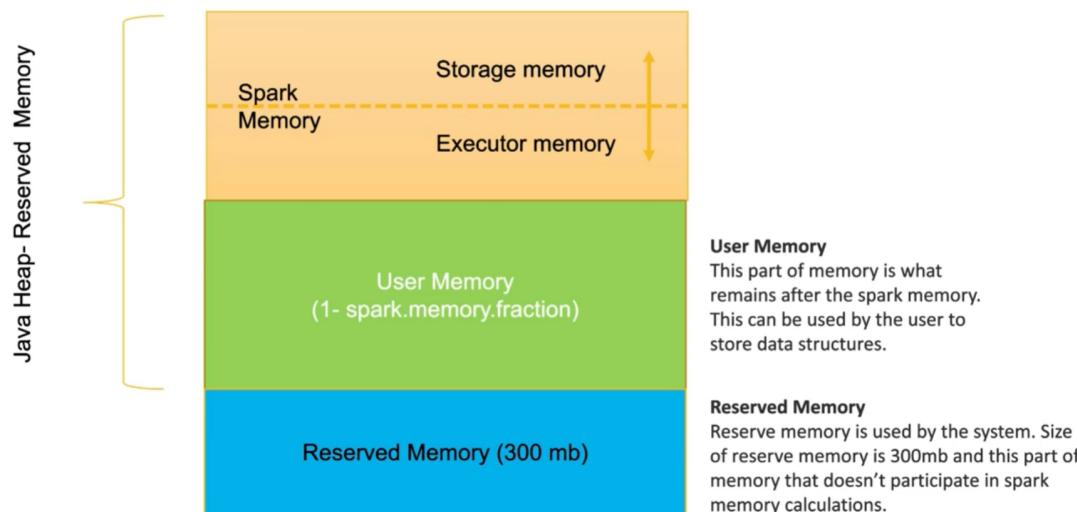
## Memory management

*why?*

→ memory critical for spark

mm = Higher performance

## Spark Memory



# Memory Types

Execution memory refers to that used for computation in shuffles, joins, sorts and aggregations

Execution Memory

Storage memory refers to that used for caching and propagating internal data across the cluster

Storage Memory



Execution and storage share a unified region

Previously memory was static  
memory allocation.

## Spark 1.6 onwards

The memory management model used by Spark 1.6 onwards is **Unified memory management**.

Instead of specifying execution and storage memory as two separate chunks, it's specified as a single unified space that they both can share.

Advantage –

- Execution and storage can use the unutilized memory
- Prevents disk spills
- Supports variety of workloads



Nowadays unified memory management is used.

- It specifies execution and storage as single chunk and only specified reserved memory.

# Determining Memory Consumption



## Memory consumption of a dataset

Create an RDD, put it into cache, and look at the "Storage" page in the web UI. The page will tell how much memory the RDD is occupying.

## Memory consumption of a particular object

Use SizeEstimator estimate method

We can see the consumption of dataset  
— create a RDD put that into cache  
go into webUI and check the  
consumption on the storage.

# Tuning Considerations



Find out the Memory consumption of a dataset/object

Serialize RDD storage

Tune Garbage collection

Level of Parallelism

Broadcasting

## RDD from slides

### Resilient Distributed Dataset

- collection of in-memory objects divided into multiple parts which resides on different machine
- Basically Architecture of Spark

- RDD exposes two types of operations

#### \* Transform operations :-

- map
- filter
- reduceByKey

#### \* Action operations :-

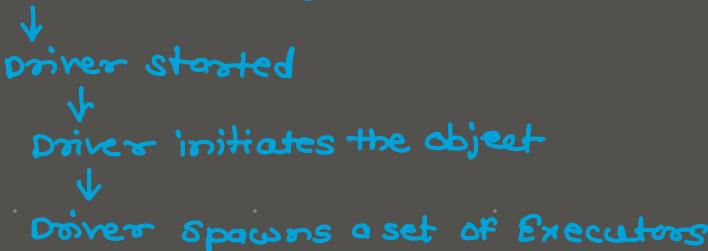
- collect
- count

RDD objects are immutable

## Deploy a Spark Program

- following event happens
  - driver process started on the same JRM
    - ↳ Driver will initiate the object to keep track of all configuration of App
  - ↳ Driver process spawns a set of Executors via a cluster manager.

Spark Object initiated by Spark Program



- Driver process is also responsible for job scheduling
- FIFO scheduler
  - fair scheduler

### Driver process job

- Construct the DAG from Spark Job
- generate the execution plan outlining the Spark Tasks to be executed.
- Invoke the task scheduler to schedule tasks to run on selected executors.

following commands for configuration of clusters

num-executors → Number of executors

executor-cores → Number of cores used per executor

executor-memory → Memory size used per executor

One of the following cluster Manager is used according to the value.

- Spark Native cluster manager ( master = spark://host:port )
- YARN ( master = yarn )
- Mesos ( master = mesos://host:port )

## Spark RDD programming

Setup of Deployment Environment.

- Create a Maven project with intelliJ
  - Add spark dependencies to POM.xml
    - Spark Core

- Spark SQL
- Spark ML
- ~ etc.

## Main types of RDD objects

- Java RDD
- JavaPair RDD

### 1. Java RDD

- Collection of partitioned distributed objects
- Created via `SparkContext` object
- Applicable Transform operations
  - `map`
  - `filter`
  - `distinct`

### • JavaPair RDD

- Collection of partitioned distributed key-value pair objects
- created by invoking the `mapToPair` operation on a JavaRDD object
- Applicable Transform operations
  - `reduceByKey`
  - `groupByKey`
  - `sortByKey`
  - `join`

## Common Action operations on both RDD types

- Return a RDD object to the driver program
  - `count()`, `collect()`, `first()`, `reduce()`
  - `foreach()`

- Persist a RDD object in a file
  - SaveAsTextFile()

## Common Transform Operations on both RDD types

- repartition(), coalesce()

*very expensive operations as they shuffle the data across many partitions hence try to minimize them.*

*sparkContext.parallelize(range(0, 20), 6) & RDD positions (0 to 4) from 6 to 14*

*odd, repartition(4)*

*coalesce() ← To only decrease the no. of partitions in efficient way*

*repartition() ← To increase decrease the no. of partition in RDD*

**Programming pattern**

- SparkContext.parallelize()
- SparkContext.textfile()

RDD automatically recover from node failures.

Spark supports two types of shared variables

① broadcast variables

↳ used to cache a value in memory on all nodes

② accumulators

↳ type of variables that can only be added to such as counters and sums.

Parallelized collections are created by calling SparkContext's parallelize method on an existing iterables or collection in your driver program.

## RDD Operations :-

Two types of operations:

① Transformations

— Which creates a new dataset from an existing one

## ② Actions

- Which returns a value to the driver program after running a computation on the dataset.

Example :-

Map is transformation that passes element through a function and returns a new RDD representing the results.

All transformation on Spark are lazy

- GroupByKey()
- Join()
- ReduceByKey()
- SortByKey()

They just remember the transformation applied to some base dataset. Transformation only computed when action requires results to be returned to the driver.

By default, each transformed RDD may be recomputed each time you run an action on it.

However, You can also persist an RDD in memory using the **persist** (or **cache**) method, in which case Spark will keep the elements around on the cluster for much faster access the next time you query it.

shuffle :- Spark mechanism for re-distributing data so that's grouped differently across partitions.

\* In Spark we can decide what size of container should be also we can choose how many container we want.

- **Spark Dynamic execution allocation**

for allocation for data and it will automatically kill itself after completion.

- Data locality is not 100% guaranteed.

- Block size in Hadoop 128MB.

In Spark More partition = more speed

Executor is JVM. 1 cores is required to manage it.

More partition = more processing power you can get.

\* In Spark If you request 10GB memory it will only give you 54%.

because 10% goes for system calls

remaining goes for JVM & garbage collection

So you only get the remaining memory 64%

— After transformation for example we filtered the data then we run coalesce() it won't show the data because it's lazy transformation which means it won't show it without using Action.

Once processing done Pipeline will be empty.

DAG - Directed acyclic graph.

SparkSQL is much more optimized = DataFrame.



**Sheesh**



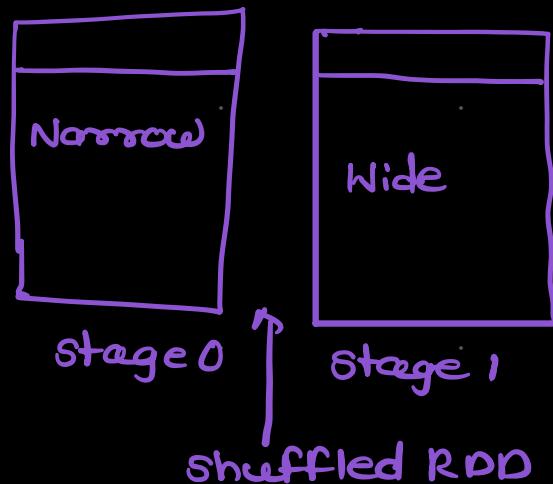
Spark is in memory processing doesn't mean it will use ram everytime. It will use ram if its available.

### \* Job, Stages and Task

- Transformation has two types

① Narrow (Where is no communication needed betn two task)  
ex. filter, map.

② Wide (Where communication betn two task required  
ex. join()  
groupByKey()



after successful completion executors will be deleted automatically but driver will still be there. Which takes around 1GB of resources so they have to be stopped manually.

pySpark2 --master local[3]

[3] = Number of cores to use

[\*] = Use whatever cores available.

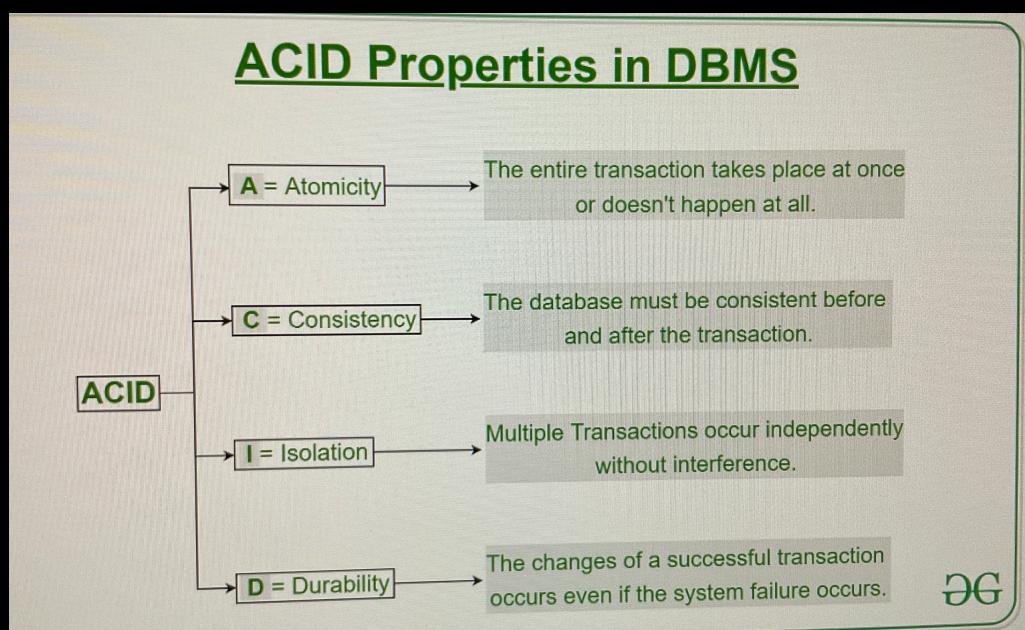


# ACID

## properties 4 MCQ

Why to use?

In order to maintain consistency in database, before and after the transaction.



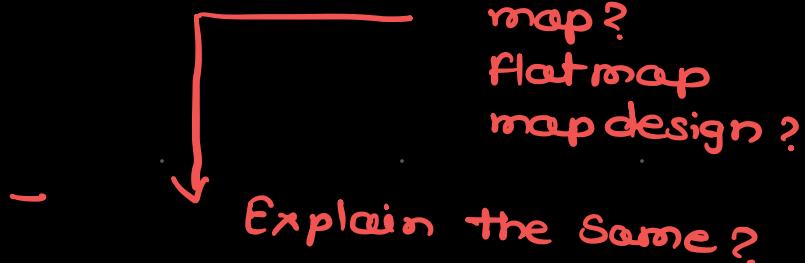
Atomicity - either the transaction takes place at once doesn't

40%. ① 20 words match the columns

- E ~ map reduce
- E ~ HDFS
- E ~ spark
- E ~ sharding
- E ~ mongod b
- F shard key

30%. ② 2 parts

1. - Which spark method used for partition
  - output txt what changed repartition
2. which method for k:v pair.



30%. ③ Numerical

Map reduce based

- 50 nodes 256GB

Yarn container how many ?

$$\# \text{ of containers} = \frac{\min(2 \times \text{cores}, 1.8 \times \text{Disk (total available RAM)})}{\text{min-container size.}}$$

=

$$\text{Reserved memory} = \text{Reserved for stack memory} + \\ \text{Reserved for HBase}$$

Ex. 12 CPU cores, 48GB RAM & 12 disks

$$\text{Reserved memory} = 6 \text{ GB for system} + 8 \text{ GB for HBase} \\ \text{Min container size} = 2 \text{ GB}$$

$$= \min \left[ 12, 1.8 \times 12, \left( \frac{48}{2} \right) \right]$$

$$\text{RAM} = 48 - 6 \text{ (removing reserved)}$$

$$\min \left( 2 \times 12, 1.8 \times 12, \left( \frac{48 - 6}{2} \right) \right)$$

$$= \min (24, 21.6, 21)$$

$$= 21 ?$$

RAM - per container ?

$$\text{RAM-per container} = \text{Max}(\text{min\_container\_size}, \frac{\text{(total available RAM)}}{\text{(containers)}})$$

$$= \text{Max}(2, \left( \frac{48-6}{21} \right))$$

$$= \max(2, 2) = 2$$

$$\# \text{ container} = 2 \times \text{cores}, 1.8 \times \text{disk}, \frac{\text{RAM-reserved}}{\text{min container}}$$

$$= \frac{48-6}{2}$$

(21) ✓

$$\text{RAM per container} \quad \frac{48-6}{21}$$

# of containers required ?

50 datanodes

256 GB

=  $50 \times 256 \text{ GB}$

= 12.5 TB

2GB container size we will require

$$\frac{256}{2} =$$

256GB size will require containers?

8/16/32

for Blaze engine

$$(n \times 2) + 3$$
$$= (50 \times 2) + 3$$

103 container require

could be wrong

Just theory from 2016

When you ask for Yarn yarn.scheduler.minimumAllocation-mb make sure to ask for exact no. for ex. 4GB if you round that off to 4.5GB Yarn will give the closest number which is 8GB.

offheap Java - 10-25%. for example garbage collection.

- So in total you need the requested on Heap Java memory plus 10-25% for offheap otherwise YARN will simply shoot down the container due to

Launcher is bigger than the memory limit.  
for this container.

Container could be 4-8-16-32

If you request anything smaller  
or anything off number you get next closest.

→ `repartition()`, `coalesce()`

Output after operation

from local = 5

parallelize = 6

Textfile = 10

we have 10 text files

using `spark.context.parallelize(range(0, 20), 6)`.

Which will distribute RDD into 6 partitions

p1

p2

p3

p4

p5

p6

→ using `repartition()`

→ `rdd1.repartition(4)`

with using this command output file will be reduced to 4 files

P1

P2

۳۹

p4

This will be done in shuffle manner which is an expensive to do.

→ impact of saveasText file → Sheffield output

1

P1: 1 2 3

p2: 4 5 6

P3: 789

P4: 10 11 12

PS: 13 14 15

PS: 1G 17

## Impact of

## separation (4)

PIC 1891215

P2: 26.10.14

P3: 35, 111c

P4: 47 13 17

$\rightarrow$  RDP coalesce()

IS used to only reduce the no. of

partitions This is optimized or improved

version of `repartition()`

- P1: PS:
- P2: PG:
- P3:
- P4:

What coalesce() do here is combine any files together based on requirement

for Here:- coalesce(4)

This will combine any two partitions into one  
to reduce from 6 to 4

= rdl1.coalesce(4)

P1:

P2:

P4: → P3 + P4

P5: P5 + P6

---

final verdict

repartition() — increase decrease + shuffle  
coalesce() → decrease but in indexed

textfile() → Read multiple files & output is  
single RDD file(string)

wholetextfiles() → Read multiple & o/p is  
single RDD file (tuples (string, string))

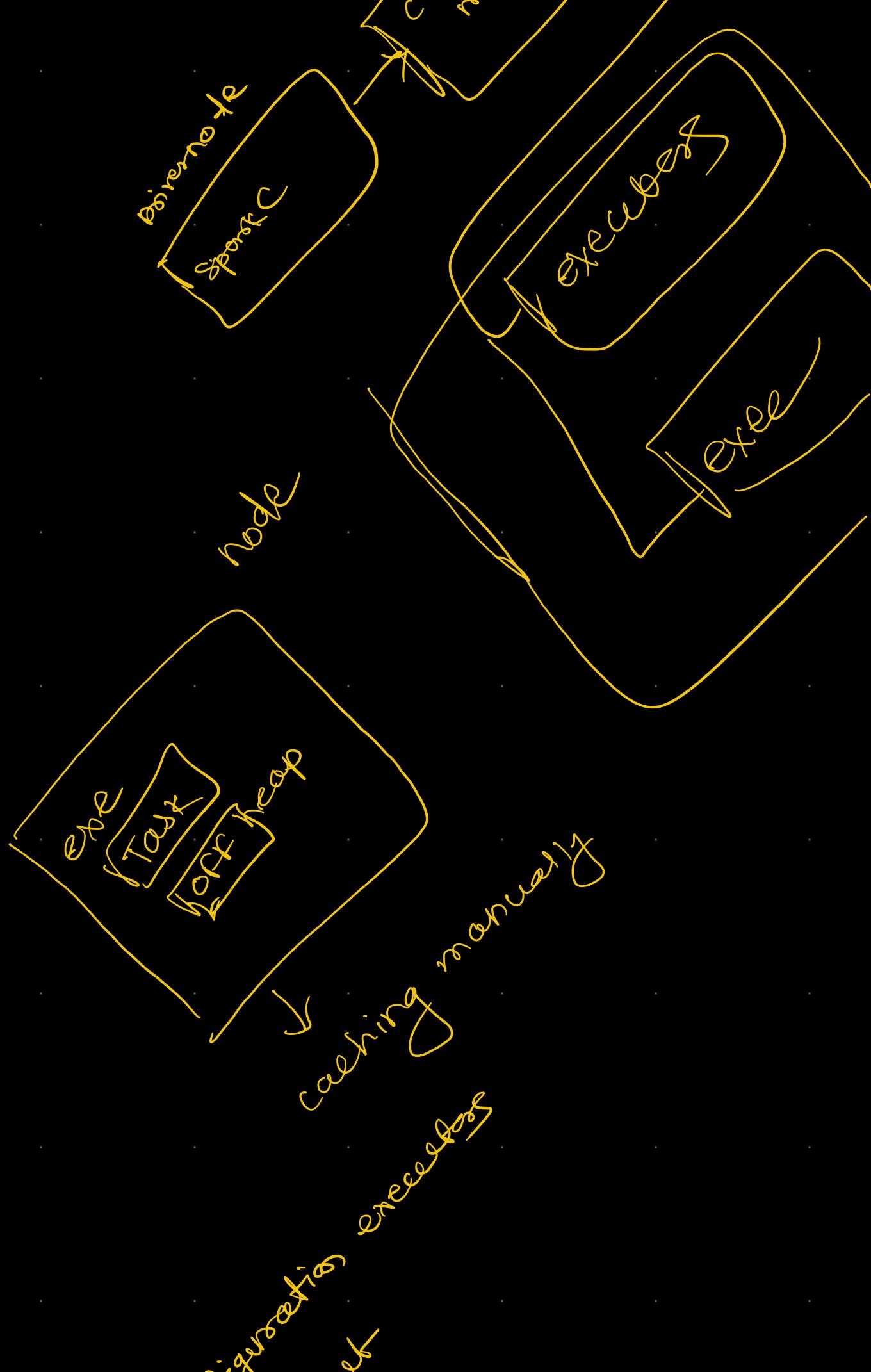
→ emptyRDD()

To create emptyRDD() function.

Tuesd 1  
we re

---

luster  
map



Spark conf  
human, exec  
executor, coher

dep( $\epsilon$ )  
~~shapefile~~  
 $P_1 \cdot O_1 \cdot e$   
 $P_2 \cdot O_2 \cdot e$   
 $P_3 \cdot O_3 \cdot e$   
 $P_4 \cdot O_4 \cdot e$

map, Topair( $\epsilon$ )  
reduce  
JavaPairRDD  
by key  
key, value

Sheffé when  
replication  
coherence  
protection  
decency  
cognit  
join  
conting

→ persistence  
coherence  
persist