10/06/20

## SPARK -

★ general purpose, in-memory, compute engine

i) compute engine -

Hadoop was providing 3 things -

1. HDFS - Storage
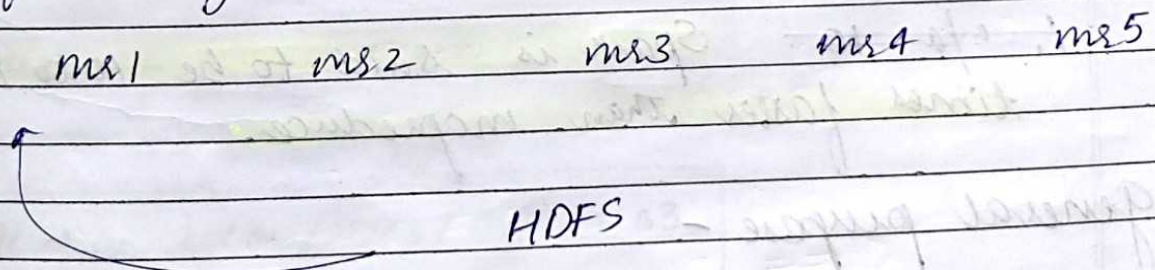2. MapReduce - Computation
3. YARN - Resource Manager

⇒ Spark is also a computation engine :
spark doesnot replace Hadoop, it is an
alternative of mapreduce.

Spark is a plug & play compute engine which needs
2 things to work with -

1. Storage - local storage, hdfs, Amazon S3
2. Resource Manager - YARN, MesOs, Kubernetes.

ii) in-memory -

For each MR job, we require 2 disk access, One is
for reading & other is for writing.

mr1       mr2       mr3       mr4       mr5

HDFS

mr runs on top of HDFS & takes i/p from HDFS
say mr1 takes i/p & writes again to HDFS.
mr2 takes i/p from mr1 & writes to HDFS.
so for every mr job; there are 2 disk access
here, we have 10 I-O disk access
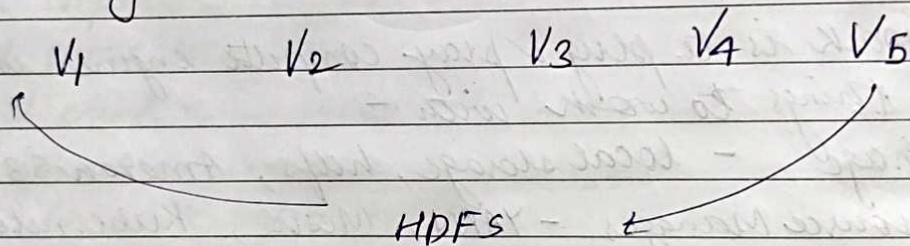for 100 MR job, - 200 I-O disk access are reqd

MapReduce has <u>high latency</u> bcz it involves more disk read & write operations than Spark.

It takes alot of time & is bottleneck for HDFS.
It is very slow bcz of disk I-O's.

⇒ In case of Spark, say it reads data from HDFS & stores it in variable V1, now for next op$^n$ stores it in V2 (V1 is in memory) V2 takes i/p from memory & store o/p in memory. Again V3 takes i/p from V2 (in-memory) and writes o/p in-memory Finally, V5 writes o/p to HDFS.
∴ Only 2 disk IO's are reqd.

$V_1$      $V_2$      $V_3$   $V_4$   $V_5$

HDFS - t

Disk IO is time consuming process.
Thus Spark optimizes it by keeping the value in-memory rather than again saving it to disk.

∴ <del>∵/a 10</del> <mark>Spark is said to be 10 to 1000 times faster than mapreduce.</mark>
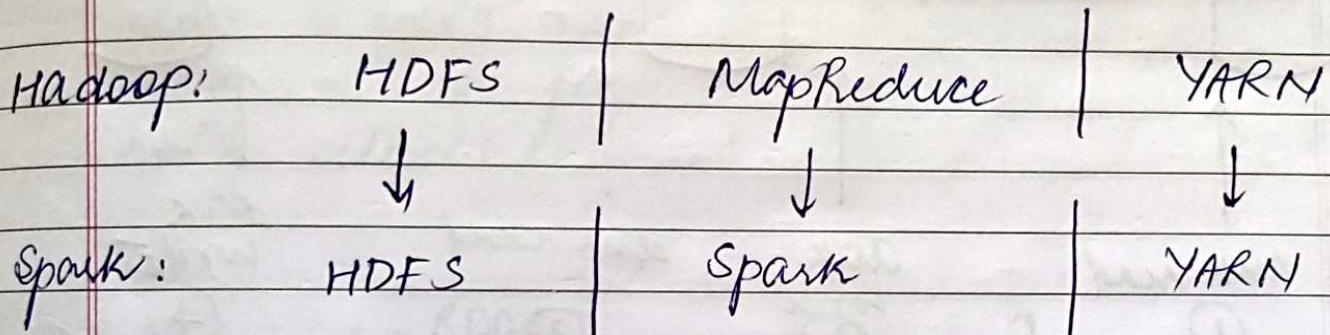
iii) general purpose -

In hadoop,
for cleaning data - Pig is used
for querying data - hive
for machine learning - Mahout
for data ingestion - sqoop.
for streaming data - storm

Also, only bound to use map & reduce only.

spark — all things are present at one place,
general purpose.
↳ Learn one style of writing code & all
things like cleaning, querying, ML,
data ingestion all these can be done with
spark!
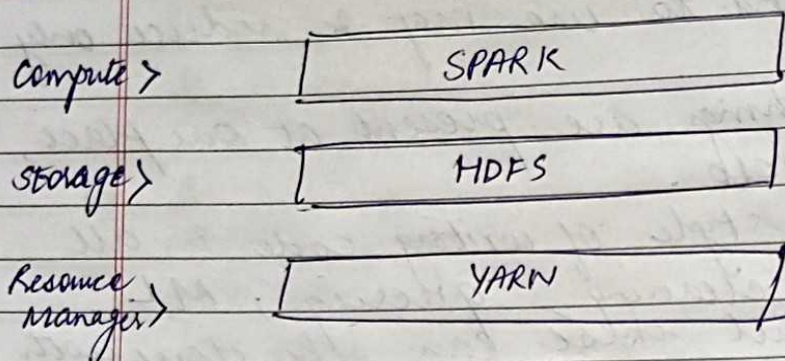↳ map reduce is also there but apart from
that many other things are there.

| Hadoop: | HDFS | MapReduce | YARN |
|---|---|---|---|
| | ↓ | ↓ | ↓ |
| Spark: | HDFS | Spark | YARN |

(compute engine)

↳ **Spark on top of Hadoop.**

| Compute → | SPARK | |
|---|---|---|
| Storage → | LOCAL / HDFS / Amazon s3 | |
| Resource Manager → | YARN / MESOS / KUBERNETES | |

↳ Hadoop Cluster.

# spark on top of Hadoop -

| | |
|---|---|
| Compute > | SPARK |
| Storage > | HDFS |
| Resource Manager > | YARN |

## Map Reduce processing:

MR jobs

iteration1       iteration 2

disk read ①    disk write ②    disk read ③    disk write ④

## Spark processing:

iteration 1 → [RAM] → [RAM] → iteration 2

memory     memory

disk read ①

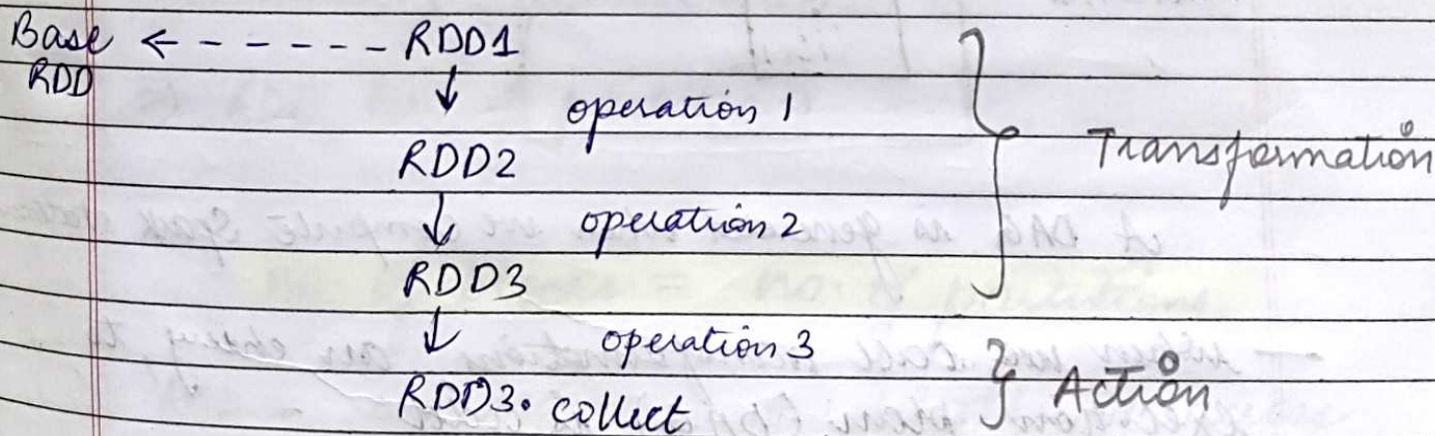initial read

disk write ② (final write)

The basic unit which holds data in Spark is called as RDD.

  Resilient Distributed Dataset.

Datatype List is collection of data, stored in 1 machine.

RDD is inmemory distributed collection (List kept in 4 diff m/c's)

Var. name
~ rdd1 = load file1 from hdfs  → = SC.textFile("—")
rdd2 = rdd1.map
rdd3 = rdd2.filter
rdd3.collect()

Base
RDD ← – – – – – – RDD1
                     ↓       operation 1        ⎫
                  RDD2                          ⎬ Transformation
                     ↓       operation 2        ⎭
                  RDD3
                     ↓       operation 3        ⎫
                  RDD3.collect                  ⎬ Action

    DAG – directed acyclic graph
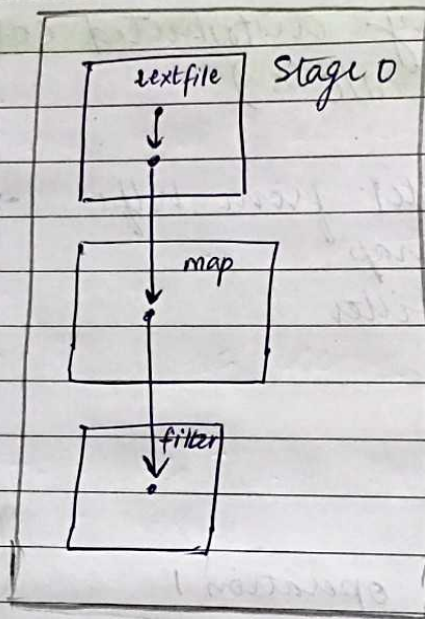↳ created at backend. (entries are made)

In Spark there are 2 kind of operations –
1. Transformations
2. Actions

Transformations are Lazy, Actions are NOT.

when we call Action, then transformations execute.

when we execute rdd1, rdd2, rdd3 = filter-

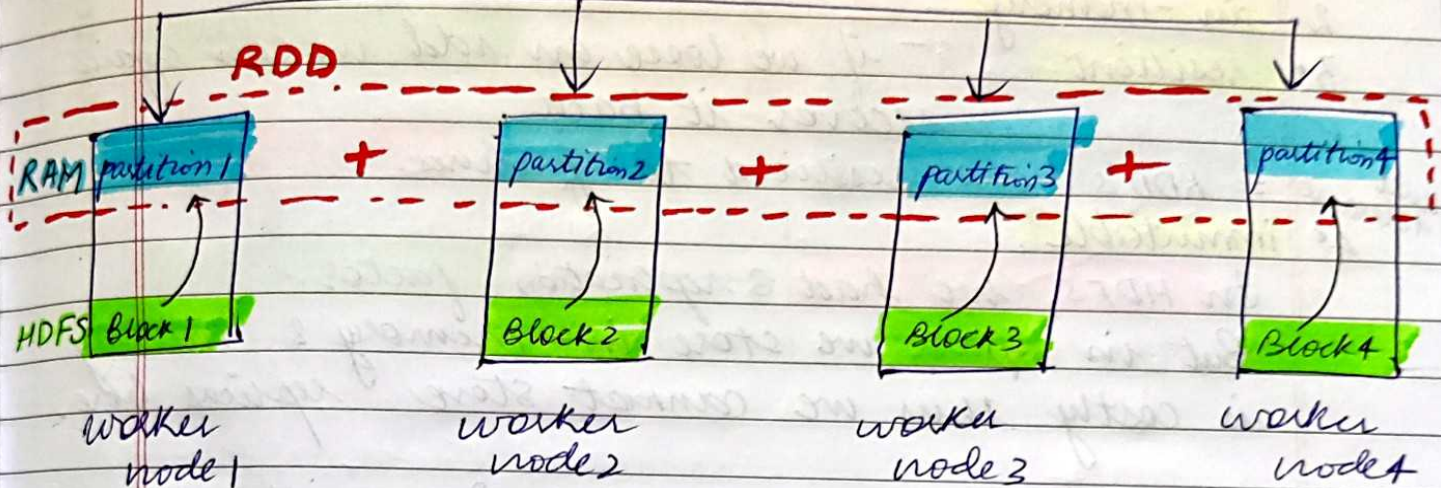no actual computation has happened. but a diagram will be created in backend. — DAG



A DAG is generated when we compute Spark statements

— when we call transformations an entry to execution plan (DAG) is add

— Execution happens when Action is encountered before that only entries are made into DAG.

4 node cluster —

All nodes are c/a worker nodes
(500MB data = blocksize 128 MB →
4 nodes)

```
                    ┌─────────────┐
                    │ driver node │
                    │             │
                    └──────┬──────┘
```

RDD

| RAM | partition 1 | + | partition 2 | + | partition 3 | + | partition 4 |
|---|---|---|---|---|---|---|---|

| HDFS | Block 1 | | Block 2 | | Block 3 | | Block 4 |

worker node 1          worker node 2          worker node 3          worker node 4

Data is loaded from HDFS (blocks) to memory (partitions)

⇒ RDD has 4 partitions

no. of blocks = no. of partitions

RDD = in-memory + distributed across machines

| Blocks | RDD |
|---|---|
| Distributed across disk (HDFS) | Distributed across memory |

Data is brought from HDFS, which are in-memory

RDD's are -

1. **Distributed**
2. **in-memory**
3. **resilient** :- if we loose an rdd we can again
recover it back.

fault tolerance = RDD's are resilient to failures.

4° **immutable**.

In HDFS, we had 3 replication factor.

But in Spark we store in-memory & memory
is costly thus we cannot store replicas here.

Then how RDD's are resilient ? we will not loose data
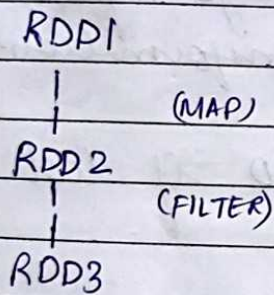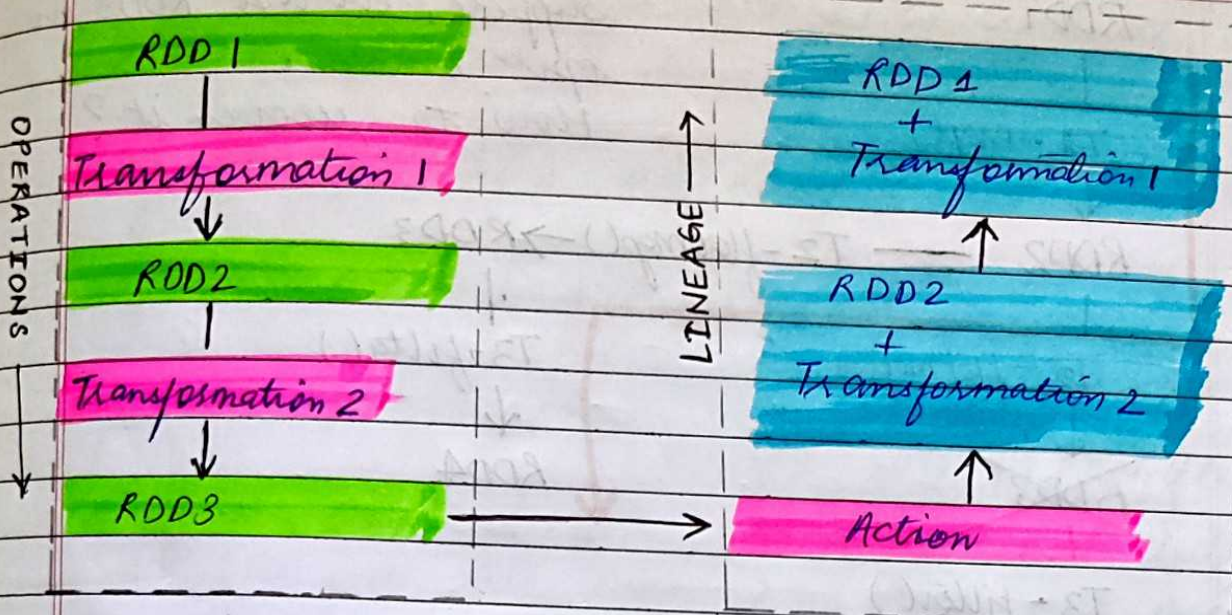
Resilient = ability to quickly recover from
failure.

= ability to recreate lost RDD's during
execution cycle

RDD provide fault tolerance through lineage graph

A lineage graph keeps a track of transformations
to be executed after an action has been called.

— When an action is encountered a path will be
created, which defines the order of execution of the
operation.

— This information is remembered in the form of lineage
graph.

**OPERATIONS**

RDD 1

↓

Transformation 1

↓

RDD 2

↓

Transformation 2

↓

RDD 3 ──────→

**LINEAGE** →

RDD 1
+
Transformation 1

↑

RDD 2
+
Transformation 2

↑

Action

---

RDD1
|
|       (MAP)
|
RDD 2
|       (FILTER)
|
RDD 3

If RDD3 is lost, then we can recover it from parent RDD i.e. RDD2.

— It will check for its parent RDD ~~it will~~ using lineage graph & it will quickly apply transformation (filter) on RDD2.

⌐RDD provides a fault tolerance.
RDD lineage graph helps recomputed any missing or damaged RDD bcz of node failure.

HDFS → fault tolerance = Resiliency → by using replication factor

RDD → we get resiliency by using Lineage graph