# Spark Execution Model And Architecture - (Kafka).

1. How can we create spark Programs?

2. How to execute spark Programs?

---

① Interactive client (spark-shell, Notebooks) {learning phase}.
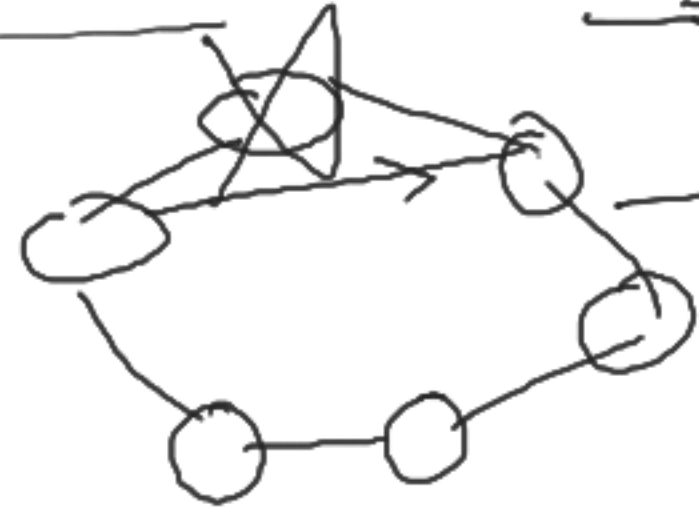
② Submit Job -        spark-submit.

created whole program → submitted to a cluster

cluster — (Hadoop) . (Multiple system)
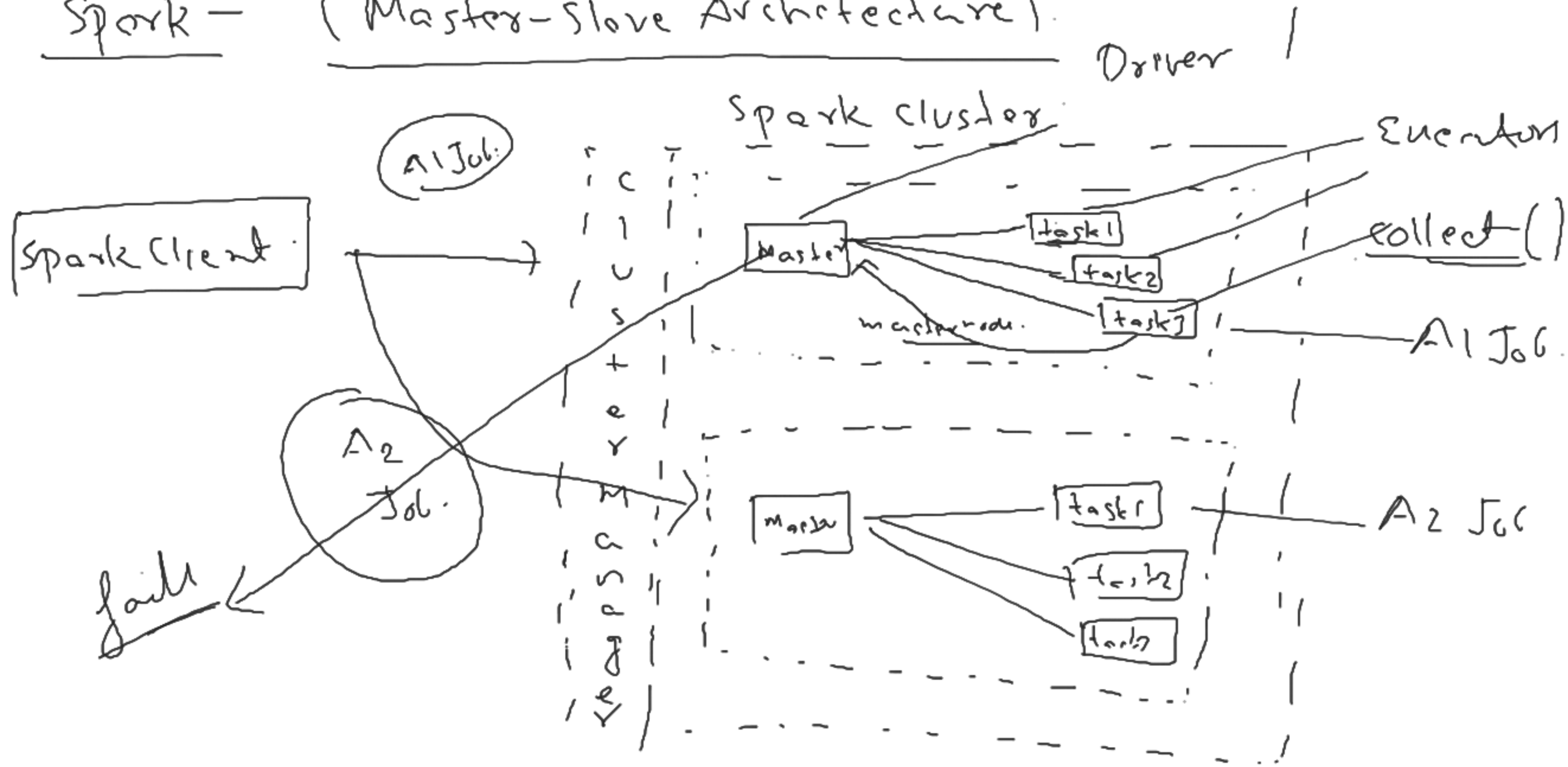
command.

① Master-Slave Architecture

w  w  w  w  w

orphan
w.

② Ring Architecture — (Cassandra / MongoDB).

latency is very high.

System Design →

# Spark — (Master-Slave Architecture)

Driver

Spark Cluster

Executor



Spark Client

A1 Job

Master

master node

task1

task2

task3

collect()

A1 Job

A2 Job

fail

Master

task1

task2

task3

A2 Job

cluster Manager

1,10,1 Ram

M

50GB → 4GB

10GB

shuffling

networking

collect()

New Student
Dept
files

0.25   0.25   1.25   1.25   1.25   1.25

6 GB

network
congestion

Bandwidth
issue

OOM

Executor OOM

latency

→ OOM

→ slowness

① collect()

② 

groupBy operator: sum()

simple python programs

voltaile   HDD

→ SDD

RAM
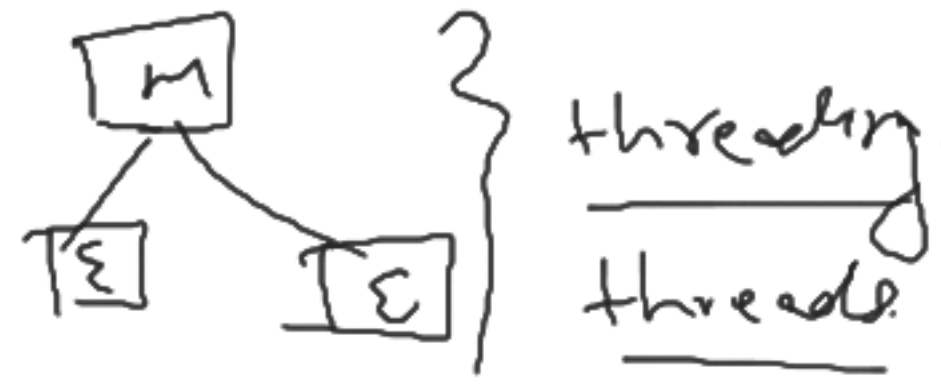
Spark Programs - (Legacy software). Python Java
Scala

① local[n] ⟶ local[2]
local[∞] ⟶ cores ⟶ ④
local[1] ⟶ one driver
one executor.

④ thread

```
[M]
[E]   [E]
```
} threading
threads

mimic the distributed system

② YARN - (EMR clust., ADB, dataproc).
Fargate)

Flink on kubernetes
Spark.
AWS
{ link (Beam).

③ Kubernetes-
↳ very interesting

single
Python program

serialization distribution

(orchestration tools).

④ Mesos - (orchestration).

⑤ Standalone - development-

① C, C++. (Evolution of technologies).      Advance Java. (7,000/-)
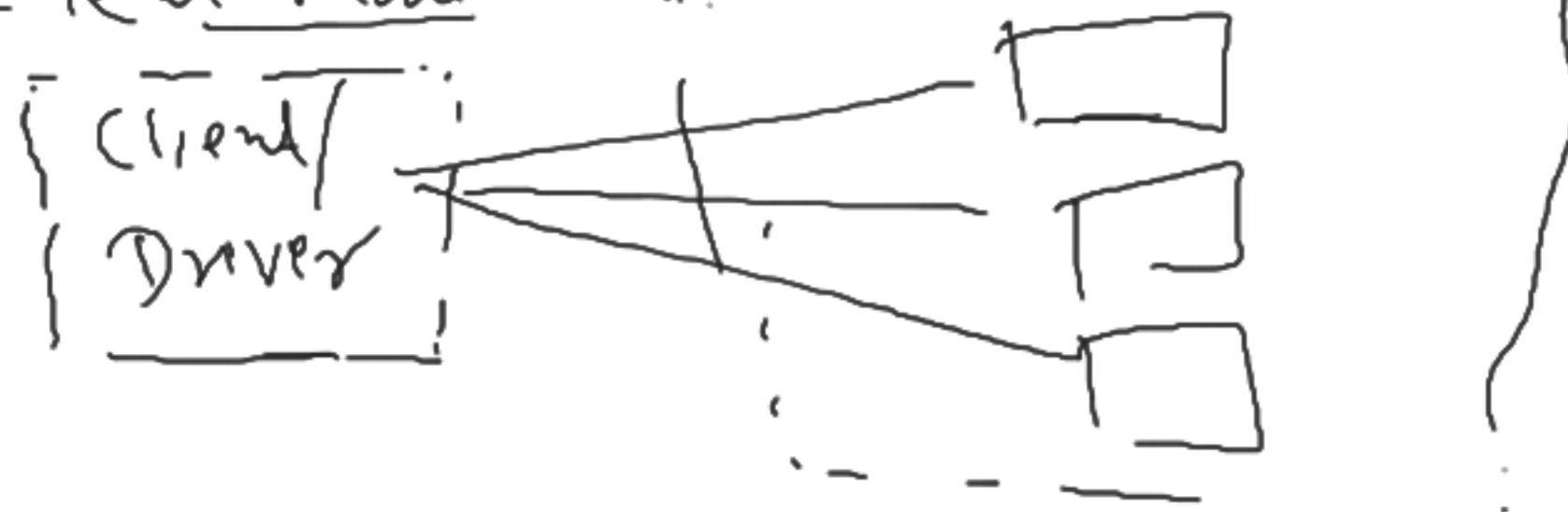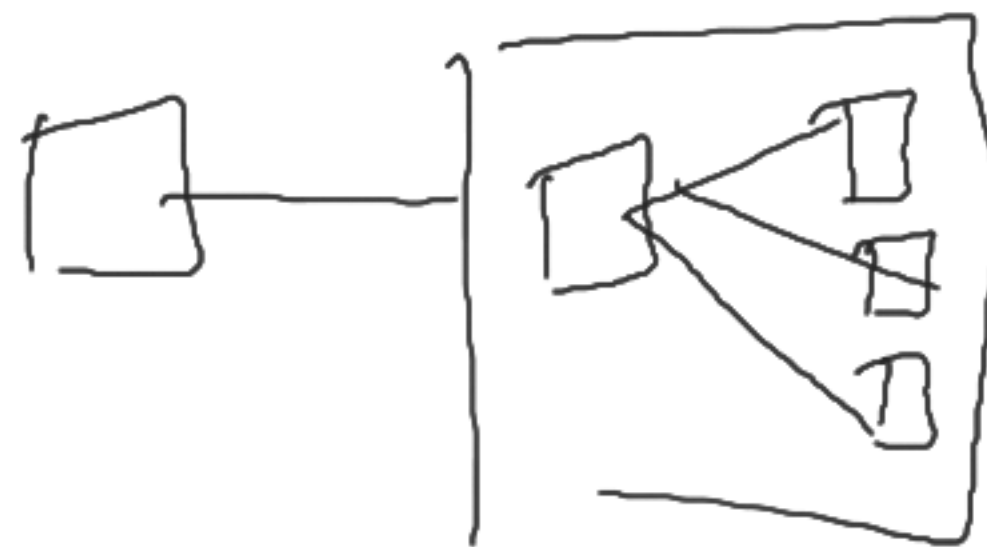                                                              J2EE.

② Spark — (Entry level Big data engineer).                → Java

      legacy. (Sqe.m)

      {Flink/Beam}

Spark with Interactive Client

      Client Mode-                                    Cluster

      {Client/
       Driver

| Cluster Manager | Execution Modes | Execution tasks. |
|---|---|---|
| 1. local[n] | client | IDE, NoteBook. → Debug mode |
| 2. YARN { on-Prem. on-cloud. | - Client | Databricks Notebook, Shell. |
| | Cluster | Spark-submit |

In memory - RAM. → (voltaile memory). → 10GB.

10,000 SR → 150cm

V / 24,000

Hadoop map-Reduce - Harddisk → knitting (slow process).

Commission ← → PMS → Computer Services. cost of

CD School → 20. Sell ~

market

8. GB. - DDR 3 → Good ~

( 40 - 50 !.)

1024 GB → 1 TB - <1000 !. →.

Technical.

battery → 4000 !

→ 2656 → 3,200

7,000 !~

1,00,00/per worker

dell charger → 1200/-

2,000/-

① Servers 20

350/cm.

Spark - R → flushed to data disk.

→ RAM.

(storage level.

① Batch form
   (Real time data).

① Data is coming in files, (historical data)

spark supports.
① Batch loads → 10,000
② Near real time streaming (micro batches)
   up-stream (Flink) (Real time streaming).

② Streaming - (Real time

water tank — main water supply

24hr

flush 10000L

5cm flow.

2.5cm. 2m

① Batch loads
② Near real time
③ real time.
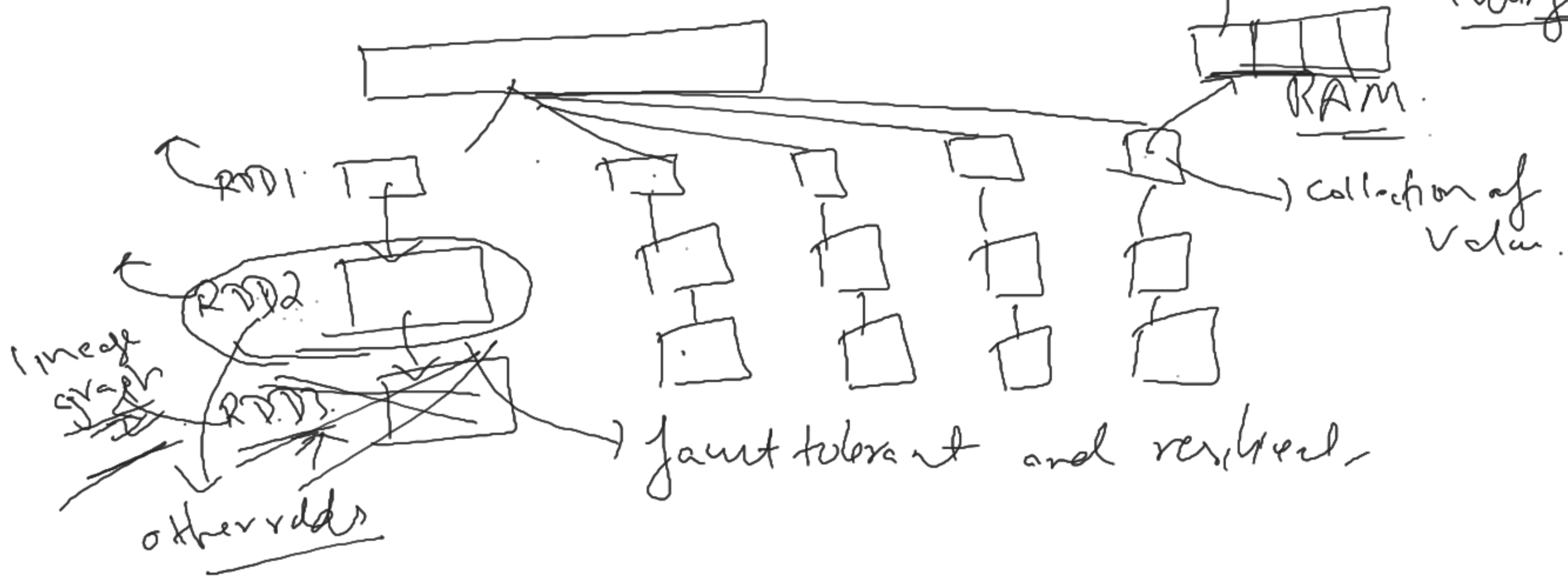
# RDD - Resilient Distributed Dataset

① Immutable - change the contents of an RDD.

Hadoop

① Replication
② Erasure coding

Column name

RAM.

RDD1:

→ Collection of values.

RDD2

lineage graph

RDD3

other rdds

→ fault tolerant and resilient.

① Transformations —

input RDD — output RDD

① wide → shuffling

② narrow → no-shuffling

until and unless an action is applied, transformation, won't be executed. (spark is lazy evaluated)

② Action.

input RDD → non-RDD.