

Agenda

- a. Real time problem statement
- b. History of data
- c. Spark Architecture
- d. How to play with Spark ?
- e. Spark Advance Concepts

Problem statement

Netflix →



Ques. 1

which Comedy Movie has the more ratings
and return the ratings?

↓ Million
Movies

(100 million)

users

① Python | Java

RDBMS (MySQL, Oracle
Movies Postgres....)

ratings ↑ 2



Table

Join



Result ✓

Request Time out

③ Hadoop (Map Reduce) ✓

↓ Time

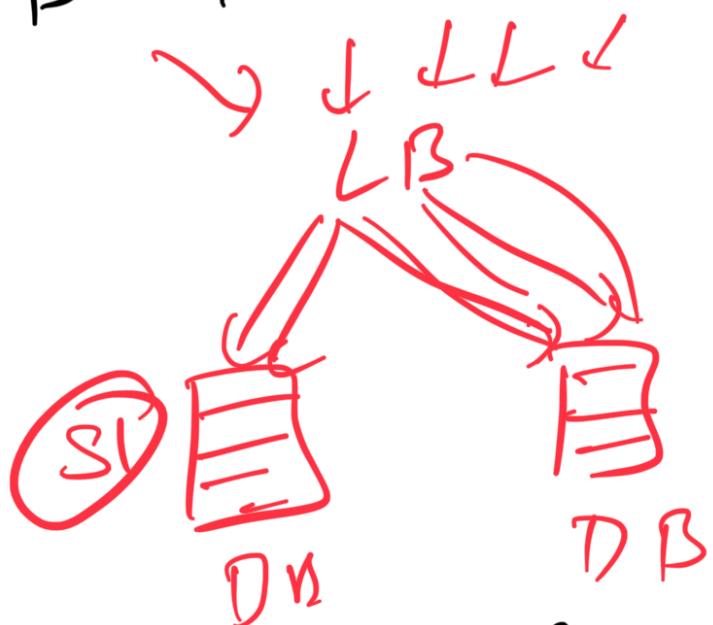
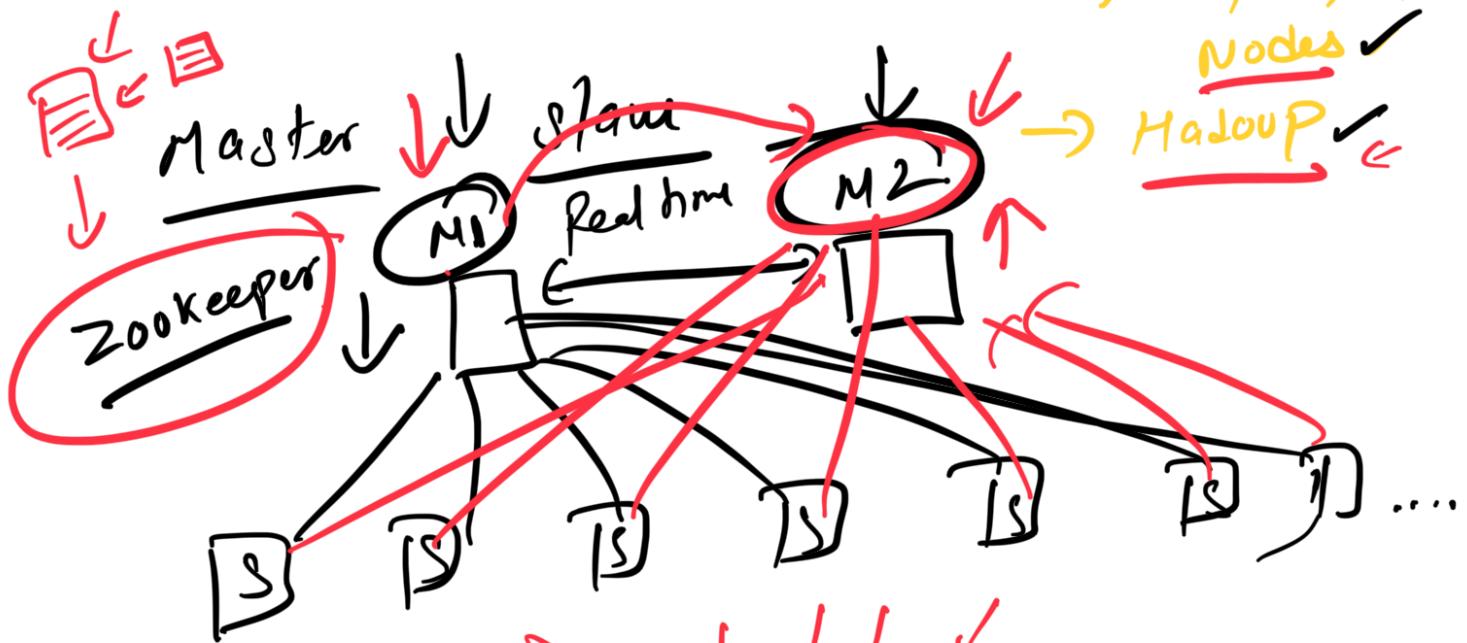
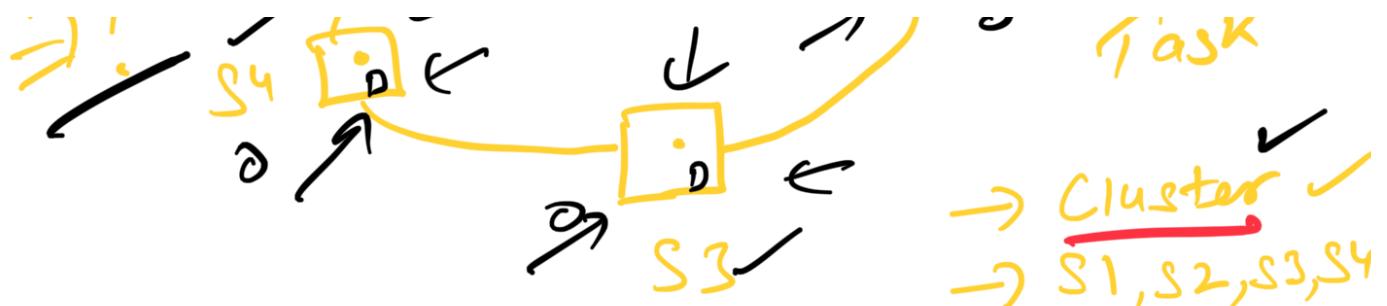
Hadoop = A Distributed System

Server 4230 M/C

↓ Go to Master S1 N/W

~9

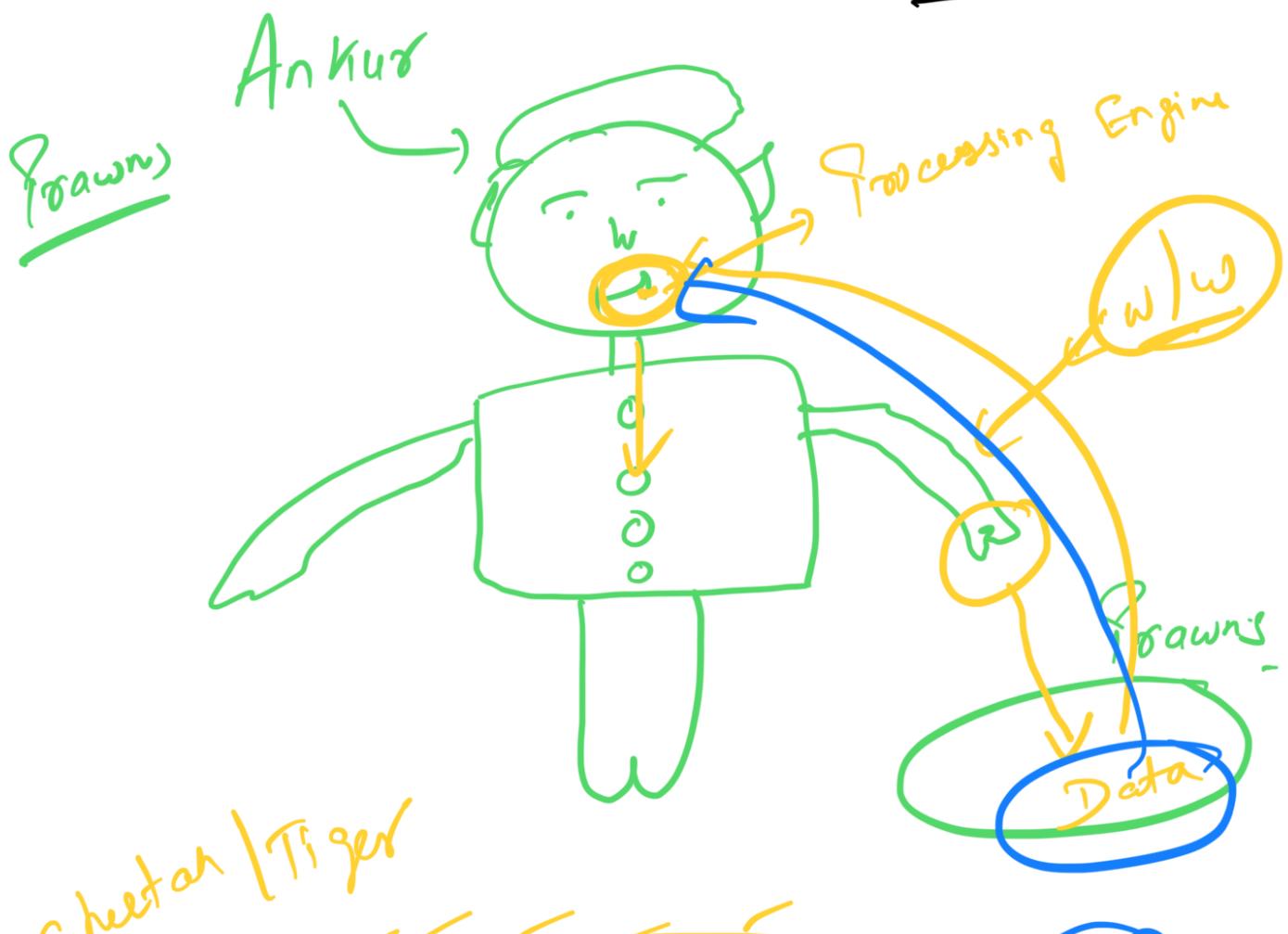


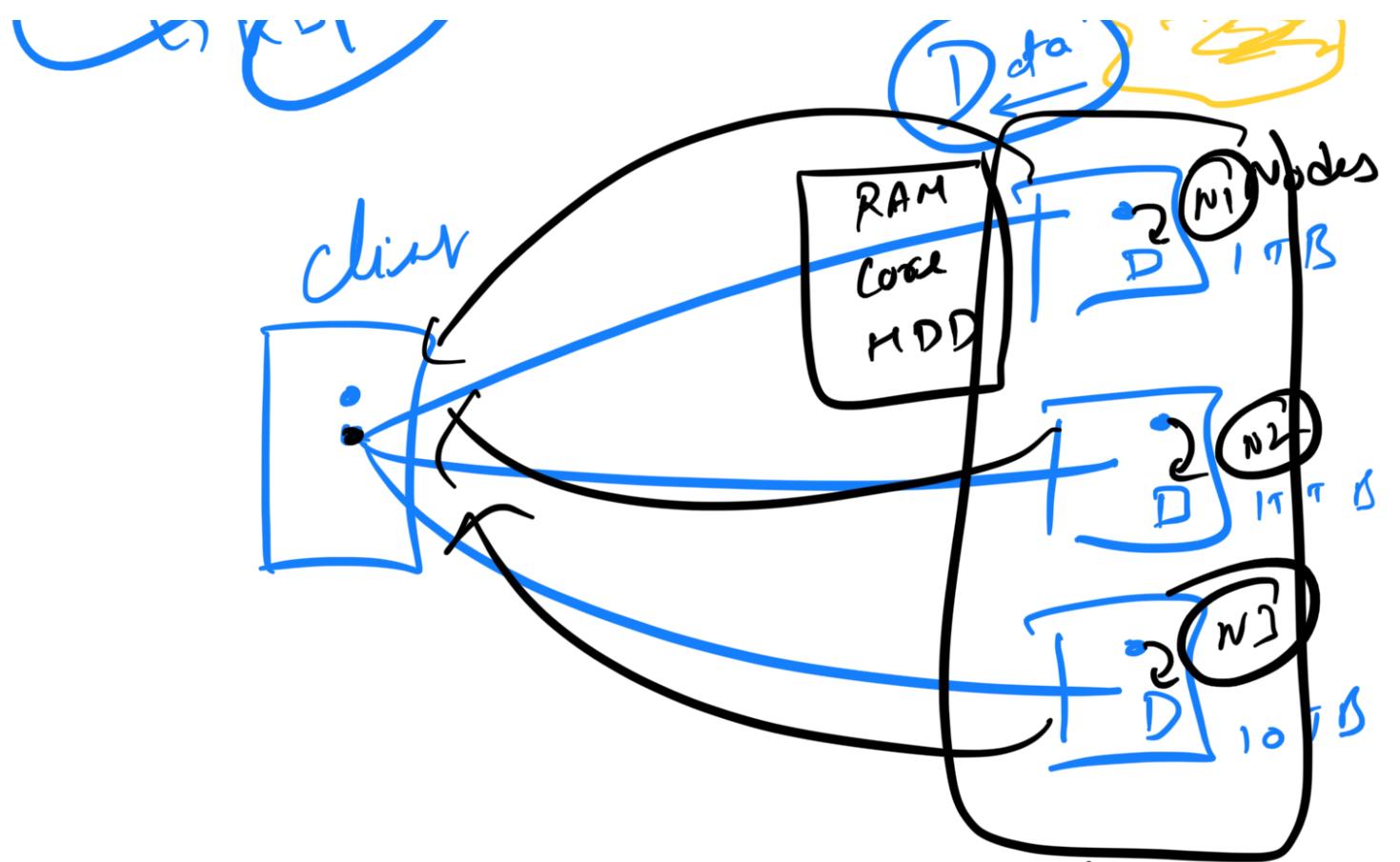


why Hadoop?

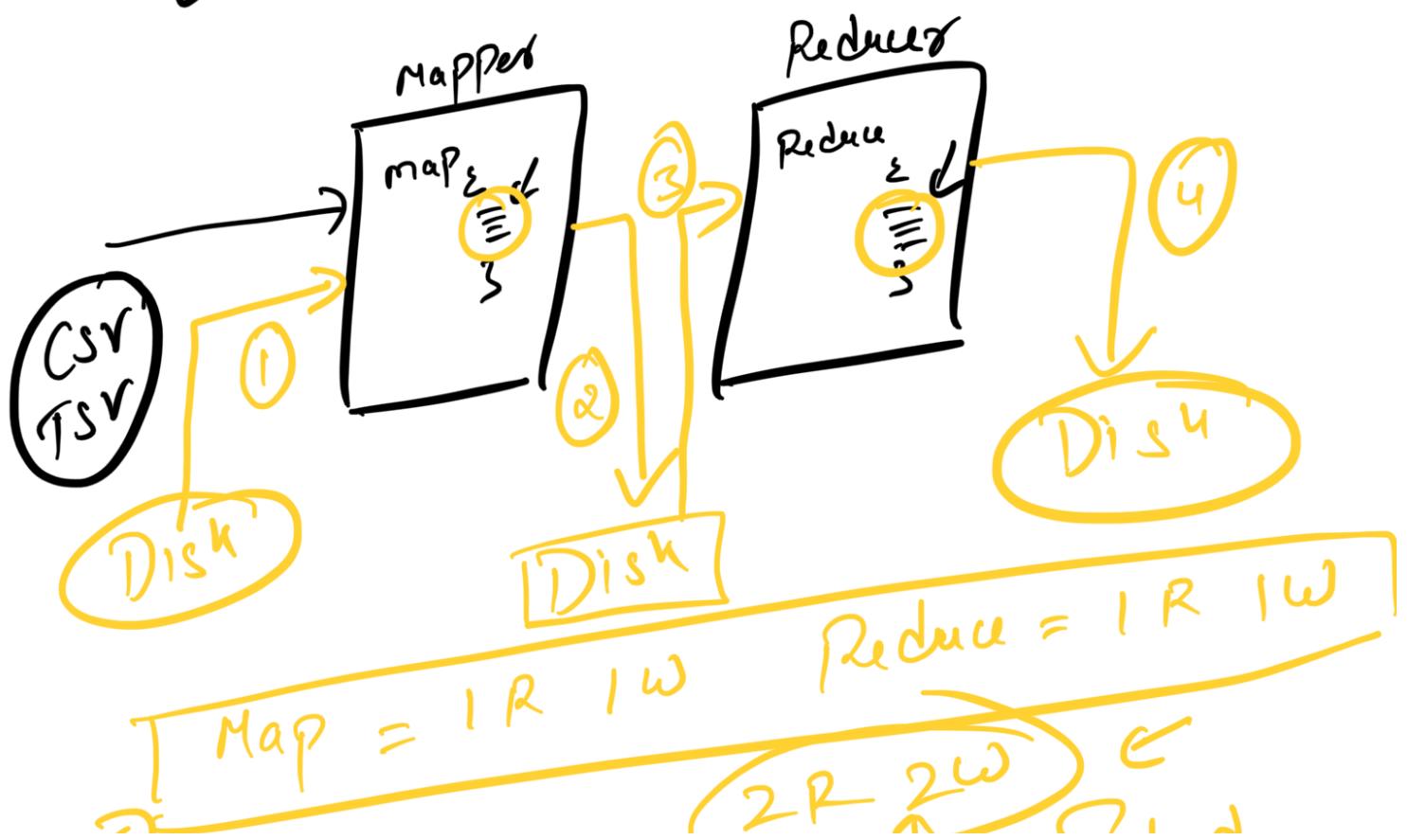
1950 → kB, Punch Cards
 1970 → MB, IBM DB
 2000 → GB,

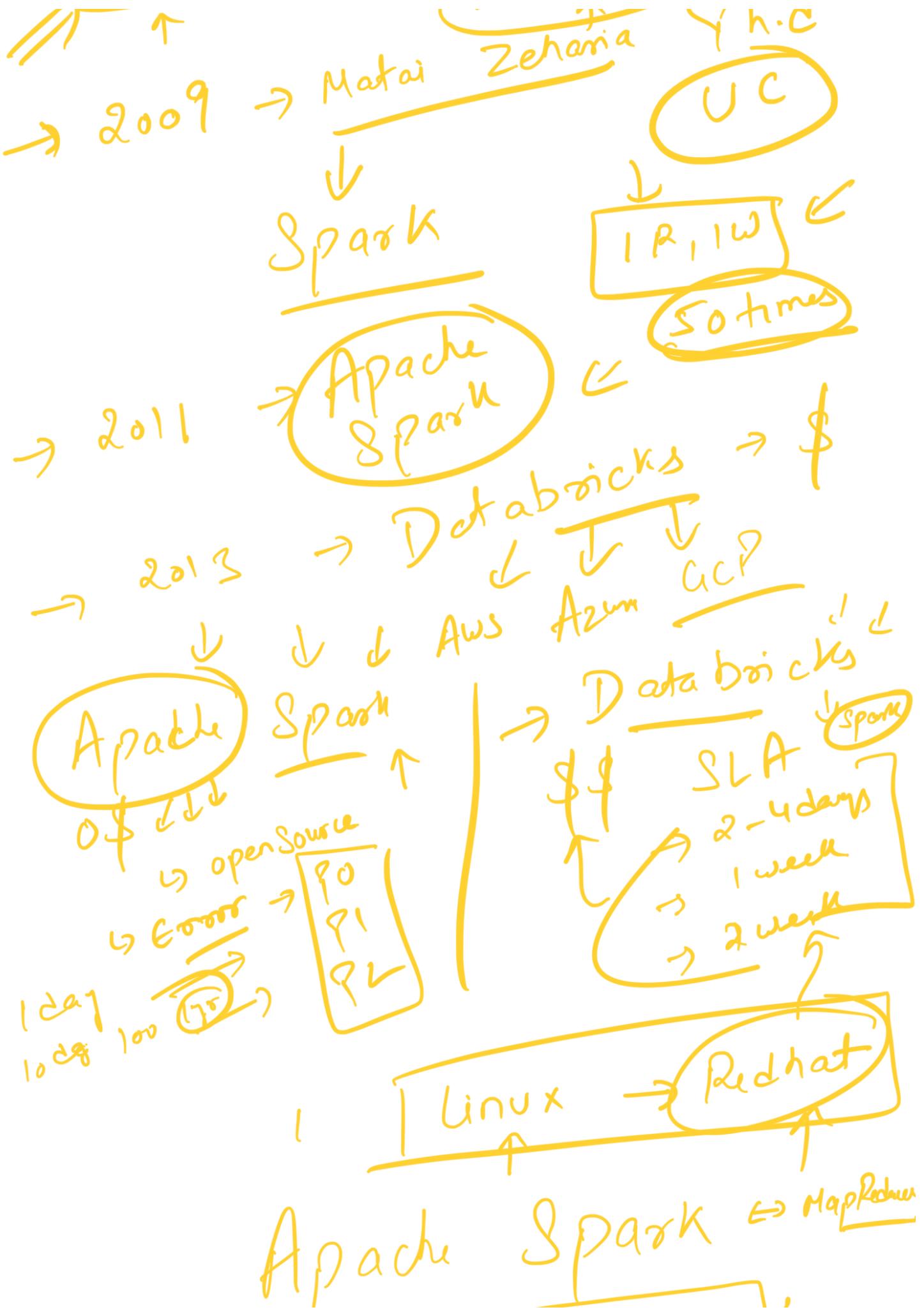






✓ Hadoop \Rightarrow Map/Reduce (Java / Python)

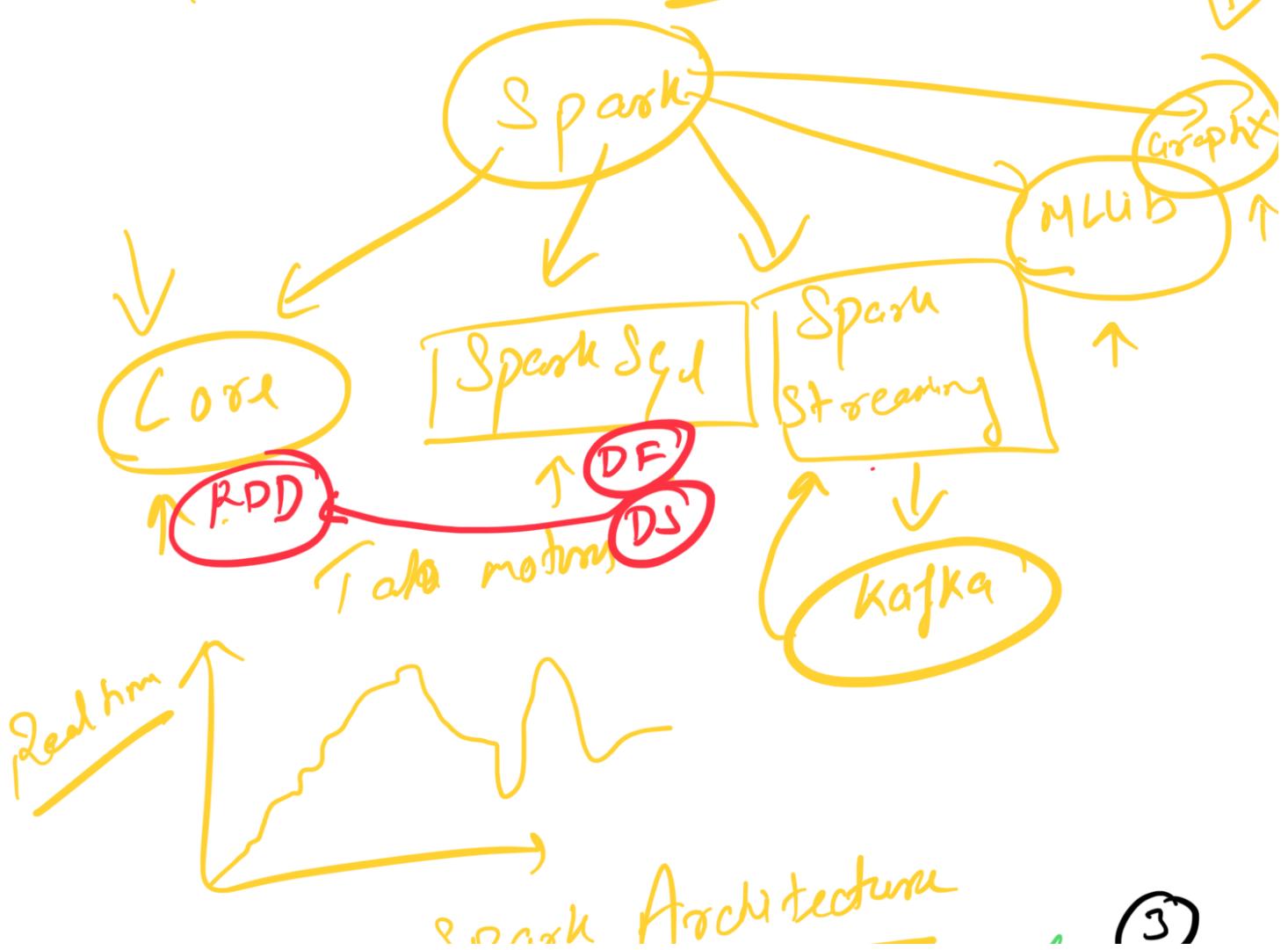




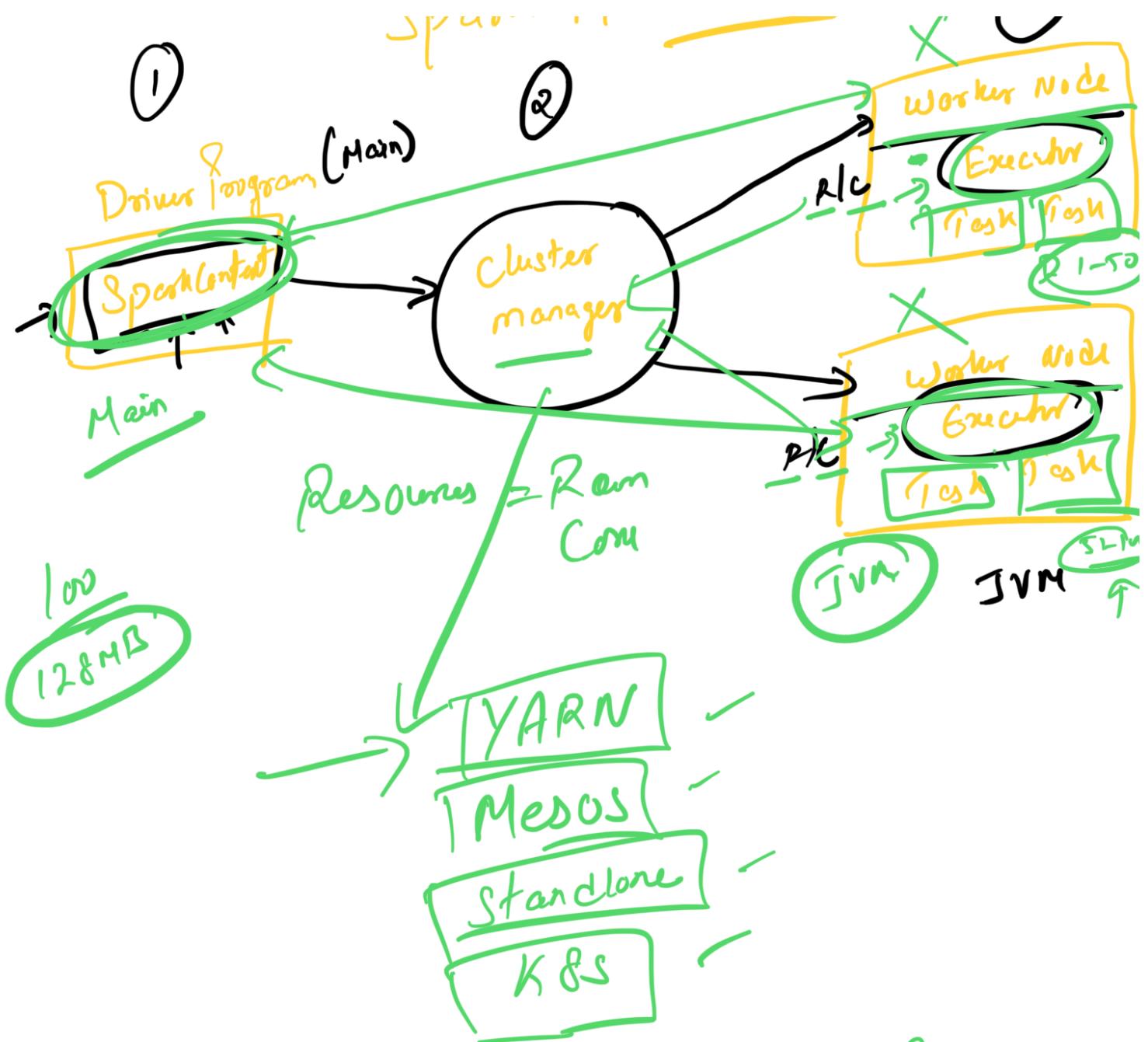
↳ Open Source ¹
 doing large scale data Processing and
 as ^{min} 100 times faster as compared to
 map Reduce.

- ① IR IW
- ② Memory Based Computations

Map Reduce → Java
 Spark → Java, Python, R, Scala



③



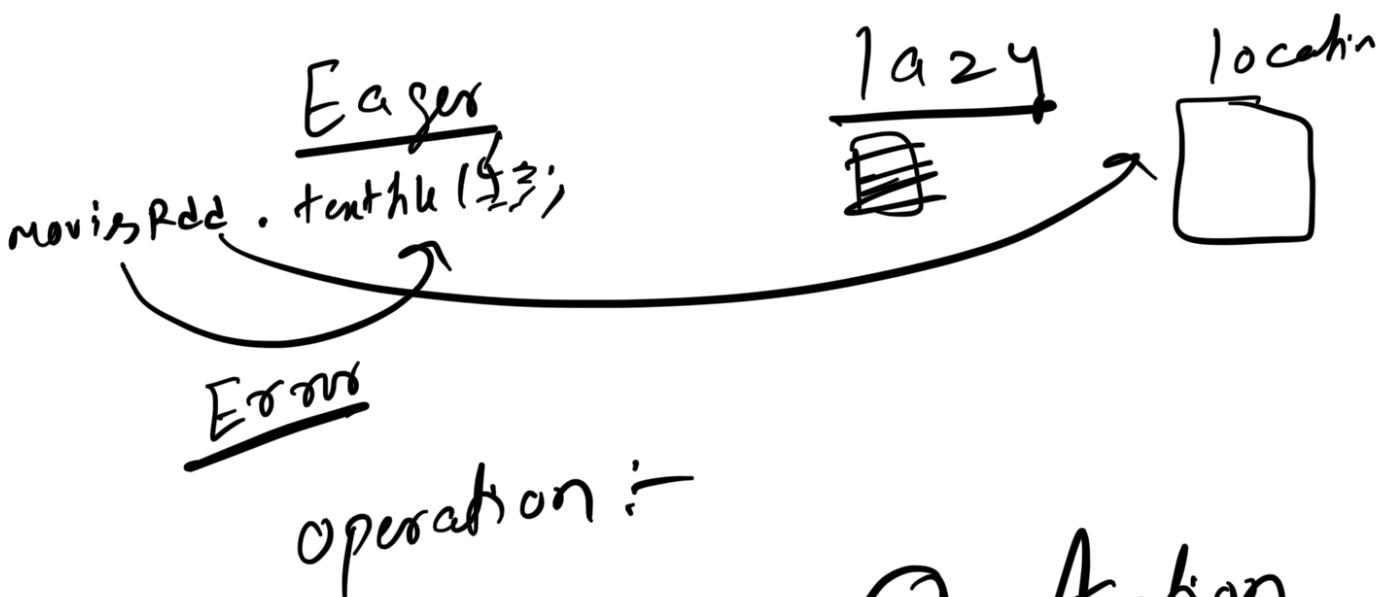
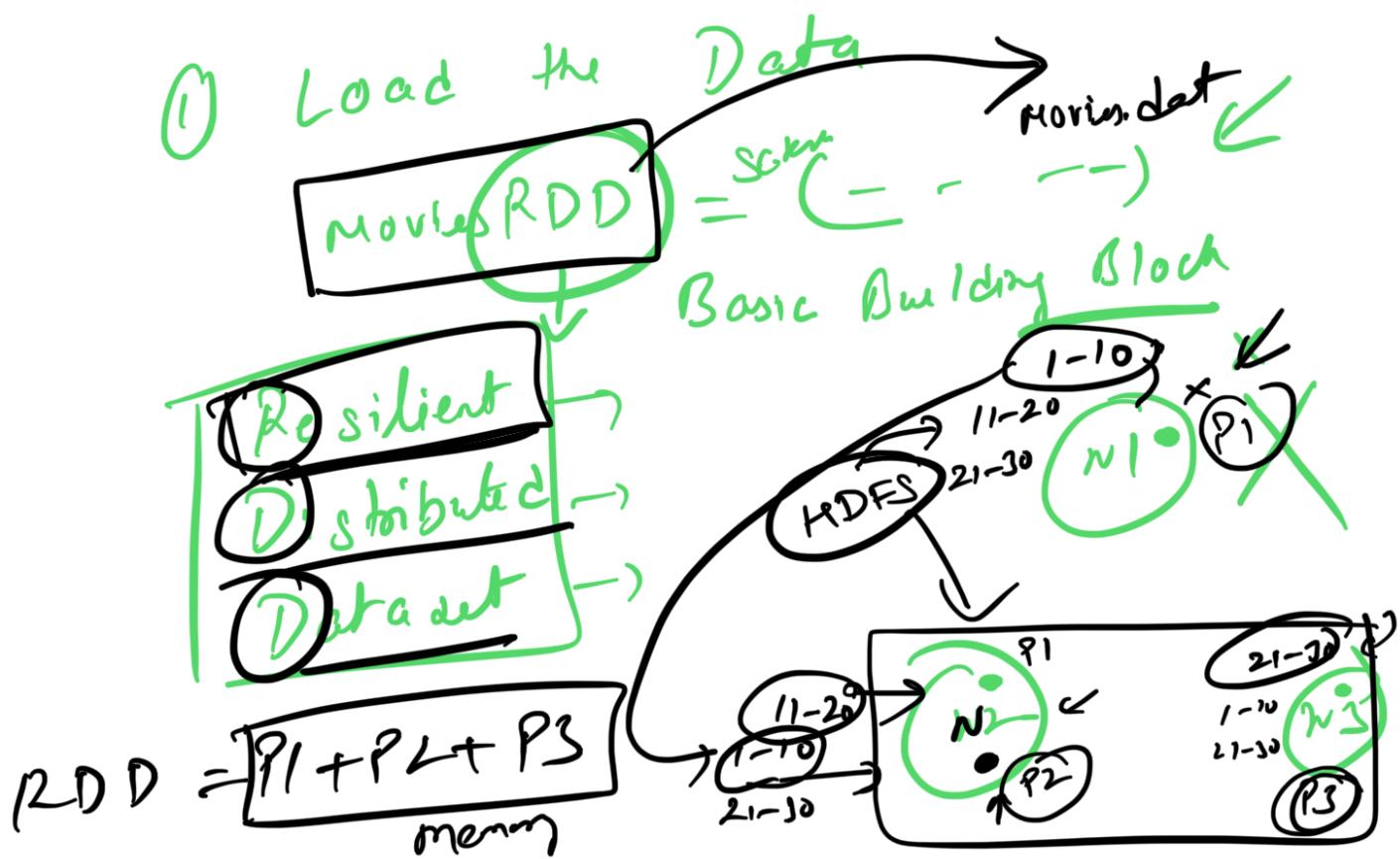
How to interact with Spark?

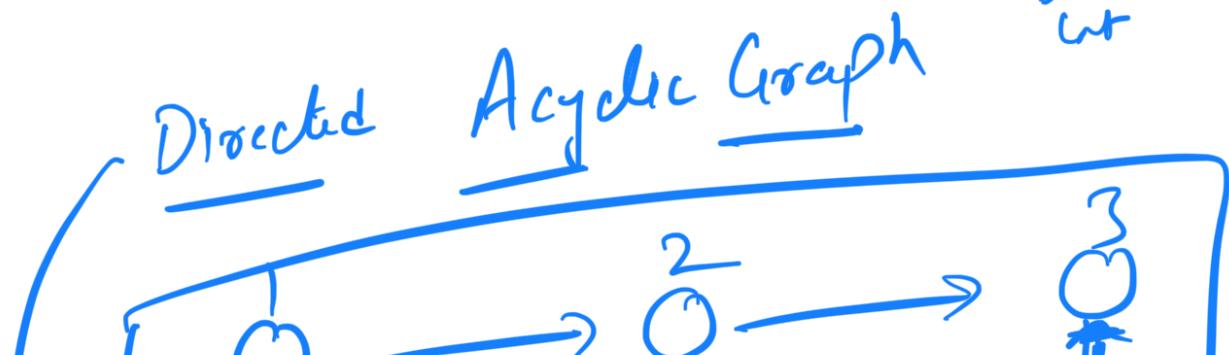
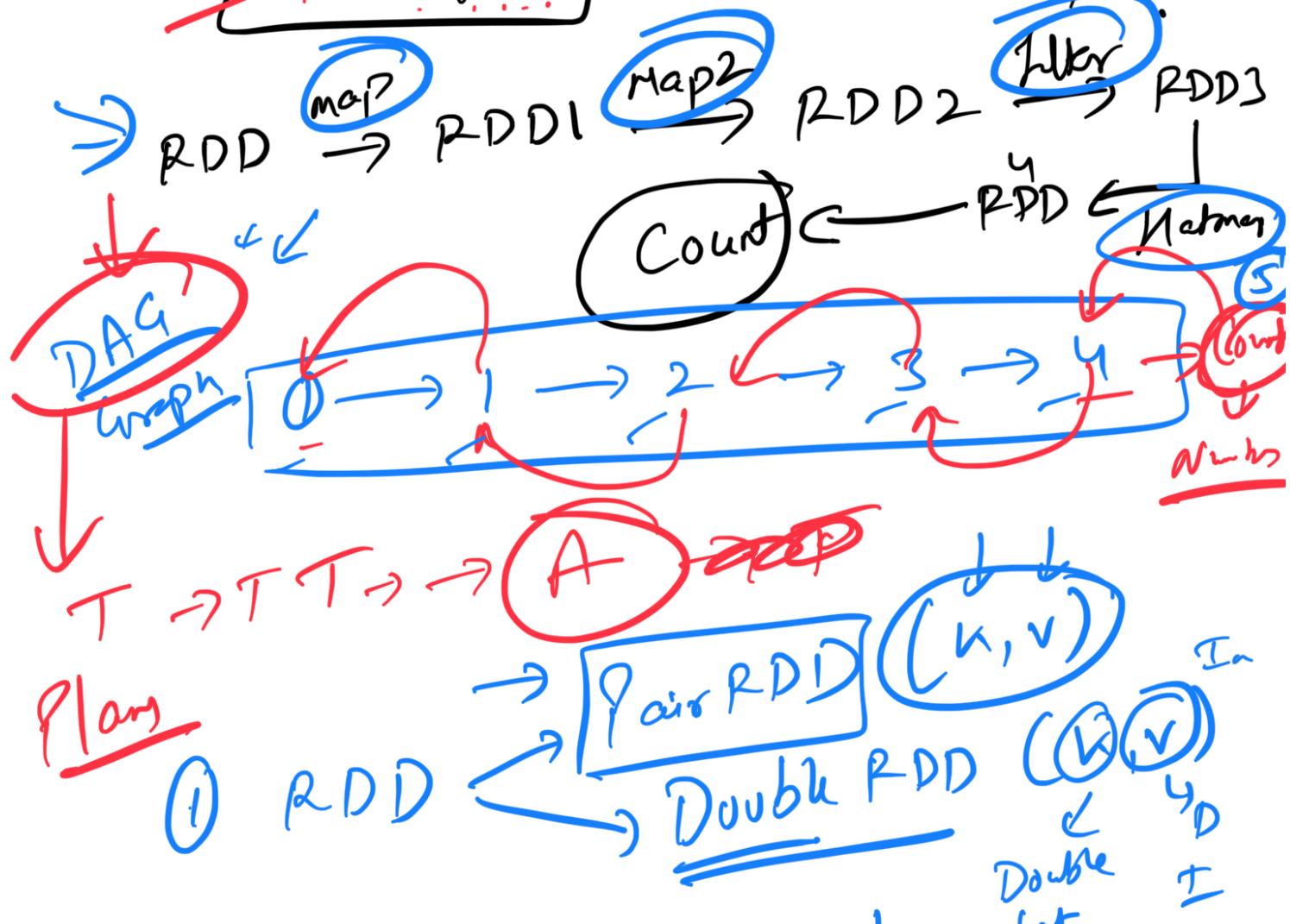
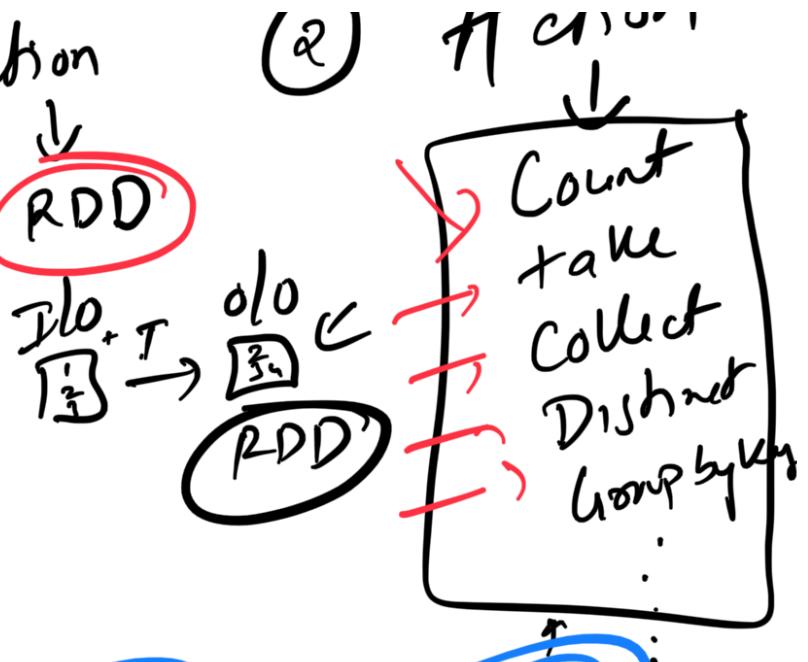
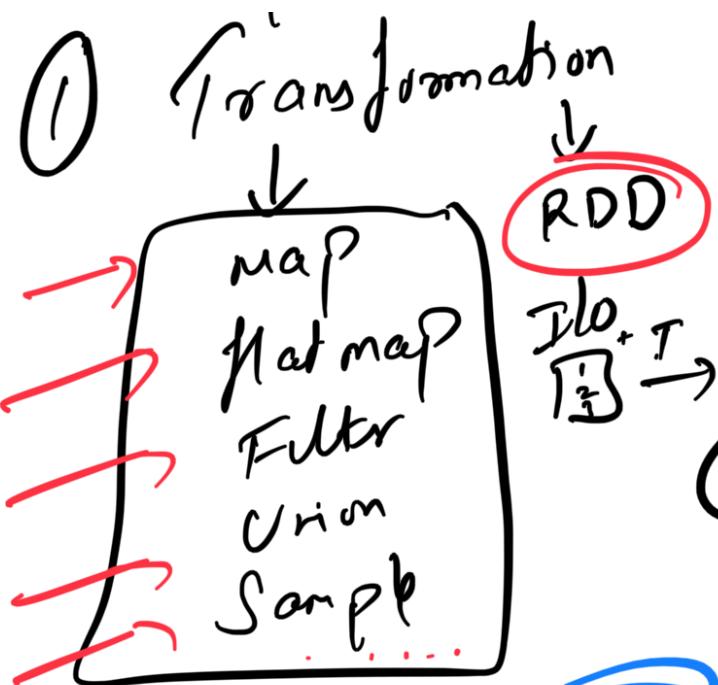
- ✓ ① **Interactive:**
 - Spark-shell → Scala
 - Pyspark → Python
 - SparkR → R
- ② **Application:** → **Spark-Submit** ↓
 - **.py**

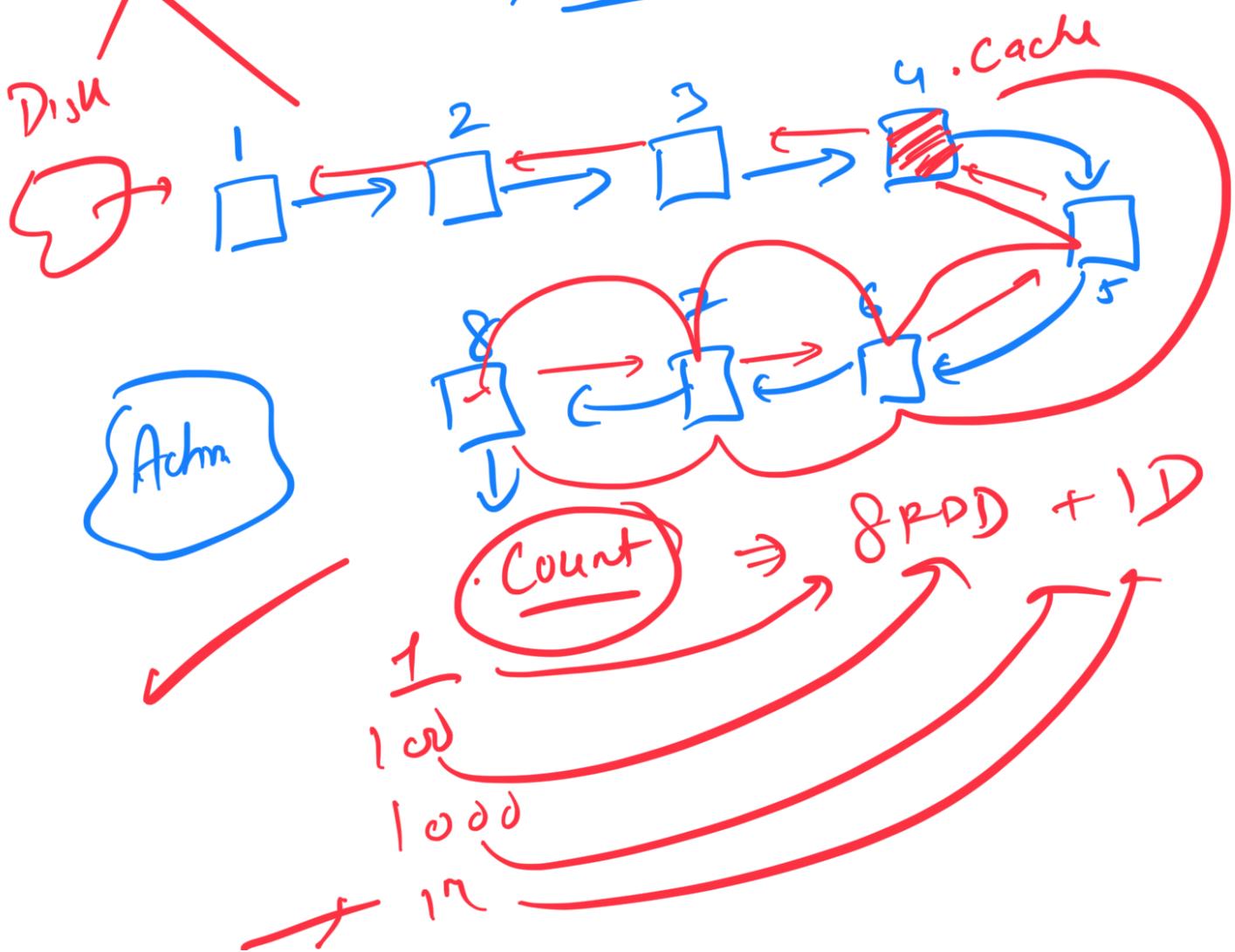
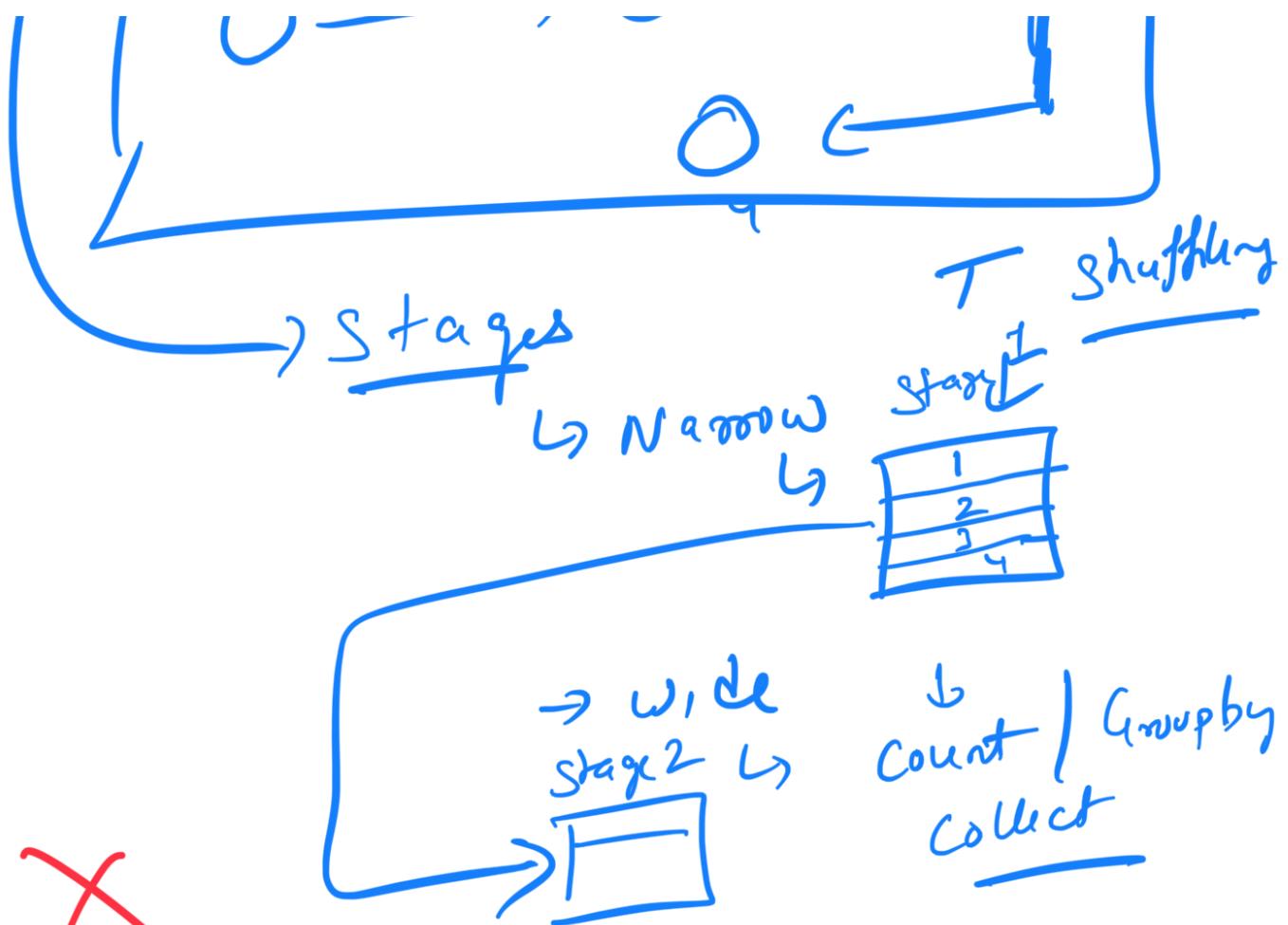
Ques.1 which Comedy movies have most
rankings and return the no of ranking:



An.1 Spark Context







.Cache()