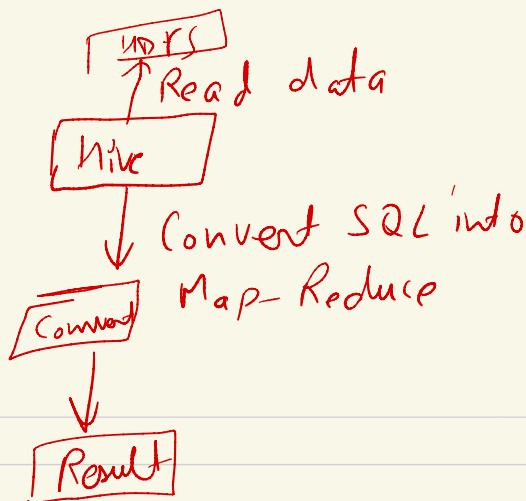



Why Hive?

- ✓ Simple to Use
- ✓ Built on top of Hadoop
- ✓ typical SQL kind of framework
 - ↳ logic will be converted in Map-Reduce code

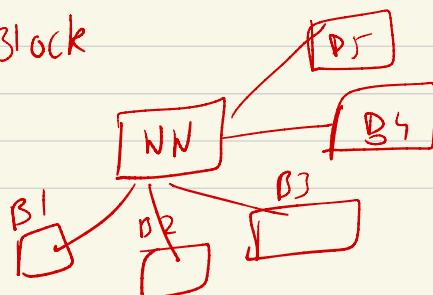


Hadoop

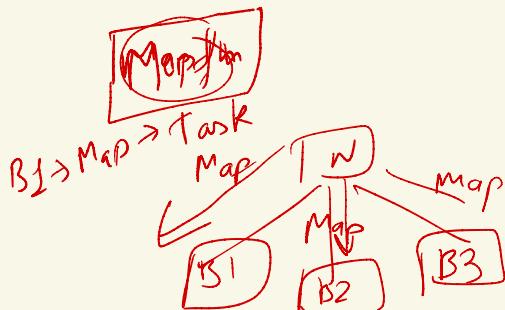
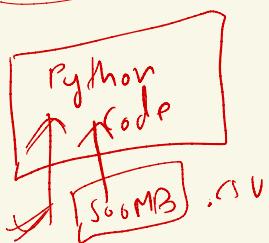
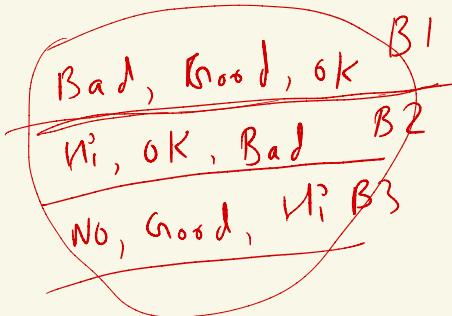
- 1) HDFS → file storage
- 2) Map-Reduce → Processing
- 3.) Yarn → Resource Manager

1TB → Block
3GB

in Local Laptop



$\boxed{ }$ \rightarrow $4 \text{ GB} + 4 \text{ GB}$
 δ external
Vertical Scaling



$M_1 \rightarrow \text{Bad: } 1, \underline{\text{Good: }} 1, \text{ Ok: } 1 \rightarrow R_1$
 $M_2 \rightarrow \text{Hi: } 1, \text{ Ok: } 1, \text{ Bad: } 1 \rightarrow R_2$
 $M_3 \rightarrow \text{No: } 1, \underline{\text{Good: }} 1, \text{ Hi: } 1 \rightarrow R_3$

{ Map \rightarrow Key, Value (output of Map)
Reduce $\sum R_1 + R_2 + R_3$

Sort - Shuffle

$\square R_{d1} \quad 0$

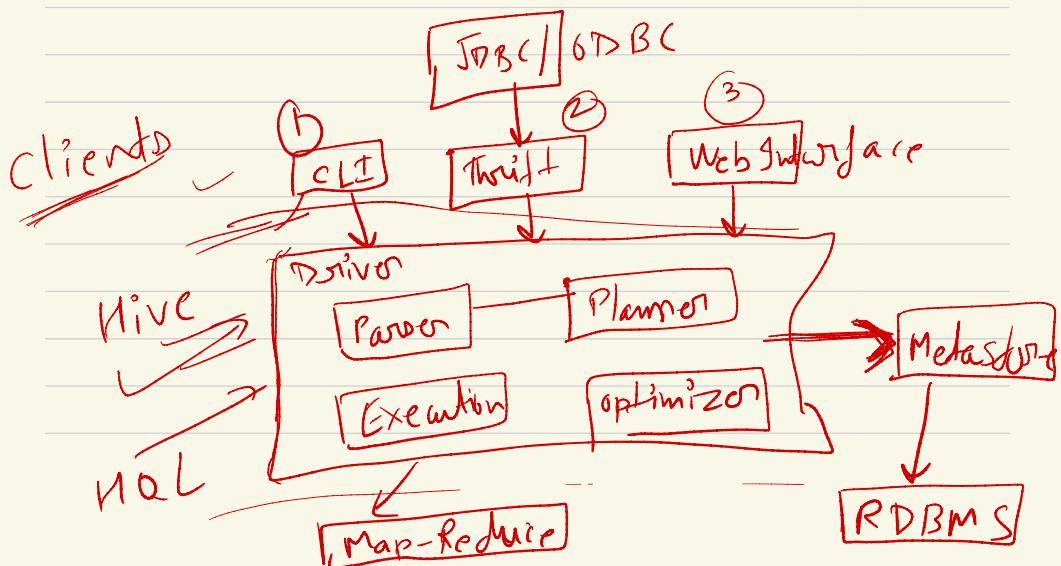
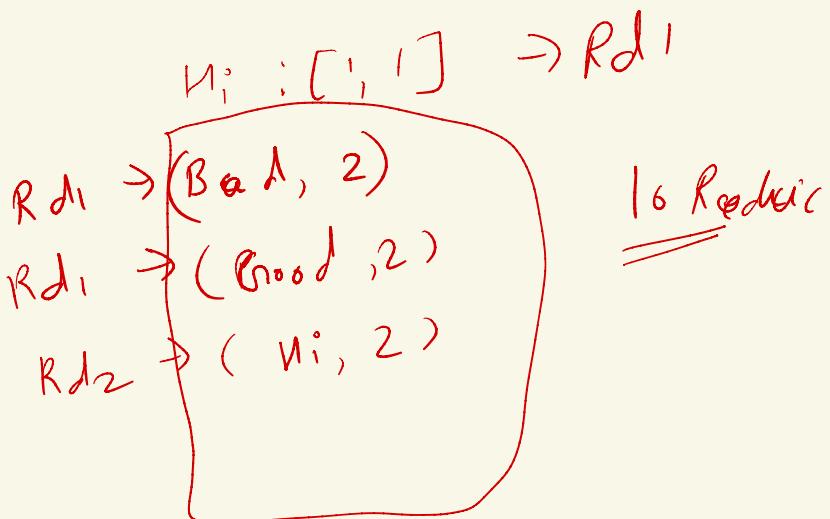
$\text{hash('Bad')} \rightarrow 1 \rightarrow \square R_{d2} \quad 1$
 $\square R_{d3} \quad 2$

Bad: [1, 1]

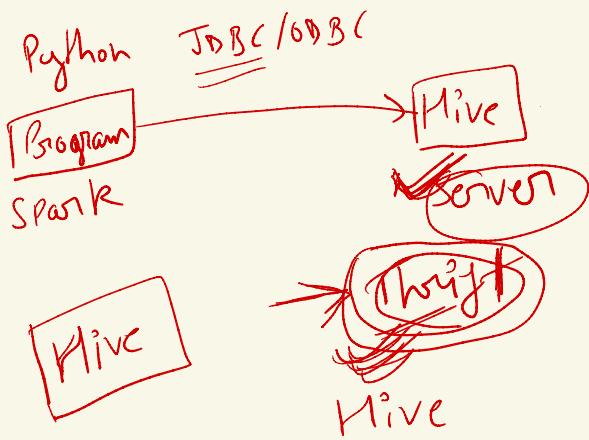
Rd2

Good: (1, 1)

hash("H") \rightarrow 0



HD FS



→ Request -1 → Driver

Driver

- ↳ controller for SQL statements
- ↳ (multiple) sessions for query
- ↳ maintain life cycle of SQL
- ↳ maintain metadata for execution
- ↳ collect output and display

Sqoop → Syntax & Semantics

↳ execution plan

↳ optimization

→ Processing

Parsing / Compilation

↳ Syntax check

↳ Execution plan

↳ Steps to get the output

↳ Raise compile time errors

Optimize

NQL → E1

↓ E2

↓ E2

Select * from Tp1

where =

group by =

(sub query)

↳ Compare execution plans

↳ Calculate Cost

↳ Execution plan

↳ A h

↳ Try to place

group by } Where

where } group by

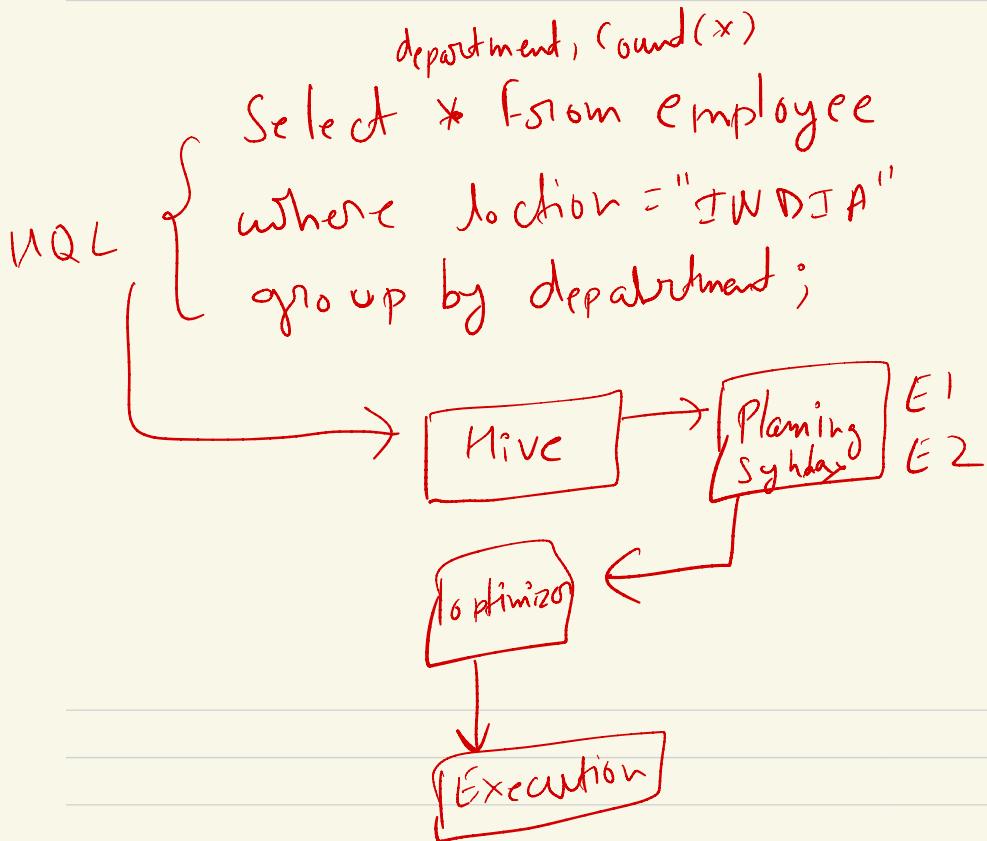
E1

E2

Cost1 Cost2

or combine them on matching together

Execution plans & PTA



employee → 2⁶⁶K

E1 → where ; group by } C1
E2 → group by , where } C2

C1 > C2

C1 < C2

$line1 = \underline{\text{"Hi, Bad, Good"}}$

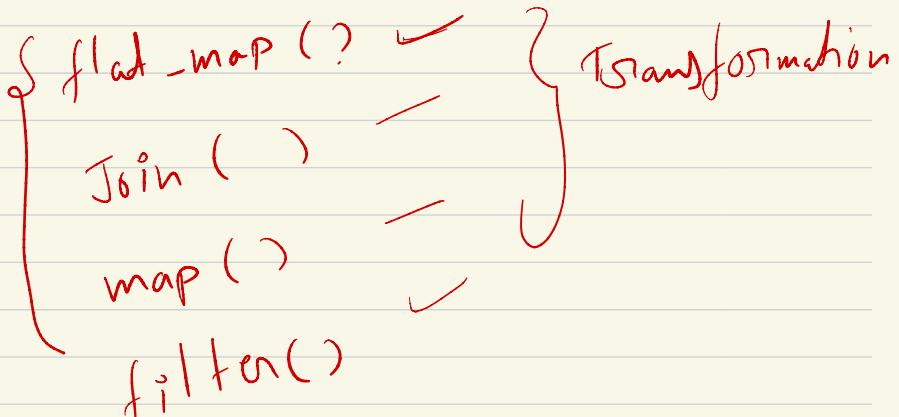
one \rightarrow many flat-map($line1 \rightarrow$)

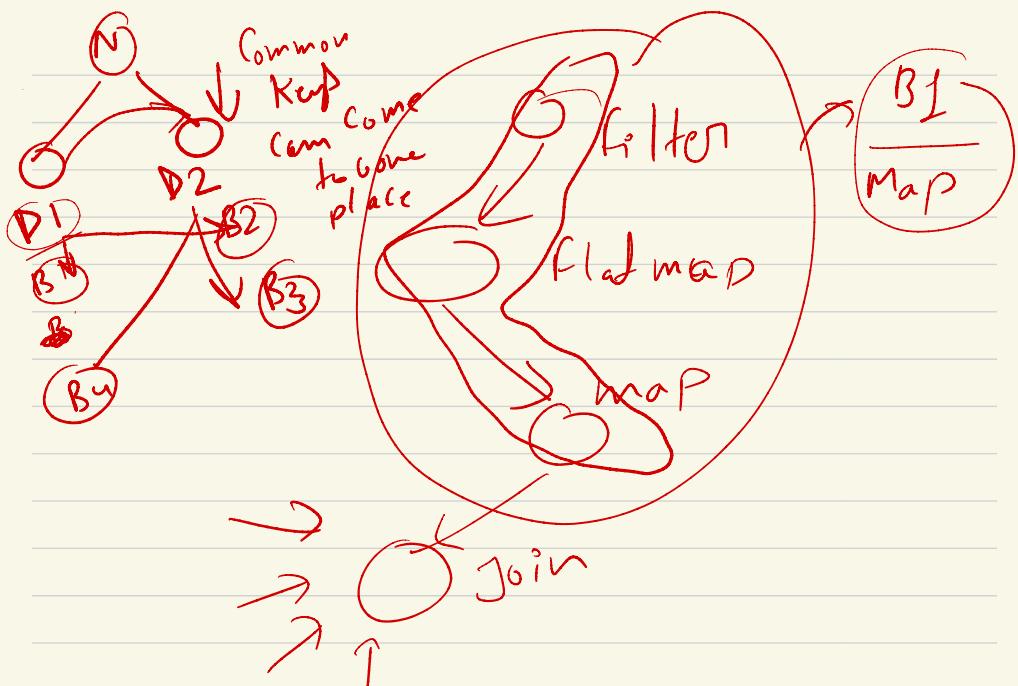
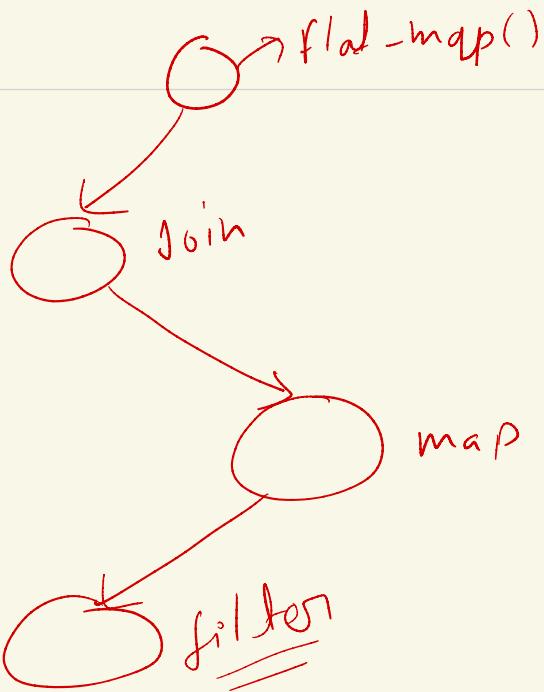
$R1 \rightarrow \{ "Hi", "Bad", "Good" \}$

map($R1 \rightarrow$)

$R2 \rightarrow Hi:1, Bad:1, Good:1$

joins, group by





map1

O D1

B1 \rightarrow id, salary

1,	1000
2,	2000
3,	5000

Map2

B2 \rightarrow id, Salary

1,	1000
3,	5000
1,	2000

B3 \rightarrow Dept-name, id

Finance, 1

IT, 2

B4 \rightarrow Dept name, id
OPS, 3

Employee \leftarrow Department

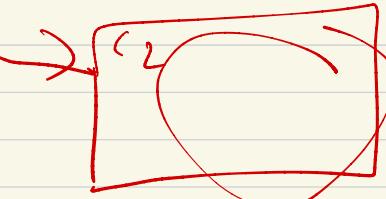
Employee

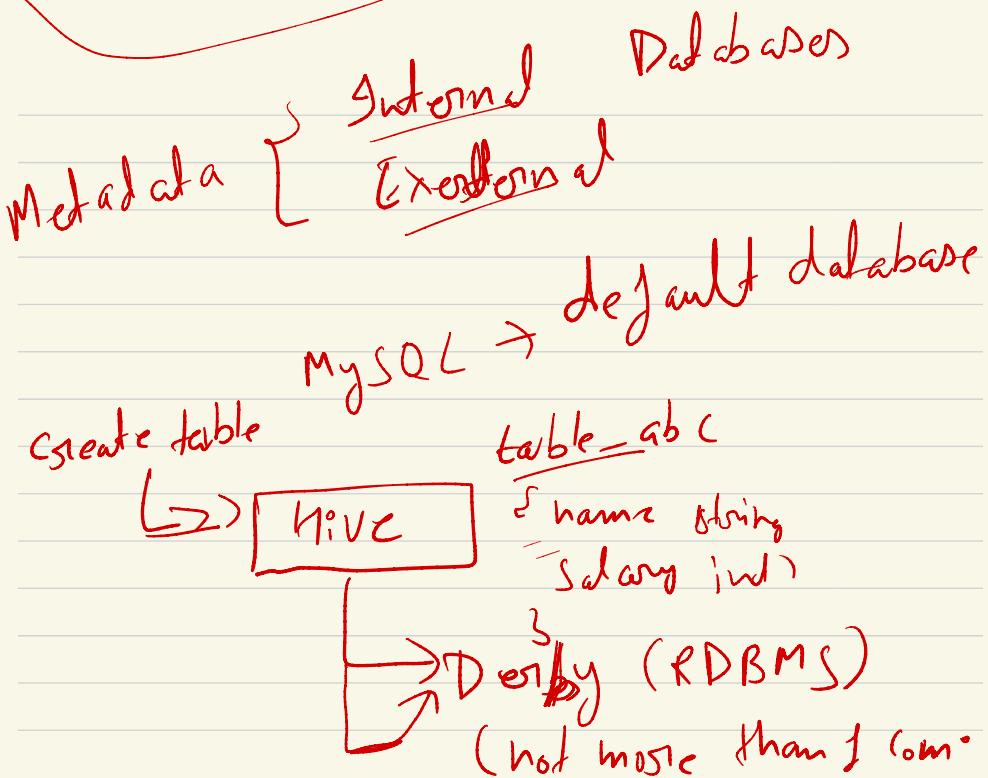
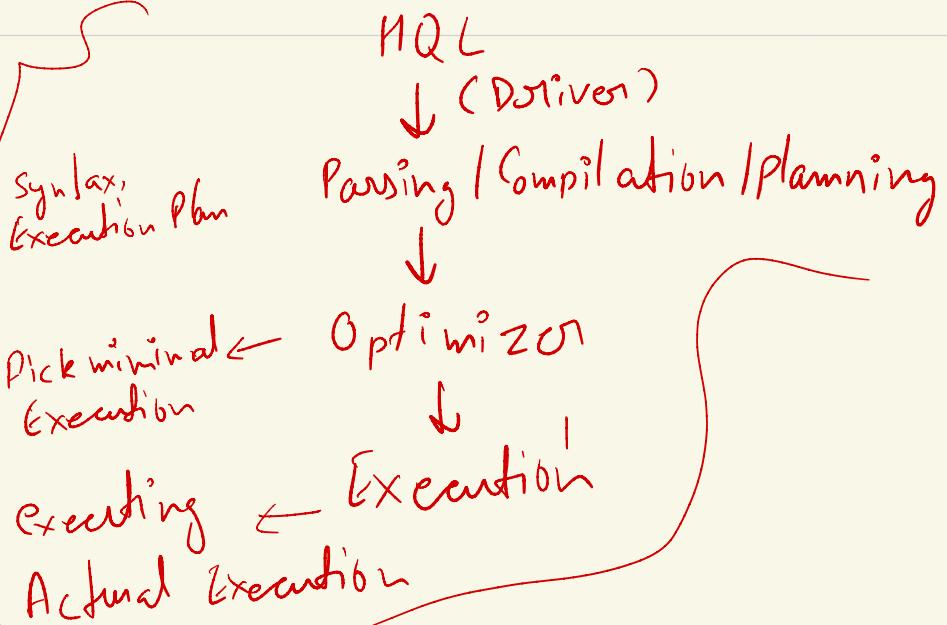
Dept-1

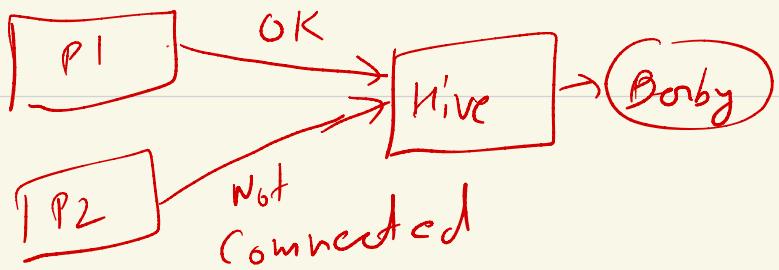
Finance, 1

1,	1000
14,	5000
1,	2000

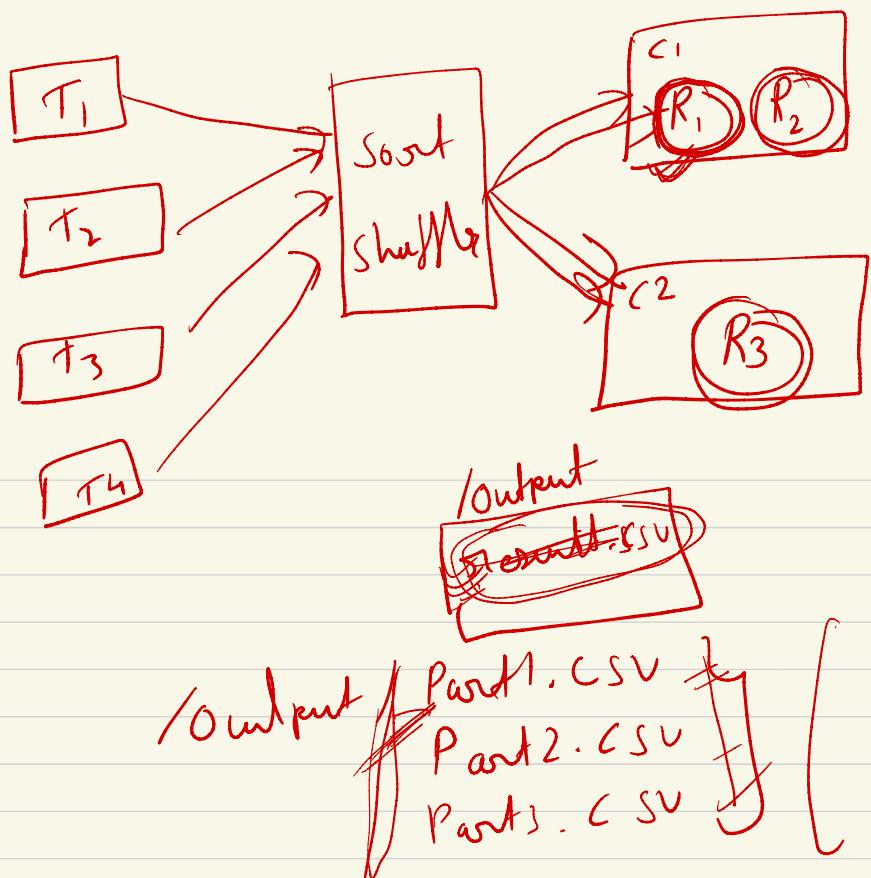
YARN

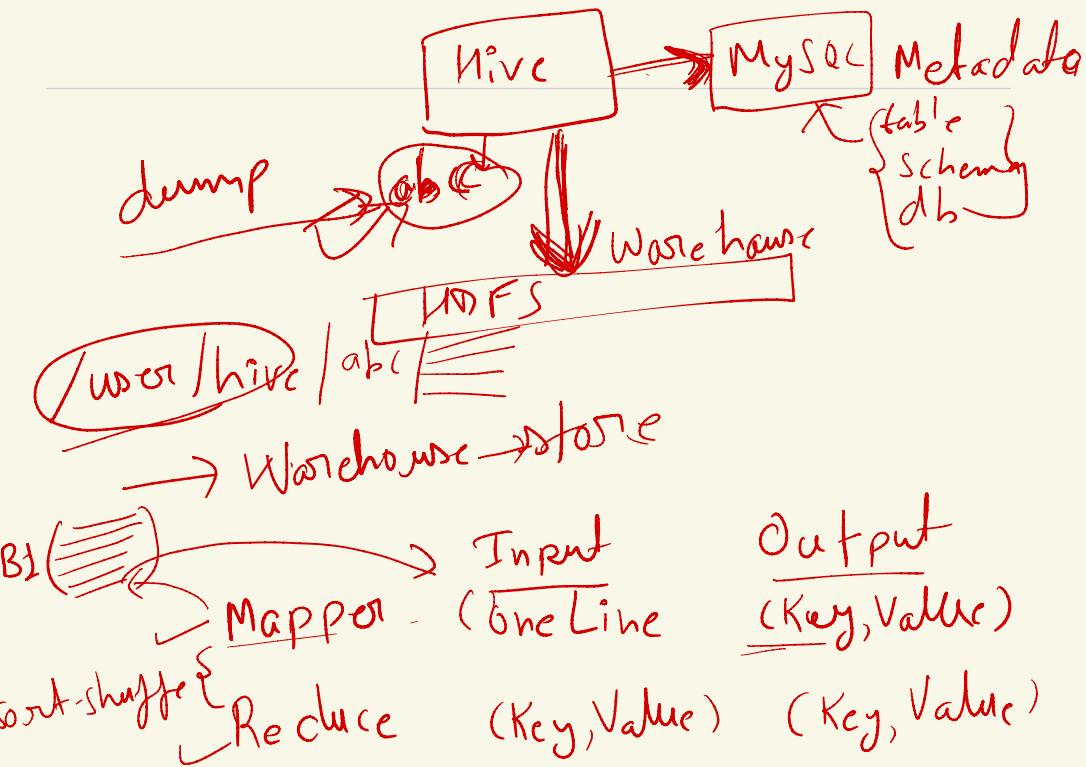






- ↳ External DB
- ↳ Metadata Back up
- ↳ More Multiple Concurrent Connection
- ↳ Expose to external use case(s)





Mapper()

(B1)

```
{
    for (0 till last line)
        {
            "Hi, Hello, Bad"
            ↗Hi : 1, Hello : 1, Bad : 1
            "OK, Hello, Bad"
            ↗OK : 1, Hello : 1, Bad : 1
        }
}
```

(1)

(2)

Map $\rightarrow (K, V)$ \rightarrow Reduce

Disk

Part2.csv

OK: 2 $\downarrow R_2$

Bad : 2

Hello : 2

Bad: [1, 1]

Hello: [1, 1]

OK: [1]

Part1.csv

Shard - shard/1c

0-10

R₁

R₂

10-20

hash(Bad) = 5

hash(Hello) = 9

hash(OK) = 15

employee

$d_1 \rightarrow \underline{\text{Read}}$

$d_2 \rightarrow \underline{\underline{\text{filter}}}(d_1)$

