

---

2018

---

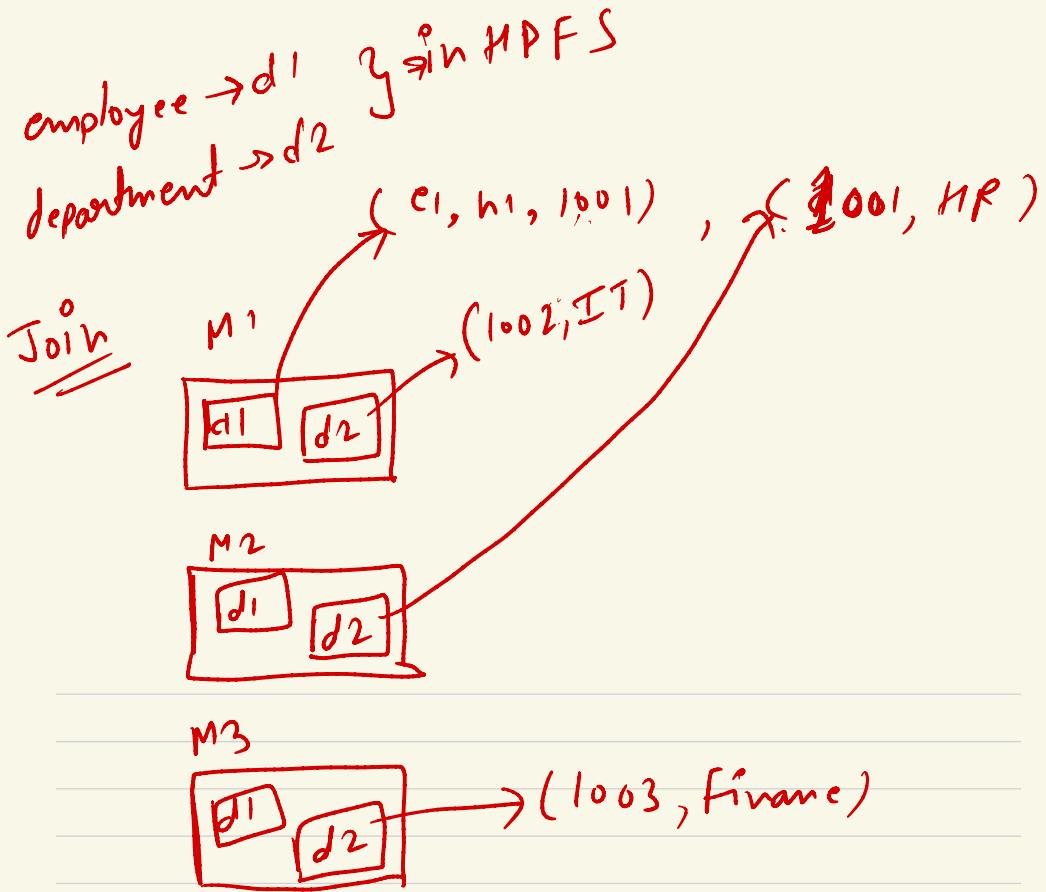
---

---

---



## Bucketing in Hive



## What is Bucketing

Basically, decomposing table datasets into more manageable parts.

## Employee Table

emp-id, emp-name, dept-id

1, S1, 100

2, S2, 103

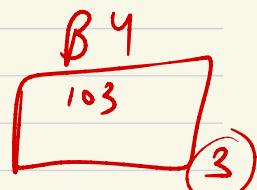
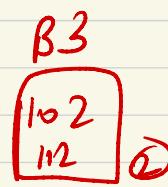
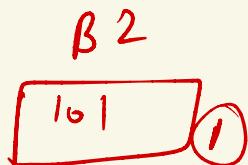
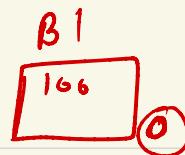
3, S3, 101

4, S5, 102

bucket columns (dept-id)  $\rightarrow$  cluster-id

buckets = 5  $\xrightarrow{\text{Key \% num-of-buckets}}$

hash-funct(key) = ?



$$① \text{hash}(100) = 100 \% 5$$

$$= 0$$

$$② \text{hash}(103) = 103 \% 5$$

$$= 3$$

$$③ \text{hash}(101) = 101 \% 5$$

$$= 1$$

$$④ \text{hash}(102) = 102 \% 5$$

Country id partition column

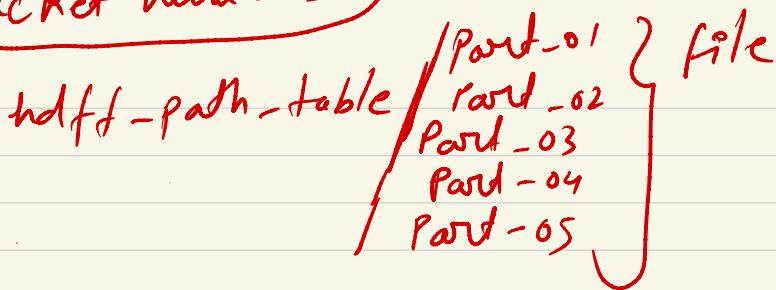
Now it looks in warehouse!

hdfs-path / Country = INDIA

hdf-path / Country = USA

x) Bucket representation for a table without partition?

(Bucket num = 5)



x) Bucket representation if a partitioned table!

hdf-path-table / country = INDIA /  $p_0$   $p_1$   $p_2$   $p_3$

hdfs-path-table / country = USA /  $p_0$   $p_1$   $p_2$   $p_3$   $p_4$

## optimized joins in Hive

### Map Join (Broadcast Join)

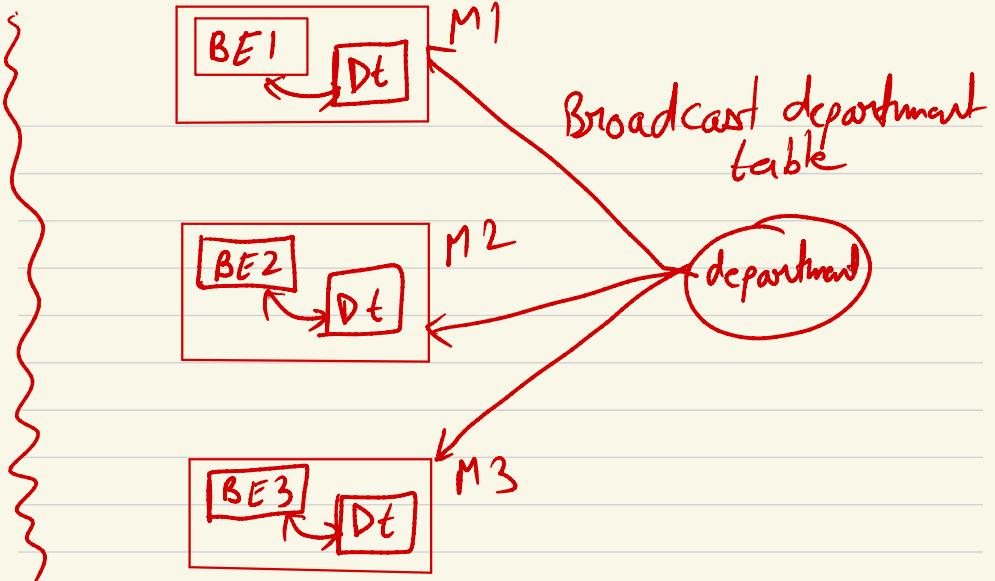
employee JOIN department

employee = 20K records

department = 10 records

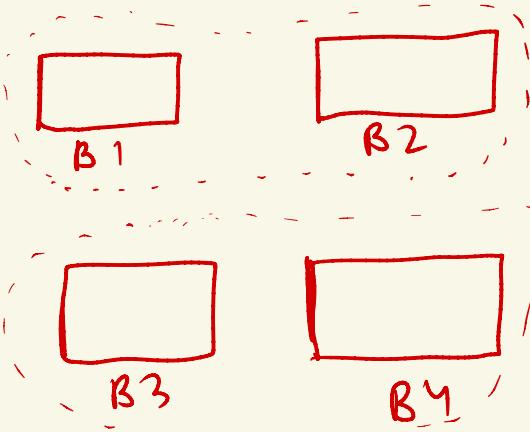
### Map Task Execution

HDFS

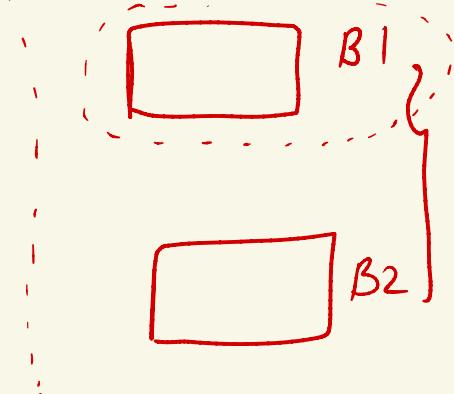


## Bucket Map Join

Employee Table Buckets  
4 buckets



Department Table  
2 only



## Condition for Bucket Map Join

- ① Tables are bigger in size
- ② Both tables are bucketized on the same join column.
- ③ Number of buckets in both of the table are multiple of each other

## City Table (2 B)

(B1)

cityid	cityN
3	NY
1	SJ

## Sales\_Table (2 B)

(cityid) (Join)

(B1)

cityid	saled
1	500
3	6000
1	400

(B2)

cityid	name
2	SF
4	LA

cityid	Saled
4	50
2	200
4	45

(B2)

## Sorted - Merge-Bucket Join

arr1 = [ ①, 5, 9, 11 ] X  
i → i → i → i

arr2 = [ ②, 4, 7, 12, 20, 21 ]  
j → j → j → j → j → j

result = [ 1, 2, 4, 5, 7, 9, 11, 12, 20, 21 ]

$O(n+m) \rightarrow$  Linear Complexity

$$L\{ \text{corr1} = [ \begin{matrix} 1, & 5, & 6, & 9, & 9, & 11, & 12, & 12 \end{matrix} ]$$

①      i → j

$$R\{ \text{corr2} = [ \begin{matrix} 3, & 4, & 6, & 9, & 9, & 10, & 11 \end{matrix} ]$$

②      j → i

$$\text{result} = [1, 3, 4, 5, (6, -, -), (9, d), (i, d)]$$