



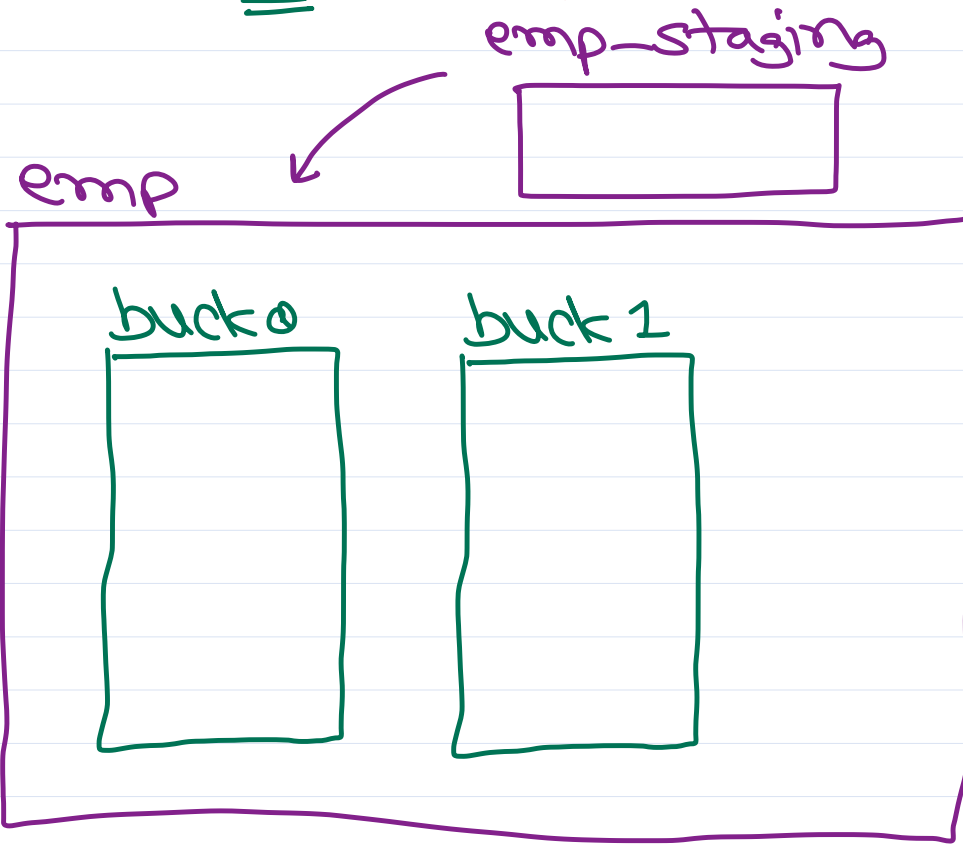
Big Data Technologies

Trainer: Mr. Nilesh Ghule.



Bucketing

emp table →
clustered by (ename)
into 2 buckets.



ename

ename
FORD \rightarrow hashCode() % 2 \rightarrow 0 \rightarrow buck0
 \rightarrow 1 \rightarrow buck1

JAMES \rightarrow hashCode() % 2 \rightarrow 0 \rightarrow bucket
 \rightarrow 1 \rightarrow bucket

KING \rightarrow hashCode() % 2

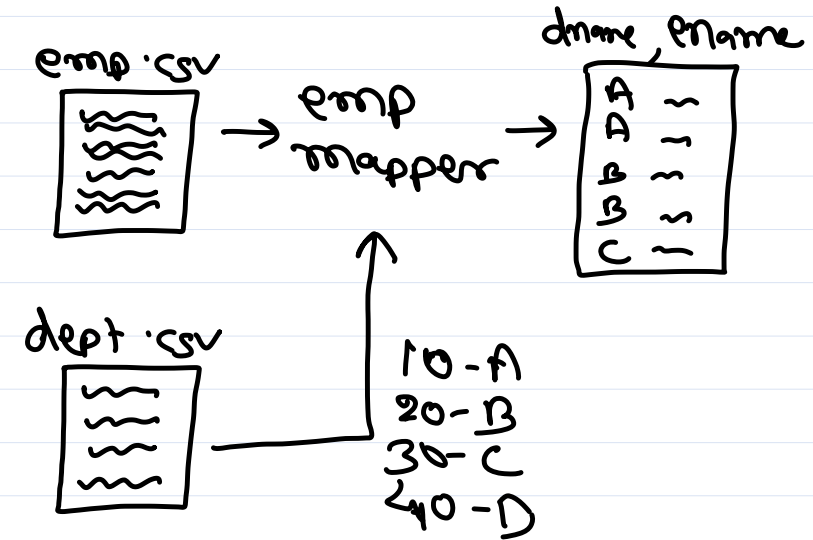
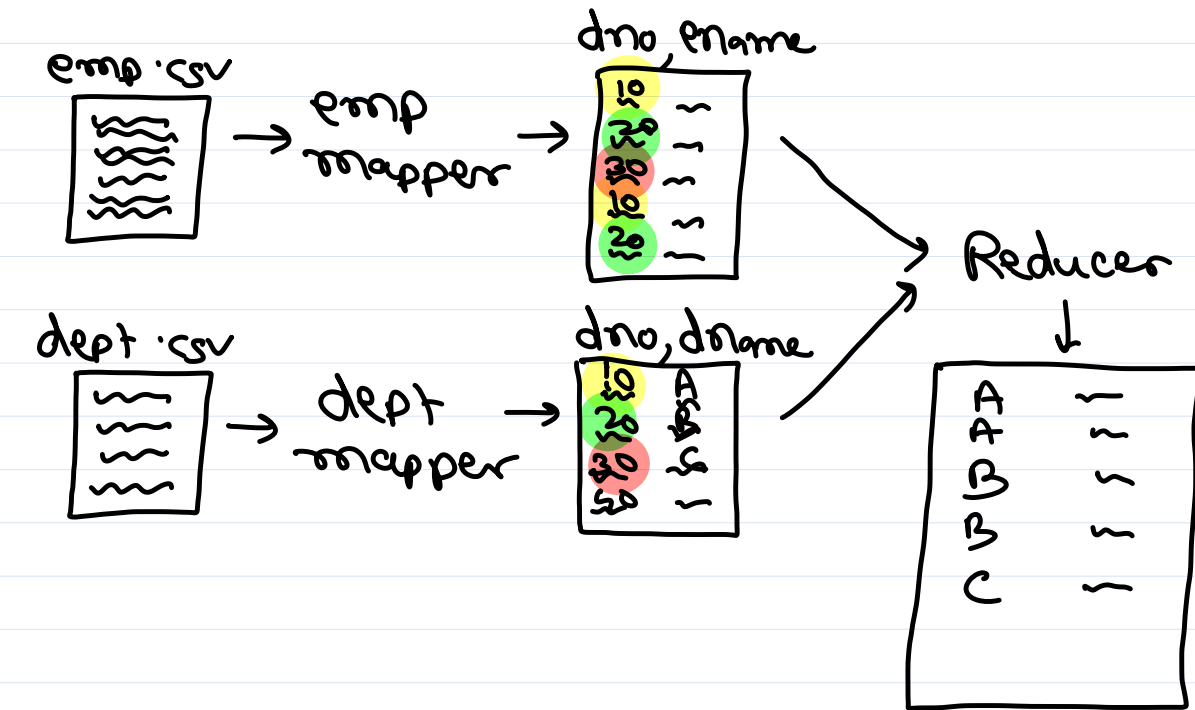
- $\nearrow 0 \rightarrow \text{bucket}$
- $\searrow 1 \rightarrow \text{bucket}$

when DML ops are performed on bucketed tables, num of reducers are set to num of buckets.

In Hive 3.x, when DML ops are performed on non-bucketed table, one bucket is considered for the table (one reducer is used).



Reduce side joins vs Map side joins





Thank you!

Nilesh Ghule <nilesh@sunbeaminfo.com>

