

Agenda

Understand how MapReduce tasks are scheduled and executed using YARN.

Explore the communication between HDFS and YARN.

Work through examples to solidify these concepts.

What is OLTP?

Why it is not suitable for Analytics

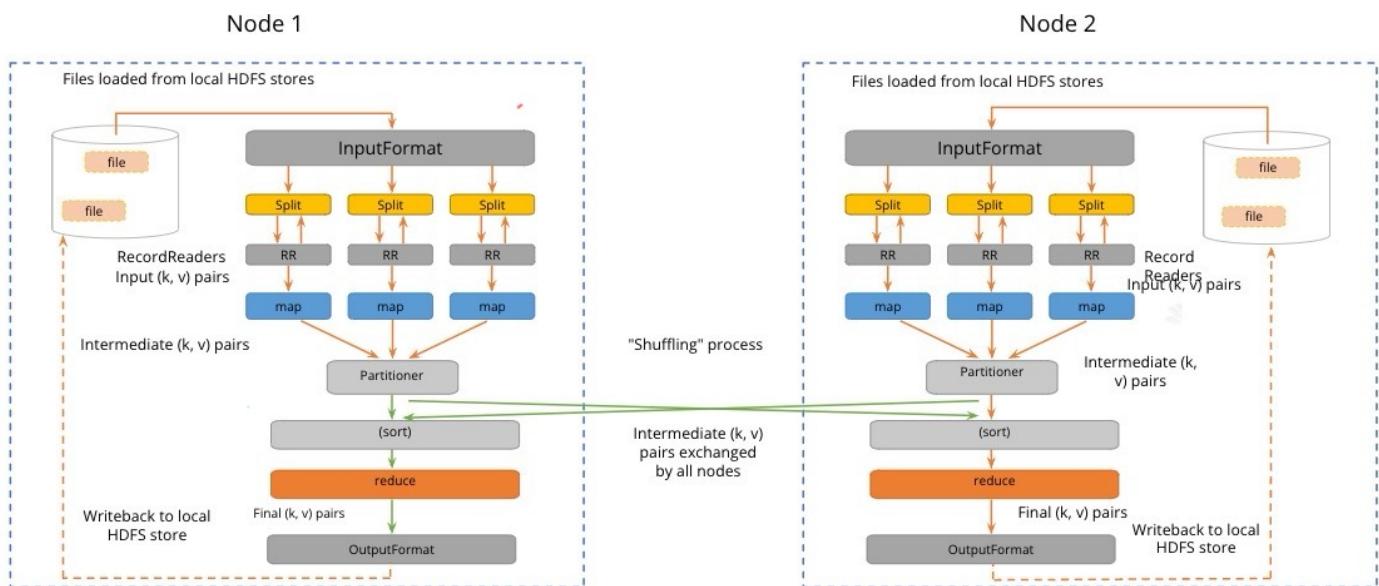
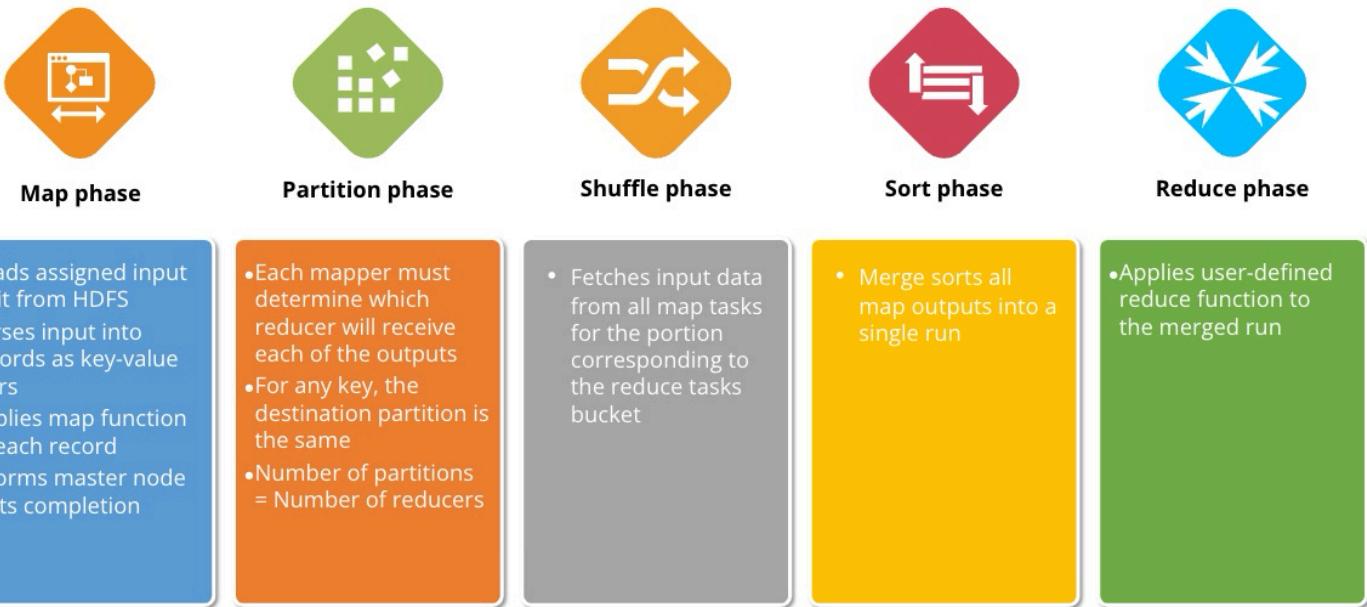
Introduction to Data warehousing

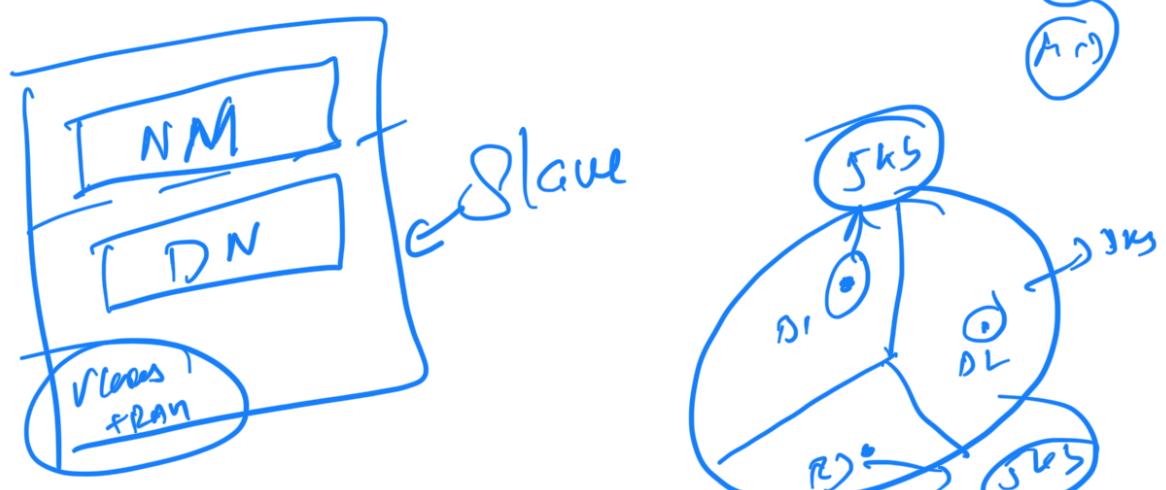
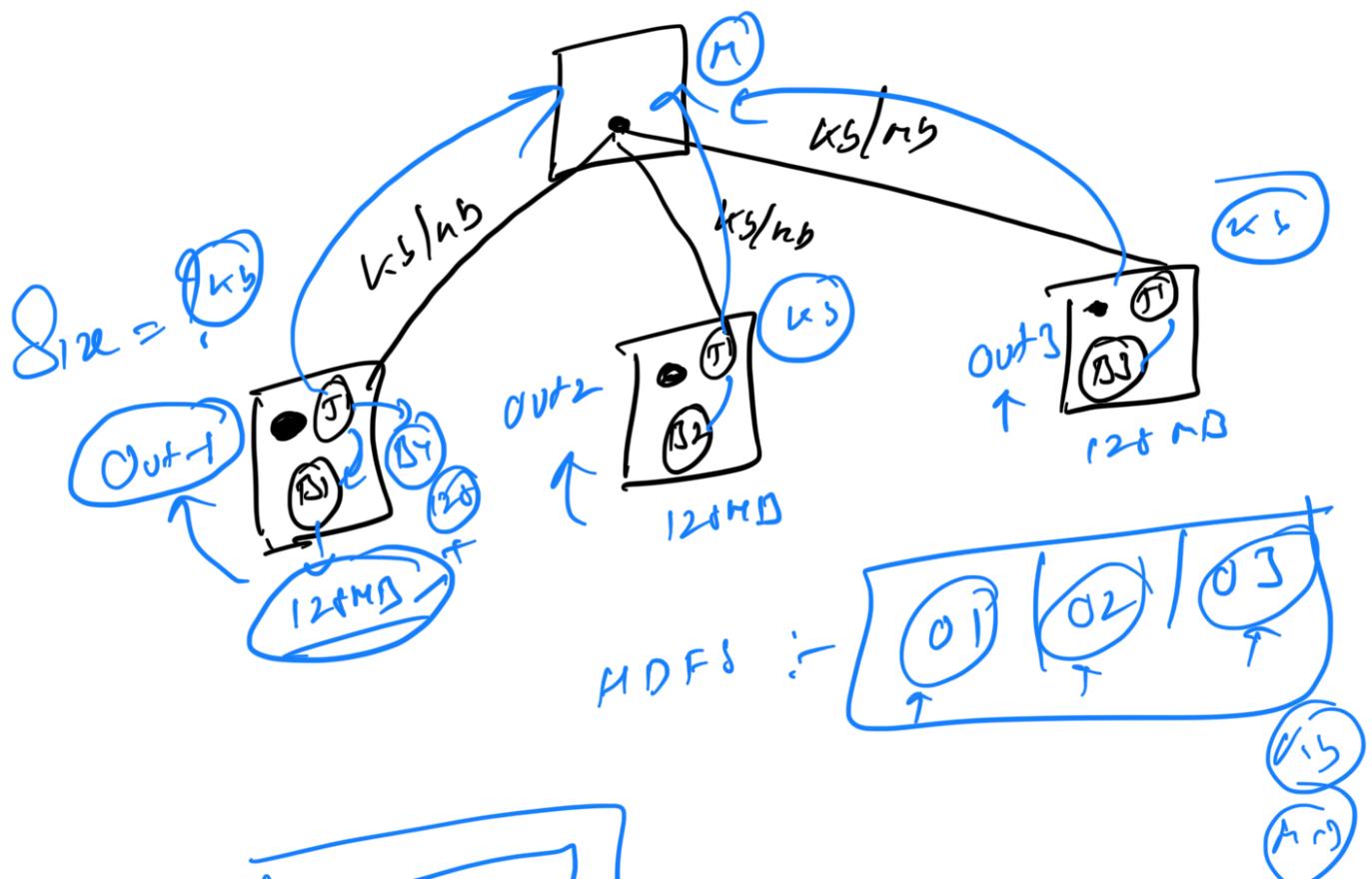
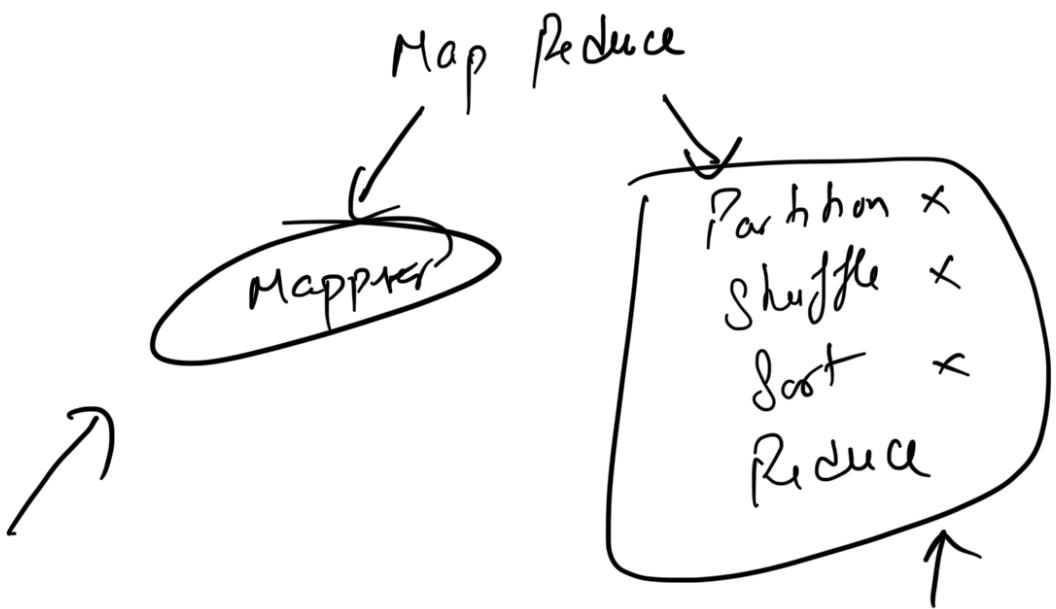
Why is data warehouse a better option?

How to design a Data Warehouse

Dimensional Modeling Approach







Input format



①

Determines how the input data is divided into chunks, known as splits. Supplies a 'Record Reader' implementation that reads the data within each split and convert into key value pairs.

②

→ Common Input format implementations:

①

TextInputFormat

default

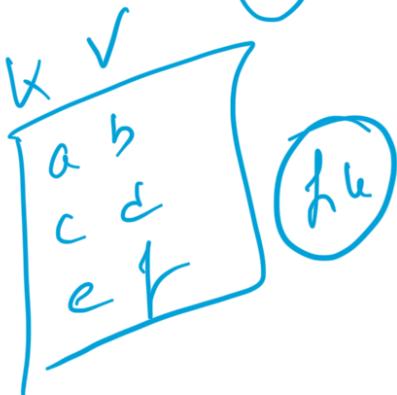
↳ K = byte offset of line
V = line of text

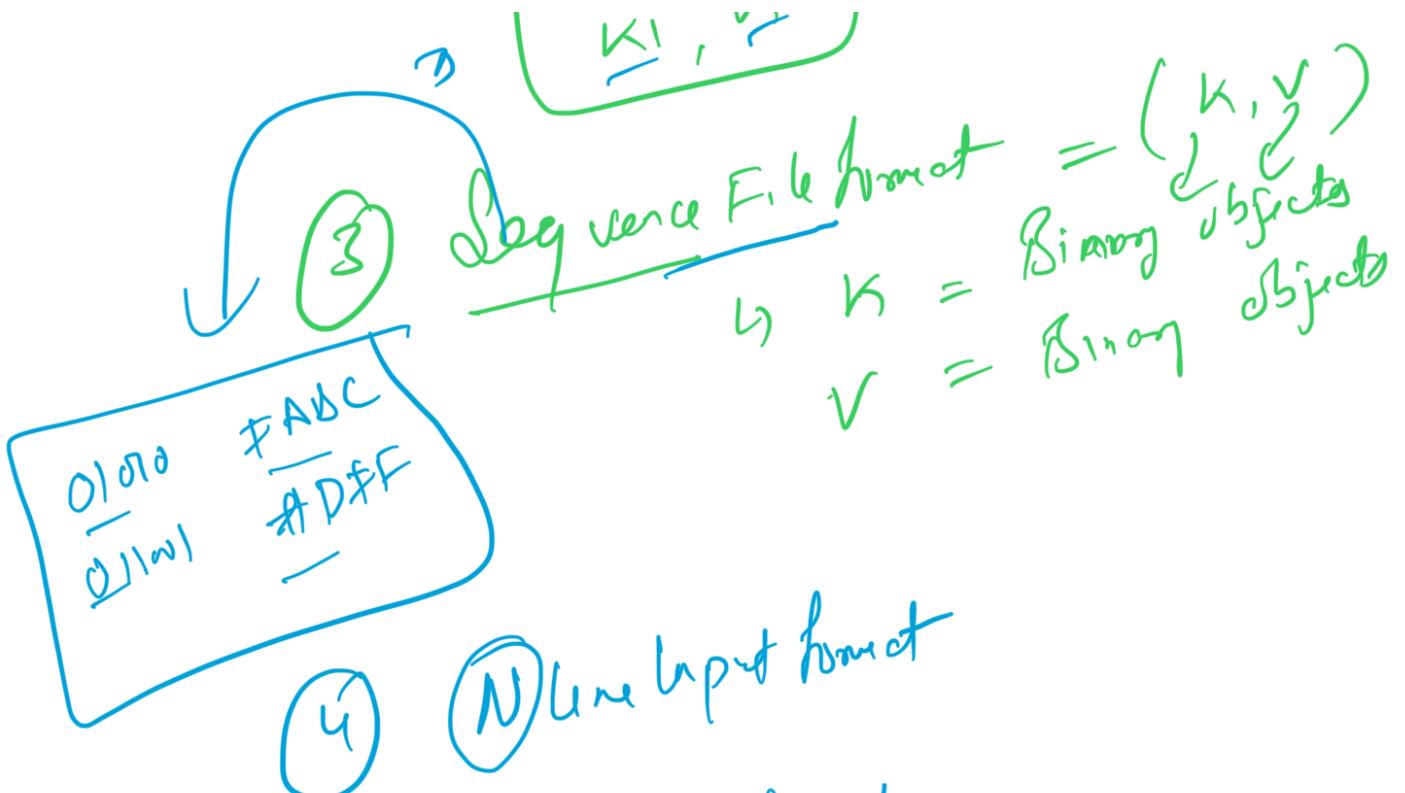
1:1

②

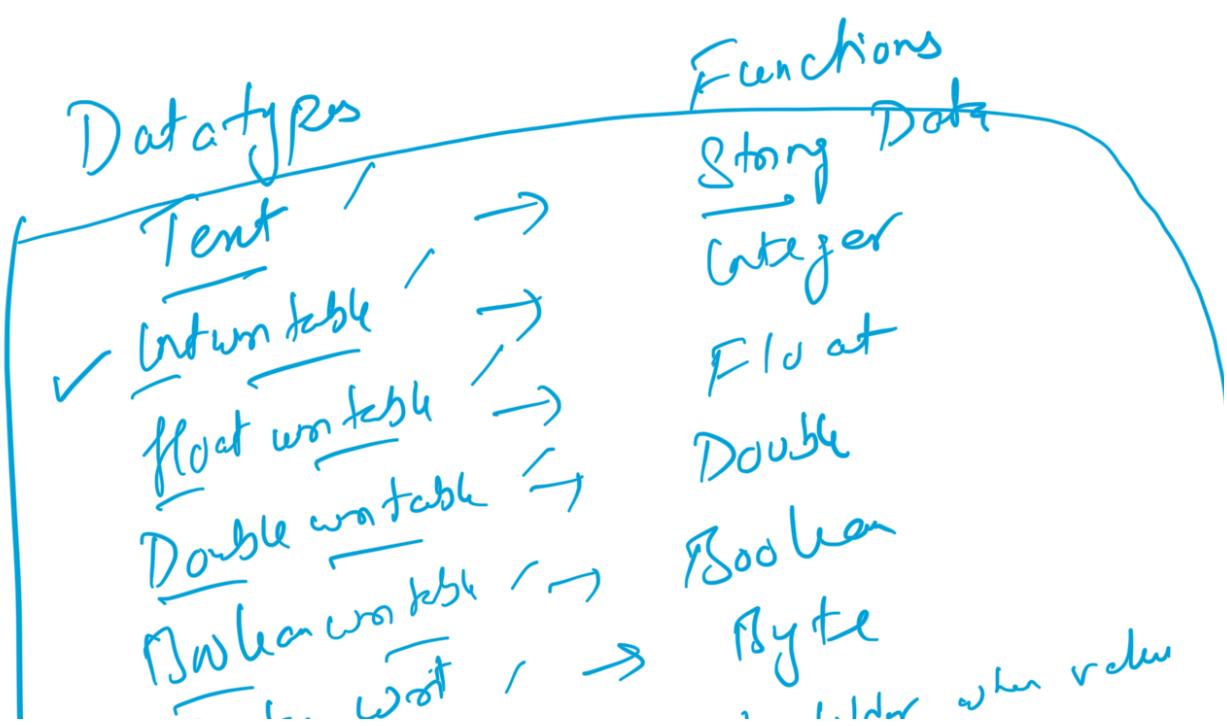
KeyValueTextInputFormat

↳ Used when line of the input file contains a (K, V) pair separated by a delimiter.

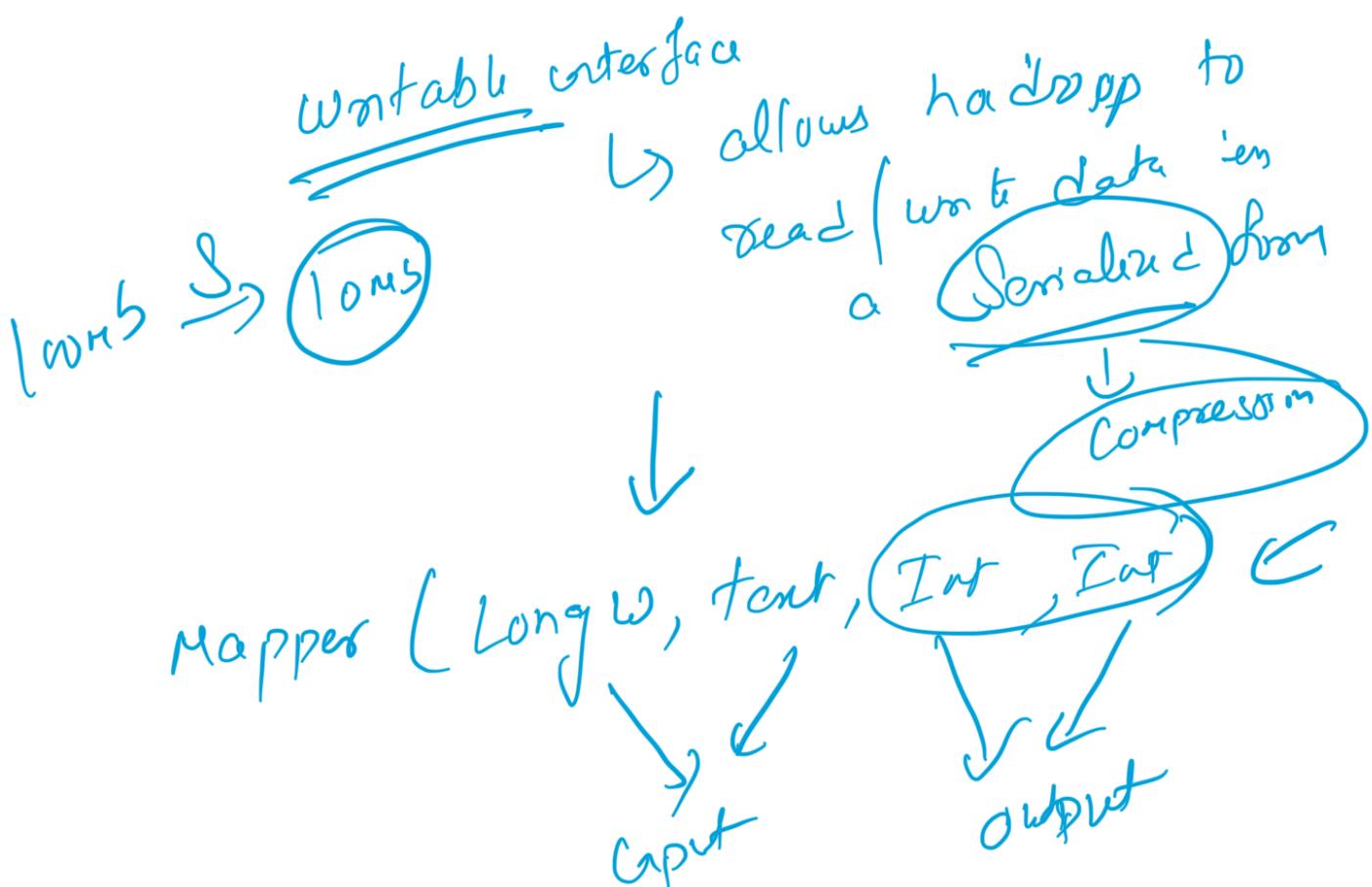
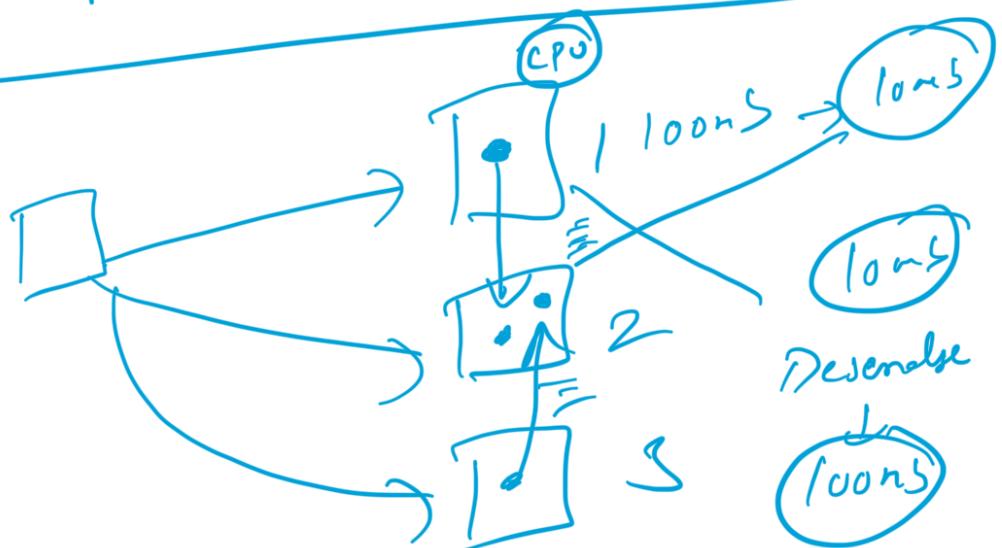




K = line number (PK)
 V = ORW

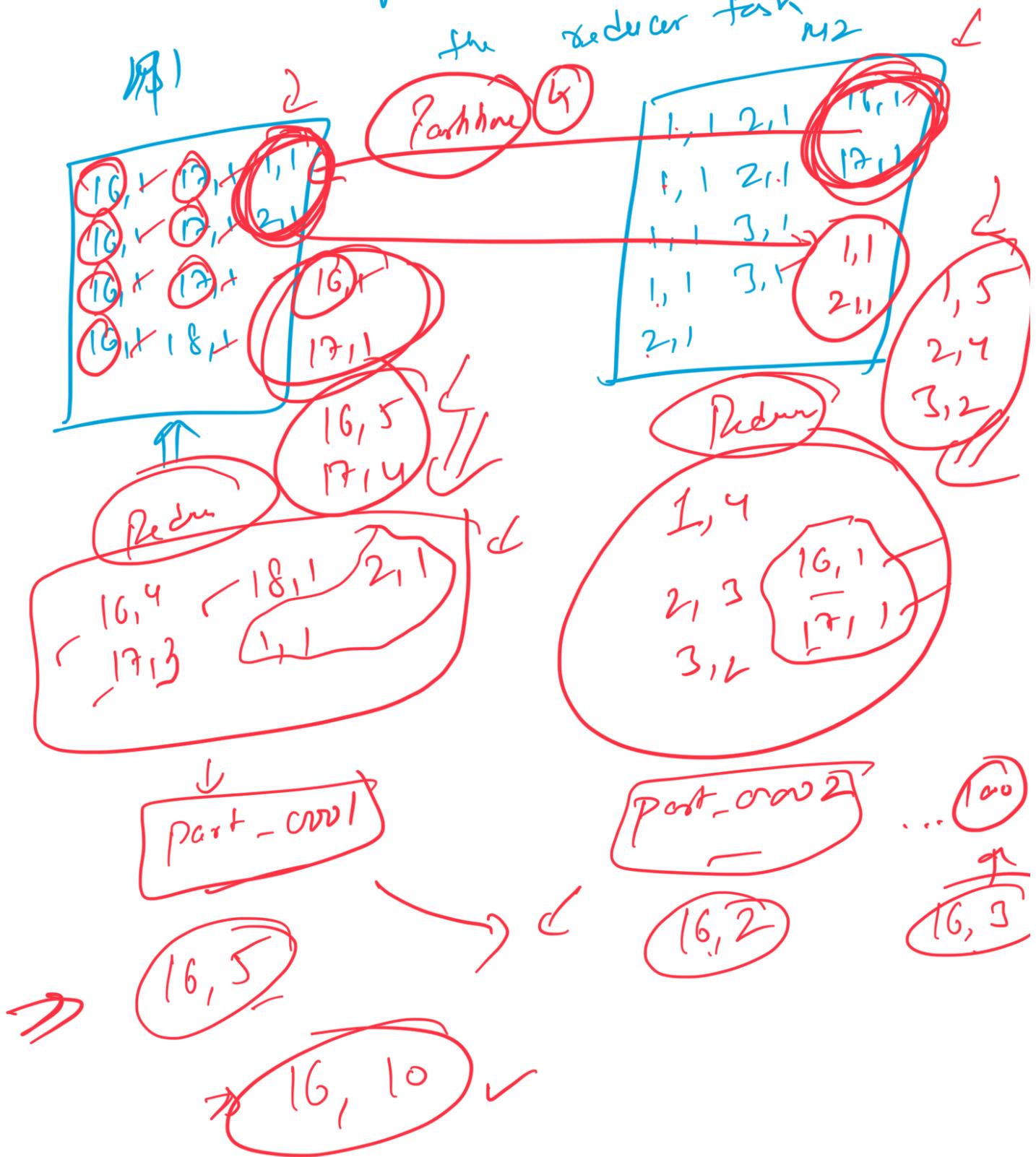


By ~~the~~
Null writable / → place
long writable is not needed.
 → Long



↳ Partitioner
 ↳ plays a crucial role in
 . the output

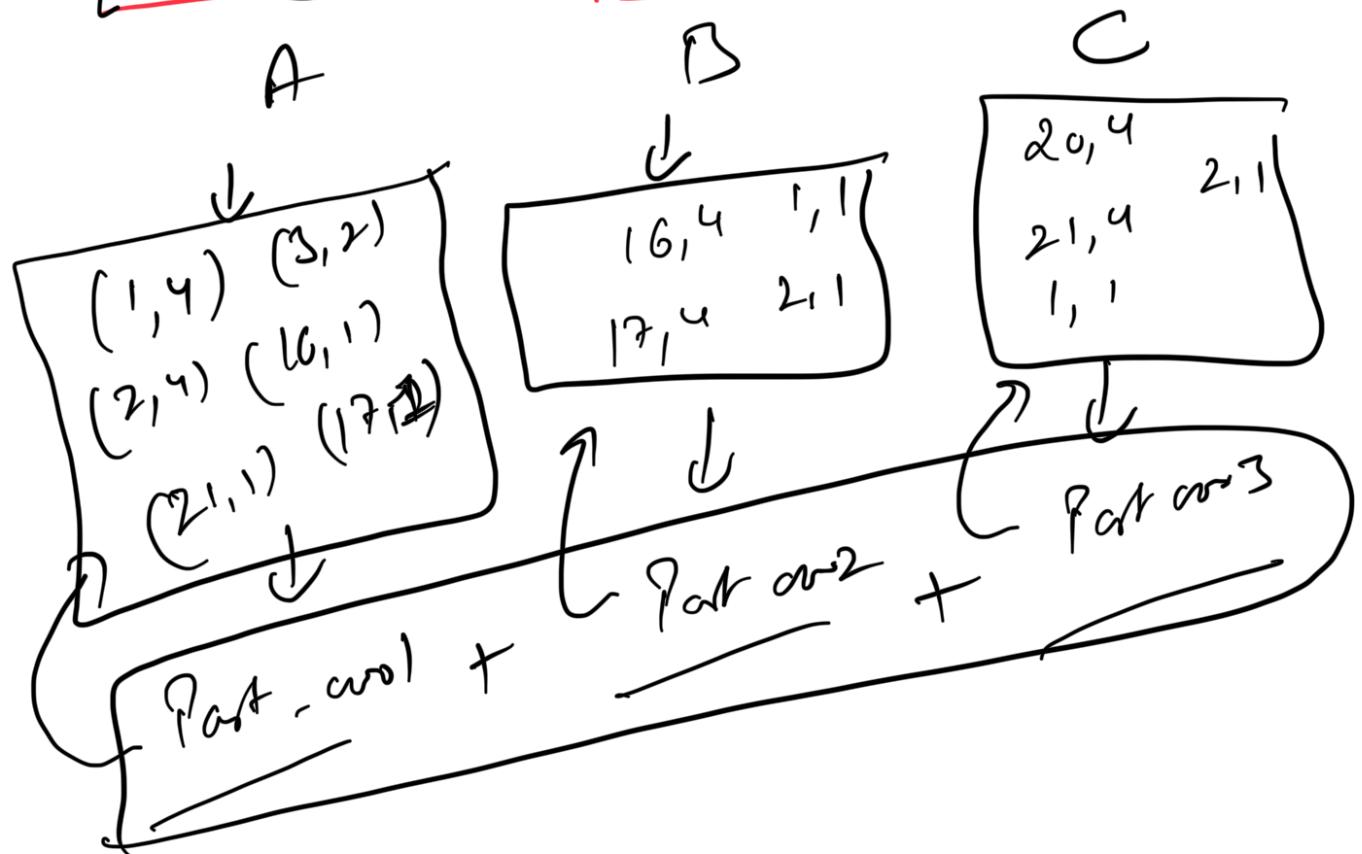
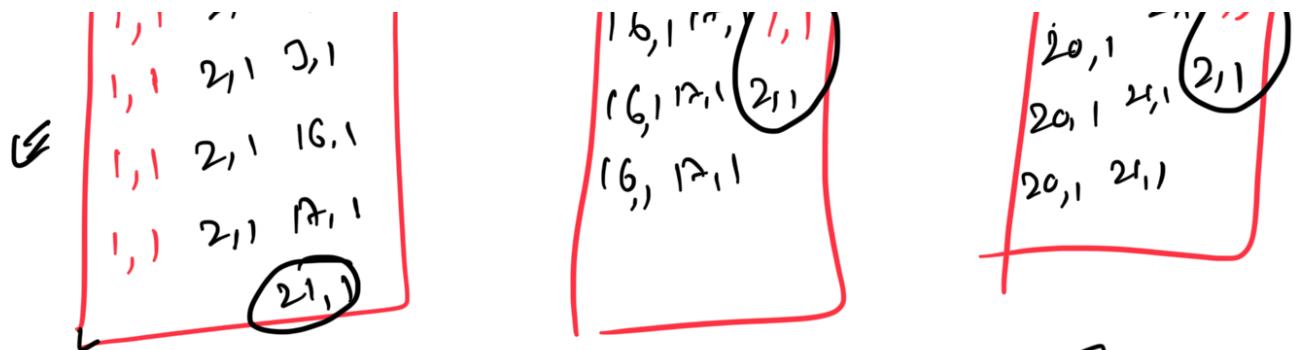
Determining how
of maps is distributed among
the reducer task M2



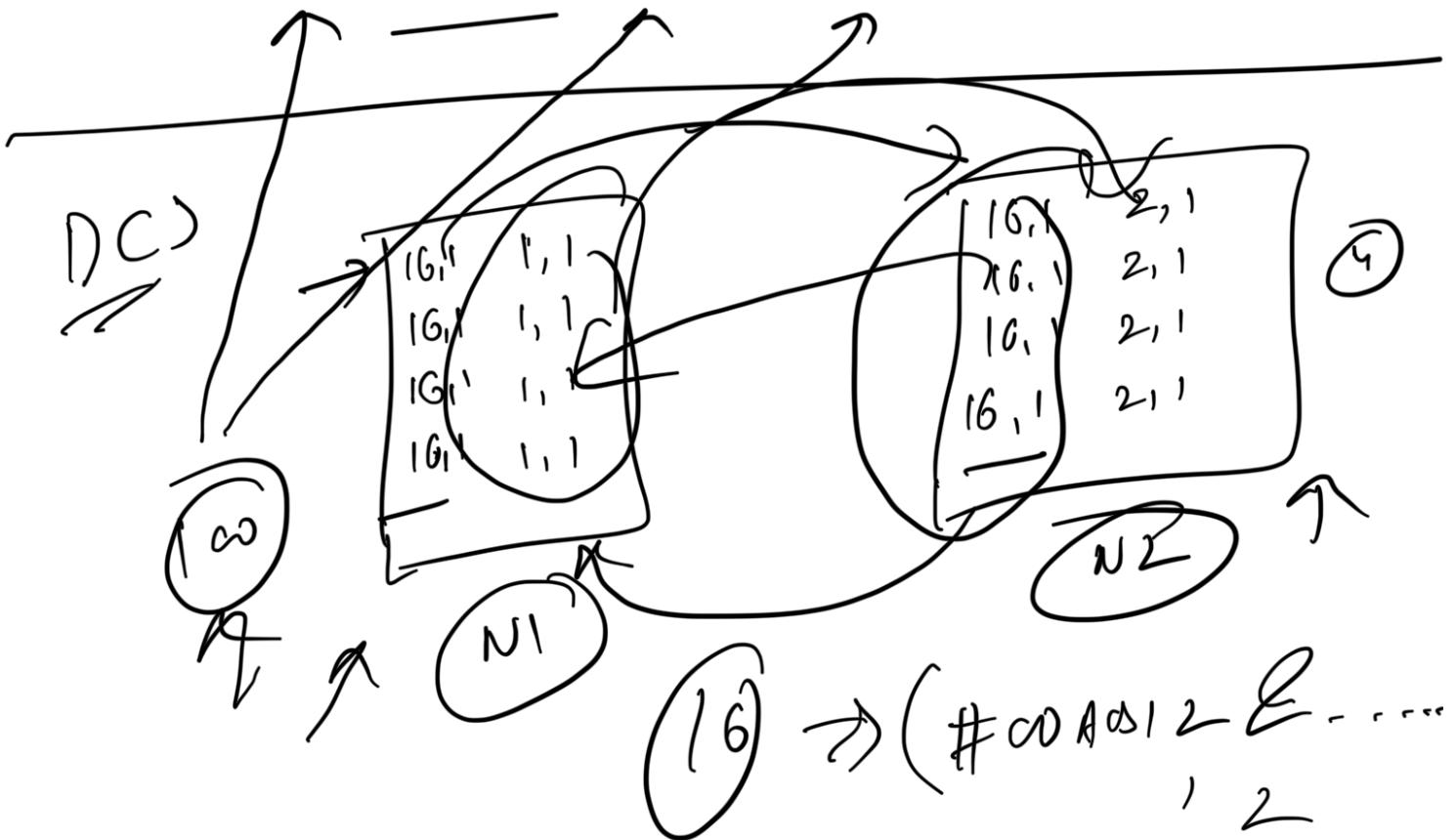
1, 1 2, 1 3, 1

1, 1 17, 1

20, 1 2, 1
2, 1 1, 1



Formula: ReducerNumber = (key.hashCode() & Integer.MAX_VALUE) % numReducers



\Rightarrow hadoop fs - getnurse

\Rightarrow File whl. Copy all w/