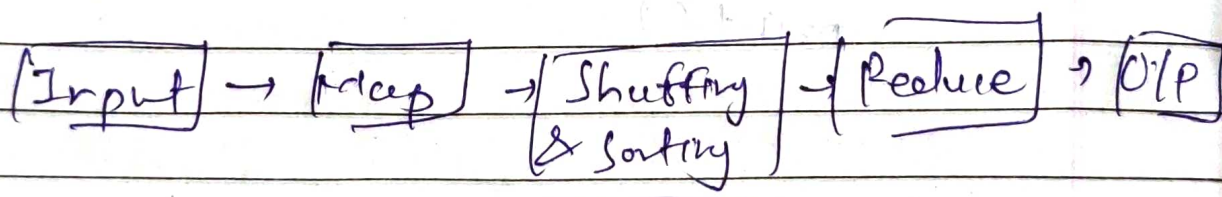


PAGE NO.:
DATE: / /

Ans-5

MapReduce is processing framework used to process data over a large no. of machines. Hadoop uses MapReduce to process the data distributed in Hadoop cluster. MapReduce is not similar to other regular processing frameworks like Hadoop, JPA, JET. All these processing frameworks are designed to use traditional sys. where the data is stored at single location like NFS, Oracle. But when processing big data, the data is located on multiple commodity machines. with helps of HDFS.



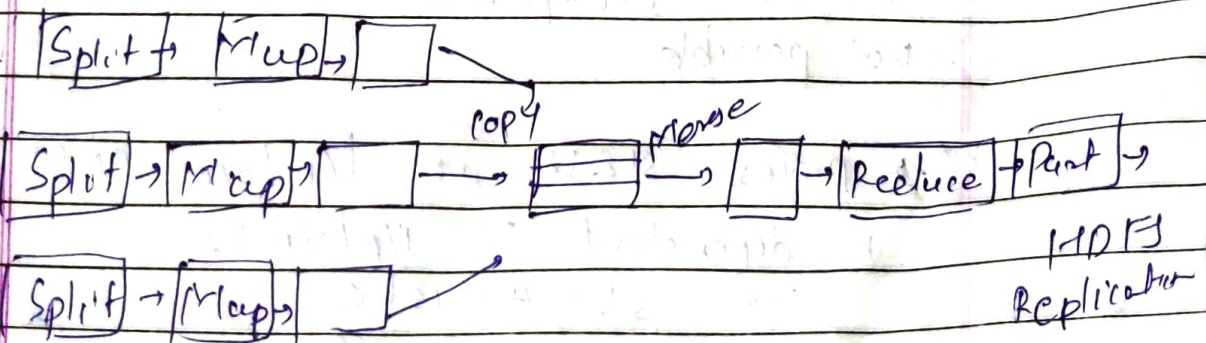
MapReduce is terminology that comes with Map phase & reducer phase. map is used for transformation while reducer is used for aggregation kind of operation. The terminology for Map & reduce is derived from some functional programming lang. like Lisp, Scala etc. MapReduce processing framework prog. comes with 3 main components i.e., our Driver code, Mapper, & Reducer.

Working of Mapper

- Mapper is initial line of code that initially interacts with HDFS data sets. Suppose if we have 120 Peter Blocks of dataset we are analysing then in that case, will be 120 Mapper program that runs in parallel on machine & produce own output known as intermediate o/p which is stored in local Disk not on HDFS. The o/p of Mapper act as i/p for Reducer which perform Sorting & aggregation o/p on data & produce final o/p

Working of Reducer

- Mapper produces o/p in form key-value pairs which works as i/p for Reducer. But before sending, this intermediate key-value pairs directly to Reducer. Some process will be done which shuffle & sort key-value pairs according to key value.
- o/p of Reducer is stored on HDFS. Reducer perform some computation + x.



MapReduce Data flow