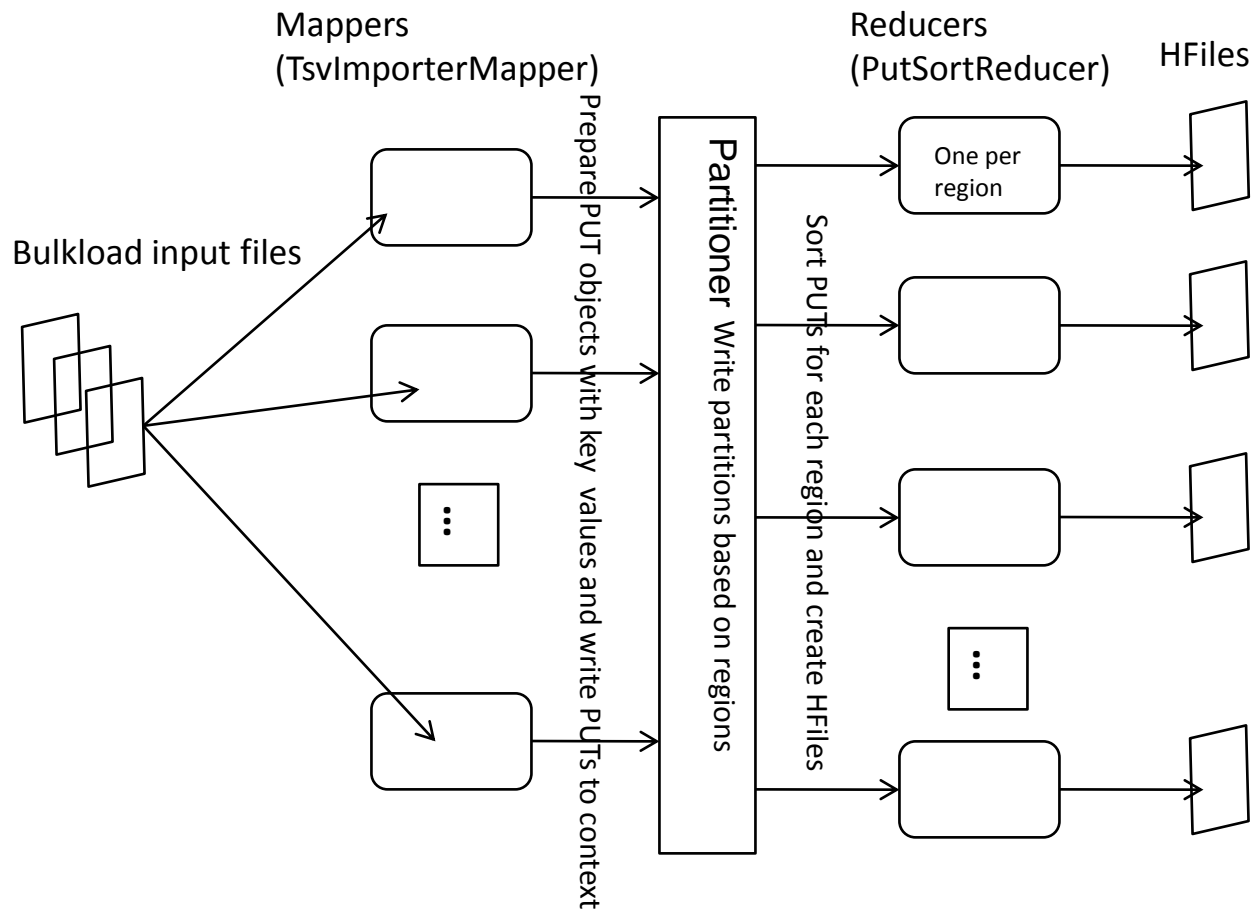


Bulkload Performance Improvement

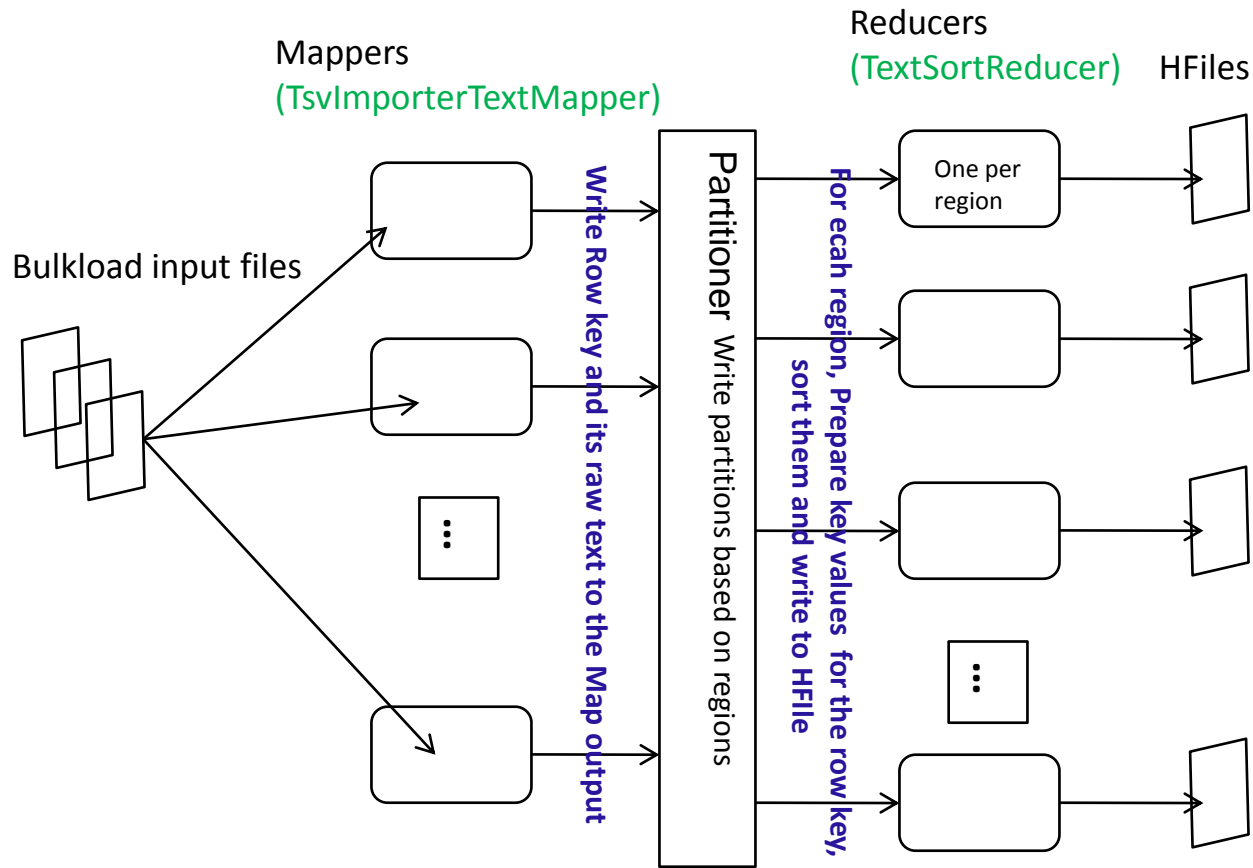
Author: Jyothi & Rajeshbabu



Bulkload with default Mapper



Bulkload with custom Mapper & Reducer



Performance Test

► Environment

◦ 3 Node cluster

- Node1 – Name node, Data Node , Resource Manager , Node Manager, HMaster, Region Server
- Node 2 – Zookeeper , Data Node, Region Server, Node Manager
- Node 3 – Data Node, Region Server, Node Manager

◦ JVM

- Default configurations

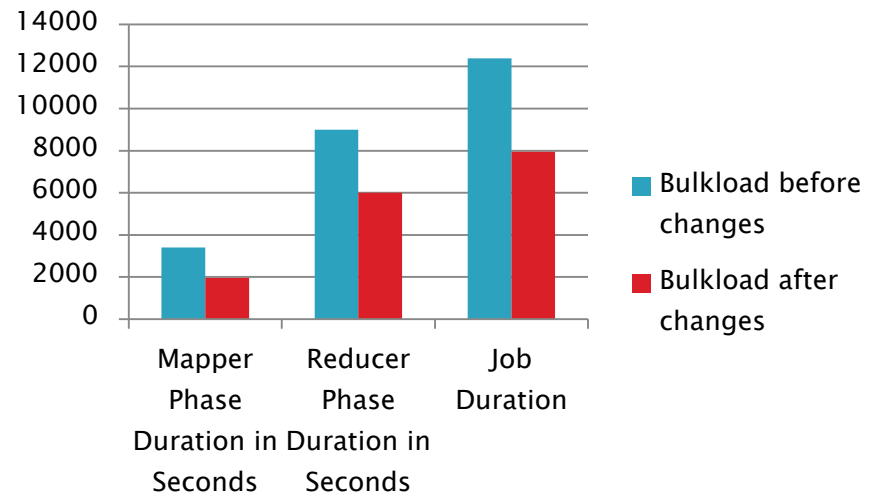
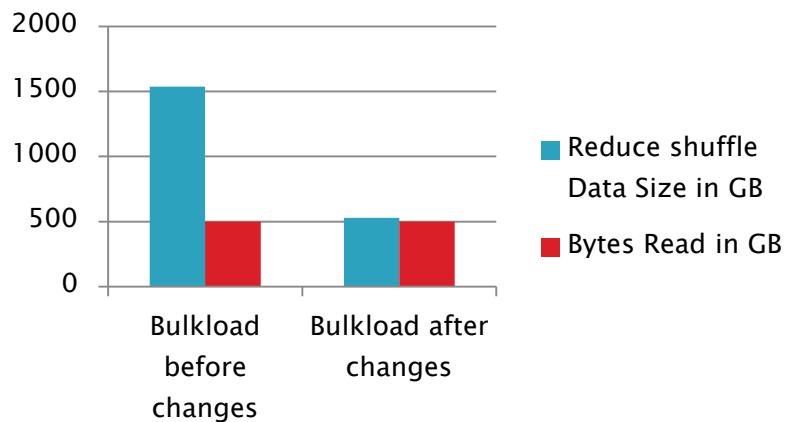
Performance Test

▶ Test Setup

- 500GB input data with 500 bytes in each line
- Put to HDFS with 64MB Block Size
- Regions & Number of Reducers – 500
- Number of Mappers – 8000
- Column Family – 1
- Columns in CF – 20

Test Report

Metric	Bulkload before changes (TsvImporterMapper & PutSortReducer)	Bulkload after changes (TsvImporterTextMapper & TextSortReducer)
Time Taken for Mapper Phase	56Min 44Sec	32 Mins 32 Sec
Time Taken for Reducer Phase	149 Mins 50 Sec	100 Mins 3 Sec
Reduce shuffle bytes	1651439196000 (1538GB)	568033518000 (529GB)
Bytes Read	536903763904	536903763904



Conclusion

- ▶ Noticed 55% improvement in throughput
- ▶ With improvement, Bytes shuffled are almost equal to input data size where as with default ones it is 3 times more than the input data size.