



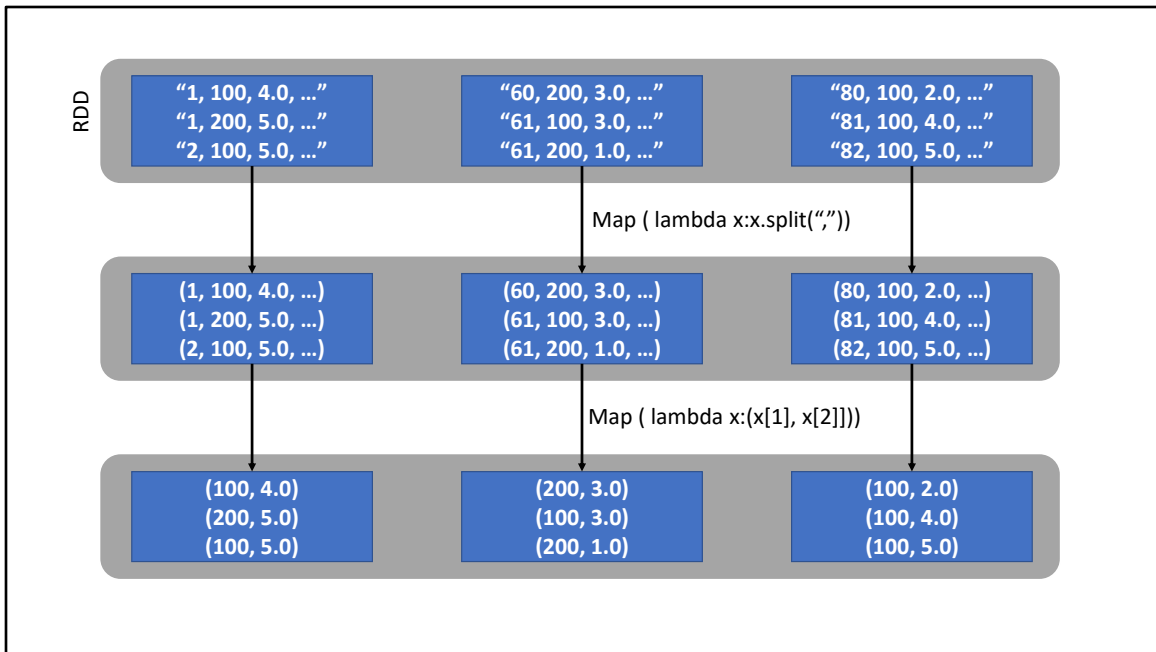
Data-Intensive Distributed Computing
CS 431/631/451/651 (Fall 2021)

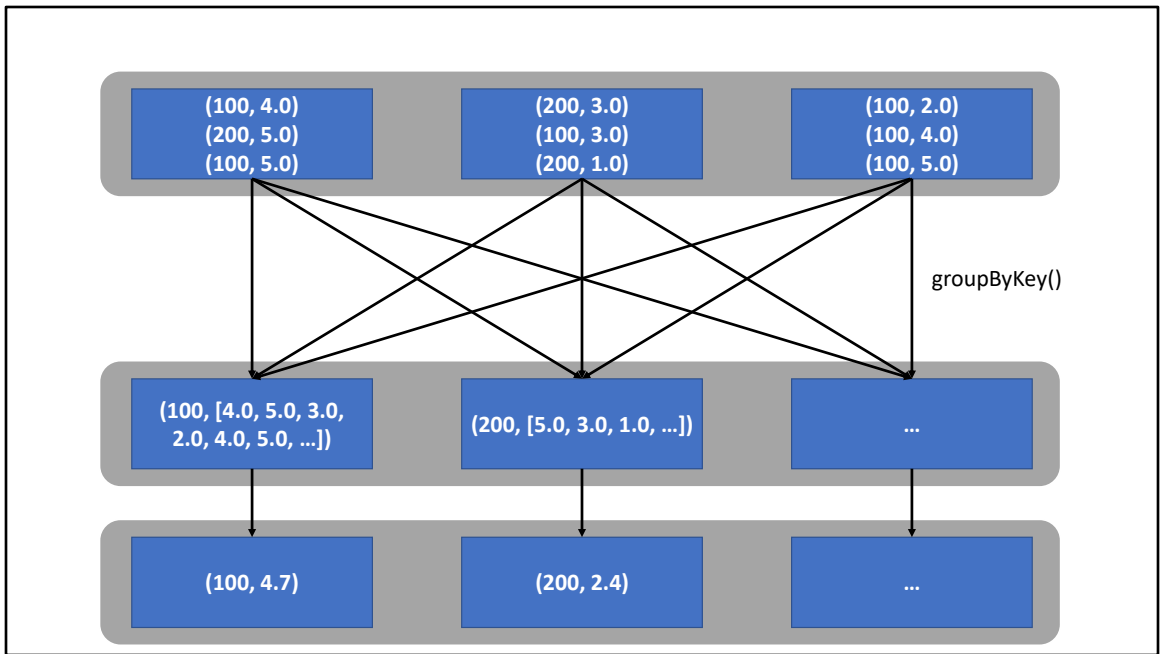
Part 3: From MapReduce to Spark (3/3)

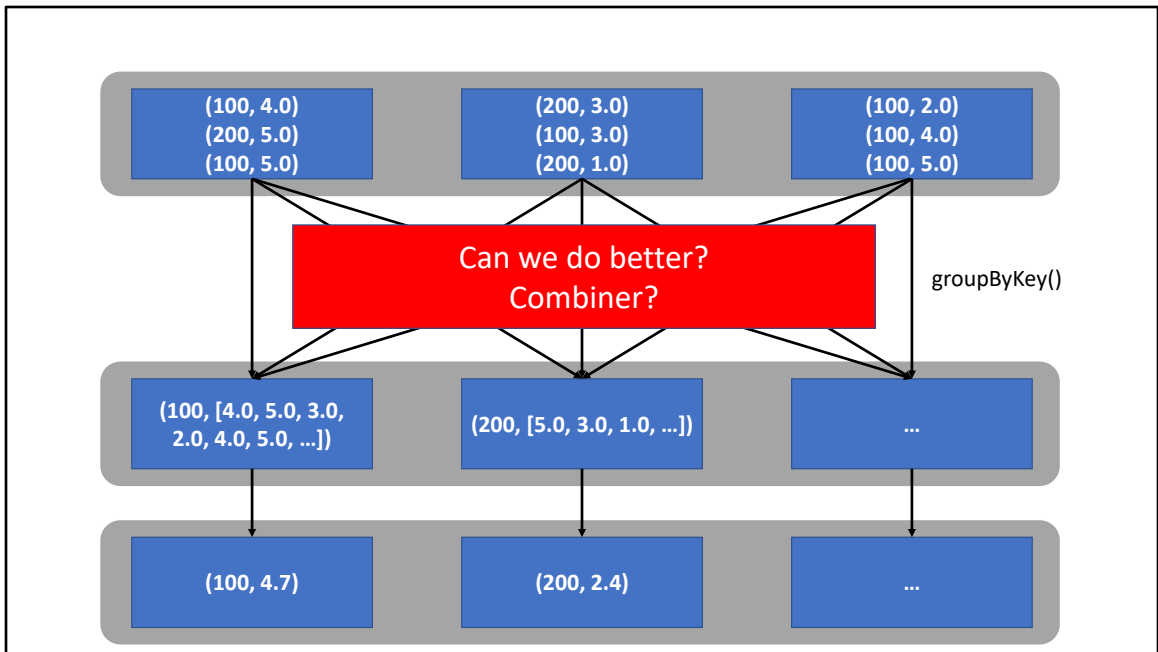
Ali Abedi

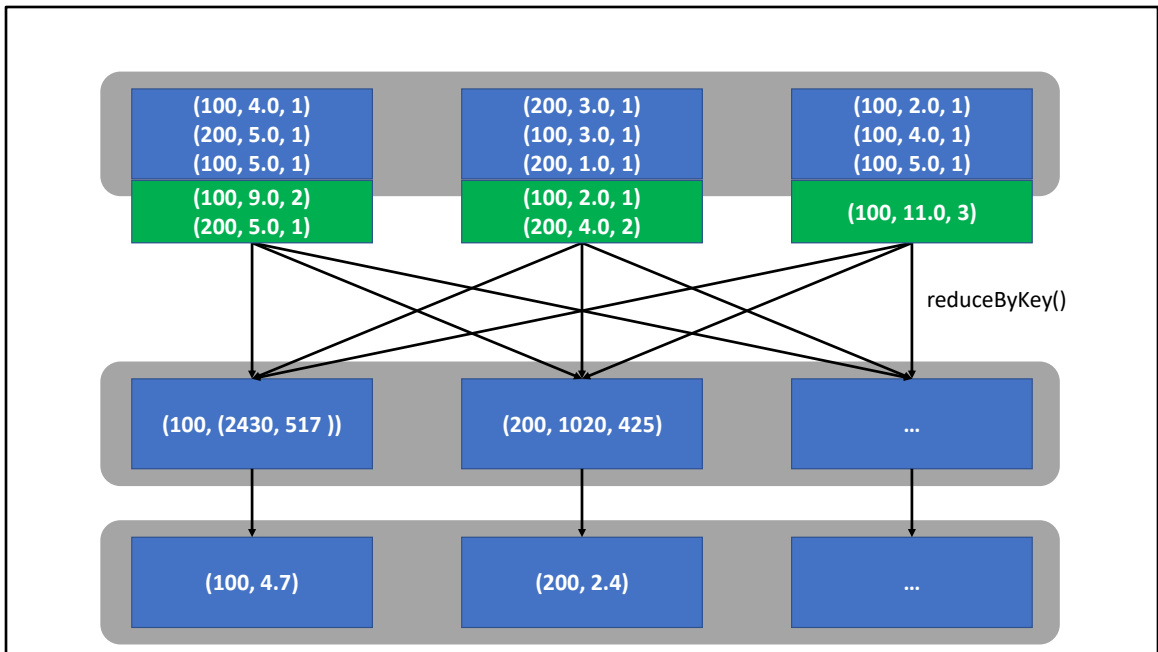
These slides are available at <https://www.student.cs.uwaterloo.ca/~cs451/>

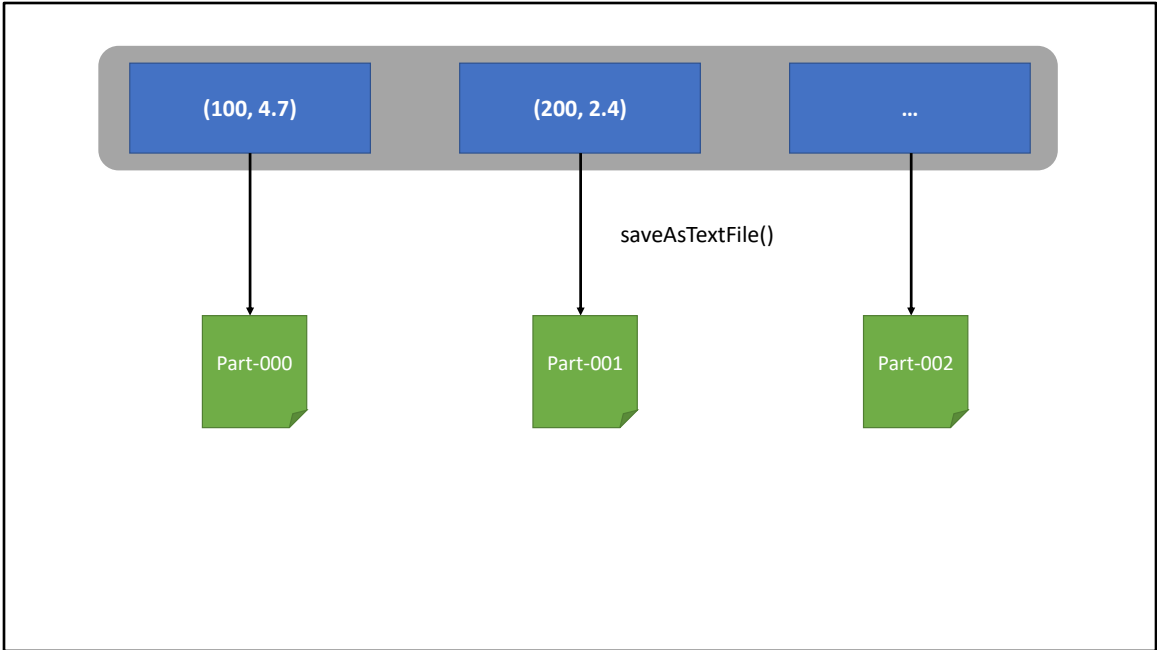
Movie rating example

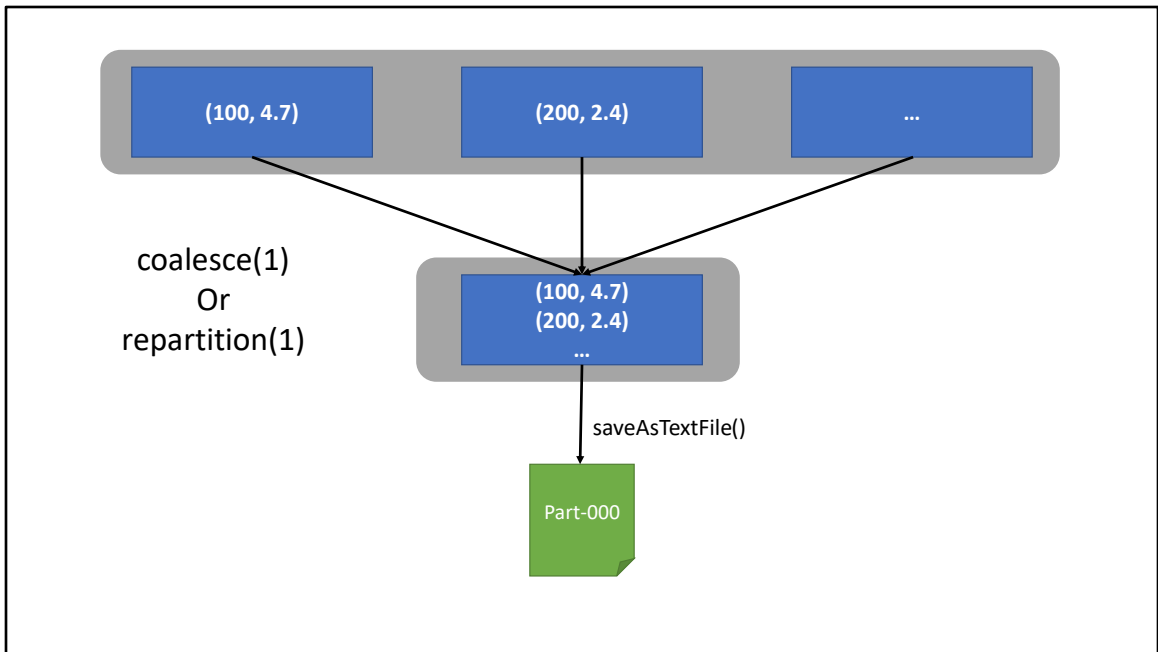












Repartition triggers shuffling but it gives more balanced partitions. It can be used to increase or decrease the number of partitions.

Coalesce can be used to only reduce the number of partitions. It avoids full shuffling so it is faster than repartition but it may give unbalanced partitions.

Other topics ...

- Cache vs Persist
- Broadcast variables vs Accumulators