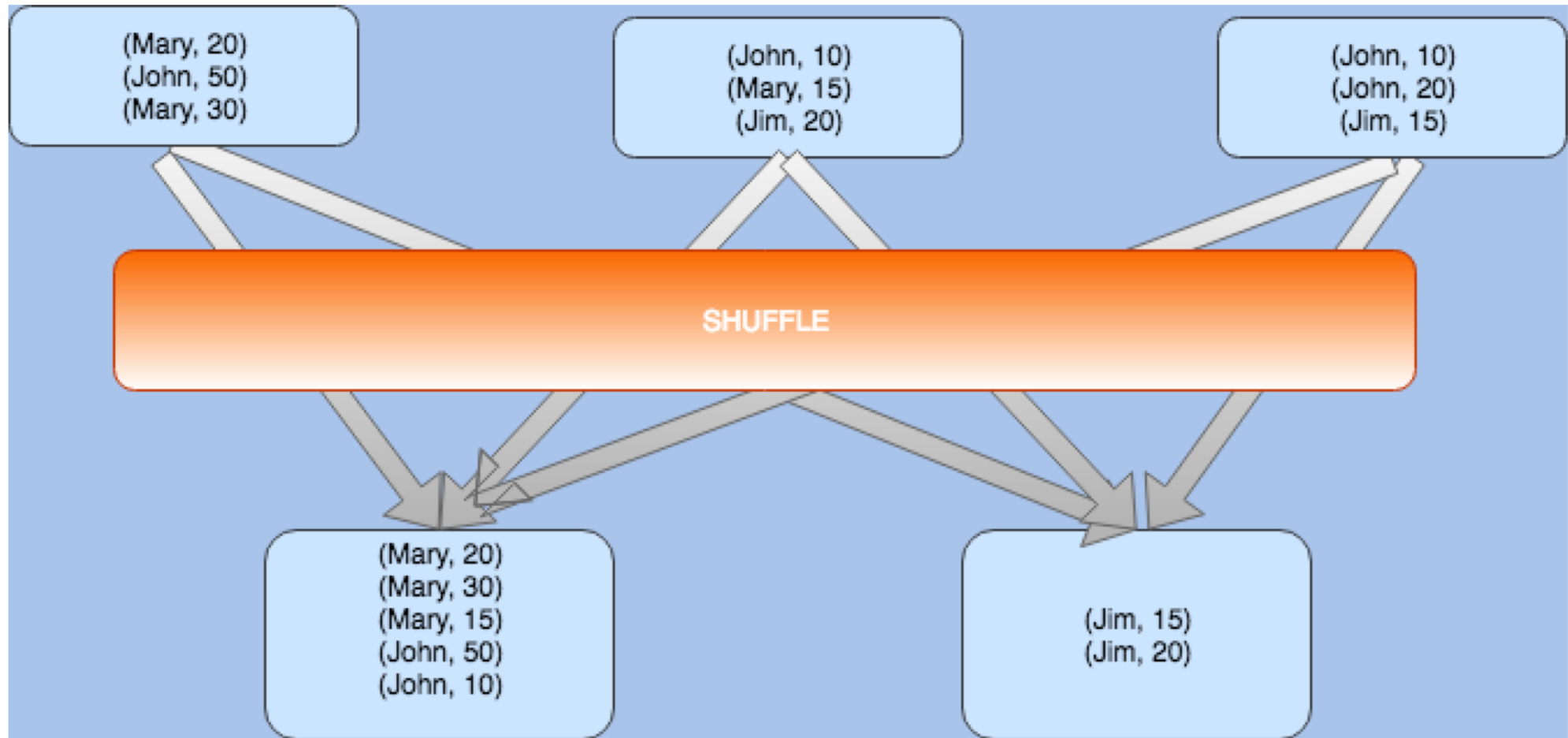


Shuffling

GroupBy Shuffle

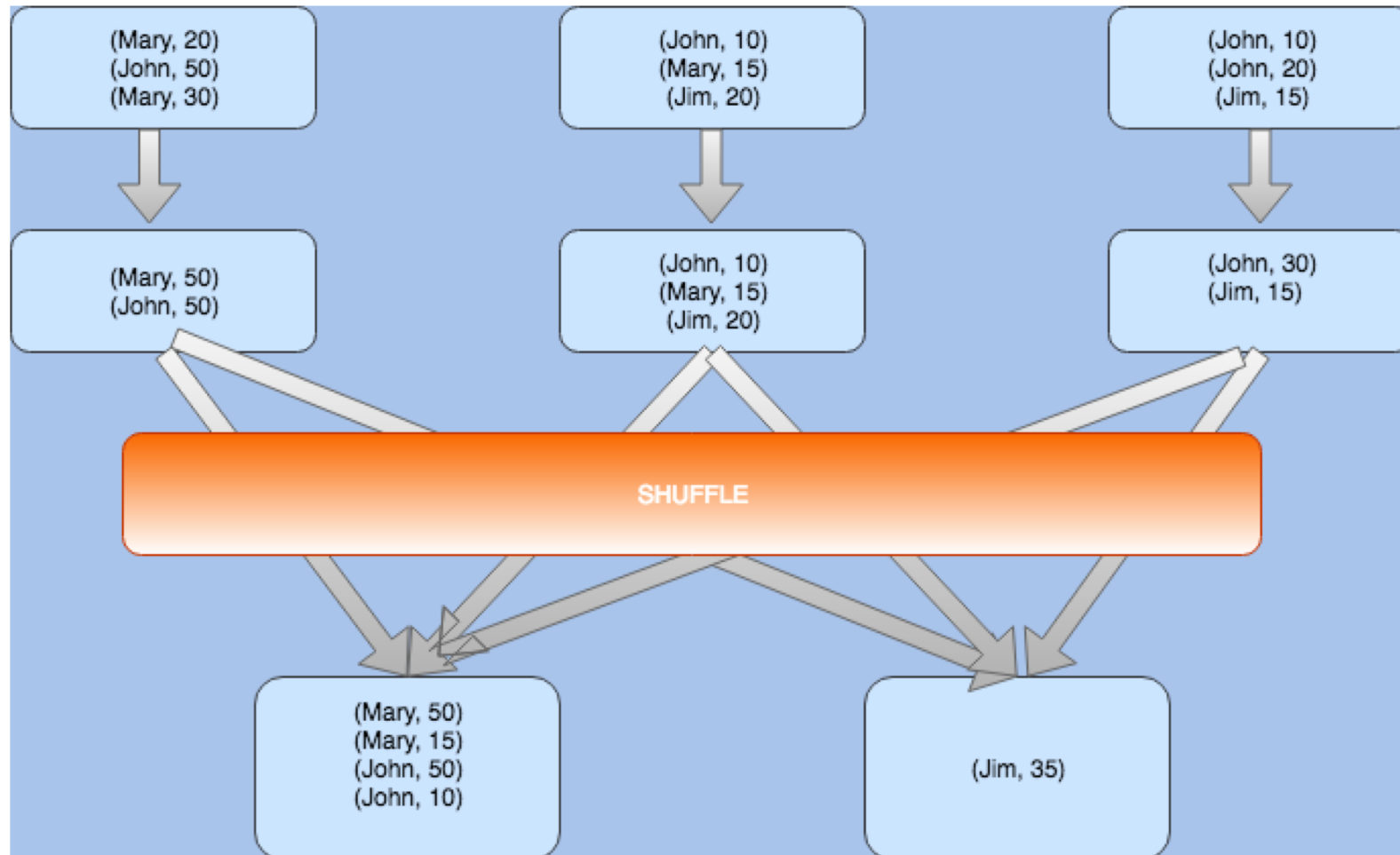
- Does not aggregate data before moving it across the network
- Increased latency

GroupBy Shuffle



ReduceBy Shuffle

- Aggregates (reduces) data on worker before shuffle
- Should be the preferred over GroupBY



Hash Partitioning

By default, Spark uses hash partitioning to determine which key-value pair should be sent to which machine.

.... How to decide what to use for buckets