# Coalesce vs Repartition

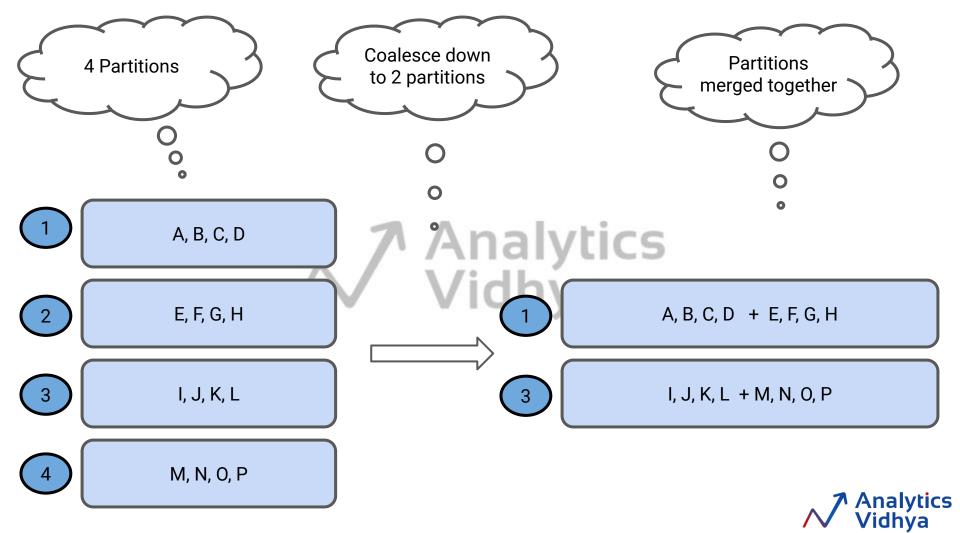# Repartition

Repartition

- To increase or decrease partitions
- Complete shuffling of data
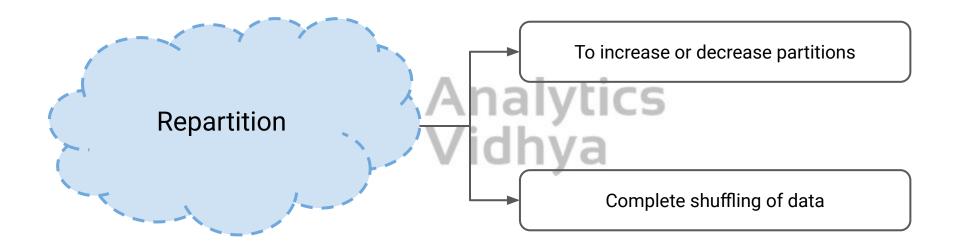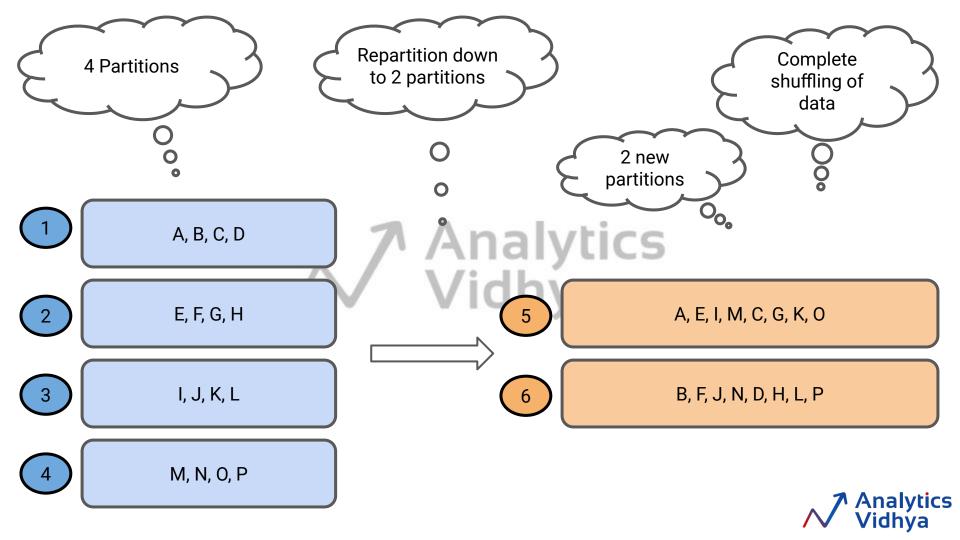
# Partition tuning: Repartition or Coalesce ?

If your dataset is **skewed**: use **repartition**

If you want **more partitions**: use **repartition**

If you want to drastically **reduce** the number of partitions(e.g. numPartitions =1): use **repartition**

If your dataset is **well balanced**(i.e. not skewed) and you want **fewer partitions,** but not dramatically fewer): use **coalesce**

Thank You