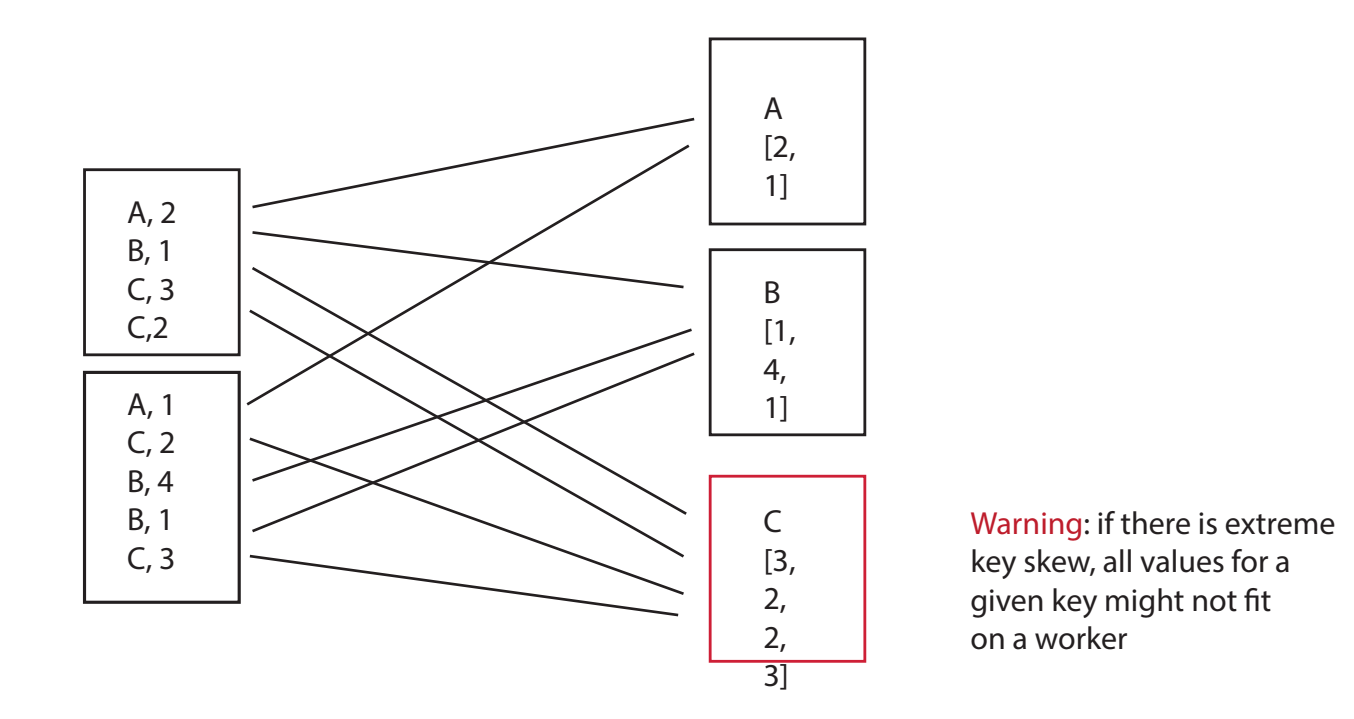


A few common aggregation methods for Spark RDDs.

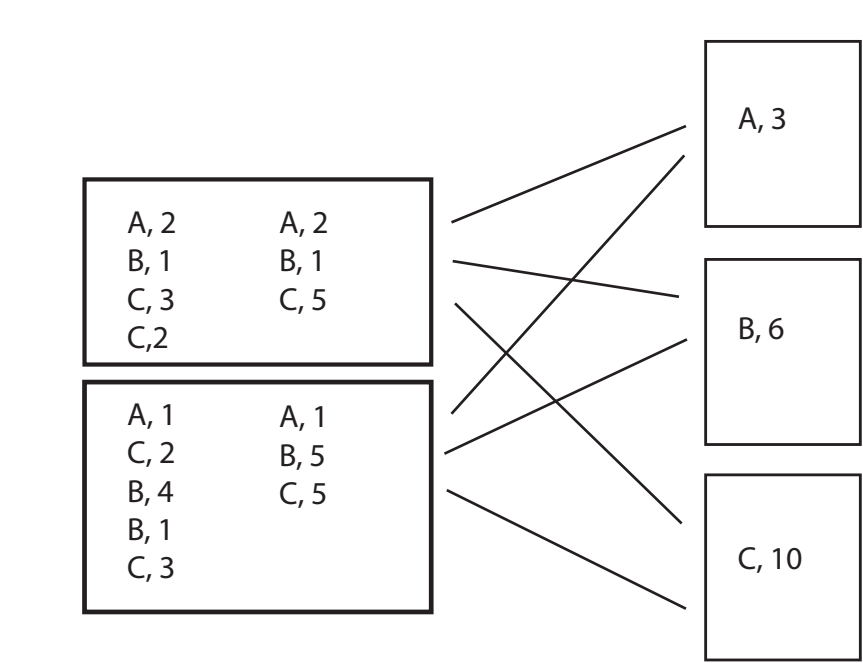
groupByKey()

https://spark.apache.org/docs/latest/api/python/_modules/pyspark/rdd.html#RDD.groupByKey



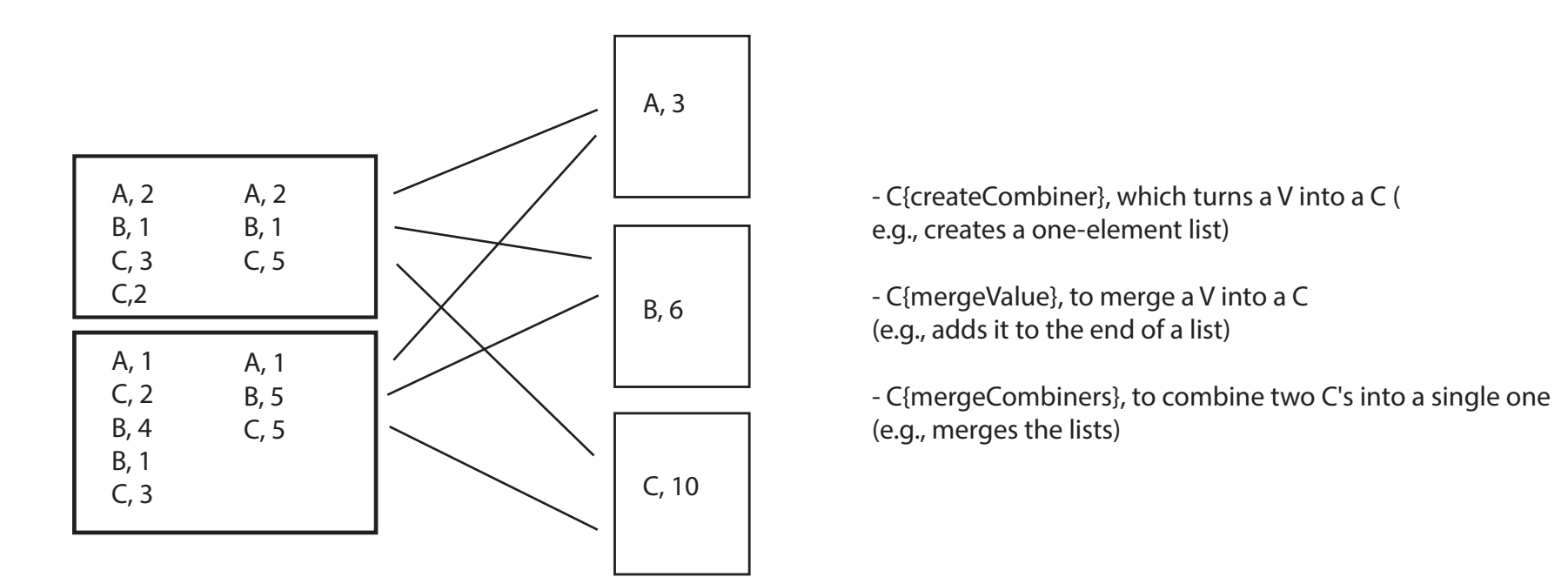
reduceByKey(Func)

https://spark.apache.org/docs/latest/api/python/_modules/pyspark/rdd.html#RDD.reduceByKey



combineByKey(createCombiner, mergeValue, mergeCombiners)

https://spark.apache.org/docs/latest/api/python/_modules/pyspark/rdd.html#RDD.combineByKey

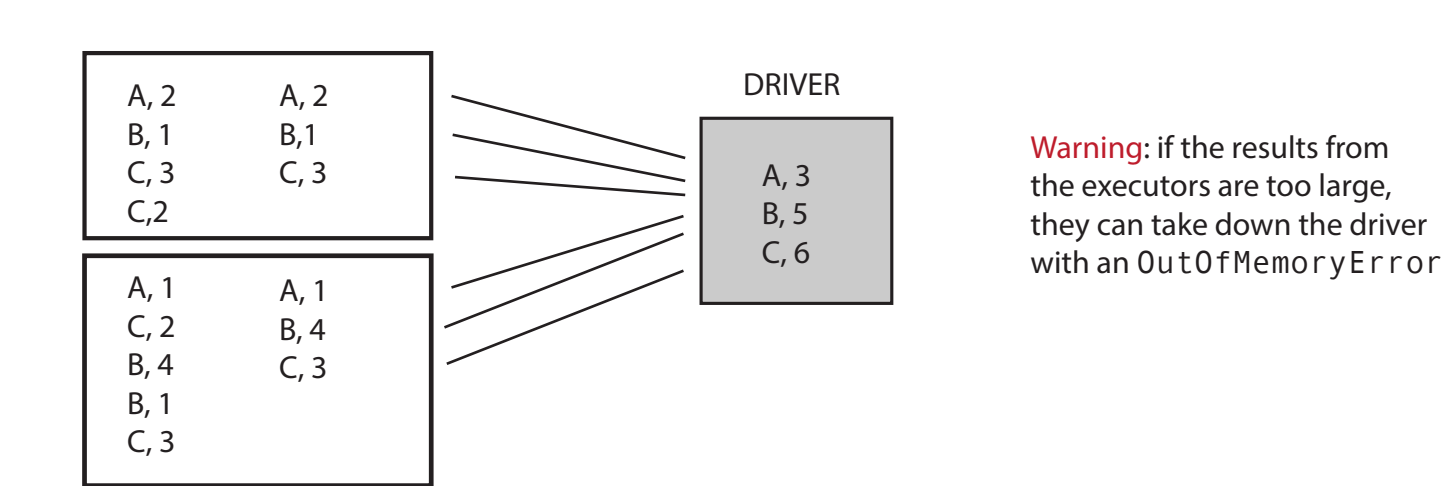


aggregateByKey(zeroValue, seqOp, combOp)

https://spark.apache.org/docs/latest/api/python/_modules/pyspark/rdd.html#RDD.aggregateByKey

start value,
seqOp -> within partition function,
combOp -> across partition function

ex: aggregateByKey(0,max,add)



foldByKey(zeroValue, Func)

https://spark.apache.org/docs/latest/api/python/_modules/pyspark/rdd.html#RDD.foldByKey

Calls combineByKey, but allows us to use a zero value which can be added to the result an arbitrary number of times, and must not change the result (eg. 0 for addition, 1 for multiplication)

treeAggregate(zeroValue, seqOp, combOp, depth)

https://spark.apache.org/docs/latest/api/python/_modules/pyspark/rdd.html#RDD.treeAggregate

Same as aggregate except it “pushes down” some of the subaggregations (creating a tree from executor to executor) before performing final aggregations on the driver.