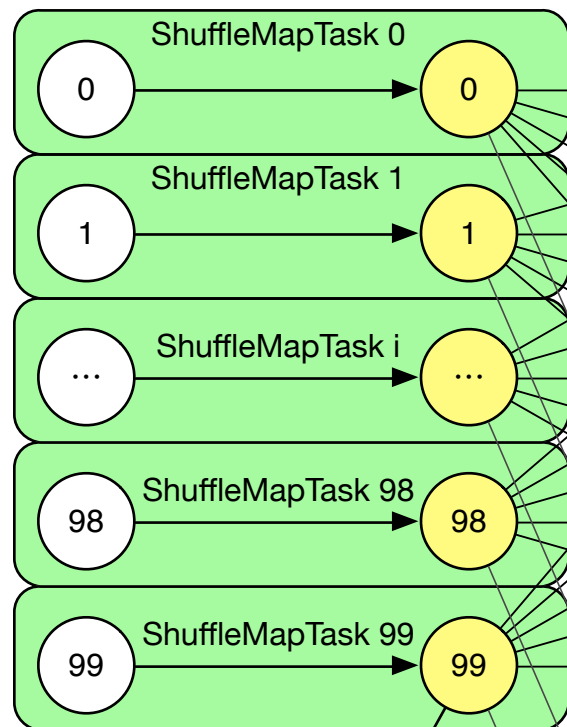


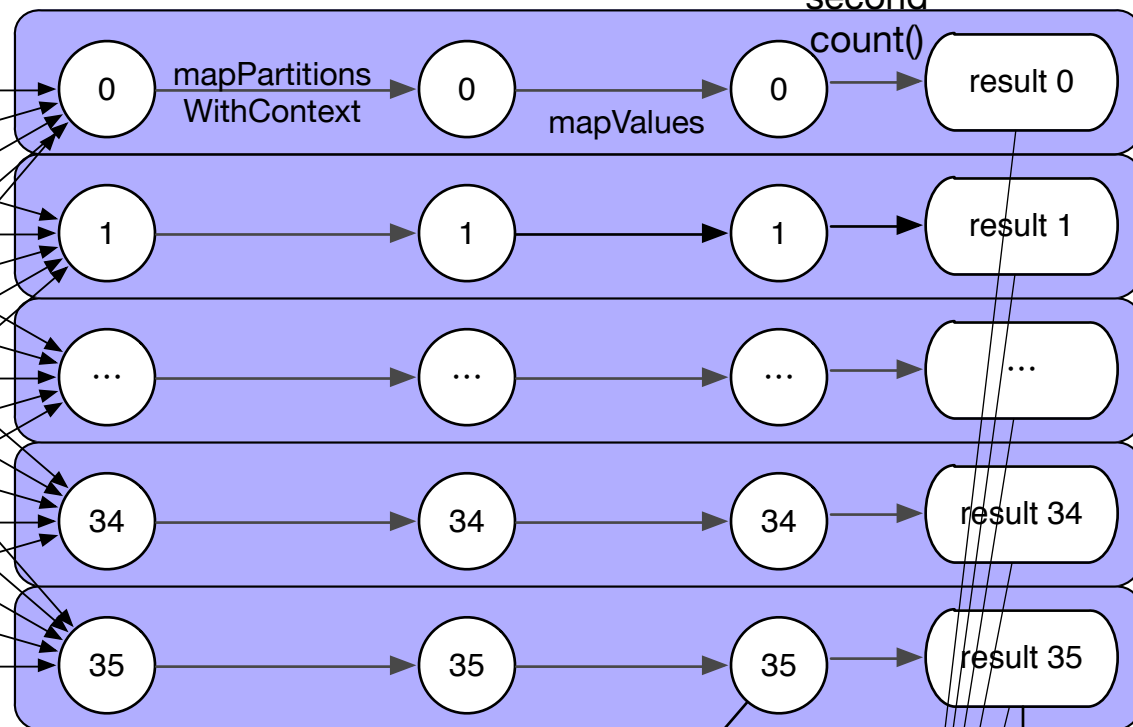
Stage 0

ParallelCollectonRDD FlatMappedRDD



Stage 2

ShuffledRDD MapPartitionsRDD MappedValuesRDD



Stage 1

Array[(Int, Array[Byte])]

(2, Byte[1000])

(2, Byte[1000])

(1, Byte[1000])

(4, Byte[1000])

(3, Byte[1000])

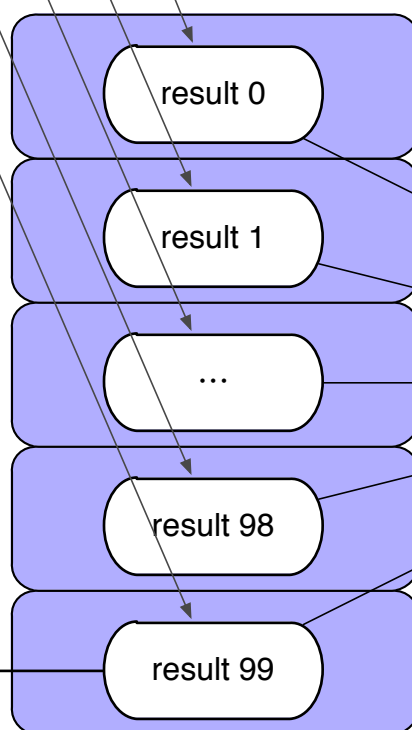
(2, Byte[1000])

(5, Byte[1000])

(1, Byte[1000])

(2, Byte[1000])

result of
count(Array[Int, Byte[]])



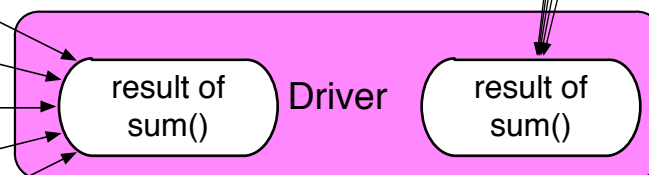
(Int, Iterable[Array[Byte]])

(1, Iterable[Array[Byte]])

(3, Iterable[Array[Byte]])

(5, Iterable[Array[Byte]])

result of
count(Array[Int,
Iterable[Array[Byte]]])



first count()

second count()

Legend:

RDD

partition:

i

cached partition:

i

assumed data
in the partition/result

ShuffleMapTask

ResultTask