



# Apache Spark

Execution Plan

# Spark Execution Plan

```
val surveyRawDF = spark.read  
  .option("header", "true")  
  .option("inferSchema", "true")  
  .csv(args(0))
```

Job 0

Job 1

```
val partitionedSurveyDF = surveyRawDF.repartition( numPartitions = 2)  
val countDF = partitionedSurveyDF.where( conditionExpr = "Age < 40")  
  .select( col = "Age", cols = "Gender", "Country", "state")  
  .groupBy( col1 = "Country")  
  .count()  
Logger.info(countDF.collect().mkString("->"))
```

Job 2

C0	C1	C2	C3
val	val	val	val
val	val	val	val
val	val	val	val

REPARTITION

C0	C1	C2	C3
val	val	val	val
val	val	val	val
val	val	val	val

WHERE

C0	C1	C2	C3
val	val	val	val
val	val	val	val
val	val	val	val

SELECT

C0	C1	C2	C3
val	val	val	val
val	val	val	val
val	val	val	val

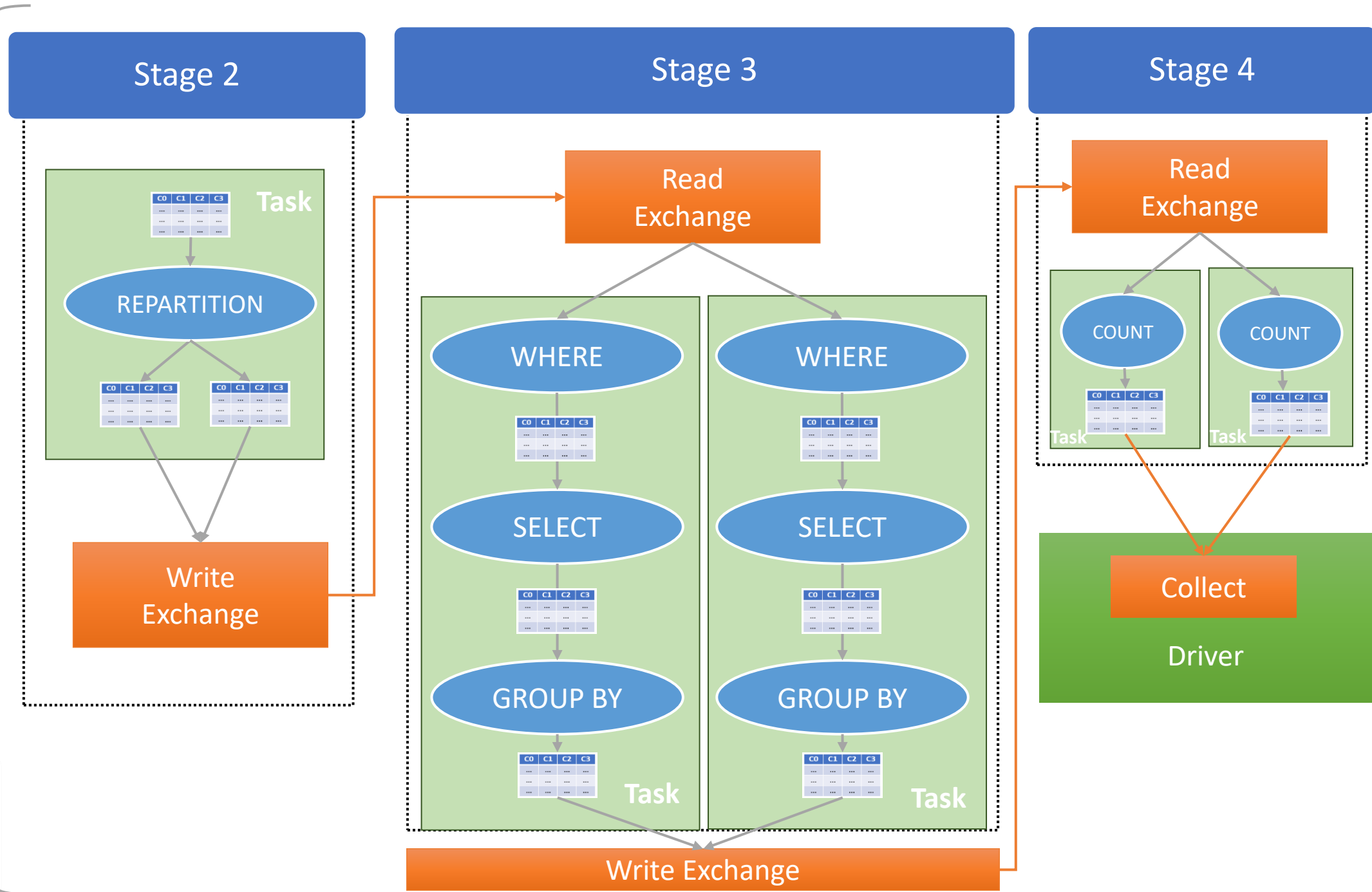
GROUP BY

C0	C1	C2	C3
val	val	val	val
val	val	val	val
val	val	val	val

COUNT

C0	C1	C2	C3
val	val	val	val
val	val	val	val
val	val	val	val

Job 2





LEARNING JOURNAL

[www.learningjournal.guru](http://www.learningjournal.guru)