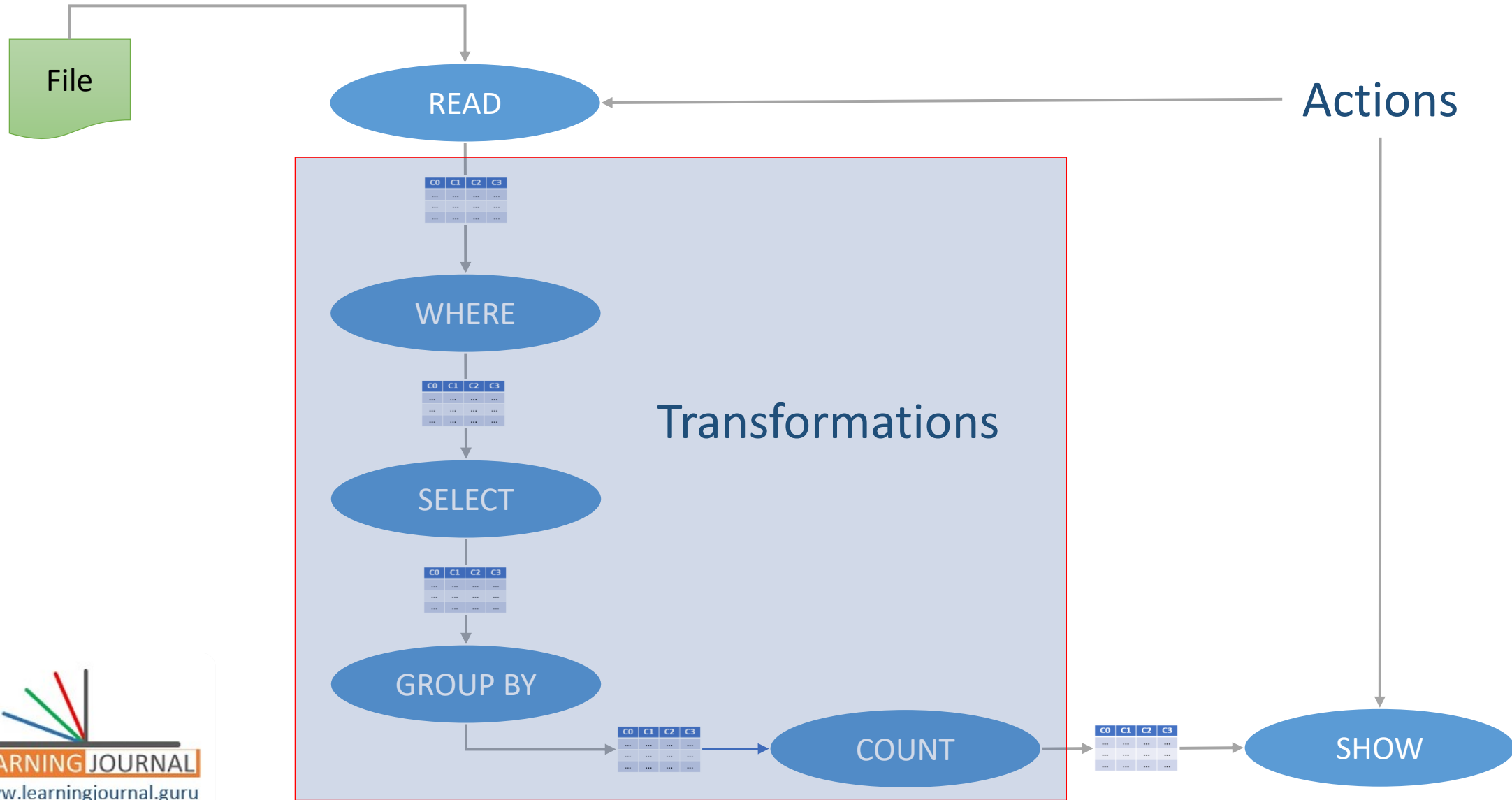




Apache Spark

Spark Jobs, Stages, and Tasks

Spark Operations

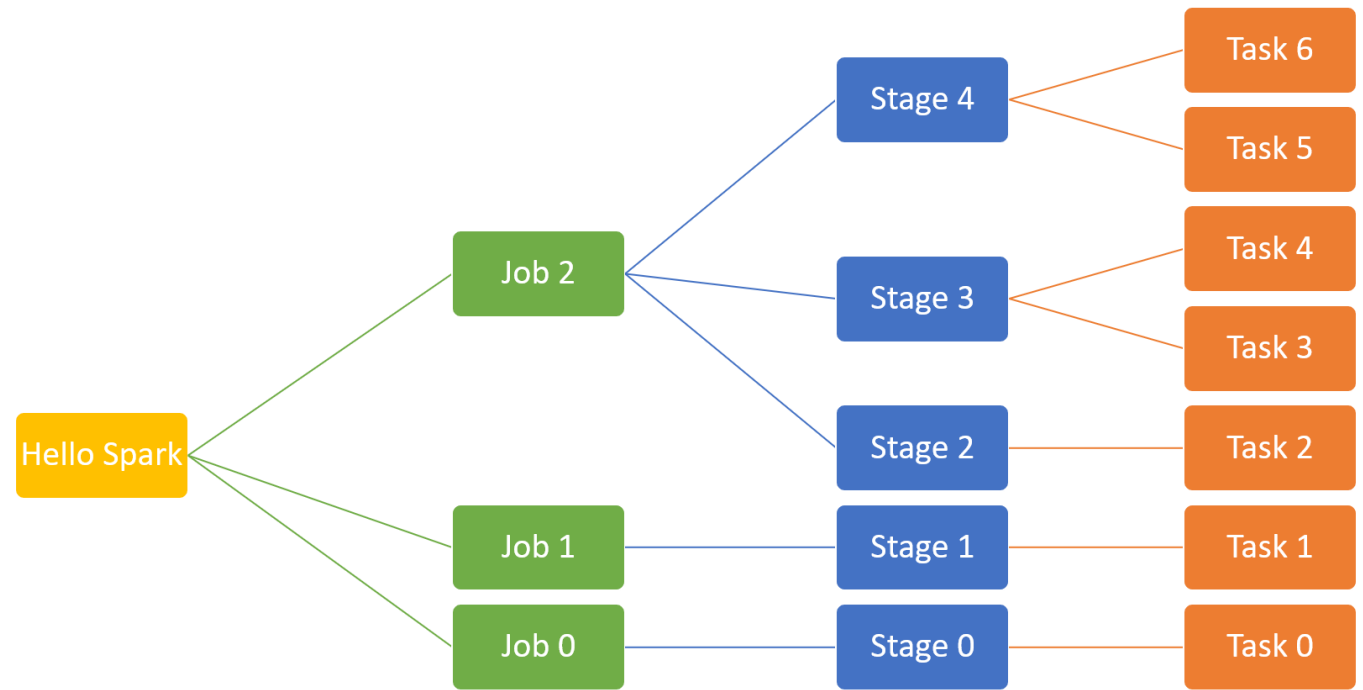


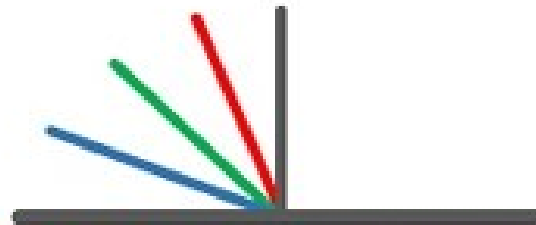
Spark Execution Plan

```
def main(args: Array[String]): Unit = {  
  if (args.length == 0) {  
    Logger.error("Usage: HelloSpark filename")  
    System.exit(status = 1)  
  }  
  Logger.info("Starting Hello Spark")  
  val spark = SparkSession.builder()  
    .config(getSparkAppConf)  
    .getOrCreate()  
  //Logger.info("spark.conf=" + spark.conf.getAll.toString())  
  
  val surveyRawDF = LoadSurveyDF(spark, args(0))  
  val partitionedSurveyDF = surveyRawDF.repartition(numPartitions = 2)  
  val countDF = countByCountry(partitionedSurveyDF)  
  
  Logger.info(countDF.collect().mkString("->"))  
  
  Logger.info("Finished Hello Spark")  
  scala.io.StdIn.readLine()  
  spark.stop()  
}
```

```
def countByCountry(surveyDF: DataFrame): DataFrame = {  
  surveyDF.where(conditionExpr = "Age < 40")  
    .select(col = "Age", cols = "Gender", "Country", "state")  
    .groupBy(col1 = "Country")  
    .count()  
}
```

: Dataset[Row]
: DataFrame
: RelationalGroupedDataset
: DataFrame





LEARNING JOURNAL

www.learningjournal.guru