

# Apache Spark 2 using Scala - Data Processing – Overview

## Apache Spark 2 – Data Processing Overview (Scala)

Spark's data processing consists of:

- **Transformations** → Lazy operations that define a new dataset from an existing one (don't execute immediately).
- **Actions** → Trigger execution and return results to the driver or save them to storage.

### 1. Spark Transformations

Transformation	Description	Scala Example
<b>map</b>	Applies a function to each element and returns a new RDD/DataFrame.	<code>val rdd2 = rdd.map(x =&gt; x * 2)</code>
<b>flatMap</b>	Like <code>map</code> , but flattens the results (one input → many outputs).	<code>val words = lines.flatMap(line =&gt; line.split(" "))</code>
<b>filter</b>	Keeps only elements satisfying a condition.	<code>val even = rdd.filter(_ % 2 == 0)</code>
<b>mapPartitions</b>	Runs a function on each partition (more efficient for batch ops).	<code>val result = rdd.mapPartitions(iter =&gt; iter.map(_ * 2))</code>
<b>mapPartitionsWithIndex</b>	Same as above but also gets partition index.	<code>val indexed = rdd.mapPartitionsWithIndex((idx, iter) =&gt; iter.map(x =&gt; (idx, x)))</code>
<b>distinct</b>	Removes duplicates.	<code>val unique = rdd.distinct()</code>
<b>sample</b>	Randomly samples elements.	<code>val sampleRDD = rdd.sample(false, 0.1)</code>
<b>union</b>	Combines two RDDs.	<code>val unionRDD = rdd1.union(rdd2)</code>
<b>intersection</b>	Elements common to both RDDs.	<code>val interRDD = rdd1.intersection(rdd2)</code>
<b>subtract</b>	Elements in one RDD but not in the other.	<code>val diff = rdd1.subtract(rdd2)</code>

<b>cartesian</b>	Cartesian product between two RDDs.	<code>val cart = rdd1.cartesian(rdd2)</code>
<b>groupBy</b>	Groups elements using a function.	<code>val grouped = rdd.groupBy(_ % 2)</code>
<b>groupByKey</b>	Groups key-value pairs by key (expensive).	<code>val grouped = rdd.groupByKey()</code>
<b>reduceByKey</b>	Merges values for each key using an operation (faster than groupByKey).	<code>val reduced = rdd.reduceByKey(_ + _)</code>
<b>aggregateByKey</b>	Aggregates values with separate seqOp and combOp.	<code>val agg = rdd.aggregateByKey(0)(_+_, _+_)</code>
<b>combineByKey</b>	Most general form of aggregation.	<code>val combined = rdd.combineByKey(v =&gt; (v,1), (c:(Int,Int),v)=&gt; (c._1+v, c._2+1), (c1:(Int,Int),c2:(Int,Int))=&gt; (c1._1+c2._1, c1._2+c2._2))</code>
<b>sortByKey</b>	Sorts key-value RDD by key.	<code>val sorted = rdd.sortByKey()</code>
<b>join</b>	Inner join between two key-value RDDs.	<code>val joined = rdd1.join(rdd2)</code>
<b>leftOuterJoin</b>	Left join between RDDs.	<code>val leftJoin = rdd1.leftOuterJoin(rdd2)</code>
<b>rightOuterJoin</b>	Right join between RDDs.	<code>val rightJoin = rdd1.rightOuterJoin(rdd2)</code>
<b>cogroup</b>	Groups values from multiple RDDs sharing the same key.	<code>val cogrouped = rdd1.cogroup(rdd2)</code>
<b>repartition</b>	Changes number of partitions (shuffle).	<code>val repart = rdd.repartition(4)</code>
<b>coalesce</b>	Reduces partitions without full shuffle.	<code>val coalesced = rdd.coalesce(2)</code>
<b>pipe</b>	Sends RDD elements to an external process.	<code>val piped = rdd.pipe("wc -l")</code>
<b>zip</b>	Combines two RDDs element-wise.	<code>val zipped = rdd1.zip(rdd2)</code>

## 2. Spark Actions

Action	Description	Scala Example
<b>collect</b>	Returns all elements to driver (careful with large data).	<code>val data = rdd.collect()</code>
<b>count</b>	Returns number of elements.	<code>val cnt = rdd.count()</code>
<b>first</b>	Returns the first element.	<code>val firstVal = rdd.first()</code>
<b>take(n)</b>	Returns first n elements.	<code>val few = rdd.take(5)</code>
<b>takeOrdered(n)</b>	Returns first n elements in order.	<code>val ordered = rdd.takeOrdered(5)</code>
<b>top(n)</b>	Returns top n elements in descending order.	<code>val topVals = rdd.top(5)</code>
<b>reduce</b>	Aggregates elements using a function.	<code>val sum = rdd.reduce(_ + _)</code>
<b>fold</b>	Same as reduce but with zero value.	<code>val sum = rdd.fold(0)(_ + _)</code>
<b>aggregate</b>	Aggregates with separate seqOp and combOp.	<code>val agg = rdd.aggregate(0)(_ + _, _ + _)</code>
<b>foreach</b>	Runs a function on each element (no return).	<code>rdd.foreach(println)</code>

<b>countByKey</b>	Returns count of values per key.	val counts = rdd.countByKey()
<b>countByValue</b>	Returns counts of each element.	val counts = rdd.countByValue()
<b>saveAsTextFile</b>	Saves RDD to text file.	rdd.saveAsTextFile("output")
<b>saveAsSequenceFile</b>	Saves RDD as Hadoop sequence file.	rdd.saveAsSequenceFile("output")
<b>saveAsObjectFile</b>	Saves RDD as serialized objects.	rdd.saveAsObjectFile("output")

## Example Program – Transformations + Actions

```
import org.apache.spark.sql.SparkSession
object SparkTransformationsActions {
  def main(args: Array[String]): Unit = {
    val spark = SparkSession.builder()
      .appName("Spark Transformations and Actions")
      .master("local[*]")
      .getOrCreate()
    val sc = spark.sparkContext
    // Sample RDD
    val rdd = sc.parallelize(List(1, 2, 3, 4, 5, 6, 2, 4))
    // Transformation: Filter even numbers and double them
    val transformedRDD = rdd.filter(_ % 2 == 0).map(_ * 2).distinct()
    // Action: Collect and print
    println("Transformed Data: " + transformedRDD.collect().mkString(", "))
    // Transformation: Create pair RDD
    val pairRDD = rdd.map(x => (x, 1))
    // Reduce by key
    val counts = pairRDD.reduceByKey(_ + _)
    // Action: Display word counts
    println("Element Counts:")
    counts.collect().foreach(println)
    spark.stop()
  }
}
```

### ✅ Key Points to Remember

- Transformations are **lazy** → nothing happens until an action is called.
- Use **reduceByKey** instead of groupByKey for performance.
- **collect()** should be avoided for large datasets.
- Partitioning and persistence can hugely improve performance.