
WORKING WITH SPARK PAIR RDD

By www.HadoopExam.com

**Note: These instructions should be used with the HadoopExam Apache Spark: Professional Trainings.
Where it is executed and you can do hands on with trainer.**

Working With Pair RDD

Example of Pair RDD: flipkart.com

✓ 10.100.0.1	→ Samsung S6
10.101.21.1	→ Samsung S6
10.25.125.3	→ iPhone 6
10.44.126.8	→ office chairs
127.45.46.3	→ Lockers
128.46.129.4	→ Kids Toys car
140.111.128.129	→ Samsung S6
✓ 10.100.0.1	→ Samsung S6
✓ 10.100.0.1	→ iPhone 6

⇒ most demanding phone — Samsung S6

⇒ Some user is comparing — iPhone 6 & Samsung S6

10.100.0.1 → is likely to buy high end phone.

⇒ Flipkart will share these search with google AdSense. Adwords.

⇒ And whenever possible google will show Samsung S6 & iPhone 6 advertisement on possible Ad space on webpage. if user is browsing. (With Buy Button from flipkart.com)

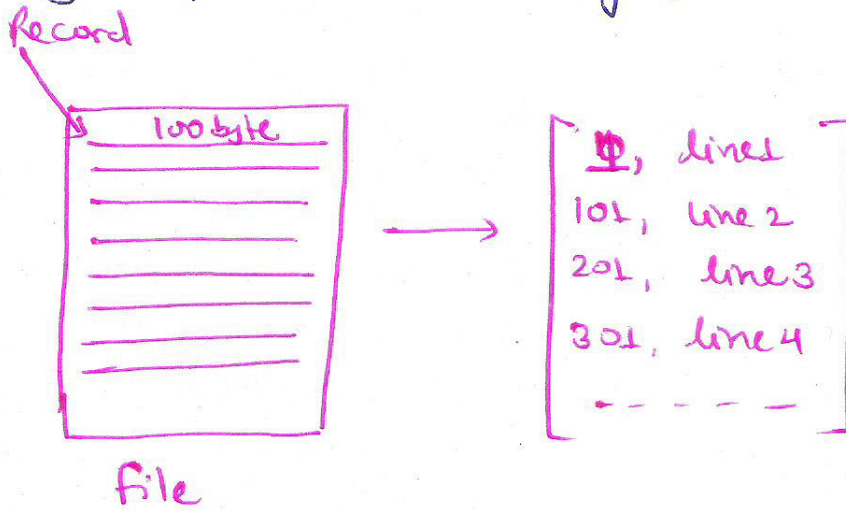
⇒ So this pair RDD has following key and Value pairs.

[ip, search keyword]

⇒ Creating Pair RDD's

⇒ In Hadoop when file is loaded, each line became a value.

⇒ Byte offset as a key for that line.



⇒ But in Spark, it is created as below.

Val

```
Val token = sc.textfile("HadoopExam.log")
Val pairs = token.map(x => x.split(' ')(0), x)
```

⇒ Regular RDD converted in a pair RDD using map() function.

Aggregation In Pair RDDs :

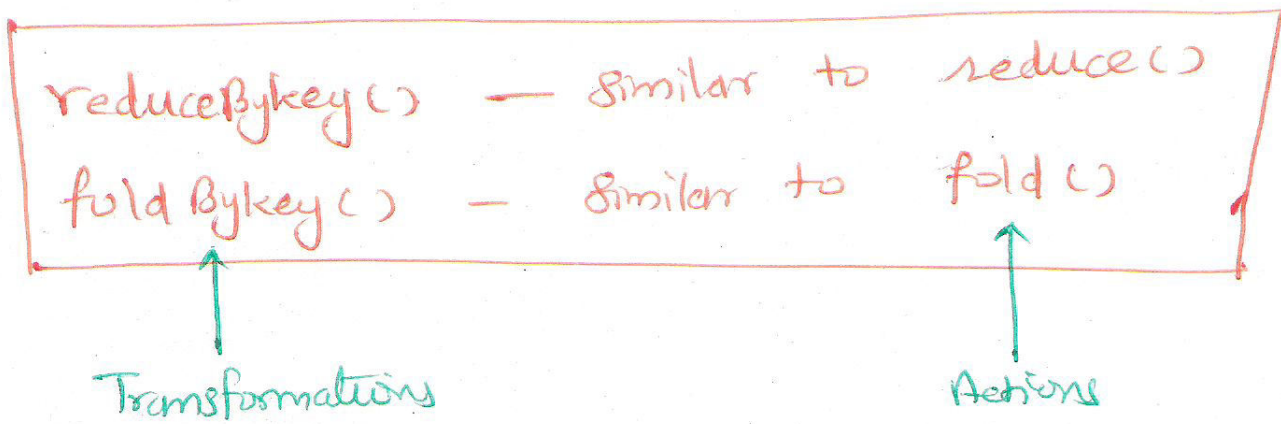
⇒ It is common to aggregate all values with the same key together.

⇒ On Regular RDD :

```
fold(), reduce(), combine()
```

All Are Actions

=> But in pair RDD, we have similar functions for transformations.



=> reduceByKey(): take a function as input and use it to combine value for some key.

e.g. (+)

=> However, it will reduce together only for some key.

=> Generally (key, value) pairs file are bigger as stored in HDFS. (GB to TB)

=> Hence, reduceByKey() is not implemented as an action.

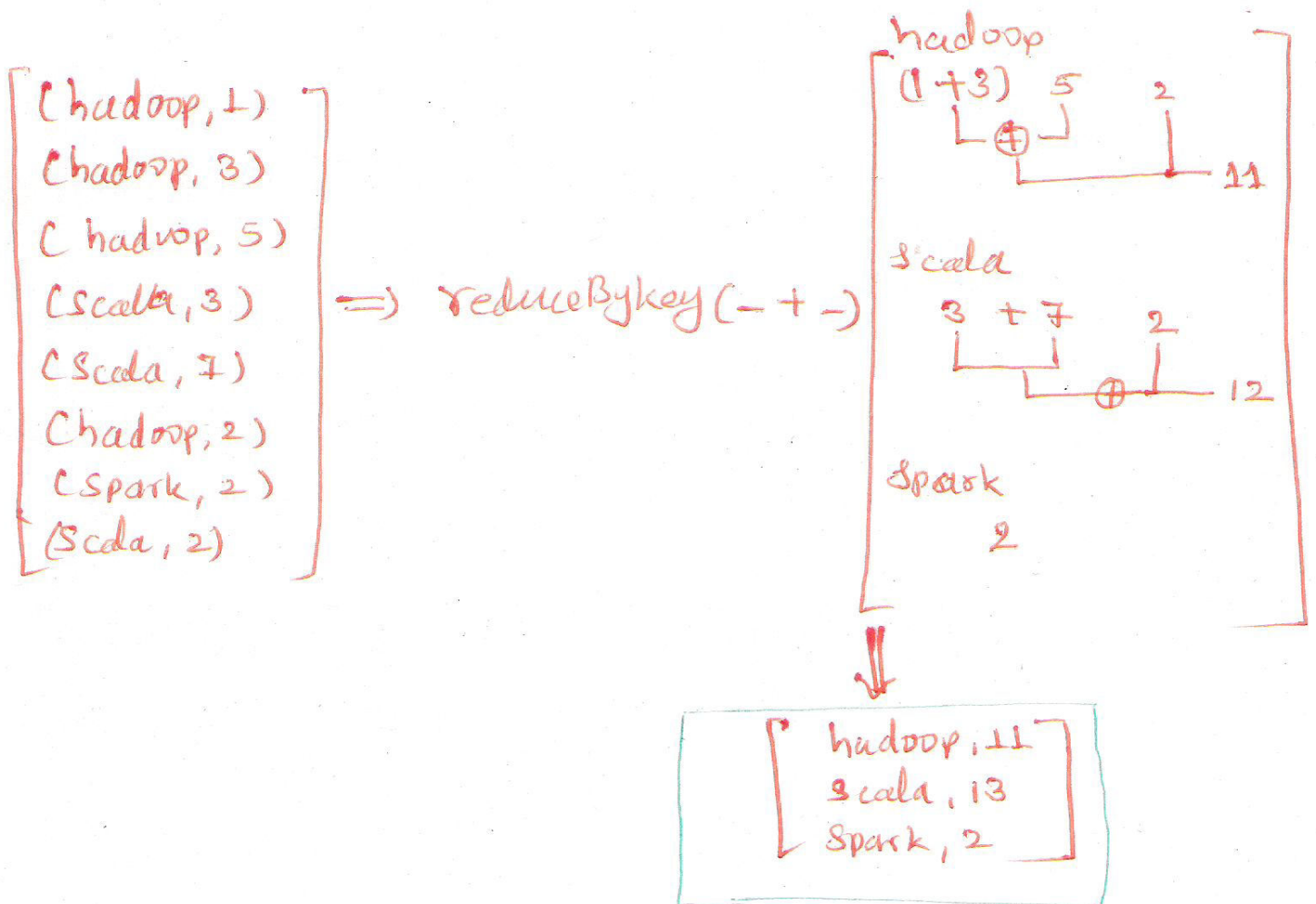
=> Because action always returns a value to the user program.

(4)

⇒ If values are very large, then it would be a memory issue or space issue on single machine.

⇒ Instead, it returns a new RDD consisting of each key and reduced values of that key.

⇒ And RDD is partitioned across cluster hence will not cause space issue.



⇒ `foldByKey()`:

It is also quite similar to `fold()` Action.

- ⇒ It requires the zero value similar to fold() function, for initialization value.
- ⇒ If you know the Hadoop framework, then there was a concept of combiner.
- ⇒ Combiner is like mini reducers, which work independently on each partition (on individual node).
- ⇒ And combines the data locally.
- ⇒ Similarly reduceByKey() and foldByKey() will automatically perform combining locally on each node/machine before computing global totals for each key.
- ⇒ In Spark user do not have to specify combiner.
- ⇒ However, you can create a custom combiner.
- ⇒ Let's see word count example.

val input = sc.textFile("hdfs://hadoopexam.log")

val words = input.flatMap(x => x.split(' '))

val results = words.map(x => x, 1)

• reduceByKey((x, y) => x + y)

=> combineByKey(): -

=> This is the underline function of many other function.

=> It is most general of the per-key aggregation function.

=> Like aggregate(), combineByKey() allows the user to return values that are not same type as our input data.

val result = input.combineByKey()

(v) => (v, 1)

(acc: (Int, Int), v) =>

acc._1 + v, acc._2 + 1)

(acc1: (Int, Int), acc2: (Int, Int))

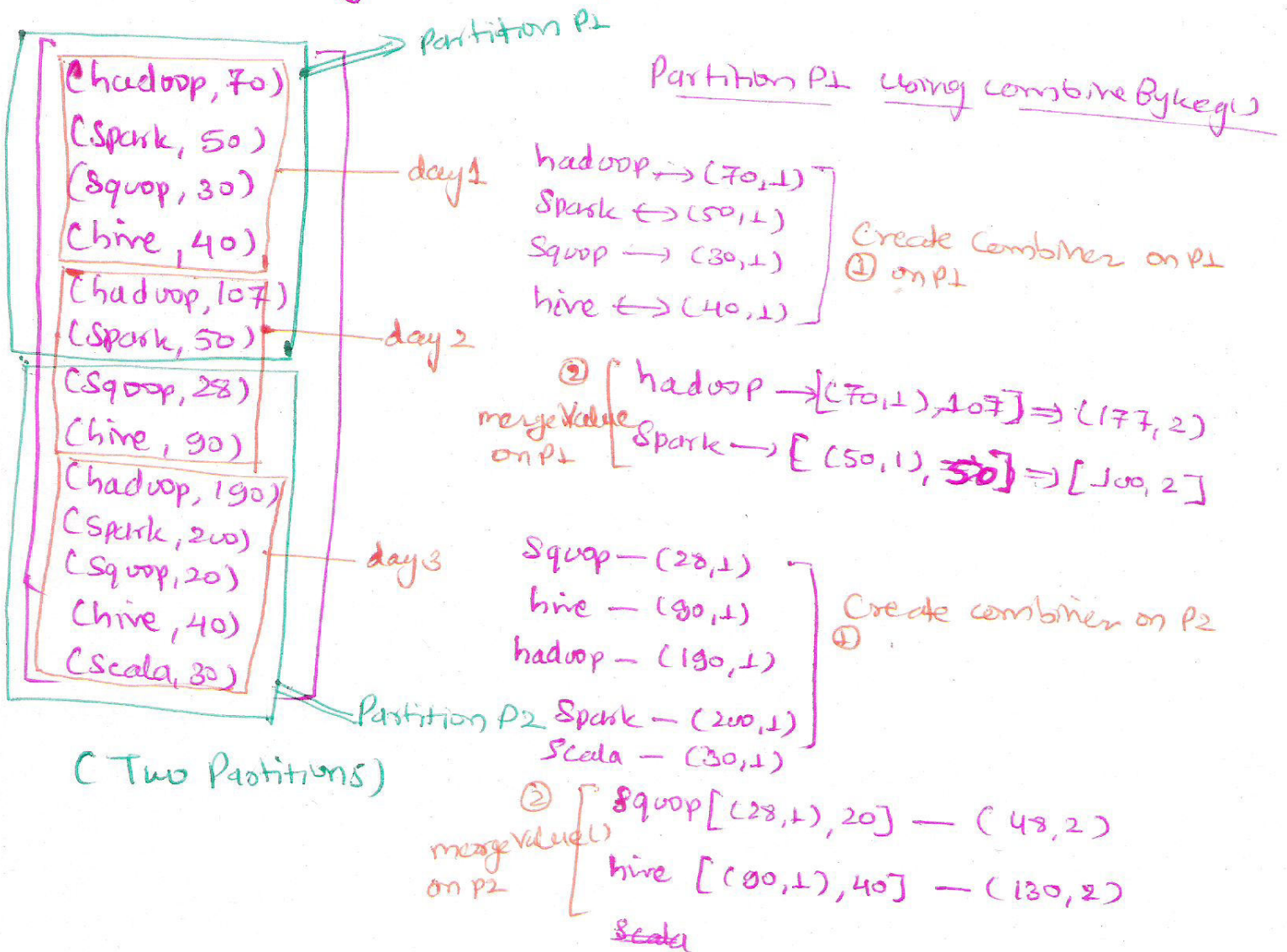
=> (acc1._1 + acc2._1, acc1._2 + acc2._2))

• map { case (key, value) => (key, value._1) }

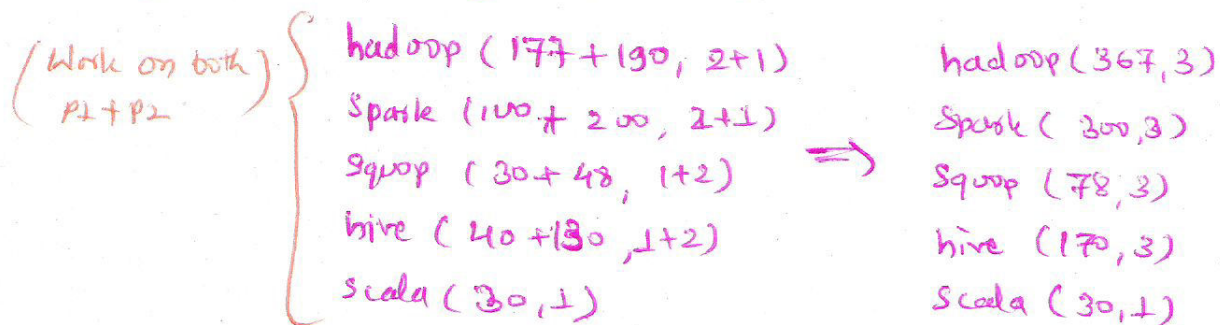
value._2

result.collectMap().map(println(-))

Example:- Average search count for each keyword during last 3 days.



③ mergeCombiner (Apply 3rd function)



⇒ Then do Average calculation
Value/Count

⇒ There are many options for combining our data by key.

⇒ Most of them are implemented on top of

`combineByKey()`

but provide a simpler interface.

WWW.HadoopExam.Com

WWW.GuideTechie.Com