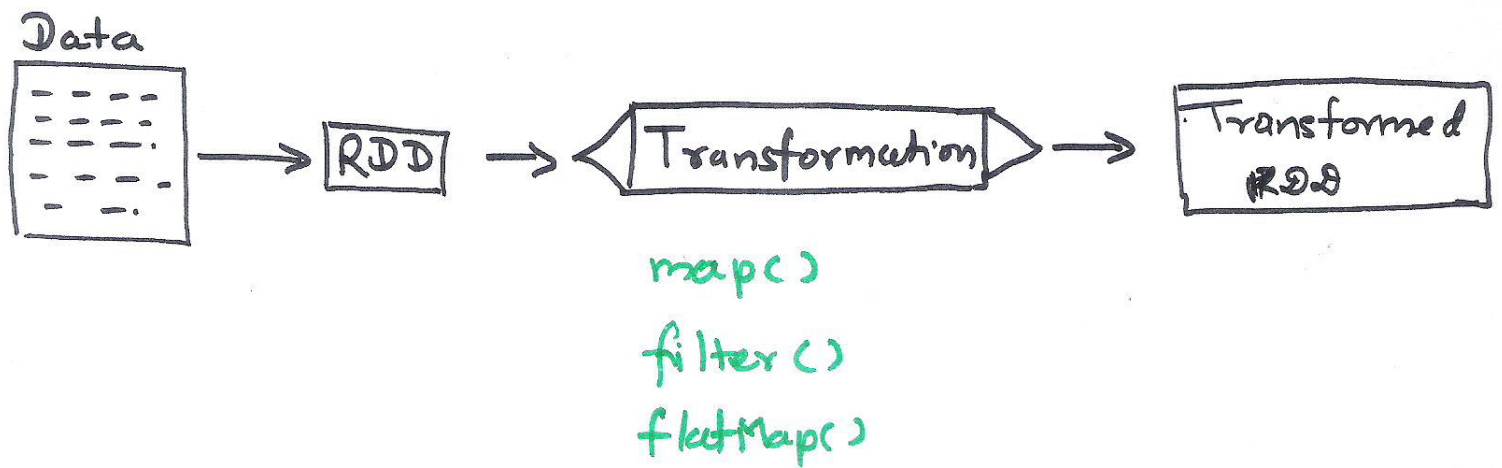

SPARK TRANSFORMATION IN DEPTH

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.**

Transformation In RDD



Element-wise Transformation: -

- Transformation works on each element
- Two most common transformation
 - map()
 - filter()

⇒ map() :

List [1, 2, 3, 4].map(?)

① What do you want to do with each element of list

② I want to double each element

③ Replace ? with function, which can double an element

List [1, 2, 3, 4].map(x ⇒ x * 2)



List [2, 4, 6, 8]

Filter :

② What do you want to do with each element of list

List [1, 2, 3, 4, 5, 6]. filter(?)

② We want to filter/select values other than 3

③ Pass the function which filter values other than 3

$$x \Rightarrow x \neq 3$$

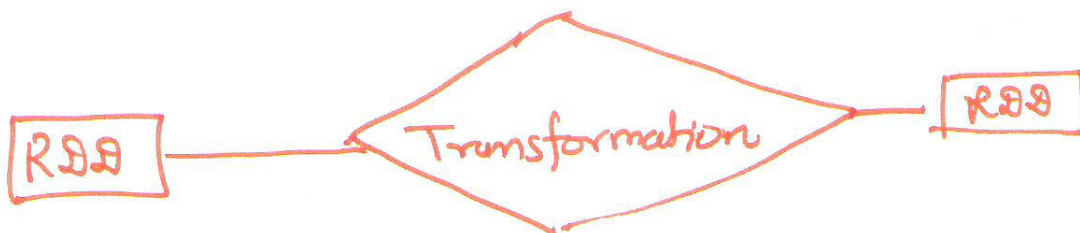
List [1, 2, 3, 4, 5, 6]. filter($x \Rightarrow x \neq 3$)

Input RDD



List [1, 2, 4, 5, 6] ← Output RDD

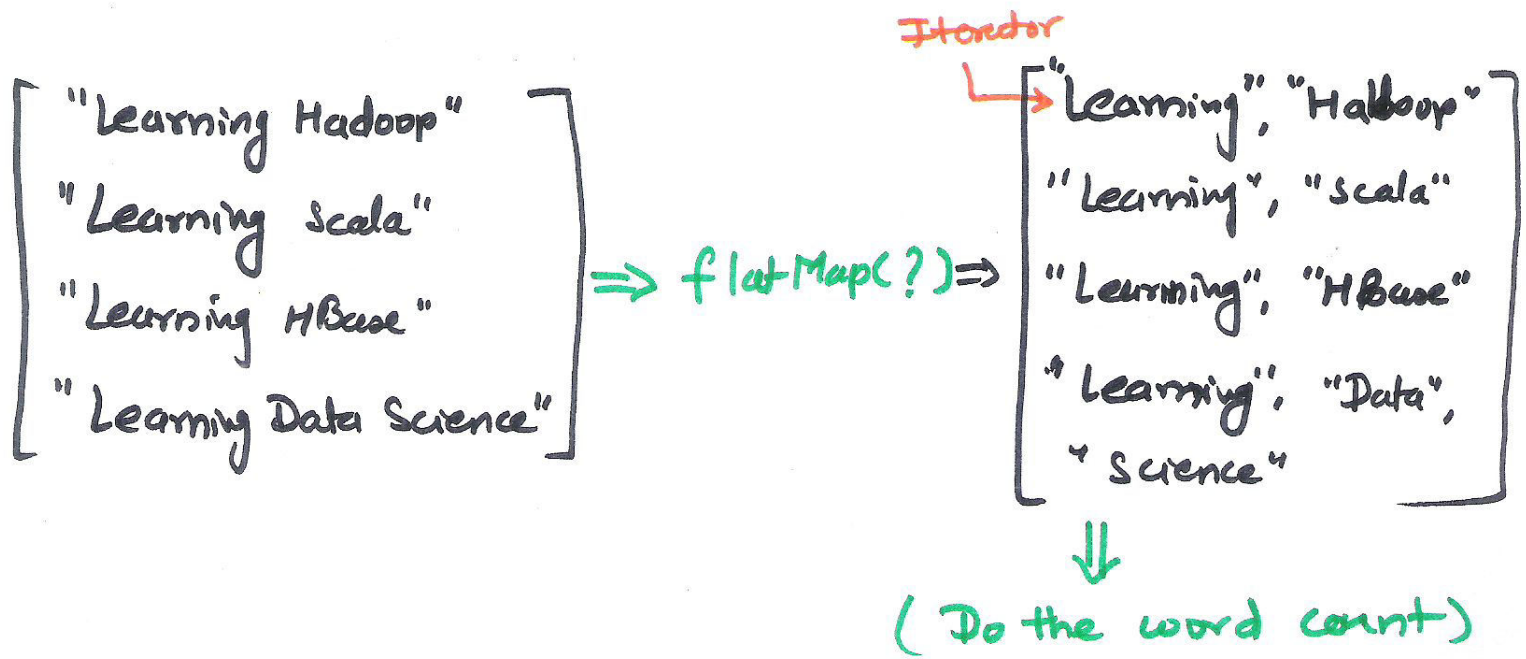
\Rightarrow map() and filter() both are the function which independently applied on each element of input RDD



\Rightarrow Return type does not have to be the same as its input type.



FlatMap: produce multiple output for each element



$\text{flatMap}(?)$ Provide the function which will help to create multiple elements (Return iterator) e.g. `split()` function

\Rightarrow The function (?) we provide to `flatMap()` is called individually for each element to input RDD

\Rightarrow This "?" function will not return individual element, it return an iterator with return values.

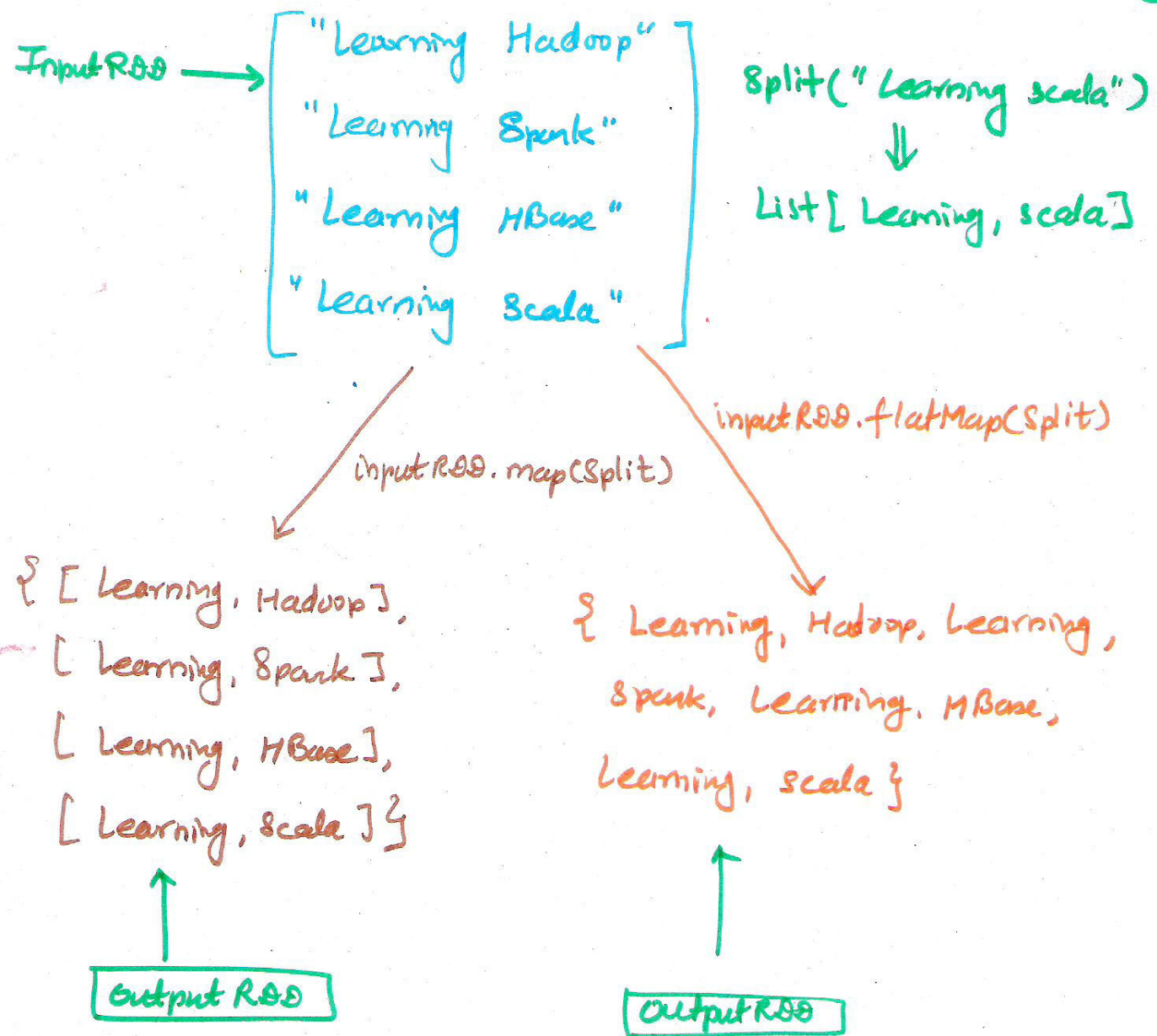
$\text{flatMap}[\text{List}(\text{Iterator1}, \text{Iterator2}, \dots, \text{IteratorN})]$

\Rightarrow Now this `flatMap` will not return list of iterators

\Rightarrow It will return all the elements combined from all iterator.

\rightarrow Given Sample Split all the words from line

④



Hands on Exercise Lab

- Map function squaring the values

```
val input = sc.parallelize(List(1, 2, 3, 4))  
val result = input.map(x => x * x)  
println(result.collect().mkString(","))
```

- Splitting lines in multiple words.

```
val lines = sc.parallelize(List("hello HadoopExam.com", "hi"))  
val words = lines.flatMap(line => line.split(" "))  
words.first() // returns "hello"
```

- Return an RDD consisting of only elements that pass the condition passed to filter().

```
val result = input.filter(x => x != 1)  
println(result.collect().mkString(","))
```

- Another example on map()

```
val l = sc.parallelize(List(1, 2, 3, 4, 5))  
  
val result = l.map(x => x*2 )  
println(result.collect().mkString(","))  
  
def f(x: Int) = if (x > 2) Some(x) else None  
  
val result = l.map(x => f(x))  
println(result.collect().mkString(","))  
  
def g(v: Int) = List(v-1, v, v+1)  
val result = l.map(x => g(x))  
println(result.collect().mkString(","))  
  
val result = l.flatMap(x => g(x))  
println(result.collect().mkString(","))
```

- Set operation on RDD

```
val rdd = sc.parallelize(List(1,2,3))  
val other = sc.parallelize(List(3,4,5))  
  
val result = rdd.union(other)
```

```
println(result.collect().mkString(","))

val result = rdd.intersection(other)
println(result.collect().mkString(","))

val result = rdd.subtract(other)
println(result.collect().mkString(","))

val result = rdd.cartesian(other)
println(result.collect().mkString(","))
```

HADOOPEXAM ALL PRODUCTS TRAINING'S

- [Hadoop BigData Professional Training \(3500INR/\\$79\)](#)
- [HBase \(NoSQL\) Professional Training \(3500INR/\\$79\)](#)
- [Apache Spark Professional Training \(3500INR/\\$79\)](#)

CLOUDERA HADOOP AND NOSQL CERTIFICATION

- [CCD410 : Hadoop Developer](#)
- [CCA50X : Hadoop Administrator](#)
- [HBase/NoSQL Developer](#)

DATABRICKSA OREILLY SPARK CERTIFICATION

- [Apache Spark Developer](#)

AWS : AMAZON WEBSERVICE CERTIFICATION

- [AWS Solution Architect : Associate](#)
- [AWS Solution Architect: Professional](#)
- [AWS Developer : Associate](#)
- [AWS Sysops Admin : Associate](#)

MICROSOFT AZURE CERTIFICATION

- [Azure 70-532](#)
- [Azure 70-533](#)

DATA SCIENCE CERTIFICATION

- [EMC E20-007](#)

EMC CERTIFICATIONS

- [EMC E20-007](#)

SAS ANALYTICS CERTIFICATION

- [SAS Base A00-211](#)
- [SAS Advanced A00-212](#)
- [SAS Analytics : A00-240](#)
- [SAS Administrator : A00-250](#)

ORACLE JAVA CERTIFICATION

- [Java 1z0-808](#)
- [Java 1z0-809](#)

ORACLE DATABASE CLOUD CERTIFICATION

- [1z0-060 \(Oracle 12c\)](#)
- [1z0-061 \(Oracle 12c\)](#)