### **RDDs and key-value pairs**

## **RDDs of key-value pairs**

- Spark supports also RDDs of key-value pairs
  - Key-value pairs in python are represented by means of python tuples
    - The first value is the key part of the pair
    - The second value is the value part of the pair
- RDDs of key-value pairs are sometimes called "pair RDDs"

## **RDDs of key-value pairs**

- RDDs of key-value pairs are characterized by specific operations
  - reduceByKey(), join(), etc.
  - These operations analyze the content of one group (key) at a time
- RDDs of key-value pairs are characterized also by the operations available for the "standard" RDDs
  - filter(), map(), reduce(), etc.

## **RDDs of key-value pairs**

- Many applications are based on RDDs of keyvalue pairs
- The operations available for RDDs of keyvalue pairs allow
  - "grouping" data by key
  - performing computation by key (i.e., by group)
- The basic idea is similar to the one of the MapReduce-based programs in Hadoop
  - But there are more operations already available

# Creating RDDs of key-value pairs

## **Creating RDDs of key-value pairs**

- RDDs of key-value pairs can be built
  - From other RDDs by applying the map() or the flatMap() transformation on other RDDs
  - From a Python in-memory collection of tuple (key-value pairs) by using the parallelize() method of the SparkContext class

## **Creating Pair RDDs**

- Key-value pairs are represented as tuples composed of two elements
  - Key
  - Value
- The standard built-in Python tuples are used

### RDDs of key-value pairs by using the Map transformation

## RDDs of key-value pairs by using the map transformation

### Goal

- Define an RDD of key-value pairs by using the map transformation
- Apply a function **f** on each element of the input RDD that **returns one tuple for each input** element
  - The new RDD of key-value pairs contains one tuple y for each element x of the "input" RDD

# RDDs of key-value pairs by using the map transformation

- Method
  - The standard map(f) transformation is used
    - The new RDD of key-value pairs contains one tuple y for each input element x of the "input" RDD
      y=f(x)

# RDDs of key-value pairs by using the map transformation: Example

- Create an RDD from a textual file containing the first names of a list of users
  - Each line of the file contains one first name
- Create an RDD of key-value pairs containing a list of pairs (first name, 1)

## RDDs of key-value pairs by using the map transformation: Example

# Read the content of the input textual file
namesRDD = sc.textFile("first\_names.txt")

# Create an RDD of key-value pairs
nameOnePairRDD = namesRDD.map(lambda name: (name, 1))

# RDDs of key-value pairs by using the map transformation: Example

# Read the content of the input textual file
namesRDD = sc.textFile("first\_names.txt")



## RDDs of key-value pairs by using the flatMap transformation

# RDDs of key-value pairs by using the flatMap transformation

### Goal

- Define an RDD of key-value pairs by using the flatMap transformation
- Apply a function **f** on each element of the input RDD that **returns a list of tuples for each input** element
  - The new PairRDD contains all the pairs obtained by applying f on each element x of the "input" RDD

## RDDs of key-value pairs by using the flatMap transformation

### Method

- The standard flatMap(f) transformation is used
  - The new RDD of key-value pairs contains the tuples returned by the execution of f on each element x of the "input" RDD
    - [y]=f(x)
      - Given a element x of the input RDD, f applied on x returns a list of pairs [y]
      - The new RDD is a "list" of pairs contains all the pairs of the returned list of pairs. It is not an RDD of lists.
    - [y] can be the empty list

# RDDs of key-value pairs by using the flatMap transformation: Example

- Create an RDD from a textual file
  - Each line of the file contains a set of words
- Create a PairRDD containing a list of pairs (word, 1)
  - One pair for each word occurring in the input document (with repetitions)

# RDDs of key-value pairs by using the flatMap transformation: Example v1

# Define the function associated with the flatMap transformation def wordsOnes(line):

```
pairs = []
for word in line.split(''):
pairs.append((word, 1))
return pairs
```

# Read the content of the input textual file linesRDD = sc.textFile("document.txt")

# Create an RDD of key-value pairs based on the input document # One pair (word,1) for each input word wordOnePairRDD = linesRDD.flatMap(wordsOnes)

# RDDs of key-value pairs by using the flatMap transformation: Example v2

# Read the content of the input textual file linesRDD = sc.textFile("document.txt")

# Create an RDD of key-value pairs based on the input document # One pair (word,1) for each input word wordOnePairRDD = linesRDD.flatMap(lambda line: \ map(lambda w: (w, 1), line.split('')))

# RDDs of key-value pairs by using the flatMap transformation: Example v2

# Read the content of the input textual file linesRDD = sc.textFile("document.txt")

# Create an RDD of key-value pairs based on the input document # One pair (word,1) for each input word wordOnePairRDD = linesRDD.flatMap(lambda line: \



### Goal

 Use the parallelize method to create an RDD of keyvalue pairs from a local python in-memory collection of tuples

#### Method

- It is based on the standard parallelize(c) method of the SparkContext class
- Each element (tuple) of the local python collection becomes a key-vaue pair of the returned RDD

- Create an RDD from a local python list containing the following key-value pairs
  - ("Paolo", 40)
  - ("Giorgio", 22)
  - ("Paolo", 35)

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

# Create the RDD of pairs from the local collection nameAgePairRDD = sc.parallelize(nameAge)



# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

# Create the RDD of pairs from the local collection
nameAgePairRDD = sc.parallelize(nameAge)

Create an RDD or key-value pairs based on the content of the local in-memory python list. The RDD is stored in the "distributed" main memory of the cluster servers

# Transformations on RDDs of key-value pairs

### Transformations on RDDs of keyvalue pairs

- All the "standard" transformations can be applied
  - Where the specified "functions" operate on tuples
- Specific transformations are available
  - E.g., reduceByKey(), groupByKey(), mapValues(), join(), ...

#### Goal

- Create a new RDD of key-value pairs where there is one pair for each distinct key k of the input RDD of key-value pairs
  - The value associated with key k in the new RDD of key-value pairs is computed by applying a function f on the values associated with k in the input RDD of key-value pairs
    - The function f must be associative and commutative
      - otherwise the result depends on how data are partitioned and analyzed
  - The data type of the new RDD of key-value pairs is the same of the "input" RDD of key-value pairs

#### Method

- The reduceByKey transformation is based on the reduceByKey(f) method of the RDD class
- A function f is passed to the reduceByKey method
  - Given the values of two input pairs, f is used to combine them in one single value
  - f is recursively invoked over the values of the pairs associated with one key at a time until the input values associated with one key are "reduced" to one single value
- The retuned RDD contains a number of key-value pairs equal to the number of distinct keys in the input key-value pair RDD

- Similarly to the reduce() action, the reduceByKey() transformation aggregate values
   However,
  - reduceByKey() is executed on RDDs of key-value pairs and returns a set of key-value pairs
    - reduce() is executed on an RDD and returns one single value (stored in a local python variable)
- And
  - reduceByKey() is a transformation
    - reduceByKey() is executed lazily and its result is stored in another RDD
  - Whereas reduce() is an action

- Shuffle
  - A shuffle operation is executed for computing the result of the reduceByKey() transformation
    - The result/value for each group/key is computed from data stored in different input partitions

### ReduceByKey transformation: Example

- Create an RDD from a local python list containing the pairs
  - ("Paolo", 40)
  - ("Giorgio", 22)
  - ("Paolo", 35)
  - The key is the first name of a user and the value is his/her age
- Create a new RDD of key-value pairs containing one pair for each name. In the returned RDD, associate each name with the age of the youngest user with that name

### ReduceByKey transformation: Example

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

# Create the RDD of pairs from the local collection nameAgePairRDD = sc.parallelize(nameAge)

# Select for each name the lowest age value youngestPairRDD=nameAgePairRDD.reduceByKey(lambda age1, age2:\ min(age1, age2))

### ReduceByKey transformation: Example

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

# Create the RDD of pairs from the local collection nameAgePairRDD = sc.parallelize(nameAge)

# Select for each name the lowest age value youngestPairRDD= nameAgePairRDD.reduceByKey(lambda age1, age2:\ min(age1, age2))

The returned RDD of key-value pairs contains one pair for each distinct input key (i.e., for each distinct name in this example)
#### Goal

- The foldByKey() has the same goal of the reduceBykey() transformation
- However, foldByKey()
  - Is characterized also by a "zero" value
  - Functions must be associative but are not required to be commutative

#### Method

- The foldByKey transformation is based on the foldByKey(zeroValue, op) method of the RDD class
- A function op is passed to the fold method
  - Given values of two input pairs, op is used to combine them in one single value
  - op is also used to combine input values with the "zero" value
  - op is recursively invoked over the values of the pairs associated with one key at a time until the input values are "reduced" to one single value
- The "zero" value is the neutral value for the used function
   op
  - i.e., "zero" combined with any value v by using op is equal to v

- Shuffle
  - A shuffle operation is executed for computing the result of the foldByKey() transformation
    - The result/value for each group/key is computed from data stored in different input partitions

### FoldByKey transformation: Example

- Create an RDD from a local python list containing the pairs
  - ("Paolo", "Message1")
  - ("Giorgio", "Message2")
  - ("Paolo", "Message3")
  - The key is the first name of a user and the value is a message published by him/her
- Create a new RDD of key-value pairs containing one pair for each name. In the returned RDD, associate each name the concatenation of its messages (preserving the order of the messages in the input RDD)

#### FoldByKey transformation: Example

# Create the local python list nameMess = [("Paolo", "Message1"), ("Giorgio", "Message2"), \ ("Paolo", "Message3")]

# Create the RDD of pairs from the local collection nameMessPairRDD = sc.parallelize(nameMess)

# Concatenate the messages of each user concatPairRDD= nameMessPairRDD.foldByKey(", lambda m1, m2:\ m1+m2)

#### Goal

- Create a new RDD of key-value pairs where there is one pair for each distinct key k of the input RDD of key-value pairs
  - The value associated with the key k in the new RDD of keyvalue pairs is computed by applying user-provided functions on the values associated with k in the input RDD of key-value pairs
    - The user-provided "function" must be associative
      - otherwise the result depends how data are partitioned and analyzed
  - The data type of the new RDD of key-value pairs can be different with respect to the data type of the "input" RDD of key-value pairs

#### Method

- The combineByKey transformation is based on the combineByKey(createCombiner, mergeValue, mergeCombiner) method of the RDD class
  - The values of the input RDD of pairs are of type V
  - The values of the returned RDD of pairs are of type U
  - The type of the keys is K for both RDDs of pairs

- The createCombiner function contains the code that is used to transform a single value (type V) of the input RDD of key-value pairs into a value of the data type (type U) of the output RDD of key-value pairs
  - It is used to transform the first value of each key in each partition to a value of type U

- The mergeValue function contains the code that is used to combine one value of type U with one value of type V
  - It is used in each partition to combine the initial values (type V) of each key with the intermediate ones (type U) of each key

- The mergeCombiner function contains the code that is used to combine two values of type U
  - It is used to combine intermediate values of each key returned by the analysis of different partitions

- combineByKey is more general than reduceByKey and foldByKey because the data types of the values of the input and the returned RDD of pairs can be different
  - For this reason, more functions must be implemented in this case

#### Shuffle

- A shuffle operation is executed for computing the result of the combineByKey() transformation
  - The result/value for each group/key is computed from data stored in different input partitions

- Create an RDD from a local python list containing the pairs
  - ("Paolo", 40)
  - ("Giorgio", 22)
  - ("Paolo", 35)
  - The key is the first name of a user and the value is his/her age
- Store the results in an output HDFS folder. The output contains one line for each name followed by the average age of the users with that name

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

# Create the RDD of pairs from the local collection nameAgePairRDD = sc.parallelize(nameAge)

# Compute the sum of ages and # the number of input pairs for each name (key) sumNumPerNamePairRDD=nameAgePairRDD.combineByKey(\ lambda inputElem: (inputElem, 1), \

lambda intermediateElem, inputElem: \
(intermediateElem[0]+inputElem, intermediateElem[1]+1),

lambda intermediateElem1, intermediateElem2: \ (intermediateElem1[0]+intermediateElem2[0], \ intermediateElem1[1]+intermediateElem2[1])

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

Given an input value (an age), it returns a tuple containing (age, 1)

lambda intermediateElem, inputElem: \
(intermediateElem[0]+inputElem, intermediateElem[1]+1),

lambda intermediateElem1, intermediateElem2: \ (intermediateElem1[0]+intermediateElem2[0], \ intermediateElem1[1]+intermediateElem2[1])

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

Given an input value (an age) and an intermediate value (sum ages, num represented values), it combines them and returns a new updated tuple (sum ages, num represented values)

# the number of input pairs for each name (key) sumNumPerNamePairRDD=nameAgePairRDD.combineByKey(\ lambda inputElem: (inputElem, 1), \

> lambda intermediateElem, inputElem: \ (intermediateElem[o]+inputElem, intermediateElem[1]+1),

lambda intermediateElem1, intermediateElem2: \ (intermediateElem1[0]+intermediateElem2[0], \ intermediateElem1[1]+intermediateElem2[1])

54

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

Given two intermediate result tuples (sum ages, num represented values), it combines them and returns a new updated tuple (sum ages, num represented values)

# the number of input pairs for each name (key) sumNumPerNamePairRDD=nameAgePairRDD.combineByKey(\ lambda inputElem: (inputElem, 1), \

> lambda intermediateElem, inputElem: \ (intermediateElem[o]+inputElem, intermediateElem[1]+1),

lambda intermediateElem1, intermediateElem2: \ (intermediateElem1[0]+intermediateElem2[0], \ intermediateElem1[1]+intermediateElem2[1])

# Compute the average for each name
avgPerNamePairRDD = \
sumNumPerNamePairRDD.map(lambda pair: (pair[0], pair[1][0]/pair[1][1]))

# Store the result in an output folder avgPerNamePairRDD.saveAsTextFile(outputPath)



#### Goal

- Create a new RDD of key-value pairs where there is one pair for each distinct key k of the input RDD of key-value pairs
  - The value associated with key k in the new RDD of keyvalue pairs is the list of values associated with k in the input RDD of key-value pairs
- Method
  - The groupByKey transformation is based on the groupByKey() method of the RDD class

- If you are grouping values per key to perform then an aggregation such as sum or average over the values of each key then groupByKey is not the right choice
  - reduceByKey, aggregateByKey or combineByKey provide better performances for associative and commutative aggregations
- groupByKey is useful if you need to apply an aggregation/compute a function that is not associative

- Shuffle
  - A shuffle operation is executed for computing the result of the groupByKey() transformation
    - Each group/key is associated with/is composed of values which are stored in different partitions of the input RDD

- Create an RDD from a local python list containing the pairs
  - ("Paolo", 40)
  - ("Giorgio", 22)
  - ("Paolo", 35)
  - The key is the first name of a user and the value is his/her age
- Store the results in an output HDFS folder. The output contains one line for each name followed by the ages of all the users with that name

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

# Create the RDD of pairs from the local collection nameAgePairRDD = sc.parallelize(nameAge)

# Create one group for each name with the list of associated ages agesPerNamePairRDD = nameAgePairRDD.groupByKey()

# Store the result in an output folder agesPerNamePairRDD.mapValues(lambda listValues: list(listValues)) .saveAsTextFile(outputPath);

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

# Create the RDD of pairs from the local collection nameAgePairRDD = sc.parallelize(nameAge)

# Create one group for each name with the list of associated ages
agesPerNamePairRDD = nameAgePairRDD.groupByKey()

In this RDD of key-value pairs each tuple is composed of
- a string (key of the pair)
- a "collection" of integers (the value of the pair) – a ResultIterable object

5))

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

This part is used to format the content of the value part of each pair before storing the result in the output folder.

-This transforms a ResultIterable object to a Python list

-Without this "map" the output will contains the pointers to ResultIterable

objects instead of a readable list of integer values

# Store the result in an output folder

agesPerNamePairRDD.mapValues(lambda listValues: list(listValues)) .saveAsTextFile(outputPath);

### **MapValues transformation**

## **MapValues transformation**

#### Goal

- Apply a function f over the value of each pair of an input RDD or key-value pairs and return a new RDD of key-value pairs
- One pair is created in the returned RDD for each input pair
  - The key of the created pair is equal to the key of the input pair
  - The value of the created pair is obtained by applying the function f on the value of the input pair
- The data type of the values of the new RDD of keyvalue pairs can be different from the data type of the values of the "input" RDD of key-value pairs
- The data type of the key is the same

## **MapValues transformation**

#### Method

- The mapValues transformation is based on the mapValues(f) method of the RDD class
- A function f is passed to the mapValues method
  - f contains the code that is applied to transform each input value into the a new value that is stored in the RDD of key-value pairs
- The retuned RDD of pairs contains a number of key-value pairs equal to the number of key-value pairs of the input RDD of pairs
  - The key part is not changed

### MapValues transformation: Example

- Create an RDD from a local python list containing the pairs
  - ("Paolo", 40)
  - ("Giorgio", 22)
  - ("Paolo", 35)
  - The key is the first name of a user and the value is his/her age
- Increase the age of each user (+1 year) and store the result in the HDFS file system
  - One output line per user

#### MapValues transformation: Example

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

# Create the RDD of pairs from the local collection nameAgePairRDD = sc.parallelize(nameAge)

# Increment age of all users
plusOnePairRDD = nameAgePairRDD.mapValues(lambda age: age+1)

# Save the result on disk plusOnePairRDD.saveAsTextFile(outputPath)

### FlatMapValues transformation

## FlatMapValues transformation

- Goal
  - Apply a function **f** over the value of each pair of an input RDD or key-value pairs and return a new RDD of key-value pairs
    - **f** returns a list of values for each input value
  - A list of pairs is inserted in the returned RDD for each input pair
    - The key of the created pairs is equal to the key of the input pair
    - The values of the created pairs are obtained by applying the function f on the value of the input pair
  - The data type of the values of the new RDD of key-value pairs can be different from the data type of the values of the "input" RDD of key-value pairs
  - The data type of the key is the same
# FlatMapValues transformation

#### Method

- The flatMapValues transformation is based on flatMapValues(f) method of the RDD class
- A function f is passed to the mapValues method
  - f contains the code that is applied to transform each input value into a set of new values that are stored in the new RDD of key-value pairs
- The keys of the input pairs are not changed

## FlatMapValues transformation: Example

- Create an RDD from a local python list containing the pairs
  - ("Sentence#1", "Sentence test")
  - ("Sentence#2", "Sentence test number 2")
  - ("Sentence#3", "Sentence test number 3")
- Select the words of each sentence and store in the HDFS file system one pair (senteceld, word) per line

#### FlatMapValues transformation: Example

# Create the RDD of pairs from the local collection sentPairRDD = sc.parallelize(sentences)

# "Extract" words from each sentence
sentIdWord = sentPairRDD.flatMapValues(lambda s: s.split(''))

# Save the result on disk sentIdWord.saveAsTextFile(outputPath)

## **Keys transformation**

# **Keys transformation: Example**

#### Goal

- Return the list of keys of the input RDD of pairs and store them in a new RDD
  - The returned RDD is not an RDD of key-value pairs
  - The returned RDD is a "standard" RDD of "single" elements
  - Duplicates keys are not removed
- Method
  - The keys transformation is based on the keys() method of the RDD class

# **Keys transformation: Example**

- Create an RDD from a local python list containing the pairs
  - ("Paolo", 40)
  - ("Giorgio", 22)
  - ("Paolo", 35)
  - The key is the first name of a user and the value is his/her age
- Store the names of the input users in an output HDFS folder. The output contains one name per line (duplicate names are removed)

# **Keys transformation: Example**

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

# Create the RDD of pairs from the local collection nameAgePairRDD = sc.parallelize(nameAge)

# Select the key part of the input RDD of key-value pairs
namesRDD = nameAgePairRDD.keys().distinct()

# Store the result in an output folder namesRDD.saveAsTextFile(outputPath);

#### **Values transformation**

# **Values transformation**

#### Goal

- Return the list of values of the input RDD of pairs and store them in a new RDD
  - The returned RDD is not an RDD of key-value pairs
  - The returned RDD is a "standard" RDD of "single" elements
  - Duplicates values are not removed
- Method
  - The values transformation is based on the values() method of the RDD class

## Values transformation: Example

- Create an RDD from a local python list containing the pairs
  - ("Paolo", 40)
  - ("Giorgio", 22)
  - ("Paolo", 22)
  - The key is the first name of a user and the value is his/her age
- Store the ages of the input users in an output HDFS folder.
  - The output contains one age per line
  - Duplicate ages/values are not removed

## Values transformation: Example

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 22)]

# Create the RDD of pairs from the local collection nameAgePairRDD = sc.parallelize(nameAge)

# Select the value part of the input RDD of key-value pairs agesRDD = nameAgePairRDD.values()

# Store the result in an output folder agesRDD.saveAsTextFile(outputPath);

#### Goal

- Return a new RDD of key-value pairs obtained by sorting, in ascending order, the pairs of the input RDD by key
  - Note that the final order is related to the default sorting function of the data type of the input keys
- The content of the new RDD of key-value pairs is the same of the input RDD but the pairs are sorted by key in the new returned RDD

#### Method

- The sortByKey transformation is based on the sortByKey() method of the RDD class
  - Pairs are sorted by key in ascending order
- The sortByKey(ascending) method of the RDD class is also available
  - This method allows specifying if the sort order is ascending or descending by means of a Boolean parameter
    - True = ascending
    - False = descending

#### Shuffle

- A shuffle operation is executed for computing the result of the sortByKey() transformation
  - Pairs from different partitions of the input RDD must be compared to sort the input pairs by key

## SortByKey transformation: Example

- Create an RDD from a local python list containing the pairs
  - ("Paolo", 40)
  - ("Giorgio", 22)
  - ("Paolo", 35)
  - The key is the first name of a user and the value is his/her age
- Sort the users by name and store the result in the HDFS file system

#### SortByKey transformation: Example

# Create the local python list nameAge = [ ("Paolo", 40), ("Giorgio", 22), ("Paolo", 35)]

# Create the RDD of pairs from the local collection nameAgePairRDD = sc.parallelize(nameAge)

# Sort by name the content of the input RDD of key-value pairs sortedNameAgePairRDD = nameAgePairRDD.sortByKey()

# Store the result in an output folder
sortedNameAgePairRDD.saveAsTextFile(outputPath);

- All the examples reported in the following tables are applied on an RDD of pairs containing the following tuples (pairs)
  - [("k1", 2), ("k3", 4), ("k3", 6)]
    - The key of each tuple is a string
    - The value of each tuple is an integer

Transformation	Purpose	Example of applied function	Result
reduceByKey(f)	Return an RDD of pairs containing one pair for each key of the "input" RDD of pairs. The value of each pair of the new RDD of pairs is obtained by combining the values of the input RDD associated with the same key. The "input" RDD of pairs and the new RDD of pairs have the same data type.	reduceByKey( lambda v1, v2: v1+v2) Sum values per key	[("k1", 2), ("k3", 10)]

Transformation	Purpose	Example of applied function	Result
foldByKey(zeroValue, op)	Similar to the reduceByKey() transformation. However, foldByKey() is characterized also by a zero value	foldByKey(o, lambda v1, v2: v1+v2) Sum values per key. The zero value is o	[("k1", 2), ("k3", 10)]

Transformation	Purpose	Example of applied function	Result
combineByKey( createCombiner, mergeValue, mergeCombiner )	Return an RDD of key-value pairs containing one pair for each key of the "input" RDD of pairs. The value of each pair of the new RDD is obtained by combining the values of the input RDD associated with the same key. The values of the "input" RDD of pairs and the values of the new (returned) RDD of pairs can be characterized by different data types.	<pre>combineByKey( lambda e: (e, 1), \ lambda c, e: (c[o]+e, c[1]+1), \ lambda c1, c2: (c1[0]+c2[0], c1[1]+c2[1]) ) Sum values by key and count the number of pairs by key in one single step</pre>	[("k1", (2,1)), ("k3", (10,2))]

Transformation	Purpose	Example of applied function	Result
groupByKey()	Return an RDD of pairs containing one pair for each key of the "input" RDD of pairs. The value of each pair of the new RDD of pairs is a "collection" containing all the values of the input RDD associated with one of the input keys.	groupByKey()	[("k1", [2]), ("k3", [4, 6])]

Transformation	Purpose	Example of applied function	Result
mapValues(f)	Apply a function over each pair of an RDD of pairs and return a new RDD of pairs. The applied function returns one pair for each pair of the "input" RDD of pairs. The function is applied only on the value part without changing the key. The values of the "input" RDD and the values of new RDD can have different data types.	mapValues( lambda v: v+1) Increment the value part by 1	[("k1", 3), ("k3", 5), ("k3", 7)]

Transformation	Purpose	Example of applied function	Result
flatMapValues(f)	Apply a function over each pair of an RDD of pairs and return a new RDD of pairs. The applied function returns a set of pairs (from o to many) for each pair of the "input" RDD of pairs. The function is applied only on the value part without changing the key. The values of the "input" RDD and the values of new RDD can have different data types.	flatMapValues(lamb da v: list(range(v,6))) for each input pair (k,v), the set of pairs (k,u) with values of u from v to 5 are returned and included in the new RDD	[("k1", 2), ("k1", 3), ("k1", 4), ("k1", 5), ("k3", 4), ("k3", 5)]

Transformation	Purpose	Example of applied function	Result
keys()	Return an RDD containing the keys of the input pairRDD	keys()	["k1", "k3", "k3"]
values()	Return an RDD containing the values of the input pairRDD	values()	[2, 4, 6]
sortByKey()	Return a PairRDD sorted by key. The "input" PairRDD and the new PairRDD have the same data type.	sortByKey() -	[("k1", 2), ("k3", 4), ("k3", 6)]

# **RDD-based programming**

# Transformations on two RDDs of key-value pairs

# Transformations on two RDDs of pairs

- Spark supports also some transformations that are applied on two RDDs of key-value pairs at the same time
  - subtractByKey, join, coGroup, etc.

# SubtractByKey transformation

# SubtractByKey transformation

#### Goal

- Create a new RDD of key-value pairs containing only the pairs of the first input RDD of pairs associated with a key that is not appearing as key in the pairs of the second input RDD or pairs
  - The data type of the new RDD of pairs is the same of the "first input" RDD of pairs
  - The two input RDD of pairs must have the same type of keys
    - The data type of the values can be different

# SubtractByKey transformation

#### Method

- The subtractByKey transformation is based on the subtractByKey(other) method of the RDD class
- The two input RDDs of pairs analyzed by subtractByKey are the one on which the method is invoked and the one passed as parameter (i.e., other)

#### SubtractByKey transformation

- Shuffle
  - A shuffle operation is executed for computing the result of the subtractByKey() transformation
    - Keys from different partitions of the two input RDDs must be compared

## SubtractByKey transformation: Example

- Create two RDDs of key-value pairs from two local python lists
  - First list Profiles of the users of a blog (username, age)
    - [("PaoloG", 40), ("Giorgio", 22), ("PaoloB", 35)]
  - Second list Banned users (username, motivation)

 [("PaoloB", "spam"), ("Giorgio", "Vandalism")]
 Create a new RDD of pairs containing only the profiles of the non-banned users

### SubtractByKey transformation: Example

# Create the first local python list profiles = [ ("PaoloG", 40), ("Giorgio", 22), ("PaoloB", 35)]

# Create the RDD of pairs from the profiles local list profilesPairRDD = sc.parallelize (profiles)

# Create the second local python list banned = [ ("PaoloB", "spam"), ("Giorgio", "Vandalism")]

# Create the RDD of pairs from the banned local list bannedPairRDD = sc.parallelize (banned)

# Select the profiles of the "good" users
selectedProfiles = profilesPairRDD.subtractByKey(bannedPairRDD)

## (Inner) Join transformation
# (Inner) Join transformation

#### Goal

- Join the key-value pairs of two RDDs of key-value pairs based on the value of the key of the pairs
  - Each pair of the input RDD of pairs is combined with all the pairs of the other RDD of pairs with the same key
  - The new RDD of key-value pairs
    - Has the same key data type of the "input" RDDs of pairs
    - Has a tuple as value (the pair of values of the two joined input pairs)
  - The two input RDDs of key-value pairs
    - Must have the same type of keys
    - But the data types of the values can be different

# (Inner) Join transformation

#### Method

- The (inner) join transformation is based on the join(other) method of the RDD class
- The two input RDDs of pairs analyzed by join are the one on which the method is invoked and the one passed as parameter (i.e., other)

# (Inner) Join transformation

### Shuffle

- A shuffle operation is executed for computing the result of the join() transformation
  - Keys from different partitions of the two input RDDs must be compared and values from different partitions must be retrieved

### (Inner) Join transformation: Example

- Create two RDDs of key-value pairs from two local python lists
  - First list List of questions (QuestionId, Text of the question)
    - [(1, "What is .. ?"), (2, "Who is ..?")]
  - Second list List of answers (QuestionId, Text of the answer)
    - [(1, "It is a car"), (1, "It is a byke"), (2, "She is Jenny")]
- Create a new RDD of pairs to associate each question with its answers
  - One pair for each possible pair question answer

### (Inner) Join transformation: Example

```
# Create the first local Python list
questions= [(1, "What is .. ?"), (2, "Who is ..?")]
```

# Create the RDD of pairs from the local list questionsPairRDD = sc.parallelize (questions)

```
# Create the second local python list
answers = [(1, "It is a car"), (1, "It is a byke"), (2, "She is Jenny")]
```

```
# Create the RDD of pairs from the local list
answersPairRDD = sc.parallelize(answers)
```

# Join questions with answers
joinPairRDD = questionsPairRDD.join(answersPairRDD)

### (Inner) Join transformation: Example

```
# Create the first local Python list
questions= [(1, "What is .. ?"), (2, "Who is ..?")]
```

# Create the RDD of pairs from the local list questionsPairRDD = sc.parallelize (questions)

The key part of the returned RDD of pairs is an integer number

The value part of the returned RDD of pairs is a tuple containing two values: (question, answer)

# Join questions with answers
joinPairRDD = questionsPairRDD.join(answersPairRDD)

- Spark supports outer joins
  - IeftOuterJoin
  - rightOuterJoin
  - fullOuterJoin
- Each of these methods joins two RDDs of pairs and returns an RDD of pairs with the results of the outer join operation
  - The join condition is again based on the key part of the input pairs

#### Goal

- Outer join the key-value pairs of two RDDs of keyvalue pairs based on the value of the key of the pairs
  - The new RDD of key-value pairs
    - Has the same key data type of the "input" RDDs of pairs
    - Has a tuple as value (the pair of values of the two outer joined input pairs)
  - The two input RDDs of key-value pairs
    - Must have the same type of keys
    - But the data types of the values can be different

#### Shuffle

- A shuffle operation is executed for computing the result of the outer join transformations
  - Keys from different partitions of the two input RDDs must be compared and values from different partitions must be retrieved

## **CoGroup transformation**

# **Cogroup transformation**

#### Goal

- Associated each key k of the two input RDDs of keyvalue pairs with
  - The list of values associated with k in the first input RDD of pairs
  - And the list of values associated with k in the second input RDD of pairs
- The new RDD of key-value pairs
  - Has the same key data type of the two "input" RDDs of pairs
  - Has a tuple as value (the two lists of values of the two input RDDs of pairs)
  - The two input RDDs of key-value pairs
    - Must have the same type of keys
    - But the data types of the values can be different

# **Cogroup transformation**

#### Method

- The cogroup transformation is based on the cogroup(other) method of the RDD class
- The two input RDDs of pairs analyzed by cogroup are the one on which the method is invoked and the one passed as parameter (i.e., other)

# **Cogroup transformation**

#### Shuffle

- A shuffle operation is executed for computing the result of the cogroup() transformation
  - Keys from different partitions of the two input RDDs must be compared and values from different partitions must be retrieved

- Create two RDDs of key-value pairs from two local python lists
  - First list List of liked movies (userId, likedMovies)
    - [(1, "Star Trek"), (1, "Forrest Gump"), (2, "Forrest Gump")]
  - Second list List of liked directors (userId, likedDirector)
    - [(1, "Woody Allen"), (2, "Quentin Tarantino"),
       (2, "Alfred Hitchcock")]

- Create a new RDD of pairs containing one pair for each userId (key) associated with
  - The list of liked movies
  - The list of liked directors

#### Inputs

- [(1, "Star Trek"), (1, "Forrest Gump"), (2, "Forrest Gump")]
- [(1, "Woody Allen"), (2, "Quentin Tarantino"),
   (2, "Alfred Hitchcock")]
- Output
  - (1, (["Star Trek", "Forrest Gump"], ["Woody Allen"]))
  - (2, (["Forrest Gump"], ["Quentin Tarantino", "Alfred Hitchcock"]))

```
# Create the first local python list
movies= [(1, "Star Trek"), (1, "Forrest Gump"), (2, "Forrest Gump")]
```

# Create the RDD of pairs from the first local list moviesPairRDD = sc.parallelize(movies)

# Create the RDD of pairs from the second local list directorsPairRDD = sc.parallelize(directors)

# Cogroup movies and directors per user cogroupPairRDD = moviesPairRDD.cogroup(directorsPairRDD)

```
# Create the first local python list
movies= [(1, "Star Trek"), (1, "Forrest Gump"), (2, "Forrest Gump")]
```

# Create the RDD of pairs from the first local list moviesPairRDD = sc.parallelize(movies)

Note that the value part of the returned tuples is a tuple containing two "lists": - The first value contains the "list" of movies (iterable) liked by a user - The second value contains the "list" of directors (iterable) liked by a user

# Cogroup movies and directors per user cogroupPairRDD = moviesPairRDD.cogroup(directorsPairRDD)

- All the examples reported in the following tables are applied on the following two RDDs of key-value pairs
  - inputRDD1: [('k1', 2), ('k3', 4), ('k3', 6)]

inputRDD2: [('k3', 9)]

Transformation	Purpose	Example	Result
subtractByKey(other)	Return a new RDD of key-value pairs. The returned pairs are those of input RDD on which the method is invoked such that the key part does not occur in the keys of the RDD that is passed as parameter. The values are not considered to take the decision.	inputRDD1. subtractByKey (inputRDD2)	[('k1',2)]
join(other)	Return a new RDD of pairs corresponding to join of the two input RDDs. The join is based on the value of the key.	inputRDD1.join (inputRDD2)	[('kȝ', (4,9)), ('kȝ', (6,9))]

Transformation	Purpose	Example	Result
cogroup(other)	For each key k in one of the two input RDDs of pairs, return a pair (k, tuple), where tuple contains: - the list (iterable) of values of the first input RDD associated with key k - the list (iterable) of values of the second input RDD associated with key k	inputRDD1. cogroup (inputRDD2)	[('k1', ([2], []) ), ('k3', ([4, 6], [9]) )]

# Actions on RDDs of key-value pairs

## Actions on RDDs of key-value pairs

- Spark supports also some specific actions on RDDs of key-value pairs
  - countByKey, collectAsMap, lookup

## **CountByKey action**

# CountByKey action

#### Goal

- The countByKey action returns a local python dictionary containing the information about the number of elements associated with each key in the input RDD of key-value pairs
  - i.e., the number of times each key occurs in the input RDD
- Pay attention to the number of distinct keys of the input RDD of pairs
- If the number of distinct keys is large, the result of the action cannot be stored in a local variable of the Driver

# **CountByKey action**

- Method
  - The countBykey action is based on the countBykey() method of the RDD class

#### Data are sent on the network to compute the final result

# CountByKey action: Example 1

- Create an RDD of pairs from the following python list
  - [("Forrest Gump", 4), ("Star Trek", 5), ("Forrest Gump", 3)]
  - Each pair contains a movie and the rating given by someone to that movie
- Compute the number of ratings for each movie

# CountByKey action: Example 1

# Create the local python list movieRating= [("Forrest Gump", 4), ("Star Trek", 5), ("Forrest Gump", 3)]

# Create the RDD of pairs from the local collection movieRatingRDD = sc.parallelize(movieRating)

# Compute the number of rating for each movie movieNumRatings = movieRatingRDD.countByKey()

# Print the result on the standard output
print(movieNumRatings)

# CountByKey action: Example 1

# Create the local python list movieRating= [("Forrest Gump", 4), ("Star Trek", 5), ("Forrest Gump", 3)]

# Create the RDD of pairs from the local collection movieRatingRDD = sc.parallelize(movieRating)

<u># Compute the number of rating for each movie</u> movieNumRatings = movieRatingRDD.countByKey()

Pay attention to the size of the returned local python dictionary (i.e., the number of distinct movies in this case).

## **CollectAsMap action**

# **CollectAsMap action**

#### Goal

- The collectAsMap action returns a local dictionary containing the same pairs of the considered input RDD of pairs
- Pay attention to the size of the returned RDD
- Data are sent on the network
- Method
  - The collectAsMap action is based on the collectAsMap() method of the RDD class

## **CollectAsMap action**

- Pay attention that the collectAsMap action returns a dictionary object
- A dictionary cannot contain duplicate keys
  - Each key can be associated with at most one value
  - If the "input" RDD of pairs contains more than one pair with the same key, only one of those pairs is stored in the returned local python dictionary
- Usually, the last one occurring in the input RDD of pairs
   Use collectAsMap only if you are sure that each key appears only once in the input RDD of key-value pairs

## CollectAsMap vs collect

- The collectAsMap() action returns a local dictionary while collect() return a list of keyvalue pairs (i.e., a list of tuples)
  - The list of pairs returned by collect() can contain more than one pair associated with the same key

# CollectAsMap action: Example 1

- Create an RDD of pairs from the following python list
  - [("User1", "Paolo"), ("User2", "Luca"), ("User3", "Daniele")]
  - Each pair contains a userId and the name of the user
- Retrieve the pairs of the created RDD of pairs and store them in a local python dictionary that is instantiated in the Driver
# CollectAsMap action: Example 1

# Create the local python list
users = [("User1", "Paolo"), ("User2", "Luca"), ("User3", "Daniele")]

#Create the RDD of pairs from the local list usersRDD = sc.parallelize(users)

# Retrieve the content of usersRDD and store it in a # local python dictionary retrievedPairs = usersRDD.collectAsMap()

# Print the result on the standard output
print(retrievedPairs)

# CollectAsMap action: Example 1

```
# Create the local python list
users = [("User1", "Paolo"), ("User2", "Luca"), ("User3", "Daniele")]
```

#Create the RDD of pairs from the local list usersRDD = sc.parallelize(users)

# Retrieve the content of usersRDD and store it in a
# local python dictionary
retrievedPairs = usersRDD.collectAsMap()

Pay attention to the size of the returned local python dictionary (i.e., the number of distinct users in this case).

#### Lookup action

# Lookup action

#### Goal

- The lookup(k) action returns a local python list containing the values of the pairs of the input RDD associated with the key k specified as parameter
- Method
  - The lookup action is based on the lookup(key) method of the RDD class

# Lookup action: Example 1

- Create an RDD of pairs from the following python list
  - [("Forrest Gump", 4), ("Star Trek", 5), ("Forrest Gump", 3)]
  - Each pair contains a movie and the rating given by someone to that movie
- Retrieve the ratings associated with the movie "Forrest Gump" and store them in a local python list in the Driver

### Lookup action: Example 1

# Create the local python list movieRating= [("Forrest Gump", 4), ("Star Trek", 5), ("Forrest Gump", 3)]

# Create the RDD of pairs from the local collection movieRatingRDD = sc.parallelize(movieRating)

# Select the ratings associated with "Forrest Gump"
movieRatings = movieRatingRDD.lookup("Forrest Gump")

# Print the result on the standard output
print(movieRatings)

#### Lookup action: Example 1

# Create the local python list movieRating= [("Forrest Gump", 4), ("Star Trek", 5), ("Forrest Gump", 3)]

# Create the RDD of pairs from the local collection movieRatingRDD = sc.parallelize(movieRating)

# Select the ratings associated with "Forrest Gump"
movieRatings = movieRatingRDD.lookup("Forrest Gump")

Pay attention to the size of the returned list (i.e., the number of ratings associated with "Forrest Gump" in this case).

# Actions on RDDs of key-value pairs: Summary

#### Actions on RDDs of key-value pairs: Summary

- All the examples reported in the following tables are applied on the following RDD of key-value pairs
  - inputRDD: [('k1', 2), ('k3', 4), ('k3', 6)]

#### Actions on RDDs of key-value pairs: Summary

Transformation	Purpose	Example	Result
countByKey()	Return a local python dictionary containing the number of elements in the input RDD for each key of the input RDD of pairs	inputRDD. countByKey()	{('k1',1), ('K3',2)}
collectAsMap()	Return a local python dictionary containing the pairs of the input RDD of pairs	inputRDD. collectAsMap()	{('k1', 2), ('k3', 6)} Or {('k1', 2), ('k3', 4)} Depending on the order of the pairs in the input RDD of pairs
lookup(key)	Return a local python list containing all the values associated with the key specified as parameter	inputRDD. lookup('k3')	[4, 6]