



University of  
**Salford**  
MANCHESTER



SCHOOL OF  
**SCIENCE, ENGINEERING  
& ENVIRONMENT**

**Program: MSc of Data Science**  
**Module: Big Data Tools and Techniques**

Week 3 - Part 2

Resilient Distributed Datasets (**RDDs**)

2025

# Learning Outcomes

1. How to Create an RDD
2. RDD Operations
3. RDD Types
4. RDD & Lazy Execution
5. Functional programming in Spark
6. Pair RDDs

# Refresher: What is RDD?

- RDD (Resilient Distributed Dataset)
  - Resilient – if data in memory is lost, it can be recreated
  - Distributed – processed across the cluster
  - Dataset – initial data can come from a file or be created programmatically
- RDDs are the fundamental unit of data in Spark
- Most Spark programming consists of performing operations on RDDs

# Resilient!!



“Life doesn’t get easier or more forgiving, we get stronger and more resilient.” — Steve Maraboli

[@resilientstories](#)

The New York Times

RESILIENCE

## What Makes Some People More Resilient Than Others

The very earliest days of our lives, and our closest relationships, can offer clues about how we cope with adversity.

Share full article

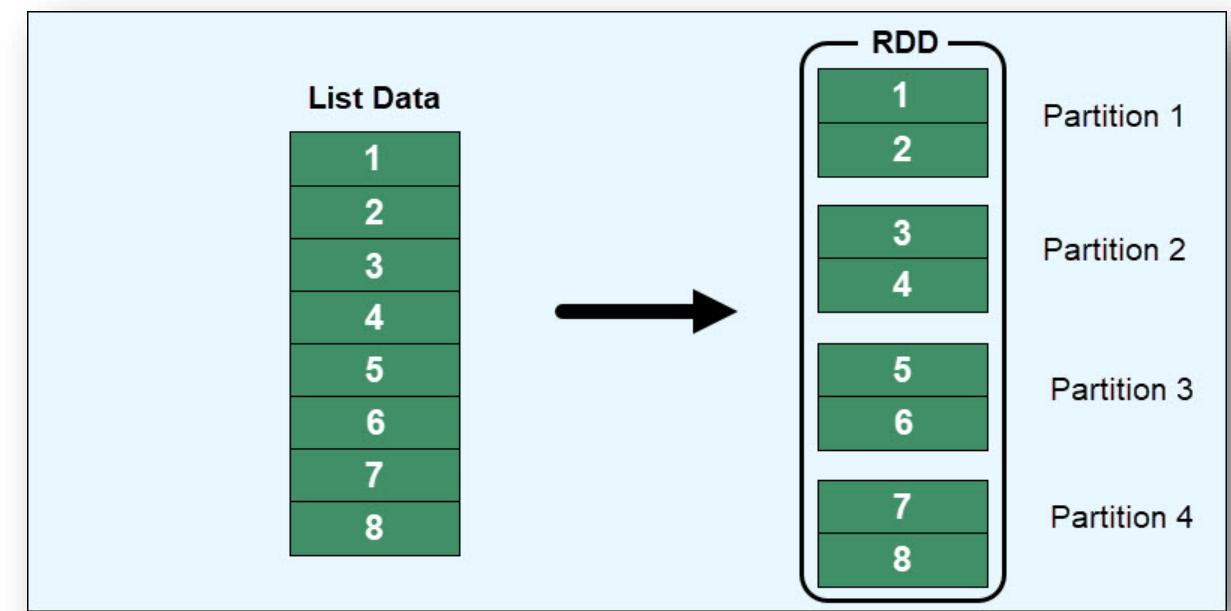


Monika Aichele

# Refresher: RDDs on a cluster

## Resilient Distributed Datasets (RDD)

- Data is partitioned in worker nodes

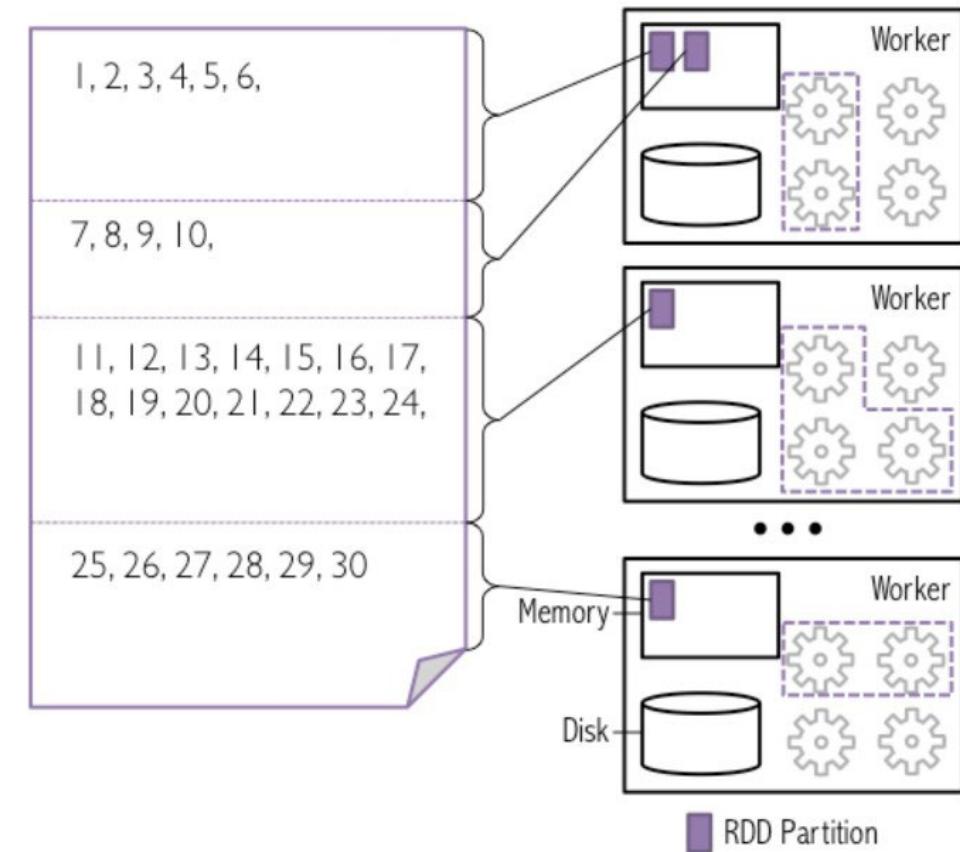


# Refresher: RDDs on a cluster

- Resilient Distributed Datasets
  - Data is partitioned across worker nodes (e.g. Data is from 1 to 30)
- Partitioning is done automatically by Spark
  - Optionally, you can control how many partitions are created

Dataset is broken into partitions

Partitions are each stored in a worker's memory



# How to Create an RDD

# Three ways to create an RDD

➤ From a file or set of files

➤ From data in memory

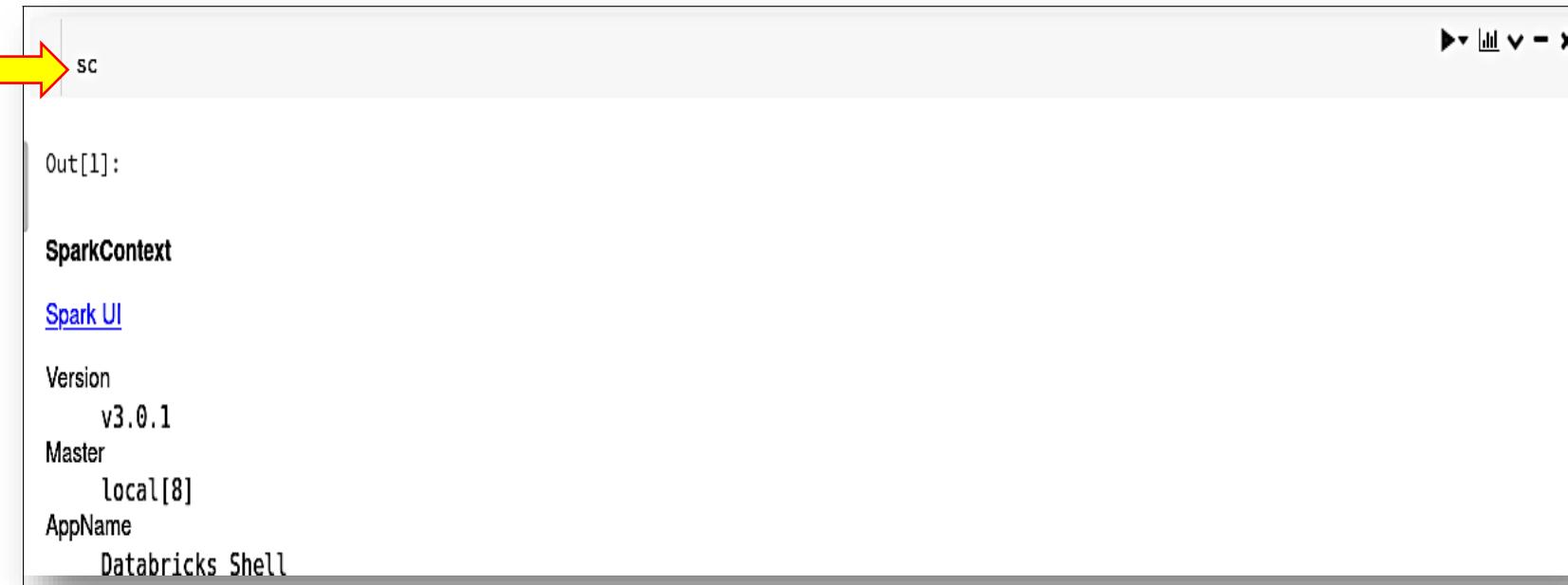
(In general, when we say "data in memory," it means that the data is stored and accessed directly from a computer's **RAM (Random Access Memory)** rather than from a storage device like a hard drive or SSD)



➤ From another RDD

# Spark Context (SC)

- Every Spark application requires a **Spark Context**
- The main entry point to the Spark API



```
sc
```

Out[1]:

```
SparkContext
Spark UI
Version
v3.0.1
Master
local[8]
AppName
Databricks Shell
```

# Example: A file based RDD

```
mydata = sc.textFile("purplecow.txt")
mydata.count()
```

File: purplecow.txt

I've never seen a purple cow.  
I hope I never see one;  
But I can tell you anyhow,  
I'd rather see than be one.



RDD: mydata

I've never seen a purple cow.
I hope I never see one;
But I can tell you anyhow,
I'd rather see than be one.

**SC:** Spark Context

```
mydata = sc.textFile("purplecow.txt")
mydata.count()
```

**textFile** is a method of  
SparkContext class that reads  
a text file and return it as an  
**RDD of Strings**.

File: purplecow.txt

```
I've never seen a purple cow.  
I hope I never see one;  
But I can tell you anyhow,  
I'd rather see than be one.
```



RDD: mydata

I've never seen a purple cow.
I hope I never see one;
But I can tell you anyhow,
I'd rather see than be one.

# Example: RDD from data in memory

```
Cmd 1
1 data = [1,2,3,4,5]
2 distData = sc.parallelize(data)
```

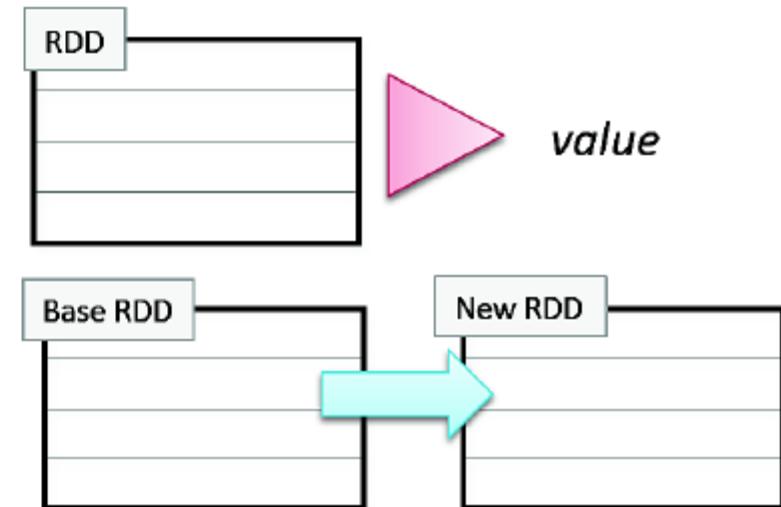


**Parallelize** is a method of SparkContext (SC) to create an RDD from an existing collection (e.g an array) which is already in the driver node.

# RDD Operations

# RDD operations

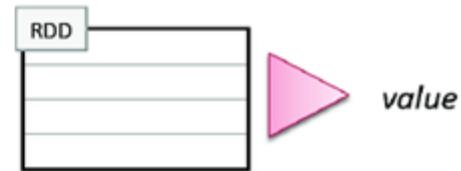
- ▶ Two types of RDD operations
  - **Actions** – return values
  - **Transformations** – define a new RDD based on the current one(s)



# RDD operations: Actions

## ► Some common actions

- ▶ `count()` – return the number of elements
- ▶ `take(n)` – return an array of the first *n* elements
- ▶ `collect()` – return an array of all elements
- ▶ `saveAsTextFile(file)` – save to text file(s)



RDD: mydata  
I've never seen a purple cow.  
I hope I never see one;  
But I can tell you anyhow,  
I'd rather see than be one.

```
1 mydata = sc.textFile("dbfs:/FileStore/tables/purplecow.txt")
2 print(mydata.count())
3 for line in mydata.take(2):
4     print(line)

▶ (2) Spark Jobs
4
I've never seen a purple cow.
I hope I never see one.
```

# Another example

```
1 data = [1,2,3,4,5]
2 distData = sc.parallelize(data)
```

Command took 0.05 seconds

Cmd 2

```
1 distData.count()
```

▶ (1) Spark Jobs

Out[11]: 5

Command took 0.51 seconds

Cmd 3

```
1 distData.collect()
```

▶ (1) Spark Jobs

Out[10]: [1, 2, 3, 4, 5]

Command took 0.24 seconds

Cmd 4

## We can create RDDs from collections instead of files

Cmd 6

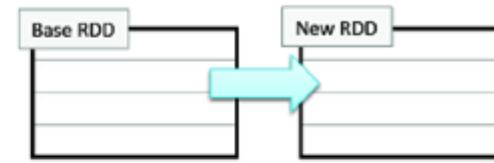
```
1 myData = ["BDTT", "ADB", "ASDV", "MLDM"]  
2  
3 myRdd = sc.parallelize(myData)  
4  
5 myRdd.take(2)
```

▶ (2) Spark Jobs

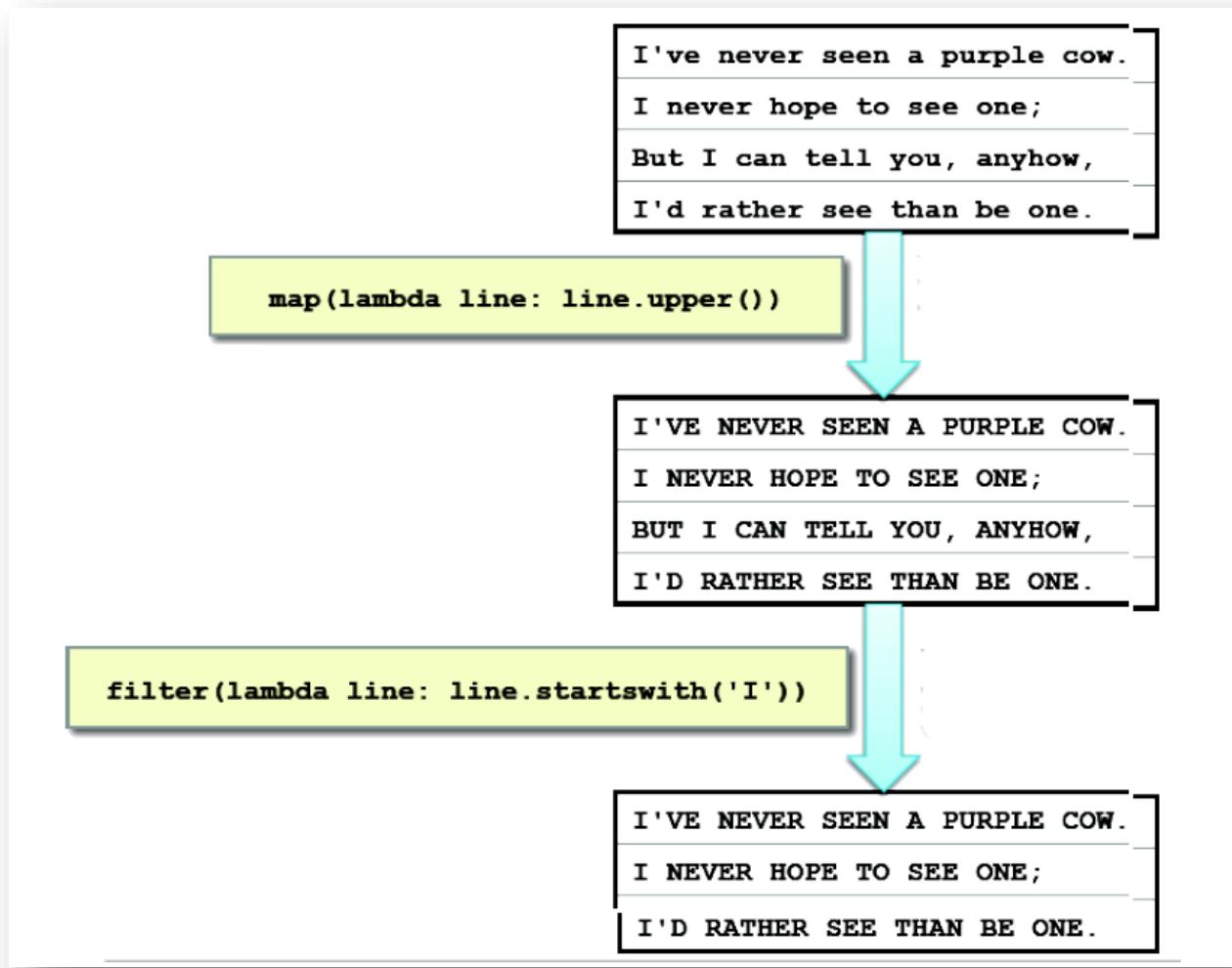
Out[19]: ['BDTT', 'ADB']

# RDD operations: Transformation

- ▶ Transformations create a new RDD from an existing one
- ▶ RDDs are immutable
  - ▶ Data in an RDD is never changed
  - ▶ Transform in sequence to modify the data as needed
- ▶ Some common transformations
  - ▶ **map(function)** – creates a new RDD by performing a function on each record in the base RDD
  - ▶ **filter(function)** – creates a new RDD by including or excluding each record in the base RDD according to a boolean function



# Example: map and filter actions



# RDD Types

# RDD Types

- ▶ RDDs can hold any type of element
  - ▶ Primitive types: integers, characters, booleans, etc.
  - ▶ Sequence types: strings, lists, arrays, tuples, dicts, etc. (including nested data types)
- ▶ Pair RDDs
  - ▶ RDDs consisting of Key-Value pairs

For a file based RDD use `sc.TextFile`

```
sc.textFile("myfile.txt")
sc.textFile("mydata/*.log")
sc.textFile("myfile1.txt,myfile2.txt")
```

Files are referenced by absolute or relative URL

- ▶ Absolute URI:
  - ▶ `file:/home/training/myfile.txt`
  - ▶ `hdfs://localhost/loudacre/myfile.txt`
  - ▶ `s3://bucket/folder/filename.csv` (AWS S3)
  - ▶ `wasb://bucket/folder/filename.csv` (Azure WASBs)
  - ▶ `gs://bucket/folder/filename.csv` (Google Cloud Storage)
  - ▶ `dbfs://folder/filename.csv` (Databricks DBFS)
- ▶ Relative URI (uses default file system): `myfile.txt`

# RDD & Lazy Execution

# RDD has a **Lazy** execution

1

Data in RDDs is not processed until an *action* is performed

File: purplecow.txt

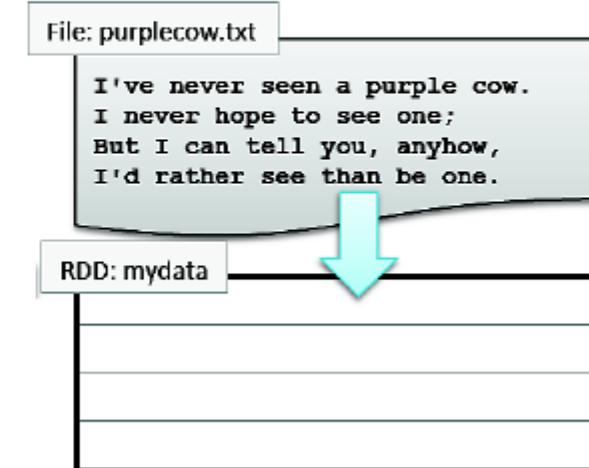
I've never seen a purple cow.  
I never hope to see one;  
But I can tell you, anyhow,  
I'd rather see than be one.

# Lazy execution

2

Data in RDDs is not processed until an *action* is performed

```
mydata = sc.textFile("purplecow.txt")
```

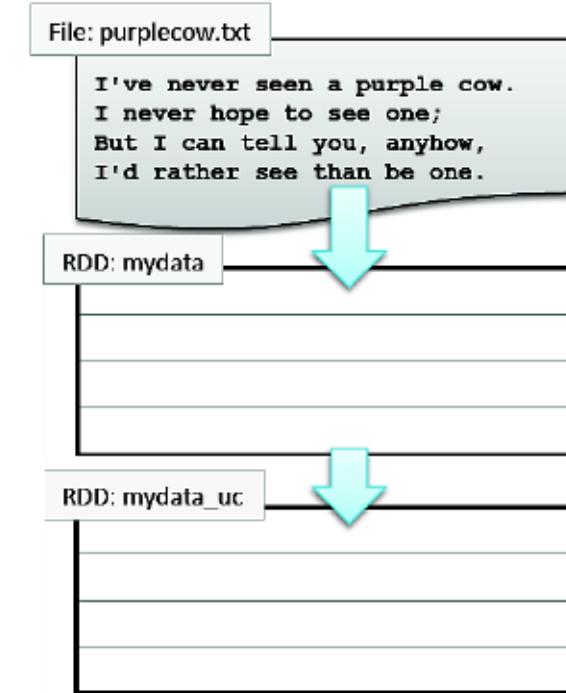


# Lazy execution

3

Data in RDDs is not processed until an *action* is performed

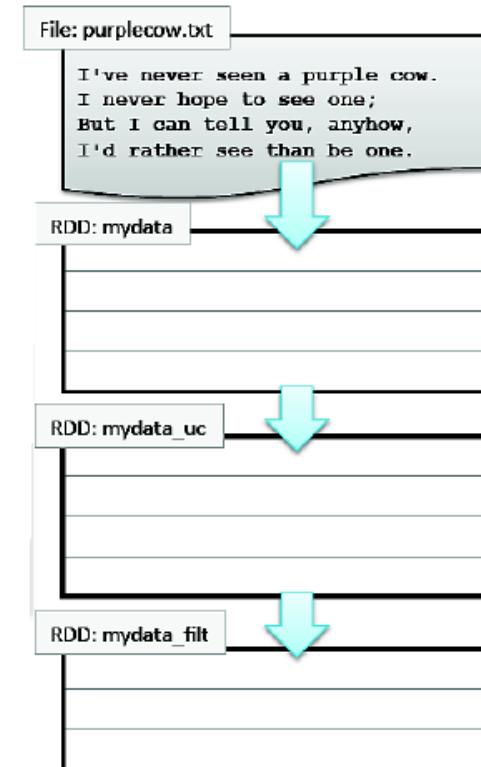
```
mydata = sc.textFile("purplecow.txt")
mydata_uc = mydata.map(lambda line:
    line.upper())
```



# Lazy execution

Data in RDDs is not processed until an *action* is performed

```
mydata = sc.textFile("purplecow.txt")
mydata_uc = mydata.map(lambda line:line.upper())
mydata_filt = mydata_uc.filter(lambda
    line: line.startswith("I"))
```

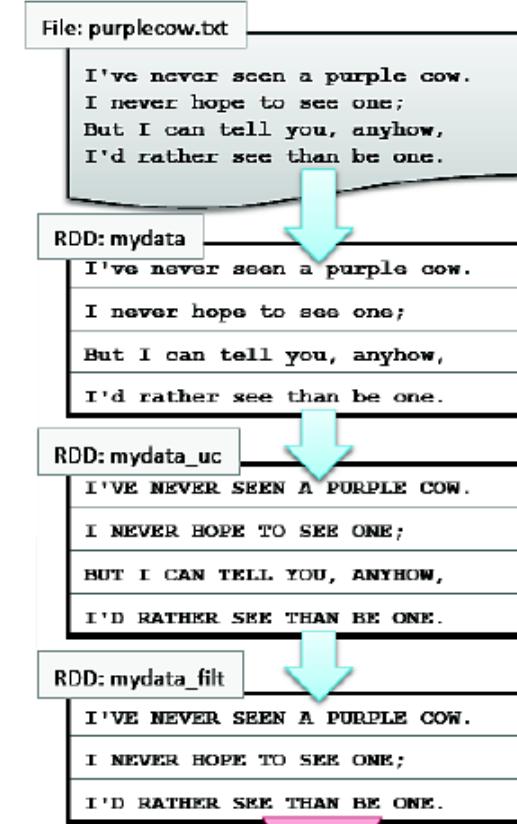


# Lazy execution

Data in RDDs is not processed until an *action* is performed

```
mydata = sc.textFile("purplecow.txt")
mydata_uc = mydata.map(lambda line:line.upper())
mydata_filt = mydata_uc.filter(lambda
    line: line.startswith("I"))
mydata_filt.count()
3
```

↑  
count() is a action



# Chaining transformation

Transformations may be chained together

```
mydata = sc.textFile("purplecow.txt")
mydata_uc = mydata.map(lambda s: s.upper())
mydata_filt = mydata_uc.filter(lambda s: s.startswith('I'))
mydata_filt.count()
3
```

is exactly equivalent to

```
sc.textFile("purplecow.txt").map(lambda line: line.upper()) \
    .filter(lambda line: line.startswith('I')).count()
3
```



Use • to chain functions (dot is a pipeline)

# Functional programming in Spark

# Functional programming in Spark

- ▶ Spark depends heavily on the concepts of functional programming
  - ▶ Functions are the fundamental unit of programming
  - ▶ Functions have input and output only
- ▶ Key concepts
  - ▶ Passing functions as input to other functions
  - ▶ Anonymous functions

purplecow.txt - Notepad

File Edit Format View Help

I never saw a purple cow.  
I never hope to see one.  
But I can tell you, anyhow,  
I'd rather see than be one!



purplecow.txt file

# Passing named functions

Cmd 4

```
1 def toUpper(s):
2     return s.upper()
3
4 mydata = sc.textFile("FileStore/tables/purplecow.txt")
5
6 mydata.map(toUpper).take(2)
```

▶ (1) Spark Jobs

Out[14]: ['I NEVER SAW A PURPLE COW.', 'I NEVER HOPE TO SEE ONE.]

# Passing anonymous functions

Cmd 5

```
1 mydata.map(lambda line : line.upper()).take(2)
```

▶ (1) Spark Jobs

```
Out[15]: ['I NEVER SAW A PURPLE COW.', 'I NEVER HOPE TO SEE ONE.]
```

Lambda functions are best for short, one-of functions

**sc.TextFile**  
**vs**  
**sc.wholeTextFiles**

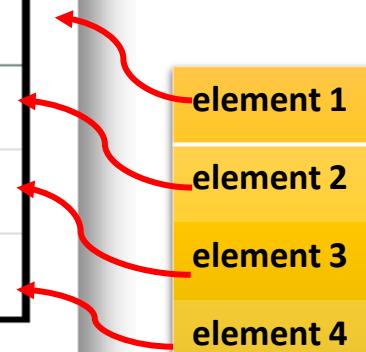
# Can Sc.TextFile works with all files?

textFile maps each line in a file to a separate RDD element

```
I've never seen a purple cow.\nI never hope to see one;\nBut I can tell you, anyhow,\nI'd rather see than be one.\n
```



```
I've never seen a purple cow.\nI never hope to see one;\nBut I can tell you, anyhow,\nI'd rather see than be one.
```



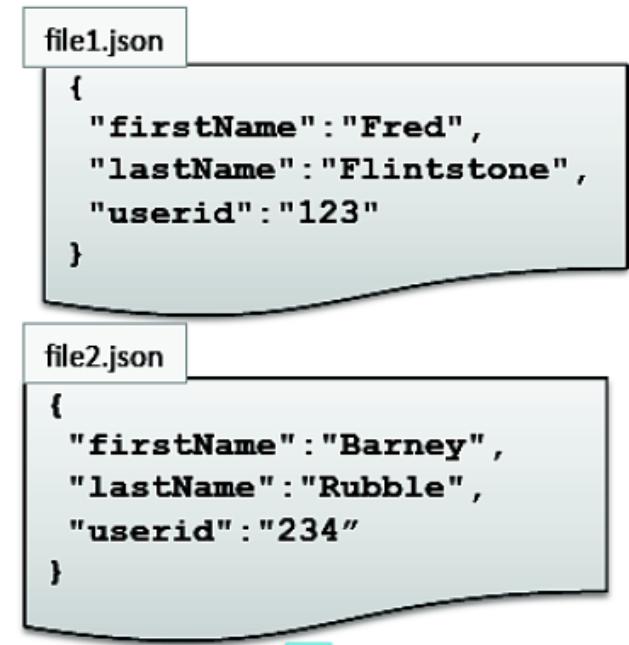
textFile only works with line-delimited files

# What about multiple files in a directory?

**sc.wholeTextFiles(directory)**

maps entire contents of each file in a directory  
to a single RDD element.

Works only for small files (elements must fit in memory)



(file1.json, {"firstName": "Fred", "lastName": "Flintstone", "userid": "123"} )
(file2.json, {"firstName": "Barney", "lastName": "Rubble", "userid": "234"} )
(file3.xml, ... )
(file4.xml, ... )

QC1.json - Notepad

File Edit Format View Help

```
{"quiz": {"sport": {"q1": {"question": "Which one is maths": {"q1": {"question": "18 + 21 = ?"}, "options":}}
```

QC2.json - Notepad

File Edit Format View Help

```
{"quiz": {"sport": {"q1": {"question": "Which one is maths": {"q1": {"question": "3 + 6 = ?"}, "options":}}
```

QC3.json - Notepad

File Edit Format View Help

```
{"quiz": {"sport": {"q1": {"question": "Which one is maths": {"q1": {"question": "5 + 7 = ?"}, "options":}}
```

We have 3 json files and we want to make one RDD from all 3 files

```
1 import json  
2  
3 myrdd1= sc.wholeTextFiles("/FileStore/tables/QC*.json")  
4  
5 myrdd1.take(3)
```

QC\* means QC1  
and QC2 and QC3

► (2) Spark Jobs

```
Out[45]: [('dbfs:/FileStore/tables/QC1.json',  
  {'quiz': {"sport": {"q1": {"question": "Which one is correct team name in NBA?", "options": ["Miami"], "\r\n" "maths": {"q1": {"question": "18 + 21 = ?", "options": ["32", "39", "41", "28"], "answer": "39"}, "\r\n" "maths": {"q1": {"question": "2 + 2 = ?", "options": ["2", "4", "6", "8"], "answer": "4"}}}}},  
  ('dbfs:/FileStore/tables/QC2.json',  
  {'quiz': {"sport": {"q1": {"question": "Which one is correct team name in NBA?", "options": ["Boston"], "\r\n" "maths": {"q1": {"question": "3 + 6 = ?", "options": ["9", "10", "11", "17"], "answer": "9"}, "\r\n" "maths": {"q1": {"question": "3 + 6 = ?", "options": ["9", "10", "11", "17"], "answer": "9"}}}}},  
  ('dbfs:/FileStore/tables/QC3.json',  
  {'quiz': {"sport": {"q1": {"question": "Which one is correct team name in NBA?", "options": ["New"], "\r\n" "maths": {"q1": {"question": "5 + 7 = ?", "options": ["10", "11", "12", "13"], "answer": "12"}, "\r\n" "maths": {"q1": {"question": "5 + 7 = ?", "options": ["10", "11", "12", "13"], "answer": "12"}}}}}])]
```

QC1.json is  
1<sup>st</sup> element of myrdd1

QC2.json is  
2<sup>nd</sup> element of myrdd1

QC3.json is 3<sup>rd</sup>  
element of myrdd1

# List of RDD Operations

(**Transformations and Actions**)

# Single RDD Transformations

- ▶ `.map(f)` – returns a new RDD formed by passing each element through a function `f`
- ▶ `.filter(f)` – returns a new RDD based on selecting elements for which the function `f` returns true
- ▶ `.flatMap(f)` – similar to `map`, but the new RDD flattens all elts
- ▶ `.sample(withReplacement, fraction, seed)` – samples a fraction of the data, with or without replacement
- ▶ `.distinct()` – returns a new RDD containing the distinct elements of the source RDD
- ▶ `.zipWithIndex()` – appends (or ZIPS) the RDD with the elt indices
- ▶ `.reduceByKey(f)` – reduces the elts of the RDD using `f` by key
- ▶ `.sortByKey(asc)` – orders (key, value) RDD by key and returns an RDD in ascending or descending order
- ▶ `.sortBy` – use provided function to sort

# Map vs flatMap

## Input RDD:

```
python  
rdd = sc.parallelize(["a,b,c", "d,e", "f"])
```

Copy Edit

## Key Differences Between map and flatMap

Feature	map	flatMap
Output elements	One output element for each input element	Zero, one, or more output elements per input
Resulting RDD	Same number of elements as input	Flattens multiple outputs into a single RDD
Use case	Use when you want a 1:1 mapping	Use when you want a 1:N mapping or flattening

## Using map:

```
python  
  
mapped_rdd = rdd.map(lambda x: x.split(","))
print(mapped_rdd.collect()) # Output: [['a', 'b', 'c'], ['d', 'e'], ['f']]
```

Copy Edit

**Explanation:** Each input string is split into a list, resulting in a nested list structure.

## Using flatMap:

```
python  
  
flat_mapped_rdd = rdd.flatMap(lambda x: x.split(","))
print(flat_mapped_rdd.collect()) # Output: ['a', 'b', 'c', 'd', 'e', 'f']
```

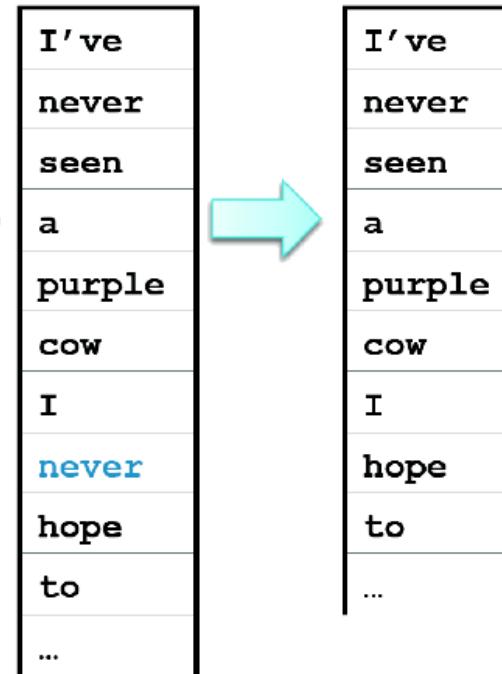
Copy Edit

**Explanation:** The lists are "flattened" into a single RDD of elements.

# Example: flatMap and distinct

```
sc.textFile(file) \
  .flatMap(lambda line: line.split()) \
  .distinct()
```

I've never seen a purple cow.  
I never hope to see one;  
But I can tell you, anyhow,  
I'd rather see than be one.



# Multi-RDD Transformations

- ▶ intersection – create a new RDD with all elements in both original RDDs
- ▶ union – add all elements of two RDDs into a single new RDD
- ▶ zip – pair each element of the first RDD with the corresponding element of the second
- ▶ join – joins two RDDs. Outer joins are supported through left outer join, right outer join, and full outer join.

# Example: subtract, zip and union

rdd1
Chicago
Boston
Paris
San Francisco
Tokyo

rdd2
San Francisco
Boston
Amsterdam
Mumbai
McMurdo Station

`rdd1.subtract(rdd2)`

Tokyo
Paris
Chicago

`rdd1.zip(rdd2)`

(Chicago, San Francisco)
(Boston, Boston)
(Paris, Amsterdam)
(San Francisco, Mumbai)
(Tokyo, McMurdo Station)

`rdd1.union(rdd2)`



Chicago
Boston
Paris
San Francisco
Tokyo
San Francisco
Boston
Amsterdam
Mumbai
McMurdo Station

# RDD Actions

- ▶ `.take()` – returns an array with the first  $n$  elements of the RDD.
- ▶ `.collect()` – returns all of the elements from the workers to the driver.
- ▶ `.reduce(f)` – aggregates the elements of an RDD by  $f$ . The  $f$  function should be commutative and associative\* so that it can be computed correctly in parallel.
- ▶ `.count()` - returns the number of elements in the RDD.

\***commutative** and **associative**, meaning that the order of operations does not affect the final result.

# RDD Actions

- ▶ `.take()` – returns an array with the first  $n$  elements of the RDD.
- ▶ `.collect()` – returns all of the elements from the workers to the driver.
- ▶ `.reduce(f)` – aggregates the elements of an RDD by  $f$ . The  $f$  function should be commutative and associative\* so that it can be computed correctly in parallel.
- ▶ `.count()` - returns the number of elements in the RDD.

\***commutative** and **associative**, meaning that the order of operations does not affect the final result.

# Some other RDD operations: Reduce(f)

## Sum of All Elements

python

```
rdd = sc.parallelize([1, 2, 3, 4, 5])
total_sum = rdd.reduce(lambda x, y: x + y)
print(total_sum) # Output: 15
```

## Finding the Maximum Value

python

```
rdd = sc.parallelize([10, 20, 5, 8, 30])
max_value = rdd.reduce(lambda x, y: x if x > y else y)
print(max_value) # Output: 30
```

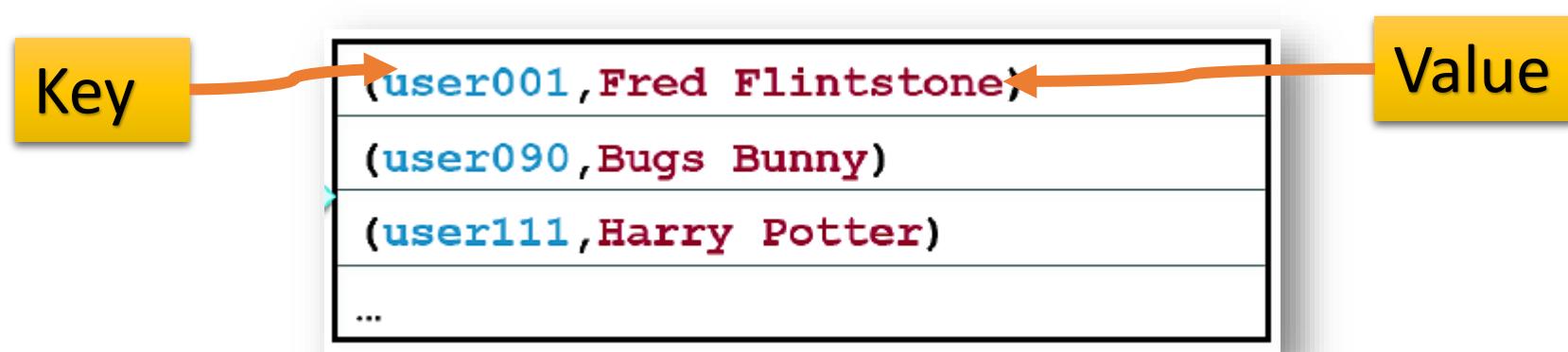
# Pair RDDs

## Pair RDDs are a special form of RDD

- ▶ Each element must be a key-value pair (a two-element tuple)
- ▶ Keys and values can be any type

Pair RDD

(key1 , value1)
(key2 , value2)
(key3 , value3)
...



# Creating pair RDDS

The first step in most workflows is to get the data into key/value form

- ▶ What should the RDD should be keyed on?
- ▶ What is the value?

Commonly used functions to create Pair RDDs

- ▶ map()
- ▶ flatMap() / flatMapValues()
- ▶ keyBy()

# Example: creating a simple pair RDD

Example: create a pair RDD from a tab-separated file

```
users = sc.textFile(file) \
    .map(lambda line: line.split('\t')) \
    .map(lambda fields: (fields[0], fields[1]))
```

user001\tFred Flintstone  
user090\tBugs Bunny  
user111\tHarry Potter  
...



(user001,Fred Flintstone)  
(user090,Bugs Bunny)  
(user111,Harry Potter)  
...

# Example 1: Keying web logs by user ID

```
sc.textFile(logfile) \
    .keyBy(lambda line: line.split(' ')[2])
```

User ID

56.38.234.188 -	<b>99788</b>	"GET /KBDOC-00157.html HTTP/1.0" ...
56.38.234.188 -	<b>99788</b>	"GET /theme.css HTTP/1.0" ...
203.146.17.59 -	<b>25254</b>	"GET /KBDOC-00230.html HTTP/1.0" ...
...		

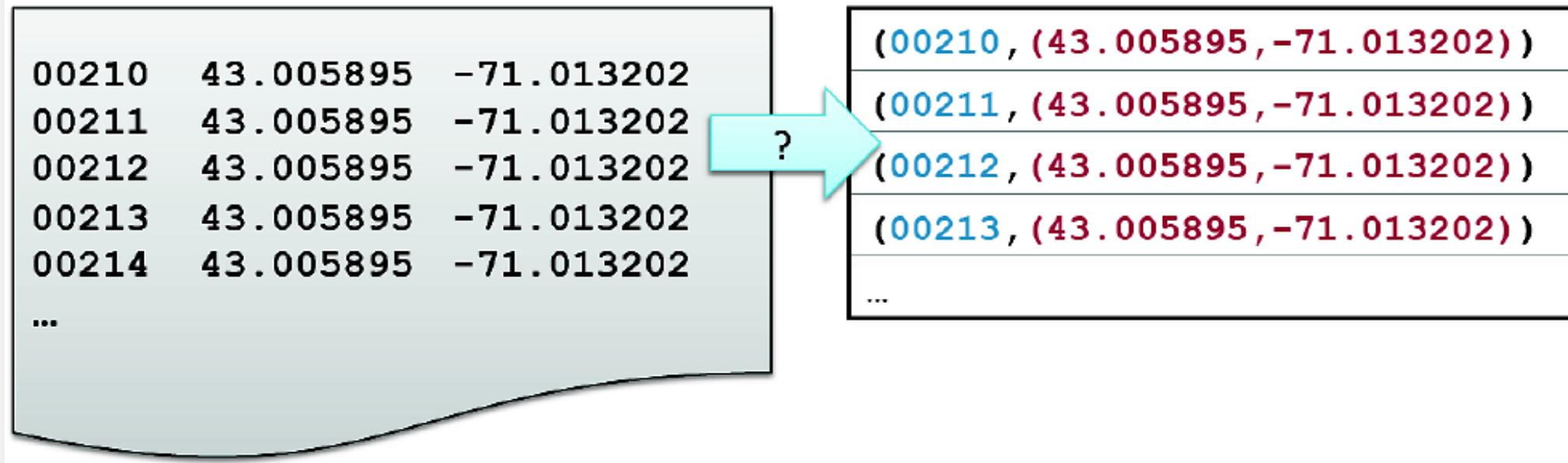
The diagram illustrates the transformation of raw log data into key-value pairs. At the top, a snippet of Scala code uses the `.keyBy` operation to extract the third field (User ID) from each log line. Below the code, the raw log entries are shown, with the User ID highlighted in a red oval. A red arrow points from this highlighted entry to the corresponding row in the resulting key-value pairs at the bottom. A large blue arrow points downwards from the raw logs to the key-value pairs, indicating the flow of data.

( <b>99788</b> , 56.38.234.188 - 99788 "GET /KBDOC-00157.html...")
( <b>99788</b> , 56.38.234.188 - 99788 "GET /theme.css...")
( <b>25254</b> , 203.146.17.59 - 25254 "GET /KBDOC-00230.html...")
...

# Question 1: Pairing complex values

How would you do this?

- ▶ Input: a list of postal codes with latitude and longitude
- ▶ Output: **postal code** (key) and **lat/long pair** (value)



# Answer 1: Pairing complex values

```
sc.textFile(file) \
    .map(lambda line: line.split()) \
    .map(lambda fields: (fields[0],(fields[1],fields[2])))
```

First, split the contents of each line  
Second, put them together in a new format

00210 43.005895 -71.013202  
01014 42.170731 -72.604842  
01062 42.324232 -72.67915  
01263 42.3929 -73.228483  
...

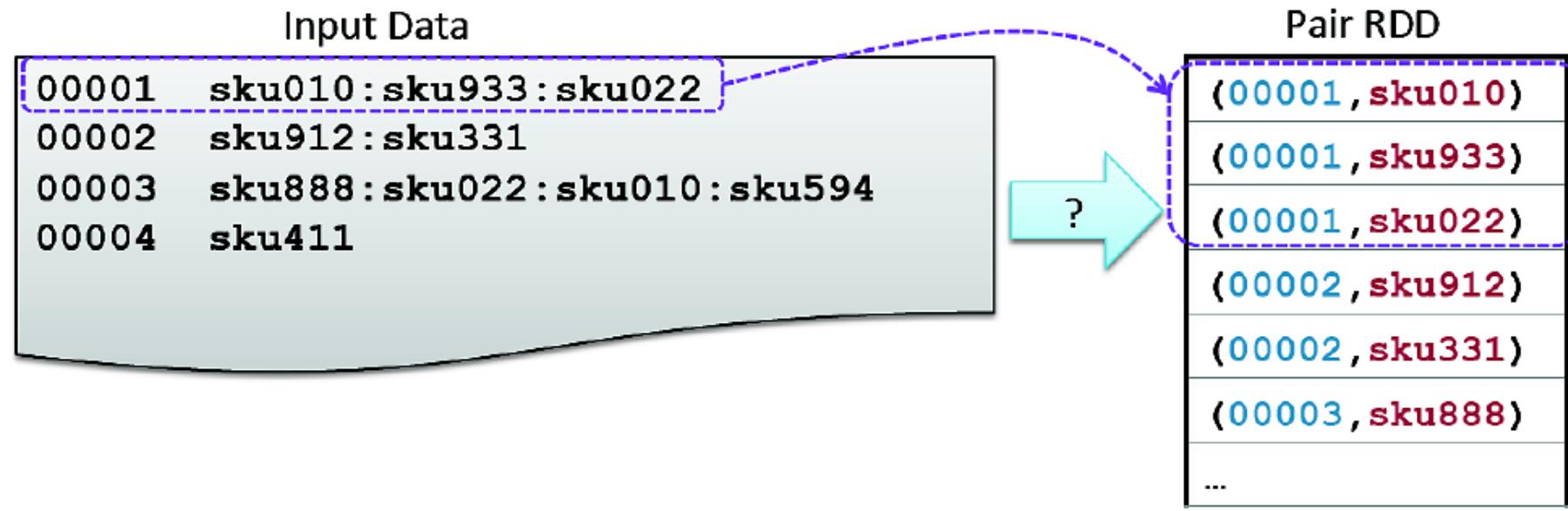
(00210, (43.005895, -71.013202))  
(01014, (42.170731, -72.604842))  
(01062, (42.324232, -72.67915))  
(01263, (42.3929, -73.228483))  
...



# Question 2: Mapping single rows to multiple pairs

How would you do this?

- ▶ Input: order numbers with a list of SKUs in the order
- ▶ Output: **order** (key) and **sku** (value)



## Answer 2: Mapping single rows to multiple pairs

```
sc.textFile(file)
```

00001 **sku010:sku933:sku022**

00002 **sku912:sku331**

00003 **sku888:sku022:sku010:sku594**

00004 **sku411**

# Answer 2: Mapping single rows to multiple pairs

2

```
sc.textFile(file)\n    .map(lambda line: line.split(' '))
```

00001 sku010:sku933:sku022

00002 sku912:sku331

00 [00001,sku010:sku933:sku022]

00 [00002,sku912:sku331]

[00003,sku888:sku022:sku010:sku594]

[00004,sku411]



Note that `split` returns  
2-element arrays, not  
pairs/tuples

## Answer 2: Mapping single rows to multiple pairs

```
sc.textFile(file)\\n    .map(lambda line: line.split(' '))\\n    .map(lambda fields: (fields[0], fields[1]))
```



Map array elements to tuples to produce a Pair RDD

## Answer 2: Mapping single rows to multiple pairs

```
sc.textFile(file) \
    .map(lambda line: line.split(' ')) \
    .map(lambda fields: (fields[0], fields[1])) \
    .flatMapValues(lambda skus: skus.split(':'))
```

00001	sku010:sku933:sku022
00002	sku912:sku331
00003	[00001,sku010:sku933:sku022]
00004	[00002,sku912:sku331]

[0]	(00001,sku010:sku933:sku022)
[0]	(00002,sku912:sku331)
	(00003,sku888:sku022:sku010:sku594)
	(00004,sku411)

(00001,sku010)
(00001,sku933)
(00001,sku022)
(00002,sku912)
(00002,sku331)
(00003,sku888)
...



# **Map-reduce programming in Spark**

- ▶ Map-reduce is a common programming model
- ▶ Easily applicable to distributed processing of large data sets

### Map-reduce in Hadoop



Hadoop MapReduce is the major implementation  
Somewhat limited

- ▶ Each job has one Map phase, one Reduce phase
- ▶ Job output is saved to files

### Map-reduce in Spark



Spark implements map-reduce with much greater flexibility

- ▶ Map and reduce functions can be interspersed
- ▶ Results can be stored in memory
- ▶ Operations can easily be chained

# Map-reduce in Spark

Map-reduce in Spark works on Pair RDDs

Map-reduce example: Word count

## Map phase

- ▶ Operates on one record at a time
- ▶ “Maps” each record to one or more new records
- ▶ e.g. map, flatMap, filter, keyBy

## Reduce phase

- ▶ Works on map output
- ▶ Consolidates multiple records
- ▶ e.g. reduceByKey, sortByKey

Input Data

```
the cat sat on the mat
the aardvark sat on the sofa
```



Result	
aardvark	1
cat	1
mat	1
on	2
sat	2
sofa	1
the	4

# Creating the initial RDD

```
counts = sc.textFile(file)
```

```
the cat sat on the  
mat  
  
the aardvark sat on  
the sofa
```

# First Mapping

```
counts = sc.textFile(file) \  
    .flatMap(lambda line: line.split())
```

the cat sat on the mat
the aardvark sat on the sofa



the
cat
sat
on
the
mat
the
aardvark
...

# Second Mapping

```
> counts = sc.textFile(file) \
    .flatMap(lambda line: line.split()) \
    .map(lambda word: (word, 1))
```

Key-Value Pairs

the cat sat on the mat
the aardvark sat on the sofa



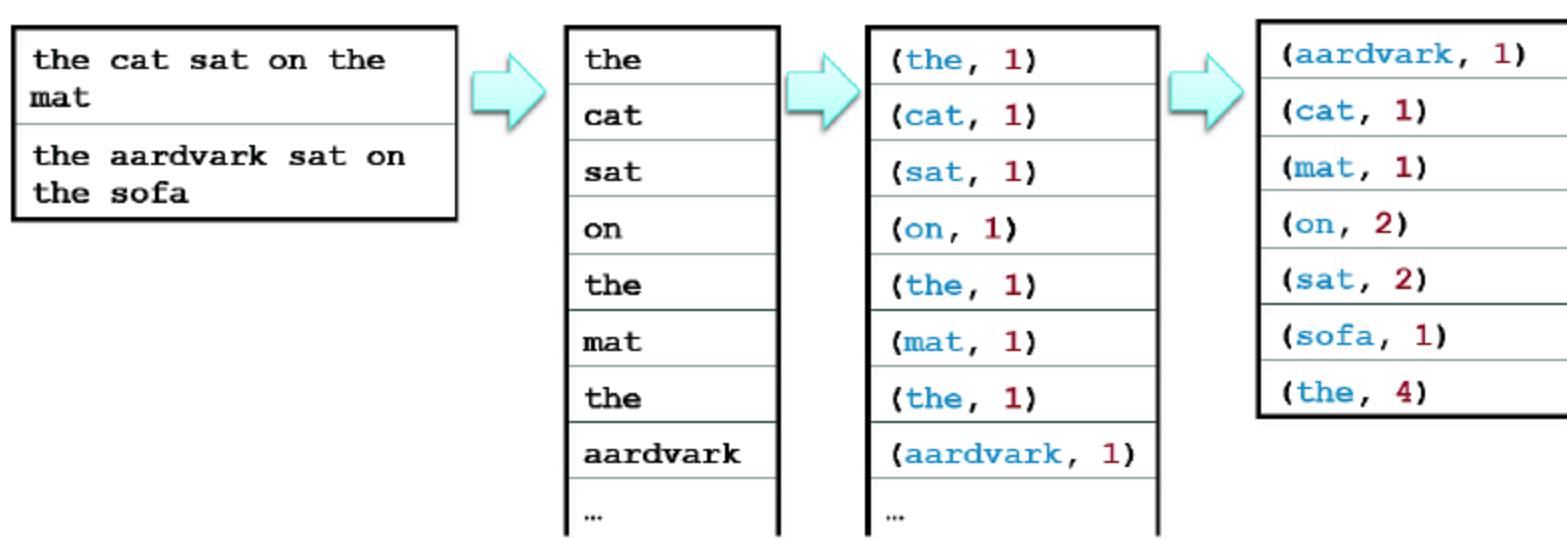
the
cat
sat
on
the
mat
the
aardvark
...



(the, 1)
(cat, 1)
(sat, 1)
(on, 1)
(the, 1)
(mat, 1)
(the, 1)
(aardvark, 1)
...

# Reducing

```
counts = sc.textFile(file) \
    .flatMap(lambda line: line.split()) \
    .map(lambda word: (word,1)) \
    .reduceByKey(lambda v1,v2: v1+v2)
```



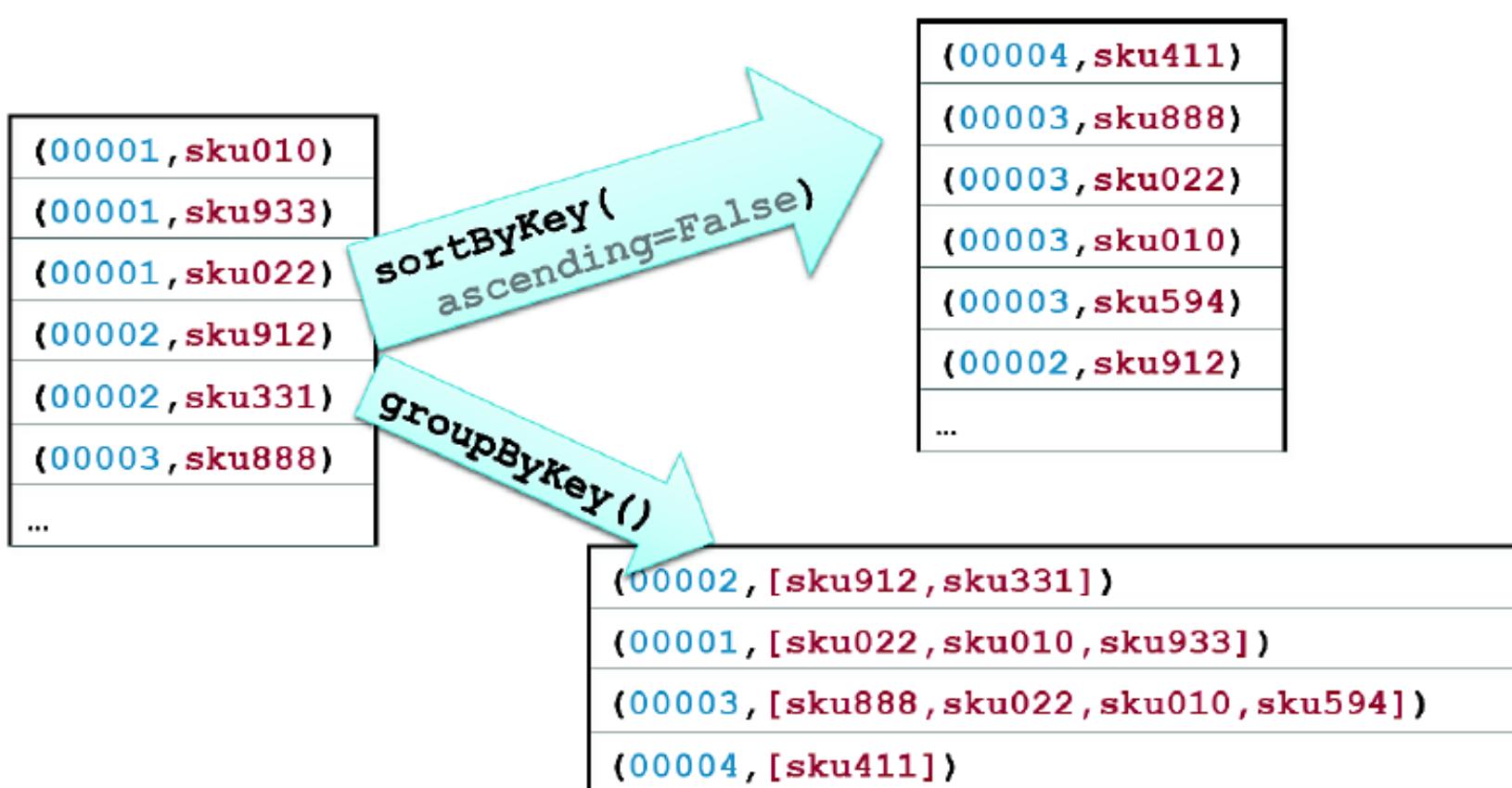
# Pair RDD Operations

# Pair RDD operations

In addition to map and reduce functions, Spark has several operations specific to Pair RDDs

- ▶ Examples
  - ▶ `reduceByKey` – merge the values for each key using a reduce function
  - ▶ `countByKey` – return a map with the count of occurrences of each key
  - ▶ `groupByKey` – group all the values for each key in an RDD
  - ▶ `sortByKey` – sort in ascending or descending order
  - ▶ `join` – return an RDD containing all pairs with matching keys from two RDDs
  - ...

# Example of pair RDD operations



# Example: Joining by key

```
movies = moviegross.join(movieyear)
```

RDD: moviegross
(Casablanca, \$3.7M)
(Star Wars, \$775M)
(Annie Hall, \$38M)
(Argo, \$232M)
...

RDD: movieyear
(Casablanca, 1942)
(Star Wars, 1977)
(Annie Hall, 1977)
(Argo, 2012)
...

(Casablanca, (\$3.7M, 1942))
(Star Wars, (\$775M, 1977))
(Annie Hall, (\$38M, 1977))
(Argo, (\$232M, 2012))
...

# Other pair operations

- ▶ `keys` – return an RDD of just the keys, without the values
- ▶ `values` – return an RDD of just the values, without keys
- ▶ `lookup(key)` – return the value(s) for a key
- ▶ `leftOuterJoin`, `rightOuterJoin`, `fullOuterJoin` – join, including keys defined in the left, right or either RDD respectively
- ▶ `mapValues`, `flatMapValues` – execute a function on just the values, keeping the key the same

## **The following offer more information on topics discussed in this lecture:**

Apache Spark Documentation

<https://spark.apache.org/docs/latest/>

RDD Programming Guide

<https://spark.apache.org/docs/latest/rdd-programming-guide.html>

Chapter 4: In-Memory Computing with Spark Book (Data Analytics with Hadoop)