

# EPL448: Data Mining on the Web – Lab 10



University of Cyprus  
Department of  
Computer Science

---

Παύλος Αντωνίου

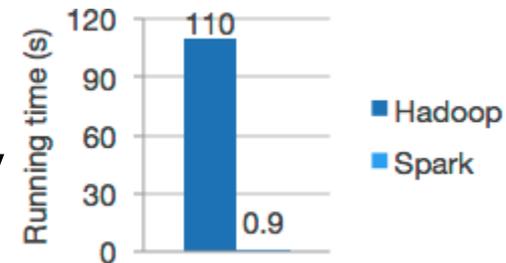
Γραφείο: B109, ΘΕΕΕ01

# Introduction to Apache Spark



- Fast and general **engine for large-scale data processing** on clusters

- claim to run programs up to 100x faster than Hadoop MapReduce for in-memory analytics, or 10x faster on disk.



- Developed in the AMPLab at UC Berkeley
  - Started in 2009
- Open-sourced in 2010 under a BSD license
- Proven scalability to over 8000 nodes in production



# Spark vs Hadoop

---



- **Spark** is an **in-memory** distributed **processing** engine
  - **Hadoop** is a framework for distributed **storage** (**HDFS**) and distributed **resource management and job scheduling** (**YARN**) and distributed **processing** (**Map/Reduce**)
  - *Spark can run with (by default) or without Hadoop components (HDFS/YARN)*
-

# Storage in Spark

---



- Distributed Storage Options:
  - Local filesystem (non distributed)
  - **Hadoop HDFS** – Great fit for batch (offline) jobs.
  - [Amazon S3](#) – For batch jobs. Commercial.
  - [Apache Cassandra](#) (DB) – Perfect for streaming data analysis (time series) and an overkill for batch jobs.
  - [Apache HBase](#) (DB)
  - [MongoDB](#) (DB)
- Cassandra vs Hbase vs MongoDB: <http://db-engines.com/en/system/Cassandra%3BHBase%3BMongoDB>

# Job Scheduling in Spark

---



- Distributed Resource Management & Job Scheduling Options:
    - **Standalone**: simple cluster manager included with Spark that makes it easy to set up a cluster
    - **Hadoop YARN**: the resource manager in Hadoop
    - **Apache Mesos**: a general cluster manager that can also run Hadoop MapReduce and service applications
  - [Hadoop vs Spark](#)
-

# Key points about Apache Spark

---



- Runs on both Windows and UNIX-like systems
  - Provides high-level APIs in Java, Scala, Python and R
  - Supports rich set of higher-level tools including:
    - [SQL DataFrames and Datasets](#) for SQL and structured data processing (e.g. csv, json files)
    - [ML](#) (DataFrame-based) and [MLlib](#) (RDD-based) **for machine learning**
    - [GraphX](#) for graph processing
    - [Spark Streaming](#) for stream processing of live data streams (e.g. sources: TCP/IP sockets, Tweets, ...)
-

# Running Apache Spark



- Apache Spark 3.0.1 is installed on your Virtual Machine
- Start Spark Shell:
  - cd /usr/local/spark/bin
  - Python Shell:
    - ./pyspark →
  - Scala Shell
    - ./spark-shell
  - R Shell
    - ./sparkR
  - [Submit an application](#) written in a file
    - ./spark-submit --master local[2] SimpleApp.py
  - Run ready-made examples
    - ./run-example <class> [params]

```
Welcome to
Spark version 3.0.0
Using Python version 3.7.8 (default, Jun 29 2020 05:46:05)
SparkSession available as 'spark'.
```

# Spark Essentials: SparkContext

---



- First thing that a Spark program does is create a ***SparkContext*** object, which tells Spark how to **access a cluster**
  - In your programs, you must use a constructor to instantiate a new ***SparkContext***
  - Then in turn ***SparkContext*** gets used to create other variables
-

# Spark Essentials: SparkSession

---



- From Spark 2.0 onwards, *SparkSession* is introduced
  - In your programs, you must use a `SparkSession.builder` to instantiate a new *SparkSession* object
  - No need to create *SparkContext*, since *SparkSession* encapsulates the same functionalities
-

# Spark Essentials: Master



- The *master* parameter determines which cluster to use

e.g. `./spark-submit --master local[2] SimpleApp.py`

| master            | description  |
|-------------------|--|
| local             | run Spark locally with one worker thread (no parallelism)                                  |
| local[K]          | run Spark locally with K worker threads (ideally set to # cores)                           |
| spark://HOST:PORT | connect to a Spark standalone cluster manager;<br>PORT depends on config (7077 by default) |
| mesos://HOST:PORT | connect to a Mesos cluster manager;<br>PORT depends on config (5050 by default)            |

# Hands on - change master param



- Through a program (via SparkContext object)
  - To create a SparkContext object (sc) you first need to build a SparkConf object that contains information about your application.

- Python

```
from pyspark import SparkContext, SparkConf
sconf = SparkConf().setAppName("App").setMaster("local[4]")
sc = SparkContext(conf=sconf)
```

- Run program

```
./bin/pyspark pythonfile.py
```

- Command line (initiate shell with 4 worker threads given that no master information defined in file):

```
./bin/pyspark --master local[4] pythonfile.py
```

# Hands on - change master param



- Through a program (via SparkSession object)

- Python

```
from pyspark.sql import SparkSession
spark =
SparkSession.builder.appName("App").master("local[4]").getOrCreate()
```

- Run program

```
./bin/pyspark pythonfile.py
```

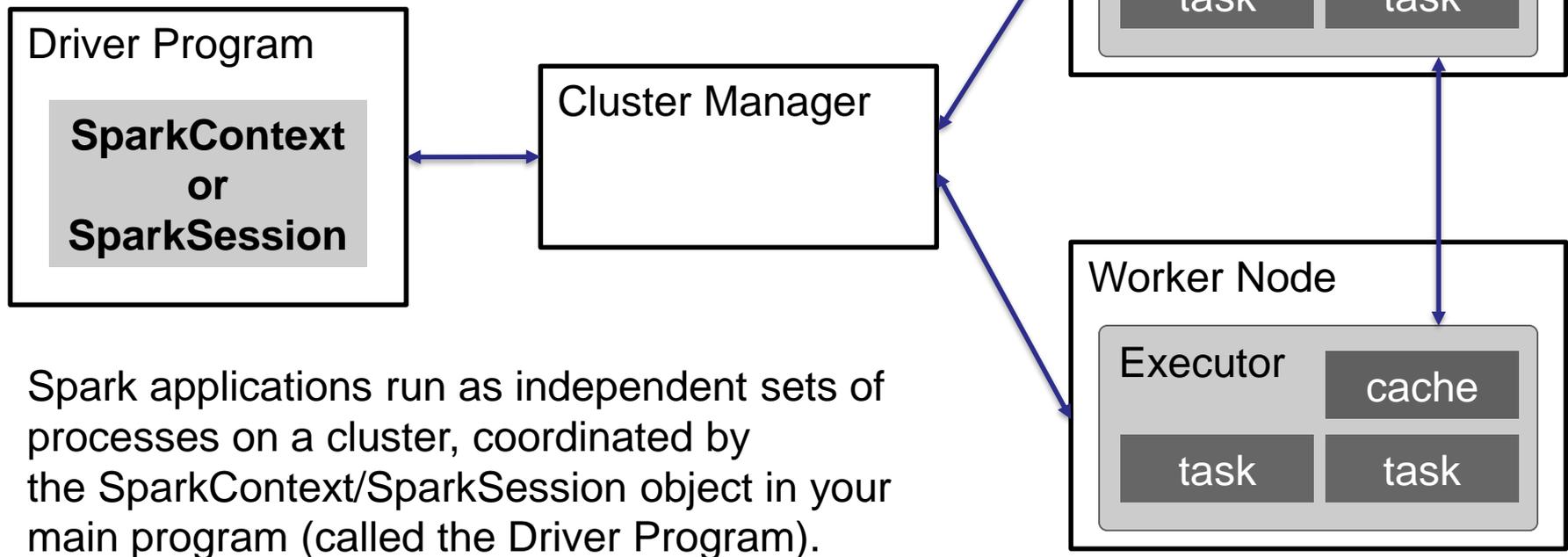
- Command line (initiate shell with 4 worker threads given that no master information defined in file):

```
./bin/pyspark --master local[4] pythonfile.py
```

# Spark Essentials: Run on cluster



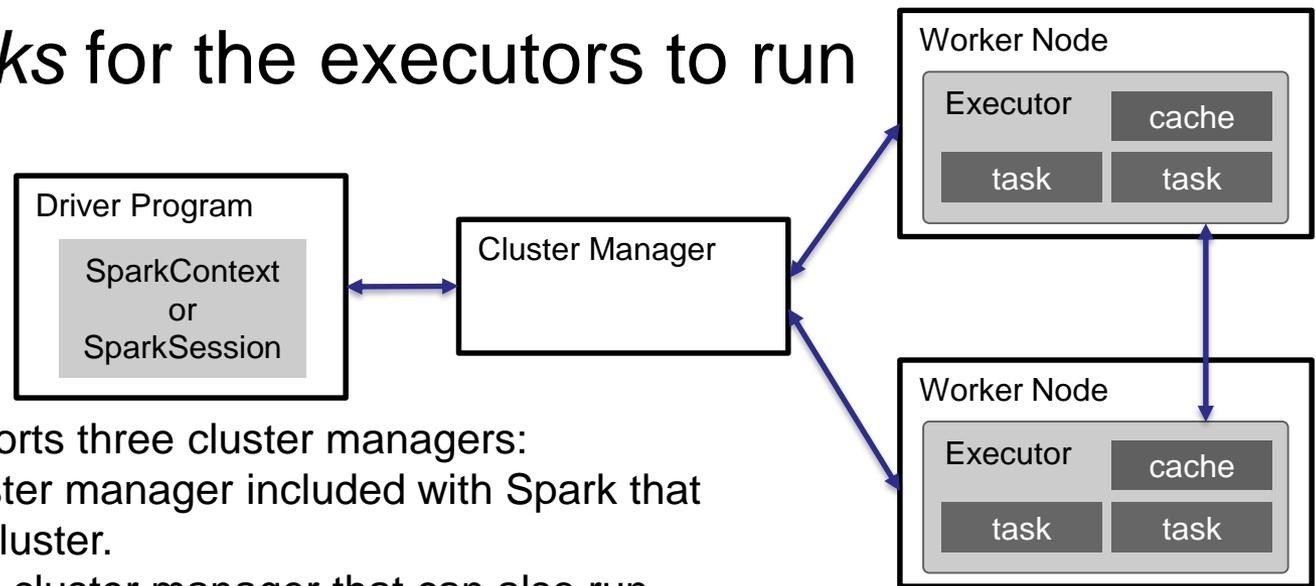
<https://spark.apache.org/docs/latest/cluster-overview.html>



# Spark Essentials: Run on clusters



1. connects to a *cluster manager* which allocates resources across applications
2. acquires *executors* on cluster nodes – worker processes to run computations and store data
3. sends *app code* to the executors
4. sends *tasks* for the executors to run



The system currently supports three cluster managers:

[Standalone](#) – a simple cluster manager included with Spark that makes it easy to set up a cluster.

[Apache Mesos](#) – a general cluster manager that can also run Hadoop MapReduce and service applications.

[Hadoop YARN](#) – the resource manager in Hadoop 2.

# Spark Web UI

<http://spark.apache.org/docs/latest/security.html#configuring-ports-for-network-security>



## Standalone mode only

| From                       | To                | Default Port | Purpose                              | Configuration Setting  | Notes  |
|----------------------------|-------------------|--------------|--------------------------------------|--|--|
| Browser                    | Standalone Master | 8080         | Web UI                               | <code>spark.master.ui.port</code> / <code>SPARK_MASTER_WEBUI_PORT</code> | Jetty-based. Standalone mode only.   |
| Browser                    | Standalone Worker | 8081         | Web UI                               | <code>spark.worker.ui.port</code> / <code>SPARK_WORKER_WEBUI_PORT</code> | Jetty-based. Standalone mode only.   |
| Driver / Standalone Worker | Standalone Master | 7077         | Submit job to cluster / Join cluster | <code>SPARK_MASTER_PORT</code>   | Set to "0" to choose a port randomly. Standalone mode only.                                      |
| External Service           | Standalone Master | 6066         | Submit job to cluster via REST API   | <code>spark.master.rest.port</code>                                      | Use <code>spark.master.rest.enabled</code> to enable/disable this service. Standalone mode only. |
| Standalone Master          | Standalone Worker | (random)     | Schedule executors                   | <code>SPARK_WORKER_PORT</code>   | Set to "0" to choose a port randomly. Standalone mode only.                                      |

Display information about running tasks, executors, and storage usage

## All cluster managers

| From                         | To                | Default Port | Purpose  | Configuration Setting                | Notes   |
|------------------------------|-------------------|--------------|--|--------------------------------------|---|
| Browser                      | Application       | 4040         | Web UI   | <code>spark.ui.port</code>           | Jetty-based                                     |
| Browser                      | History Server    | 18080        | Web UI   | <code>spark.history.ui.port</code>   | Jetty-based                                     |
| Executor / Standalone Master | Driver            | (random)     | Connect to application / Notify executor state changes | <code>spark.driver.port</code>       | Set to "0" to choose a port randomly.           |
| Executor / Driver            | Executor / Driver | (random)     | Block Manager port                                     | <code>spark.blockManager.port</code> | Raw socket via <code>ServerSocketChannel</code> |

# Spark Essentials

---



- Spark application consists of a *driver program* that runs the user's main function and executes various *parallel operations* on a cluster
- Main data abstractions to handle set of data:
  - *Resilient Distributed Datasets (RDDs)*
  - *DataFrames*
    - <https://spark.apache.org/docs/latest/sql-programming-guide.html>
  - *Datasets*

# RDDs



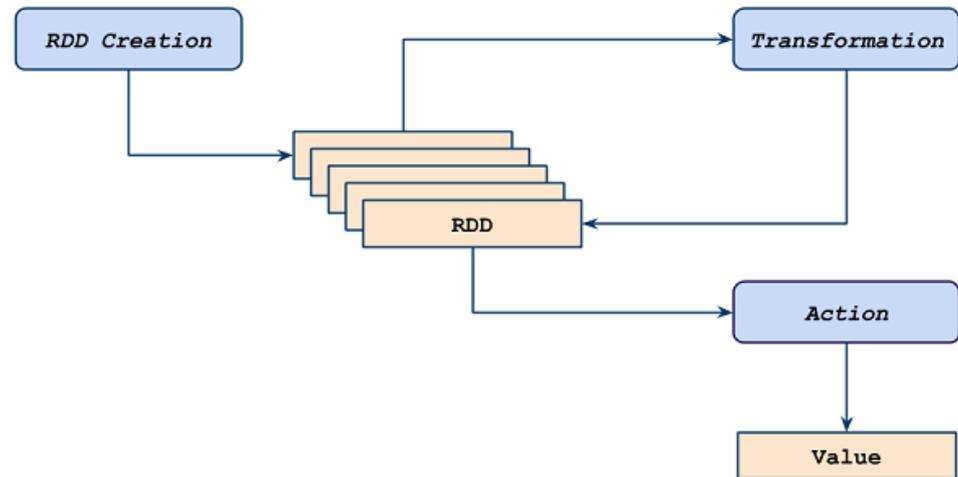
- Read-only distributed **collection of elements (Java objects)** that can be run in parallel on many nodes
  - RDDs are immutable primarily for high-speed gains
- Provided by Spark v1.0
- RDDs are created by loading data:
  - from distributed data stores (like HDFS, HBase, Cassandra, or any other data sources supported by Hadoop),
  - by parallelizing Scala collections or Python iterables lists
  - by reading files stored in the local file system
- Each record in the RDD can be divided into logical parts and then executed on different nodes of the cluster

# RDDs: Transformations & Actions



- Two types of *operations* on RDDs:

## **transformations & actions**



- Transformations create a new RDD from an existing one. **But they are lazy**: do not compute their results immediately
  - example transformations: `map()`, `filter()`
- **transformations** only **computed** when an **action** **requires** a **result** to be returned to driver
  - transformed RDDs **recomputed** each time an action runs on them. Example actions: `reduce()`, `collect()`, `count()`

# RDDs: Transformations & Actions

---



- RDD can be persisted into storage in memory or disk
- A full list with all the supported transformations and actions can be found on the following links

<http://spark.apache.org/docs/latest/rdd-programming-guide.html#transformations>

<http://spark.apache.org/docs/latest/rdd-programming-guide.html#actions>

---

# RDD Cons

---



- RDDs contain Java objects => suffer from both Garbage Collection and Java serialization\* issues which are expensive operations when the data grows.
- Unluckily, Spark does not offer any built-in optimization to speed up this kind of processes
- Thus, Spark introduced DataFrames

(\*) RDDs are sometimes **stored in serialized form** (memory or disk) to decrease memory usage (see [here](#)). Also, the **shuffle is Spark's mechanism for re-distributing data so that it's grouped differently across partitions**. This typically involves copying data across executors and machines, making the shuffle a complex and costly operation since it **involves** disk I/O, **data serialization**, and network I/O.

---

# DataFrames

---



- Included in Spark v1.3 under Spark SQL module
  - Spark SQL is a Spark module for structured data processing which allows the usage of SQL queries
  - DataFrame API is available in Scala, Java, [Python](#), & [R](#)
- Dataset organized into named columns
  - At a conceptual level, it is equivalent to a table in a relational database or to a Pandas' dataframe in Python
- DataFrames can be constructed from different sources:
  - structured data files (e.g. csv, json, xml)
  - Hive tables (stores structured data on top of Hadoop)
  - tables from external databases (MySQL, Cassandra, ...)
  - existing RDDs

# DataFrames



- Compared to RDDs, DataFrames offer a higher level of abstraction
  - Can be treated as tables
  - Spark SQL provides APIs to run SQL queries on DataFrames with a simple SQL-like syntax.
- Catalyst query optimizer & Tungsten serialization optimizer
  - speed up the execution time of Spark jobs in terms of CPU and memory efficiency, considerably reducing the thrash policies of the Garbage Collector
- DataFrames are not type safe: type checked only at runtime
  - For example, if you accidentally select the wrong column when writing a query with DataFrames, the compiler does not complain and the error is only detected when the application is run
  - Leads to tedious and time-consuming app development
- Thus, Spark introduced Datasets

# Datasets



- Included in Spark v1.6 under Spark SQL module
  - Dataset API is available in [Scala](#) and [Java](#)
- Distributed collection of data which combines the benefits of RDDs & the power of Spark SQL engine
  - constructed from JVM objects by reading data from external sources & then manipulated using transformations
  - good choice when dealing with structured or semi-structured data
- Like RDDs, Datasets are compile-time type safe
  - objects' types are inferred at compile-time which pragmatically saves development time
- Like DataFrames, Datasets involve optimizers with gains in memory use and program execution
  - suffer from Garbage Collector overhead when object serialization is necessary

# Hands on – RDDs



- From the **python** “>>>” interpreter, let’s **create a collection** (list) holding the numbers 1 to 5:  

```
data = [1, 2, 3, 4, 5]
```
- then **create an RDD** (parallelized collection) based on that data that can be operated on in parallel:  

```
distData = sc.parallelize(data)
```
- finally use a **filter transformation** to select values less than 3 and then **action collect** RDD contents back to driver

```
distData.filter(lambda s: s<3).collect()
```

- Filter transformation returns a new dataset formed by selecting those elements of the **source** on which **function** returns true

# Hands on – RDDs

---



- Lambda: anonymous functions on runtime
    - Normal function definition: `def f (x): return x**2`
      - call: `f(8)`
    - Anonymous function: `g = lambda x: x**2`
      - call: `g(8)`
-

# Hands on – pyspark



```
from pyspark import SparkContext, SparkConf
sconf = SparkConf().setAppName("SimpleApp").setMaster("local")
sc = SparkContext(conf=sconf)
# file is a list of lines from a file located on HDFS
file =
sc.textFile("hdfs://localhost:54310/user/csdeptucy/input/unixd
ict.txt")
# lineLength is the result of a map transformation, i.e. a new
RDD. Function len() is applied to each element of file list
(i.e. each line) Not immediately computed, due to laziness.
lineLengths = file.map(lambda line: len(line))
# reduce is an action. At this point Spark breaks the
computation into tasks to run on separate machines; each
machine runs both its part of the map and a local reduction,
returning only its answer to the driver program.
# a is the previous aggregate result and n is the current item
totalLength = lineLengths.reduce(lambda a, n: a + n)
print("The result is : ",totalLength)
```

**RUN APP on SPARK USING: ./spark-submit SimpleApp.py**

# Hands on – pyspark

---



- map() transformation
    - applies the given function **on every element** of the RDD  
=> returns new RDD representing the results
    - `x = [1, 2, 3, 4, 5]` # python list
    - `par_x = sc.parallelize(x)` # distributed RDD
    - `result = par_x.map(lambda i : i**2)` # new RDD
    - `print(result.collect())` → [1, 4, 9, 16, 25]
  - reduce() action
    - **aggregates** all elements of RDD using a given function and returns the final result to the driver program
-

# Hands on – pyspark



```
from pyspark import SparkContext, SparkConf
sconf = SparkConf().setAppName("SimpleApp").setMaster("local")
sc = SparkContext(conf=sconf)
```

```
# file is a list of lines from a file located on HDFS
```

```
file
```

```
sc.
```

```
textFile
```

```
# lines
```

```
lines
```

```
map
```

```
reduce
```

```
# result
```

```
collect
```

```
print
```

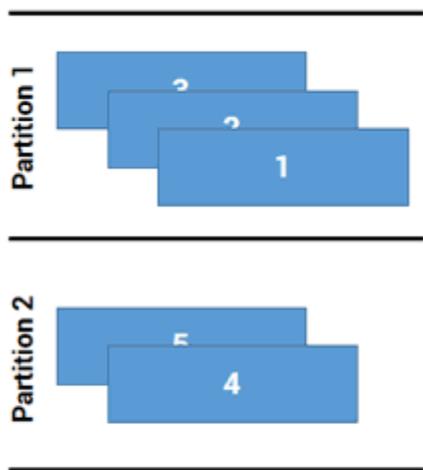
```
ret
```

```
# a is the previous aggregate result and b is the current line
```

```
totalLength = lineLengths.reduce(lambda a, n: a + n)
```

```
print("The result is : ", totalLength)
```

**RUN APP on SPARK USING:** `./spark-submit --master local SimpleApp.py`



BACK TO BASICS



reduce ( f(x) )

$f(x) = ((a, n) \Rightarrow (a+n))$

ixd

on

l,

# Hands on – pyspark



```
from pyspark import SparkContext, SparkConf
sconf = SparkConf().setAppName("SimpleApp").setMaster("local")
sc = SparkContext(conf=sconf)
# file is a list of lines from a file located on HDFS
file =
sc.textFile("hdfs://localhost:54310/user/csdeptucy/input/unixd
ict.txt")
# lineLength is the result of a map transformation. Function
len() is
Not immediate
lineLengths.persist() OR
# reduce lineLengths.cache()
computational
machine
returning
# a is the previous aggregate result and b is the current line
totalLength = lineLengths.reduce(lambda a, n: a + n)
print("The result is : ",totalLength)
```

If we also wanted to use lineLengths again later, we could add (before reduce):

```
lineLengths.persist() OR
lineLengths.cache()
```

which would cause lineLengths to be saved in memory after the first time it is computed.

ch  
ction,

**RUN APP on SPARK USING:** ./spark-submit --master local SimpleApp.py

# Spark Essentials: Persistence

---



- Spark can persist (or cache) RDDs, DataFrames or Datasets in memory across operations
    - `cache()` : use only default storage level `MEMORY_ONLY`
    - `persist()` : specify storage level (see next slide)
  - Each node stores in memory any slices of it that it computes and reuses them in other actions on that dataset – often making future actions more than 10x faster
  - The cache is fault-tolerant: if any partition of a dataset is lost, it will automatically be recomputed using the transformations that originally created it
-

# Spark Essentials: Persistence



| <i>transformation</i>                                  | <i>description</i>  |
|--|---|
| <b>MEMORY_ONLY</b>                                     | Store RDD as deserialized Java objects in the JVM. If the RDD does not fit in memory, some partitions will not be cached and will be recomputed on the fly each time they're needed. This is the default level. |
| <b>MEMORY_AND_DISK</b>                                 | Store RDD as deserialized Java objects in the JVM. If the RDD does not fit in memory, store the partitions that don't fit on disk, and read them from there when they're needed.                                |
| <b>MEMORY_ONLY_SER</b>                                 | Store RDD as serialized Java objects (one byte array per partition). This is generally more space-efficient than deserialized objects, especially when using a fast serializer, but more CPU-intensive to read. |
| <b>MEMORY_AND_DISK_SER</b>                             | Similar to MEMORY_ONLY_SER, but spill partitions that don't fit in memory to disk instead of recomputing them on the fly each time they're needed.  |
| <b>DISK_ONLY</b>                                       | Store the RDD partitions only on disk.  |
| <b>MEMORY_ONLY_2,</b><br><b>MEMORY_AND_DISK_2, etc</b> | Same as the levels above, but replicate each partition on two cluster nodes.  |

# Spark Essentials: Persistence



- **How to choose Persistence:**

- If your RDDs fit comfortably with the default storage level (**MEMORY\_ONLY**), leave them that way. This is the most CPU-efficient option, allowing operations on the RDDs to run as fast as possible.
- If not, try using **MEMORY\_ONLY\_SER** and selecting a fast serialization library to make the objects much more space-efficient, but still reasonably fast to access.
- Don't spill to disk unless the functions that computed your datasets are expensive, or they filter a large amount of the data. Otherwise, re-computing a partition may be as fast as reading it from disk.
- Use the replicated storage levels if you want fast fault recovery (e.g. if using Spark to serve requests from a web application). All the storage levels provide full fault tolerance by re-computing lost data, but the replicated ones let you continue running tasks on the RDD without waiting to re-compute a lost partition.

# More on transformation & actions

---



- flatMap() transformation
    - same as map but instead of returning just one element per element **returns a sequence per element** (which can be empty) – flattens the results
  - reduceByKey() transformation
    - operation on key-value pairs
      - in Python, key-value pairs can be implemented as tuples
    - **aggregates all values having the same key** using a given function
    - returns a distributed dataset (RDD)
-

# RDD Word Count Example



- Now lets run word count example:

```
file = sc.textFile("file:///usr/local/spark/bin/pg4300.txt")
words = file.flatMap(lambda line: line.split(" "))
             .map(lambda word: (word, 1))
             .reduceByKey(lambda a, b: a + b)
words.saveAsTextFile("output") # output folder on local fs
```

- of course `pg4300.txt` could be located in HDFS as well:

```
file =
sc.textFile("hdfs://localhost:54310/user/csdeptucy/input/pg4
300.txt")
words = file.flatMap(lambda line: line.split(" "))
             .map(lambda word: (word, 1))
             .reduceByKey(lambda a, b: a + b)
words.saveAsTextFile("hdfs://...")
```

# Hands on – DataFrames



```
from pyspark.sql import SparkSession
from pyspark.sql.types import StructType, StructField, StringType, IntegerType

# create Spark Session
spark = SparkSession.builder.appName("DataFrameApp").master("local").getOrCreate()

# Load JSON dataset
df = spark.read.load("file:///usr/local/spark/bin/people.json", format="json")

print(df.schema)
# print all people
df.show()
# print only name column
df.select("name").show()
# print people with age > 23
df.filter(df.age > 23).show()
# count people with the same age
df.groupBy("age").count().show()
```

| age | id   | name   |
|-----|------|--------|
| 25  | 1201 | george |
| 28  | 1202 | john   |
| 39  | 1203 | paul   |
| 23  | 1204 | andrew |
| 23  | 1205 | thomas |

| name   |
|--------|
| george |
| john   |
| paul   |
| andrew |
| thomas |

[{"id": "1201", "name": "george", "age": 25},  
{"id": "1202", "name": "john", "age": 28},  
{"id": "1203", "name": "paul", "age": 39},  
{"id": "1204", "name": "andrew", "age": 23},  
{"id": "1205", "name": "thomas", "age": 23}]

| age | id   | name   |
|-----|------|--------|
| 25  | 1201 | george |
| 28  | 1202 | john   |
| 39  | 1203 | paul   |

| age | count |
|-----|-------|
| 39  | 1     |
| 25  | 1     |
| 28  | 1     |
| 23  | 2     |

# Machine Learning on Spark

---



- **ML and MLlib** are Apache Spark's scalable machine learning libraries
    - ML is Dataframe-based library
    - MLlib is RDD-based library
  - Write applications quickly in Java, Scala, Python, and R
  - Run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk
  - Runs on any Hadoop data source (e.g. HDFS, HBase, or local files), making it easy to plug into Hadoop workflows
-

# What is MLlib?

---



- Classification: logistic regression, **linear support vector machine (SVM)**, naive Bayes ...
- Clustering: **K-means**, Gaussian mixtures ...
- Regression: generalized linear regression, survival regression,...
- Decomposition: **singular value decomposition (SVD)**, **principal component analysis (PCA)**
- Decision trees, random forests, and gradient-boosted trees
- Recommendation: alternating least squares (ALS)
  - based on collaborative filtering
- Topic modeling: latent dirichlet allocation (LDA)
- Frequent itemsets, association rules, and sequential pattern mining

# Why MLlib?



- scikit-learn?
  - Algorithms:
    - Classification: SVM, nearest neighbors, random forest, ...
    - Clustering: k-means, spectral clustering, ...
    - Regression: support vector regression (SVR), ridge regression, Lasso, logistic regression, ...
    - Decomposition: PCA, non-negative matrix factorization (NMF), independent component analysis (ICA), ...
- Mahout?
  - Algorithms
    - Java/Scala library
    - Distributed => Scalable
    - Core algorithms on top of Hadoop Map/Reduce => Runs slow on large datasets
    - Can run on Spark
    - Classification: logistic regression, naive Bayes, random forest, ...
    - Clustering: k-means, fuzzy k-means, ...
    - Recommendation: ALS, ...
    - Decomposition: PCA, SVD, randomized SVD, ...

# K-means (RDD-based)



- Study the file [kmeans-rdd.py](#)
  - Input file: /usr/local/spark/data/mllib/kmeans\_data.txt

```
0.0 0.0 0.0
0.1 0.1 0.1
0.2 0.2 0.2
9.0 9.0 9.0
9.1 9.1 9.1
9.2 9.2 9.2
```

- cd /usr/local/spark/bin
- wget <https://www.cs.ucy.ac.cy/courses/EPL448/1abs/LAB09/kmeans-rdd.py>
- ./spark-submit kmeans-rdd.py

# K-means (DataFrame-based)



- Study the file [kmeans-dataframe.py](#)
  - Input file: /usr/local/spark/data/mllib/kmeans\_data.txt

```
0.0 0.0 0.0
0.1 0.1 0.1
0.2 0.2 0.2
9.0 9.0 9.0
9.1 9.1 9.1
9.2 9.2 9.2
```

- cd /usr/local/spark/bin
- wget <https://www.cs.ucy.ac.cy/courses/EPL448/1abs/LAB09/kmeans-dataframe.py>
- ./spark-submit kmeans-dataframe.py

# What is Spark Streaming?



- Enables scalable, high-throughput, fault-tolerant **stream processing of live data streams**
- Data ingestion from many sources like Kafka, Kinesis, TCP sockets, Twitter, ...
- Data processing using complex algorithms expressed with functions like map, reduce, join, ...
- Processed data can be pushed out to filesystems, databases, and live dashboards



# How Spark Streaming works?



- Internally Spark Streaming receives live input data streams and divides the data into batches, which are then processed by the Spark engine to generate the final stream of results in batches.



DStream: a continuous stream of data.

DStream is represented as a sequence of RDDs.

# Spark Examples: Spark Streaming



```
from pyspark import SparkContext
from pyspark.streaming import StreamingContext
# Create a local StreamingContext with two working threads and batch
interval of 1 second
sconf = SparkContext("local[2]", "NetworkWordCount")
ssc = StreamingContext(sconf, 5)
# Create a DStream to connect to hostname:port, like localhost:9999
# lines DStream represents the stream of data that will be received
from the data server. Each record in this DStream is a line of text.
lines = ssc.socketTextStream("localhost", 9999)
# Split each line into words
words = lines.flatMap(lambda line: line.split(" "))
# Count each word in each batch
pairs = words.map(lambda word: (word, 1))
wordCounts = pairs.reduceByKey(lambda x, y: x + y)
# Print the first ten elements of each RDD generated in this DStream
to the console
wordCounts.pprint()
ssc.start() # Start the computation
ssc.awaitTermination() # Wait for the computation to terminate
```

See next slides

[Detailed explanation here.](#)

# Spark Examples: Spark Streaming



```
# Firstly, in one terminal run Netcat (http://nc110.sourceforge.net)
# to generate a data stream on the localhost:9999 TCP socket
```

```
$ nc -lk 9999
```

```
hello world
```

```
hi there fred
```

```
what a nice world there
```

```
# In another terminal run the NetworkWordCount example
# expecting a data stream on the localhost:9999 TCP socket
```

```
$ /usr/local/spark/bin/spark-submit network-word-count.py
```

```
# TERMINAL 1:
# Running N
# etcat
```

```
$ nc -lk 99
99
```

```
hello world
```

```
...
```

Scala Java **Python**

```
# TERMINAL 2: RUNNING network_wordcount.py
```

```
$ ./bin/spark-submit examples/src/main/python/streaming/network_wordcount.p
y localhost 9999
```

```
...
```

```
-----
Time: 2014-10-14 15:25:21
-----
```

```
(hello,1)
```

```
(world,1)
```

```
...
```

# Spark Examples: Spark Streaming



- Input DStreams represent the stream of input data received from streaming sources
- `lines` was an input DStream as it represented the stream of data received from the netcat server
- Every input DStream (except file stream) is associated with a **Receiver** object which receives the data from a source and stores it in Spark's memory for processing.
- Spark worker/executor is a long-running task occupying one of the cores allocated to the Spark Streaming application.
  - Important: Spark Streaming application needs to be allocated enough cores (or threads, if running locally) to process received data, as well as to run the receiver(s).
  - When running a Spark Streaming program locally, do not use “local” or “local[1]” as the master URL. Either of these means that only one thread will be used for running tasks locally. If you are using an input DStream based on a receiver (e.g. sockets, Kafka, etc.), then the single thread will be used to run the receiver, leaving no thread for processing the received data.