*"The goal is to turn data into information, and information into insight"*

Carly Fiorina, HP CEO, 2000s

# Hands-on H5
# Dataflow Programming

CS205: Computing Foundations for Computational Science
Dr. Ignacio M. Llorente
Spring Term 2020

# Before We Start
## Where We Are

**Computing Foundations for Computational and Data Science**

How to use modern computing platforms in solving scientific problems

Intro: Large-Scale Computational and Data Science

A. Parallel Processing Fundamentals

B. Parallel Computing

**C. Parallel Data Processing**

    C1. Batch Data Processing

    **C2. Dataflow Processing**

    C3. Stream Data Processing

    C4. Complex Parallel Processing Workflows

Wrap-Up: Advanced Topics

# CS205: Contents

## APPLICATION SOFTWARE

| Application Parallelism | | Program Design |
|---|---|---|

**Application Software**

**BIG COMPUTE**

| OpenACC | Optimization | | Spark |
|---|---|---|---|
| OpenMP | MPI | | Map-Reduce |

**Programming Model**

**BIG DATA**

| Slurm | **Platform** | Yarn |
|---|---|---|

**Architecture**



| Cloud Computing | | Computing Cluster |
|---|---|---|

HARVARD
School of Engineering and Applied Sciences

IACS
INSTITUTE FOR APPLIED COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

# Before We Start
## Where We Are

Concepts ➡ Platform ➡ Programming

## Week 9: **Batch Data Processing** => MapReduce

| 3/23 | 3/24 | 3/25 | 3/26 | 3/27 |
|------|------|------|------|------|
|      | **Lecture C1** | Lab I8 | Hands-on H4 |      |
|      | Batch Data Processing | MapReduce Hadoop Cluster | MapReduce Programming |      |
|      | (Quiz & Reading) |      |      |      |

## Week 10: **Dataflow Processing => Spark**

| 3/30 | 3/31 | 4/1 | 4/2 | 4/3 |
|------|------|------|------|------|
|      | Lecture C2 | Lab I9 | Hands-on H5 |      |
|      | Dataflow Processing | Spark Single Node | Spark Programming |      |
|      | (Quiz & Reading) |      |      |      |

# Context
## The Spark Programming Model



The Fundamental Data Structure - Resilient Distributed Dataset
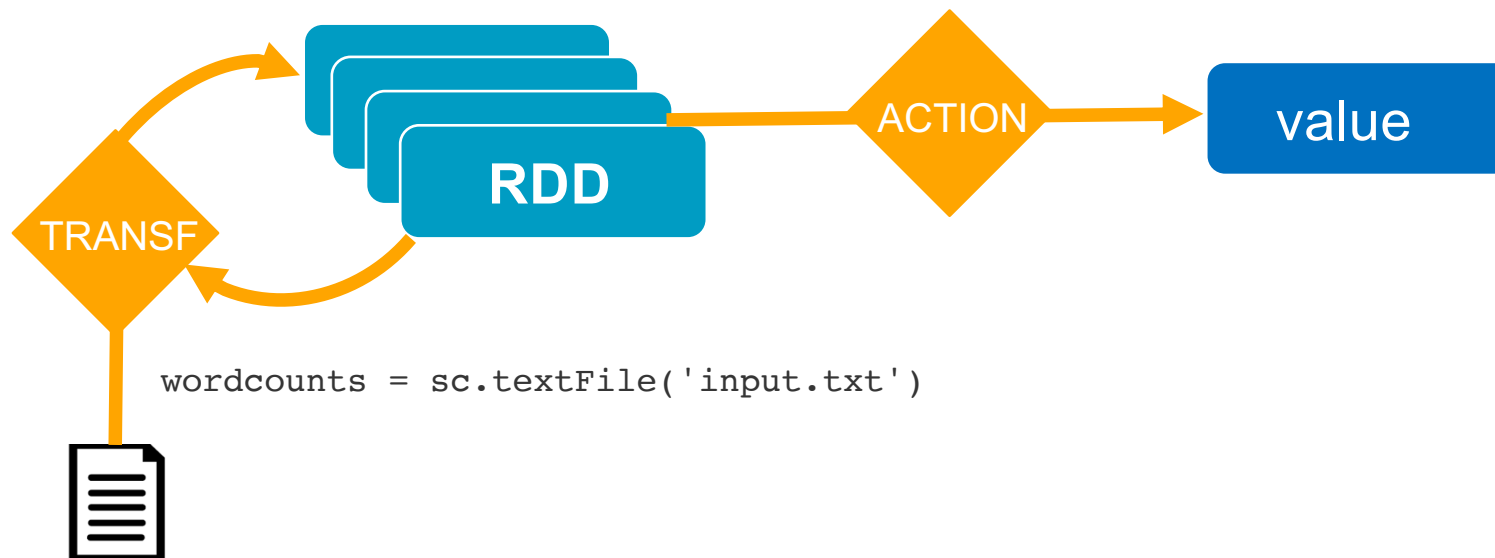
- **Resilient:** Fault-tolerant
- **Distributed:** Multiple-node
- **Dataset:** Collection of partitioned data organized in records

(the, 7)   (at, 1)    (cloud, 76)
(od, 4)    (bok, 8)   (data, 5)
(spark, 1) (home, 1)  (set, 34)

**Operations: Transformations and Actions**

```
.filter(lambda line: "spark" in line)          .count()
```

RDD

ACTION → value

TRANSF

```
wordcounts = sc.textFile('input.txt')
```

# Hands-on Examples
## Requirements

1. Unix-like shell (Linux, Mac OS or Windows/Cygwin)

2. Python installed

3. Installation of Spark (see guide "Install Spark in Local Mode")

# Roadmap
## Dataflow Programming

PySpark

Resilient Distributed Datasets

Distributed Collections

Transformations

Actions

Caching and Persistence

Pipelining

HARVARD
School of Engineering
and Applied Sciences

IACS
INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

# PySpark
## Interactive Shell for Python

---

**Spark Interactive Shell for Python**

- Easiest way to try Spark.

- Responsible for linking the python API to the spark core and initializing the Spark context

- Runs in local mode on 1 thread by default, but can control through `MASTER` environment variable

---

```
Python 2.7.12 (default, Nov 19 2016, 06:48:10)
[GCC 5.4.0 20160609] on linux2
Type "help", "copyright", "credits" or "license" for more information.
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 2.2.0
      /_/

Using Python version 2.7.12 (default, Nov 19 2016 06:48:10)
SparkSession available as 'spark'.
>>>
```

HARVARD
School of Engineering
and Applied Sciences

IACS
INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

# Resilient Distributed Datasets
## RDD Creation

**Different Ways to Create RDDs**

- Text File

```
>>> logFile = '/var/log/syslog'
>>> textRDD = sc.textFile(logFile) // create RDD
>>> textRDD.count() // RDD => result
>>> linesWithRoot = textRDD.filter(lambda line: 'root' in line) // RDD => RDD
>>> linesWithRoot.take(9) // RDD => result
```

- HDFS : Data residing on a distributed file system

```
>>> sc.textFile( "hdfs://namenode:9000/path/file"
```

- New Defined RDD

```
>>> data = [1, 2, 3, 4, 5]
>>> distData = sc.parallelize(data) // create distributed collection
```

- From Other RDD

```
>>> distDataS = distData.map(lambda x: x * x) // RDD => RDD
```
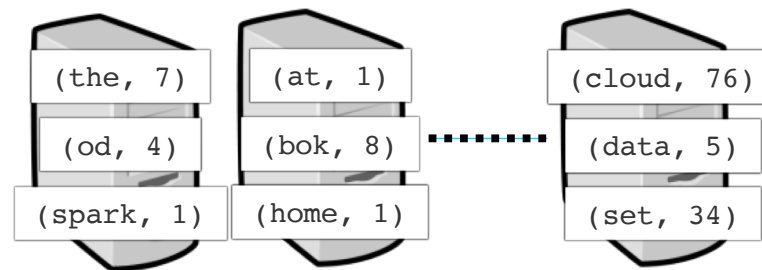
# Distributed Collections
## Parallel Processing

**RDD Partitions**

- A "parallelized" data set where the elements are copied across the nodes of a distributed system to form a distributed collection that can be computed in parallel

```
>>> data = [1, 2, 3, 4, 5]
>>> sc.defaultParallelism // default number of partitions
>>> distData = sc.parallelize(data) // create a distributed collection
>>> distDataP = sc.parallelize(data, 3) // slice the data set into 3
partitions, 3 way parallelism
>>> distDataP.count() // do some 'statistics'
>>> distDataP.getNumPartitions() // give number of partitions
>>> distDataP.reduce(lambda x, y : x + y) // more 'statistics'
```

```
(the, 7)        (at, 1)          (cloud, 76)
(od, 4)         (bok, 8)  ......  (data, 5)
(spark, 1)      (home, 1)        (set, 34)
```

HARVARD
School of Engineering
and Applied Sciences

IACS
INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

# Transformations
## Create a New RDD from an Existing One

**map()**

• Reads one element at a time

• Takes one value, creates a new value

```
>>> rdd = sc.parallelize([1, 2, 3, 4])
>>> rdd.map(lambda x: x * 2).collect()
Out[1]: [2, 4, 6, 8]
```

**flatMap()**

• Produce multiple elements for each input element

```
>>> rdd = sc.parallelize([1, 2, 3])
>>> rdd.map(lambda x: [x, x * 2]).collect()
Out[1]: [[1, 2], [2, 4], [3, 6]]
>>> rdd.flatMap(lambda x: [x, x * 2]). collect()
Out[2]: [1, 2, 2, 4, 3, 6]
```

# Transformations
## Create a New RDD from an Existing One

---

**filter()**

- Reads one element at a time
- Evaluates each element
- Returns the elements that pass the filter()

```
>>> rdd = sc.parallelize([1, 2, 3, 4])
>>> rdd.filter(lambda x: x % 2 == 0).collect()
Out[1]: [2, 4]
```

---

**Key-value Operations**

```
>>> pets = sc.parallelize([('cat', 1), ('dog', 1), ('cat', 2)])
>>> pets. reduceByKey(lambda x, y: x + y).collect() // => [('cat', 3), ('dog', 1)]
>>> pets.groupByKey().collect() // => [('cat', Seq(1, 2)), ('dog', Seq(1)]
>>> pets.sortByKey().collect() // => [('cat', 1), ('cat', 2), ('dog', 1)]
```

---

# Transformations
## Create a New RDD from an Existing One

### union()

- Merges two RDDs together

```
>>> john_smith = [('physics',85),('maths',75),('chemistry',95)]
>>> paul_adams = [('physics',65),('maths',45),('chemistry',85)]
>>> john = sc.parallelize(john_smith)
>>> paul = sc.parallelize(paul_adams)
>>> john.union(paul).collect()
```

### join()

- Joins two RDDs based on a common key

```
>>> Subject_wise_john = john.join(paul)
>>> Subject_wise_john.collect()
```

# Transformations
## Create a New RDD from an Existing One

**intersection()**

- Gives you the common terms or objects from the two RDDS

```
>>> techs = ['sachin', 'abhay', 'michael', 'rahane', 'david', 'ross',
'raj', 'rahul', 'hussy', 'steven', 'sourav']

>>> managers = ['rahul', 'abhay', 'laxman', 'bill', 'steve ']

>>> techRDD = sc.parallelize(techs)

>>> managersRDD = sc.parallelize(managers)

>>> managertechs = techRDD.intersection(managersRDD)

>>> managertechs.collect()
```

# Transformations
## Create a New RDD from an Existing One

**distinct()**

• Gets rid of any ambiguities

```
>>> best_screenplay = ["movie10","movie4","movie6","movie7","movie3"]
>>> best_story = ["movie9","movie4","movie6","movie5","movie1"]
>>> best_direction = ["movie10","movie4","movie7","movie12","movie8"]
>>> story_rdd = sc.parallelize(best_story)
>>> direction_rdd = sc.parallelize(best_direction)
>>> screen_rdd = sc.parallelize(best_screenplay)
>>> total_nomination_rdd = story_rdd.union(direction_rdd).union(screen_rdd)
>>> total_nomination_rdd.distinct().collect()
```

HARVARD
School of Engineering
and Applied Sciences

IACS
INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

# Actions
## Compute a Result Based on an Existing RDD

**count()**

```
>>> rdd = sc.parallelize([1, 2, 3, 4])
>>> rdd.count()
Out[1]: 4
```

**collect()**

- `collect()` retrieves the entire RDD
- Useful for inspecting small datasets locally and for unit tests

```
>>> rdd = sc.parallelize([1, 2, 3])
>>> rdd.collect()
[1, 2, 3]
```

# Actions
## Compute a Result Based on an Existing RDD

**take(), first(), top(), takeSample()**

- `take(n)` returns n elements from an RDD

- `takeSample()` - more suitable for taking a sample

- Use `takeOrdered(), top(n)` for ordered return

**takeOrdered()**
```
>>> rdd = sc.parallelize([5, 1, 3, 2])
>>> rdd.takeOrdered(4)
Out[1]: [1, 2, 3, 5]
>>> rdd.takeOrdered(4, lambda n: -n)
Out[2]: [5, 3, 2, 1]
```

**reduce()**

- Takes two elements of the same type and returns one new element
```
>>> rdd = sc.parallelize([1, 2, 3])
>>> rdd.reduce(lambda x, y: x * y)
Out[1]: 6
```

HARVARD
School of Engineering
and Applied Sciences

IACS
INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

Lecture H5. Spark Programming
CS205: Computing Foundations for Computational Science

Dr. Ignacio M. Llorente
18

# Caching and Persistence
## Efficiency

**Transformations are lazy!**

•A transformed RDD is only executed when actions run on it

```
>>> pyLines = lines.filter(lambda line: 'Python' in line)
>>> pyLines.first()
```

•No need for Spark to load all the lines containing "Python" into memory!

**An Example**

```
>>> textFile = sc.textFile("/user/emp.txt")
```

•It does nothing. It creates an RDD that says "we will need to load this file". The file is not loaded at this point.

•RDD operations that require observing the contents of the data cannot be lazy (**these are called actions**). An example is `RDD.count`

•So if you write `textFile.count`, at this point the file will be read, the lines will be counted, and the count will be returned.

•What if you call `textFile.count` again? The same thing: the file will be read and counted again. **Nothing is stored. An RDD is not data.**

HARVARD
School of Engineering and Applied Sciences

IACS
INSTITUTE FOR APPLIED COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

# Caching and Persistence
## Efficiency

**Caching**

- Decreases the computation time **by almost 100X** when compared to other distributed computation frameworks like hadoop mapreduce

```
>>> textFile = sc.textFile("/user/emp.txt")
```

```
>>> textFile.cache
```

- It does nothing. `RDD.cache` is also a lazy operation. The file is still not read. But now the RDD says *"read this file and then cache the contents"*. If you then run `textFile.count` the first time, the file will be loaded, cached, and counted. If you call `textFile.count` a second time, the operation will use the cache. It will just take the data from the cache and count the lines.

- `cache` is like `persist(MEMORY_ONLY)`

# Caching and Persistence
## Efficiency

**Persistence**

- RDDs are recomputed for every action, can be expensive and can also cause data to be read from disk again!

- RDDs can be cached for reuse, `rdd.persist()`

```
>>> lines = sc.textFile("README.md")
>>> lines.count()
>>> pythonLines = lines.filter(lambda line : "Python" in line)
>>> pythonLines.count()
```

Causes Spark to reload lines from disk!!!

```
>>> lines = sc.textFile("README.md")
>>> lines.persist() # Spark keeps lines in RAM
>>> lines.count()
>>> pythonLines = lines.filter(lambda line : "Python" in line)
>>> pythonLines.count()
```

Spark will avoid reloading lines every time it is used

# Pipelining
## Defining a Workflow

**Building a Pipeline of Operations**

```
>>> lines = sc.textFile("README.md")
>>> lines.map(…).filter(…).count(…)
>>> (lines
        .map(…)
        .filter(…)
        .count(…))
```

# Pipelining
## The WordCount Example with Spark

### A Pipeline of Transformations

```
wordcounts = sc.textFile('input.txt')
```

> 'The Project Gutenberg EBook of Moby Dick; or The Whale, by Herman'
> 'Melville. This eBook is for the use of anyone anywhere at no cost and'

```
.map(lambda x: x.replace(',',' ').replace('.',' '). lower())
```

> 'the project gutenberg eBook of moby dick or the whale by herman'
> 'melville this eBook is for the use of anyone anywhere at no cost and'

```
.flatMap(lambda x: x.split())
```

> 'the' 'project' 'gutenberg' 'eBook' 'of' 'moby' 'dick' 'or' 'the 'whale'
> 'by' 'herman' 'melville' 'this' 'eBook' 'is' 'for' 'the' 'use' 'of'

```
.map(lambda x: (x, 1))
```

> '(the, 1)' '(project ,1)' '(gutenberg, 1)' '(eBook, 1)' '(of, 1)' '(moby
> , 1)' '(dick, 1)' '(or, 1)' '(the, 1)' '(whale, 1)' '(by, 1)'

```
.reduceByKey(lambda x,y:x+y)
```

> '(the, 11)' '(project ,10)' '(gutenberg, 9)' '(eBook, 37)' '(of, 15)'
> '(moby , 5)' '(dick, 7)' '(or, 9)' '(the, 9)' '(whale, 123)' '(by, 98)'

# DataFrames
## Processing of Tabular Data

```
Year   Type Size

2018   A    120

2018   A    200

2019   B    300

2020   C    150
```

•Create RDD

```
>>> rdd = sc.parallelize([(2018,"A", 120),(2018,"A",
200),(2019,"B", 300),(2020,"C",150)])
>>> rdd.collect()
```

How to extract some of the columns, for example Year and Type?

How to select a group of rows, for example those from Year 2019?

How to aggregate rows, for example count by Type?

HARVARD
School of Engineering
and Applied Sciences

IACS INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

# DataFrames
## Processing of Tabular Data

```
>>> rdd = sc.parallelize([(2018,"A", 120),(2018,"A", 200),(2019,"B",
300),(2020,"C",150)])


# Convert RDD into DataFrame (you can read directly from JSON, CSV and XML)
>>> df = rdd.toDF(["year","type","size"])


# Displays the content of the DataFrame to stdout
>>> df.show()


>>> df.printSchema()


>>> df.select("type").show()


>>> df.select(df['year'], df['size'] + 1).show()


>>> df.filter(df['year'] > 2019).show()


>>> df.groupBy("type").count().show()
```

# Parallel Execution
## Single Node

**Compute PI Number with Spark**

```
from pyspark import SparkConf, SparkContext

import string


conf = SparkConf().setMaster('local[2]').setAppName('Pi')

sc = SparkContext(conf = conf)


N = 10000000

delta_x = 1.0 / N

print sc.parallelize( xrange (N),4 ).map( lambda i: (i +0.5) *
delta_x ).map( lambda x: 4 / (1 + x **2) ).reduce ( lambda a, b:
a+b) * delta_x
```

Execute with different number of partitions and threads, and compare number of tasks and execution time

HARVARD
School of Engineering and Applied Sciences

IACS
INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

# Next Steps

- Get ready for next **lecture:**
  C3. Stream Data Processing

- Get ready for next's **lab session**
  I10. Spark Cluster on AWS

# Questions
## Spark Programming

http://piazza.com/harvard/spring2020/cs205/home