

DĚLAT  
DOBRÝ SOFTWARE  
NÁS BAVÍ

# PROFINIT

## Spark

Jan Hučín

December 4th, 2019

# Agenda

1. What is it (for)
2. How to learn it
3. How to work with it
4. How it works
  - logical / technical level
  - transformations, actions, caching
  - examples
  - architecture, sources

Later:

- › Spark SQL
- › Spark streaming



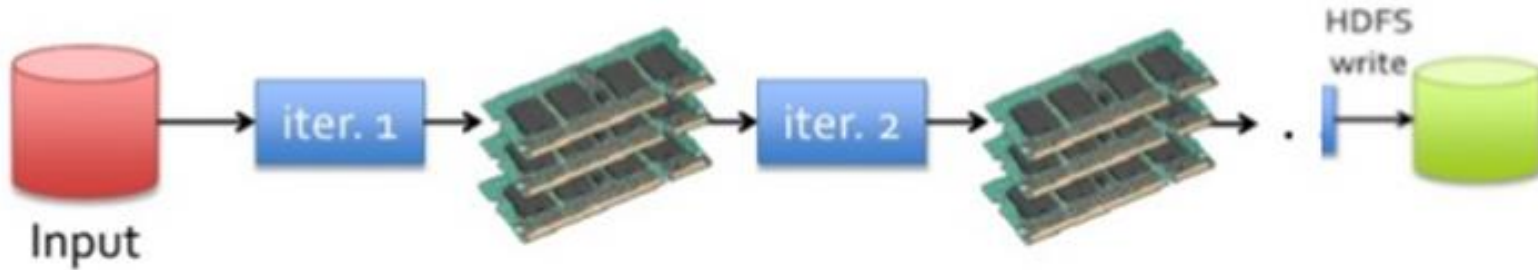
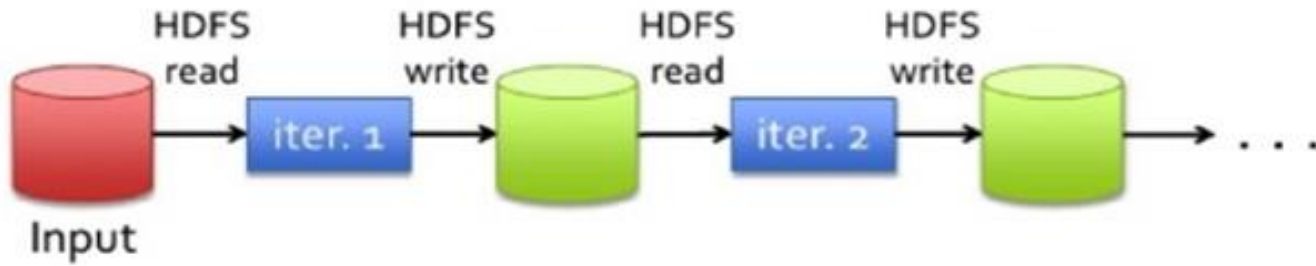
The background of the slide is a dark gray field filled with a complex, overlapping pattern of translucent, light gray polygons. These polygons vary in size and orientation, creating a sense of depth and movement, similar to a low-poly 3D environment or a crystalline structure.

# What is Spark

## (For) what is Spark

- › framework for distributed computations
- › improved map-reduce – up to 2 orders faster
  - **in-memory processing** – less I/O operations, good for iterations and data analysis
  - **operation optimisation** before processing
  - enhanced by **SQL-like** commands
- › API for Scala, Java, Python, R
- › on Hadoop (using HDFS, YARN) or standalone
- › written in Scala
- › processing in JVM
- › the most active (2017) opensource Big Data project

# Spark vs. map-reduce



## Suitable tasks

- › big enough, but not extreme
- › able to be run in parallel
- › iterative
- › not resolvable by traditional technologies

E. g.

1. sophisticated client features (risk score)
2. complicated SQL tasks for DWH
3. compute once (night), use many (day)
4. graph/network analysis
5. text-mining

## Not suitable tasks

- › too small
- › with extreme memory demands
- › custom-tailored for other technologies (SQL, Java)
- › unable to run in parallel
- › strictly real-time

E. g.

1. small data modelling
2. median computation, random skipping in the file
3. JOIN of really big tables

# How to learn it

- › <http://spark.apache.org>
- › basics of Python | Scala | Java | R
- › **self-practice**
- › advice of experienced, StackOverflow etc.



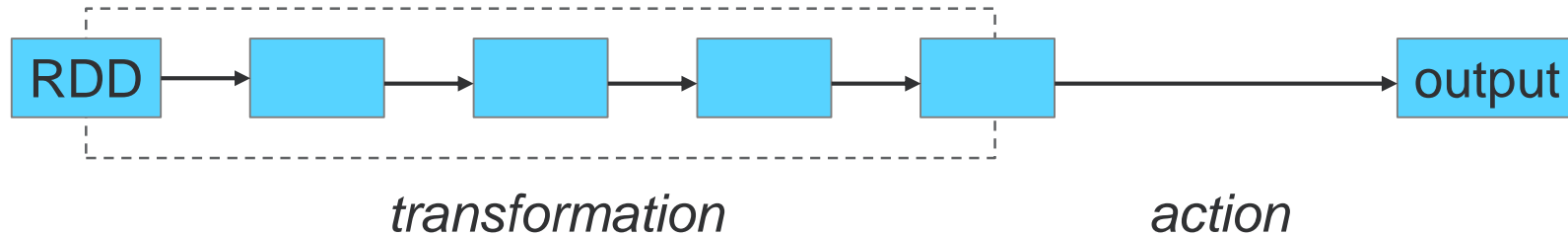
## How to work with it

- › interactively
  - command line (shell for Python and Scala)
  - Zeppelin/Jupyter notebook
- › batch / application
  - compiled .jar file
  - Python script

The background of the slide is a dark gray field filled with a complex, overlapping pattern of translucent, light gray polygons. These polygons vary in size and orientation, creating a sense of depth and movement, similar to a low-poly 3D environment or a crystalline structure.

# How Spark works

## Logical level (high)

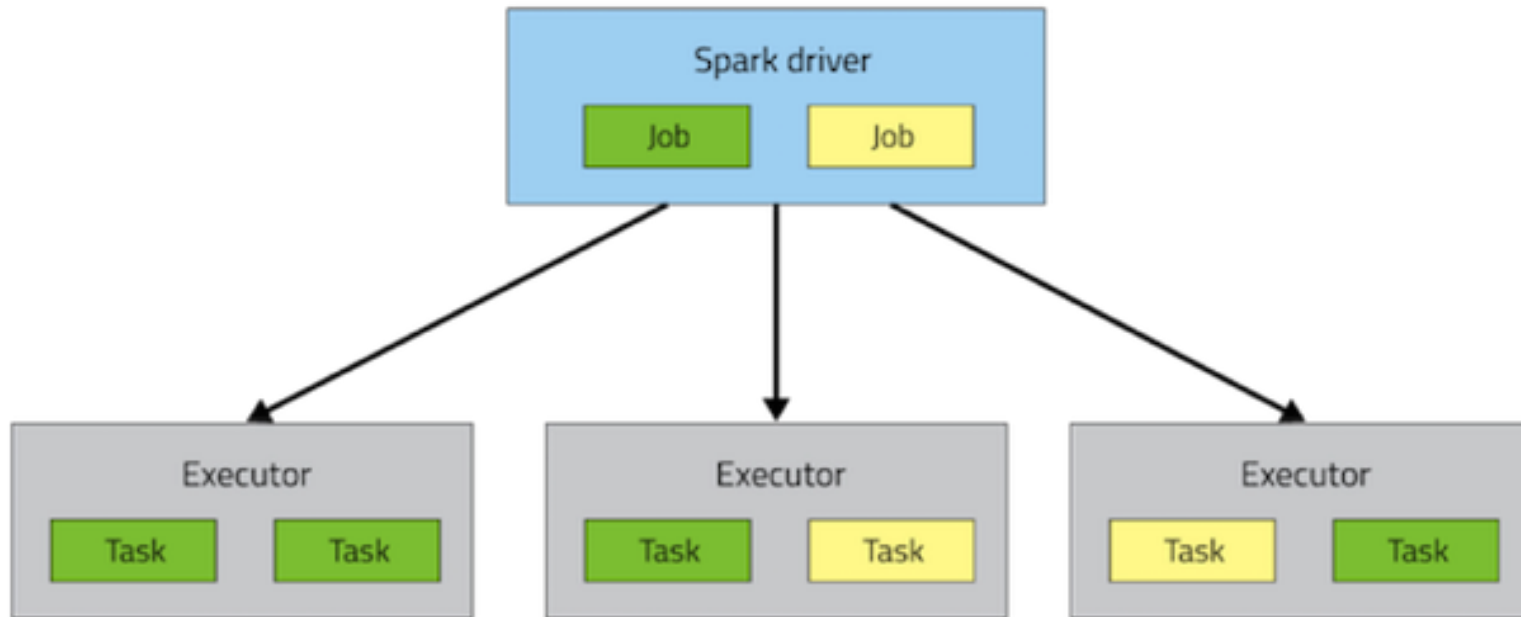


- › series of transformations finished by action
- › transformation are **planned** and **optimized**, but **not made**
- › **lazy evaluation**: action starts all process

## What is RDD?

- › resilient distributed dataset
- › item collections (line of a text file, files in a direktory, data matrix, etc.)
- › must be divisible to parts – Spark does the division itself

## Technical level (mid)



- › creates JVM on nodes (=executors)
- › application → jobs; job → tasks
- › task (and data) distribution to nodes
- › process control
- › **more later – Spark architecture**

# Spark RDD – transformations


RDD1  $\Rightarrow$  RDD2, item by item („row by row“)

- › **map** (item  $\Rightarrow$  transforming function  $\Rightarrow$  new item)
- › **flatMap** (item  $\Rightarrow$  transforming function  $\Rightarrow$  0 až N new items)
- › **filter**, **distinct** (only items meeting condition / unique items)
- › **join** (joining other RDD by key)
- › **union**, **intersection**
- › **groupByKey**, **reduceByKey** (items agregation by key)
- › ... and many others

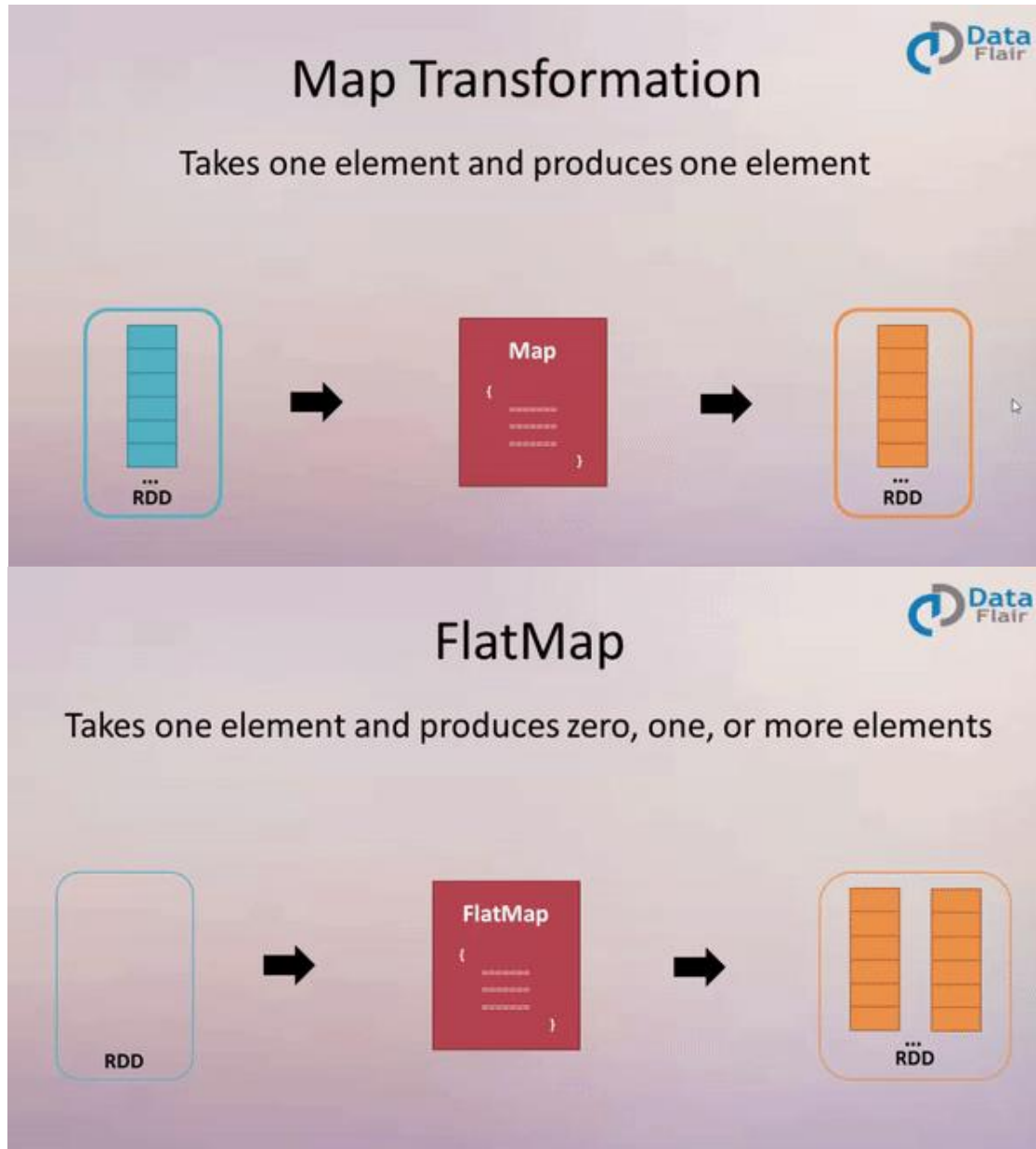
## How to get „key“?

- › first member of *tuple* = key
- › result of transformation: word  $\Rightarrow$  (word, 1)

"tuple"  
(Scala, Python)



# map a flatMap



## Example 1 – word count

- › Task: count frequency of words in document
- › Input: text file divided to lines
- › Output: RDD with items (word, frequency)
- › Workflow:
  - import text file as RDD
  - lines transformation: line  $\Rightarrow$  words  $\Rightarrow$  tuples (word, 1)
  - grouping items with same key, summing ones = count

# How to launch the interactive shell

**pyspark** (Python) | **spark-shell** (Scala)

- › from Linux shell in console
- › Spark on master node (local) or on YARN:
  - `pyspark --master local`
  - `pyspark --master yarn`
- › creates important objects: `sc` (SparkContext), `sqlContext`
- › many parameters – later
- › closed by `exit()`



## Example 1 – word count

- › Task: count frequency of words in document
- › Input: text file divided to lines
- › Workflow:
  - › import text file as RDD  
`lines = sc.textFile("/user/pascepel/bible.txt")`
  - › lines transformation: line  $\Rightarrow$  words  
`words = lines.flatMap(lambda line: line.split(" "))`
  - › lines transformation: words  $\Rightarrow$  tuples (word, 1)  
`pairs = words.map(lambda word: (word, 1))`
  - › grouping items with same key, summing ones = count  
`counts = pairs.reduceByKey(lambda a, b: a + b)`

to be  
or  
not to  
be

to  
be  
or  
not  
to  
be

(to, 1)  
(be, 1)  
(or, 1)  
(not, 1)  
(to, 1)  
(be, 1)

(to, 2)  
(be, 2)  
(or, 1)  
(not, 1)

# Why does it not count anything?

Because we have done no action so far.

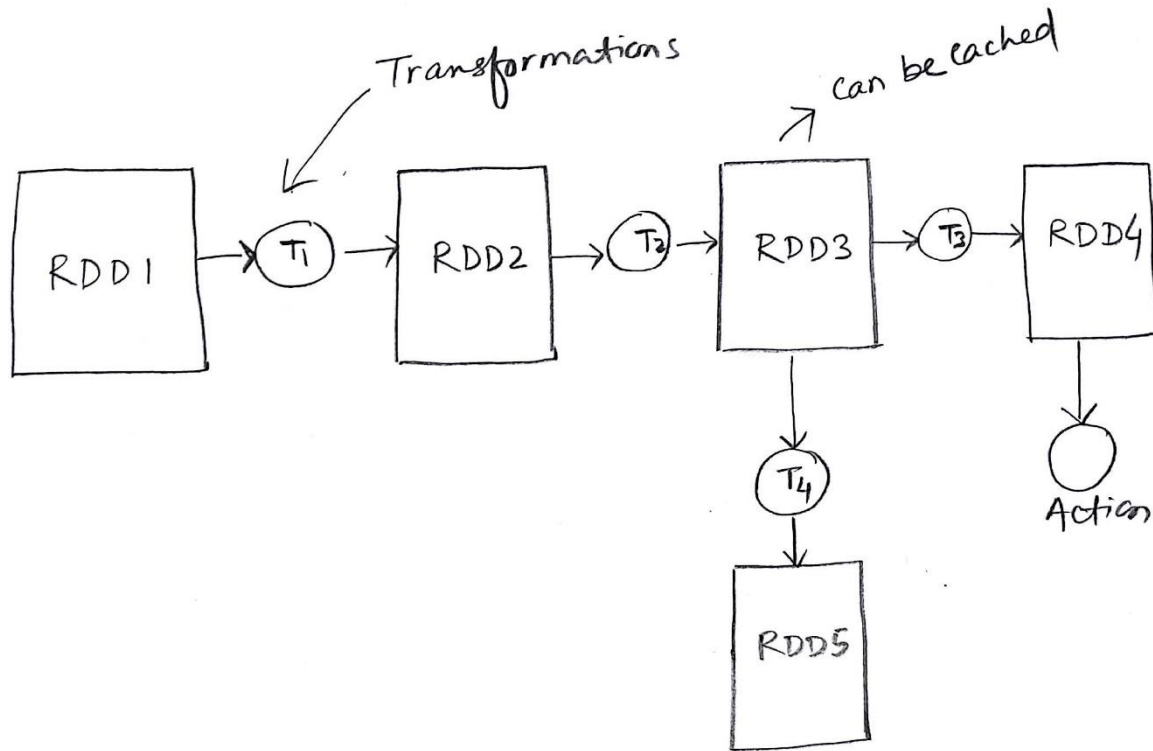
## Spark RDD – actions

- › **count**
  - › **take** (gets first  $n$  RDD items)
  - › **collect** (gets all items in RDD)
  - › **reduce** (aggregation of all RDD items with function provided)
  - › **saveAsTextFile**
  - › ... and others
- 
- › Action starts all chain from the beginning!
    - After run, all results are „forgotten“.
    - **If we don't want this, we have to cache some RDD.**

# Caching

- › Caching: RDD is not forgotten but saved into memory / on disk.
- › **Caching methods:**
  - **cache** (try to save in memory)
  - **persist** (more general – can control serialization, disk/memory)
  - **unpersist** (release RDD from memory/disk)
- › **Caching types:**
  - MEMORY\_ONLY
  - MEMORY\_AND\_DISK
  - MEMORY\_ONLY\_SER
  - MEMORY\_AND\_DISK\_SER
- › SER = serialization – less memory, high computation demands
  - Useful only for Java and Scala; Python makes serialization always
- › Caching is not an action!

# Spark program as a graph



## Example 2 – EXIF image processing

- › Task: to extract EXIF data from image files
- › Input: directory with JPG files – RDD items
- › Workflow:
  - reading binary code of a file
  - EXIF extraction (Python package PIL)
- › Result of transformation:  
RDD with items (filename, DateTime, ExposureTime, ...)
- › Output – saving data to text file

## Example 2 – EXIF image processing

Task: to extract EXIF data from image files

Input: directory with JPG files – RDD items

```
imgs = sc.binaryFiles("/user/pascepel/data/images/")
```

Workflow:

- › (filename, binary code) → (filename, dict EXIF tags)  
we use a function *get\_exif*  

```
imgs2 = imgs.map(lambda im: (im[0], get_exif(im[1])))
```
- › next operations...  

```
imgs3 = imgs2.map(...transformation function...)
```

Output: saving to a textfile

```
imgs3.saveAsTextFile('/user/pascepel/data/images_stat')
```



# Other Spark RDD operations



# Essential Core & Intermediate Spark Operations



## TRANSFORMATIONS

### General

- map
- filter
- flatMap
- mapPartitions
- mapPartitionsWithIndex
- groupBy
- sortBy

### Math / Statistical

- sample
- randomSplit

### Set Theory / Relational

- union
- intersection
- subtract
- distinct
- cartesian
- zip

### Data Structure / I/O

- keyBy
- zipWithIndex
- zipWithUniqueID
- zipPartitions
- coalesce
- repartition
- repartitionAndSortWithinPartitions
- pipe



## ACTIONS

- reduce
- collect
- aggregate
- fold
- first
- take
- foreach
- top
- treeAggregate
- treeReduce
- foreachPartition
- collectAsMap

- count
- takeSample
- max
- min
- sum
- histogram
- mean
- variance
- stdev
- sampleVariance
- countApprox
- countApproxDistinct

- takeOrdered

- saveAsTextFile
- saveAsSequenceFile
- saveAsObjectFile
- saveAsHadoopDataset
- saveAsHadoopFile
- saveAsNewAPIHadoopDataset
- saveAsNewAPIHadoopFile

# Essential Core & Intermediate PairRDD Operations



## General

- flatMapValues
- groupByKey
- reduceByKey
- reduceByKeyLocally
- foldByKey
- aggregateByKey
- sortByKey
- combineByKey

## Math / Statistical

- sampleByKey

## Set Theory / Relational

- cogroup (=groupWith)
- join
- subtractByKey
- fullOuterJoin
- leftOuterJoin
- rightOuterJoin

## Data Structure

- partitionBy

- 
- keys
  - values

- countByKey
- countByValue
- countByValueApprox
- countApproxDistinctByKey
- countApproxDistinctByKey
- countByKeyApprox
- sampleByKeyExact





# Spark architecture

# Important terms 1

## › Application master

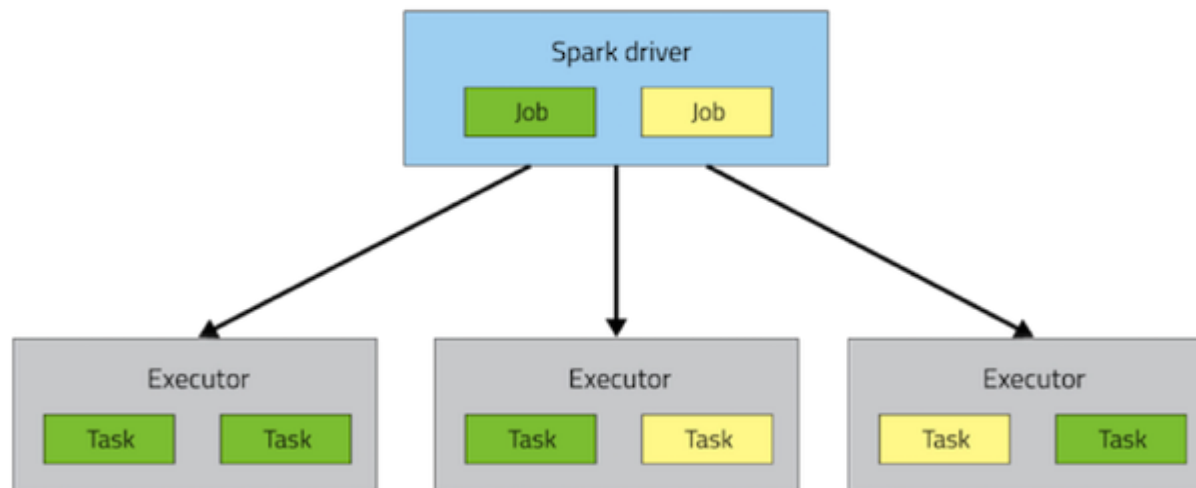
- process for negotiation of sources

## › Driver

- main process
- workflow planning
- distributing work to executors

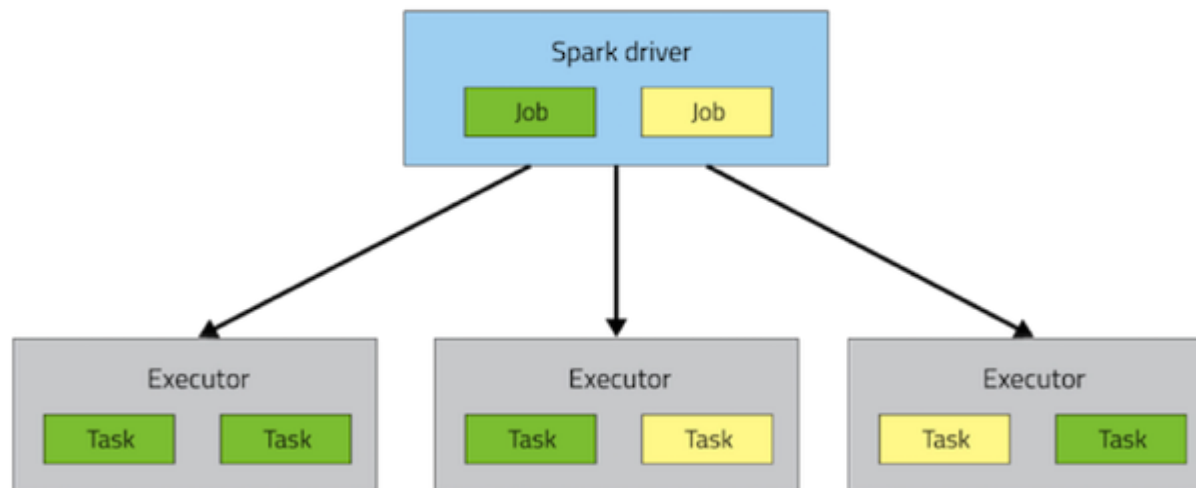
## › Executor

- process running on a node – JVM
- doing tasks (possibly in parallel if it has multiple cores)

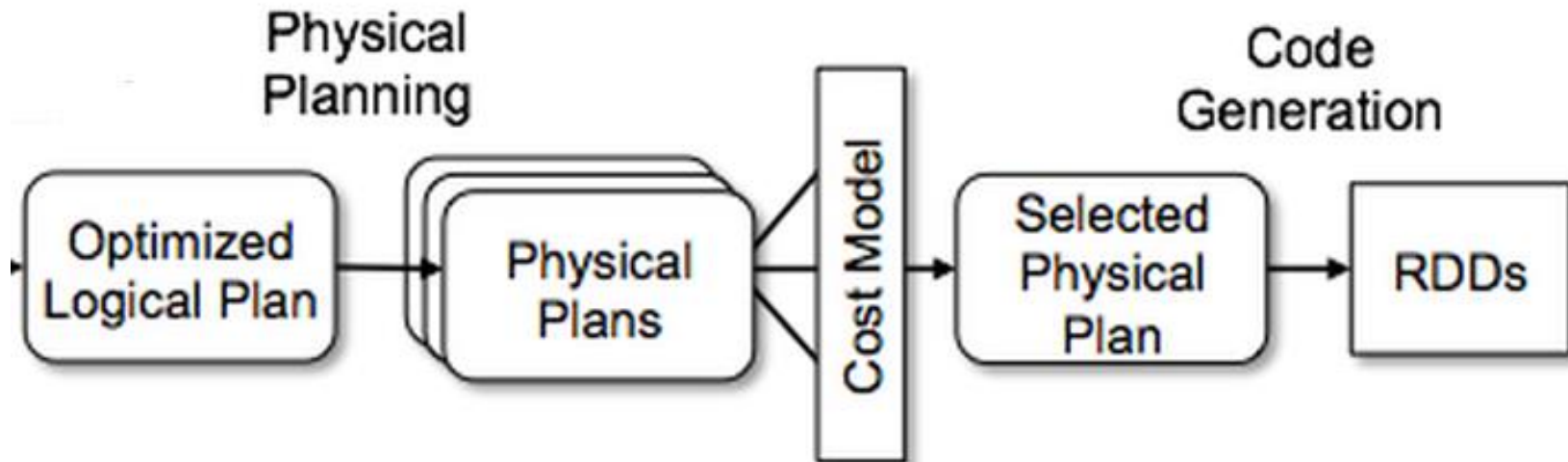
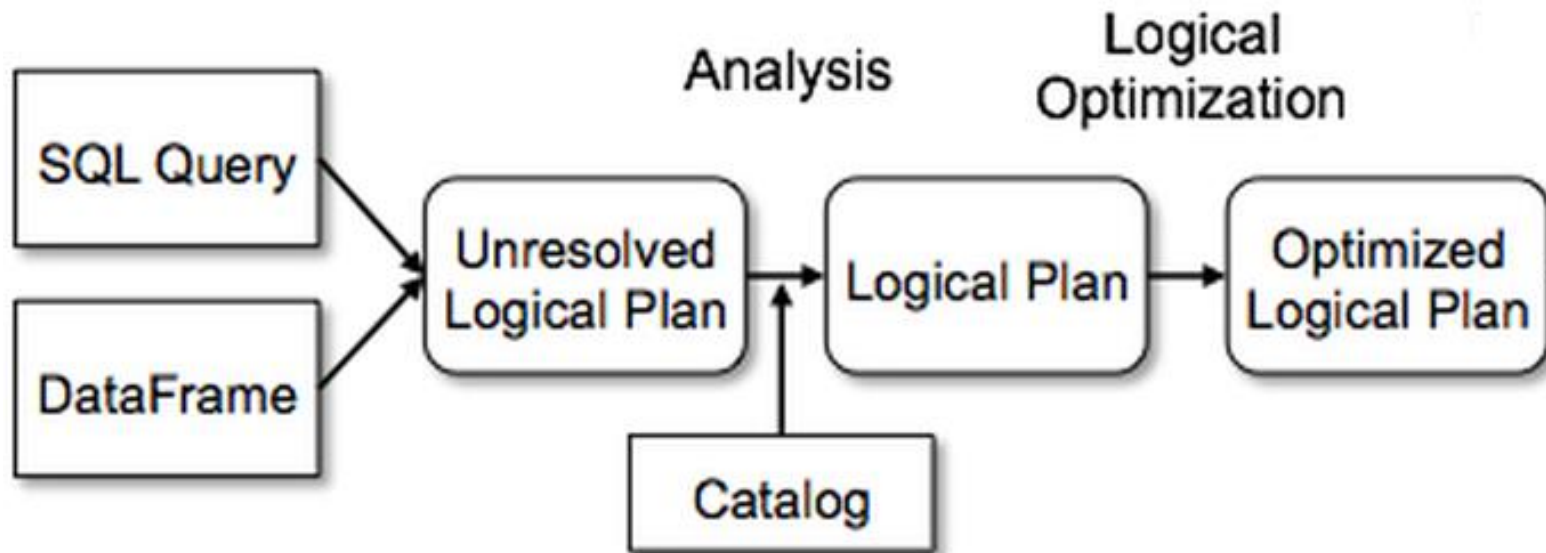


## Important terms 2

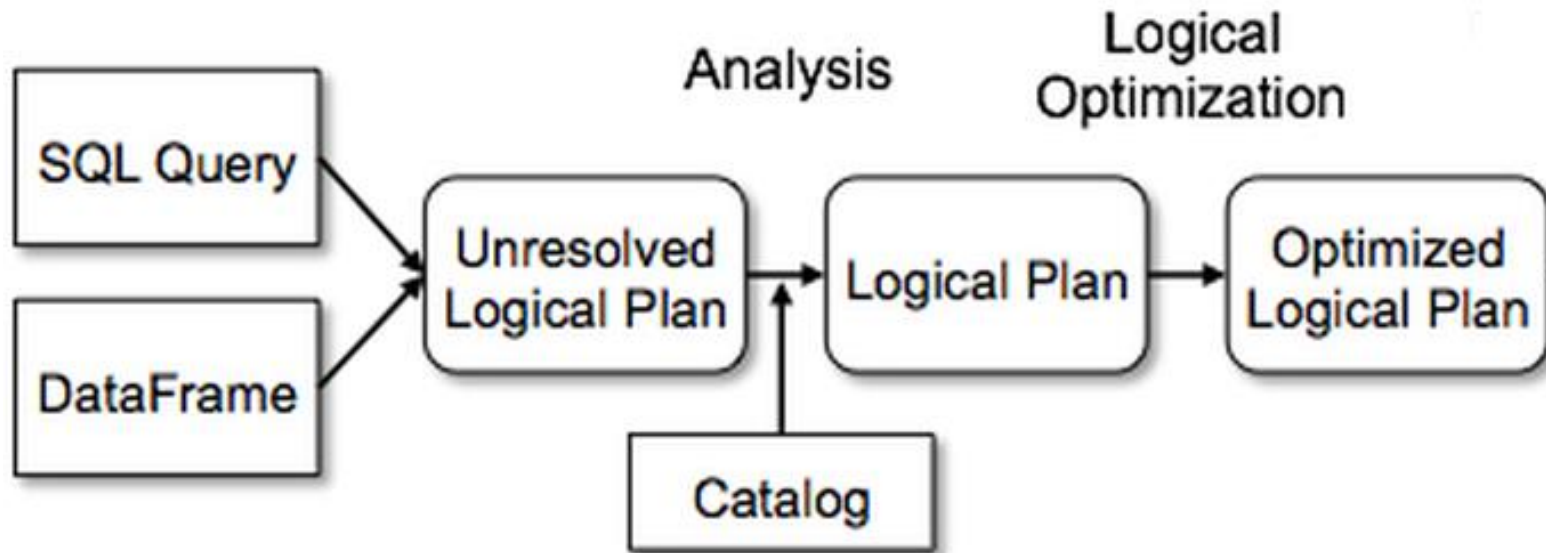
- › **Job**
  - part of application, driven by driver
- › **Stage**
  - set of transformations which can be done without a shuffle
- › **Task**
  - unit of work which is carried out by executor on some pice of data



# Planning and optimization

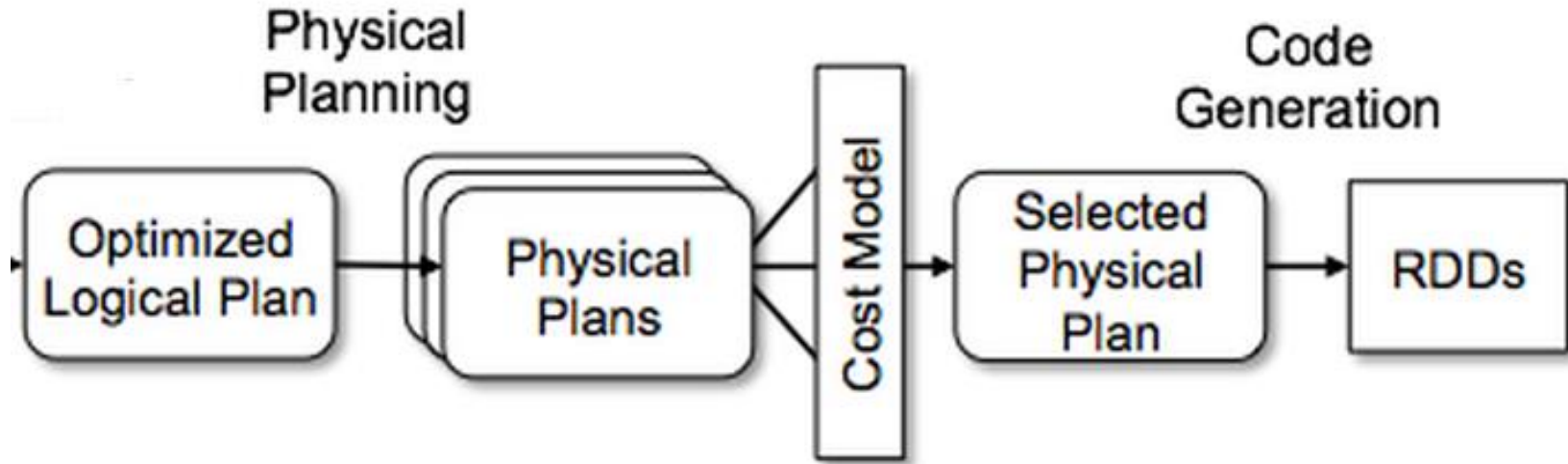


# Planning and optimization



- 
- › typecasting
  - › rearranging of transformations (e. g. swap FILTER and JOIN)
  - › JOIN type selection, exploiting clustering, partitions, data skew
  - › etc.

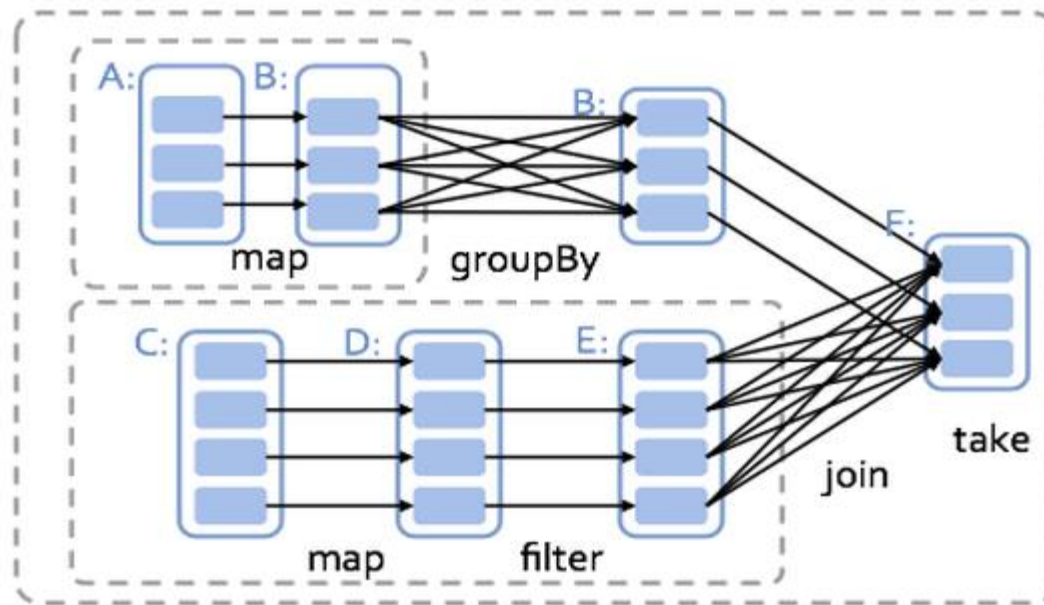
# Planning and optimization



- 
- › partitioning and data distribution
  - › translating of transformations and actions to JVM language
  - › etc.



# Example of graphic planning



- › DAG – directed acyclic graph
- › describes computation flow
- › dependencies (X must be done before Y)
- › optimisation with respect to dependencies

## Data partitioning

- › **partition** – piece of data processed in one task
- › by default 1 partition = 1 HDFS block = 1 task = 1 core
- › partition processed on the node where is saved
- › more partitions  $\Rightarrow$  more tasks  $\Rightarrow$  higher parallelism  $\Rightarrow$  smaller partition  $\Rightarrow$  lower efficiency  $\Rightarrow$  higher overhead
- › ... and the other way around

### Is it possible to set? And how?

- › data input: `sc.textFile(file, partition_number)`
- › at run: `coalesce, repartition, partitionBy`
- › Every repartitioning causes a shuffle!

# Launching and configuration

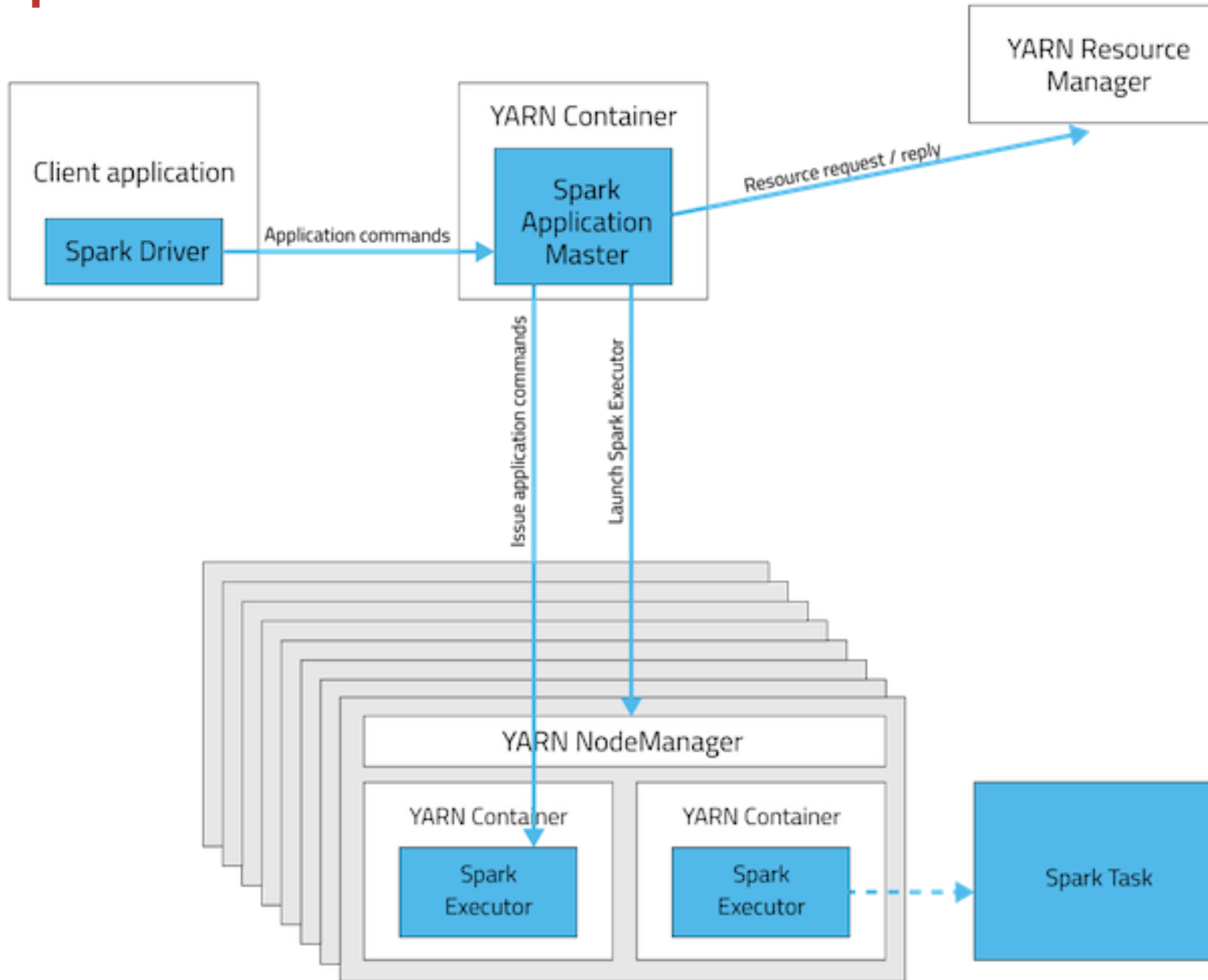
# Spark launching

`pyspark` | `spark-shell` | `spark-submit --param value`

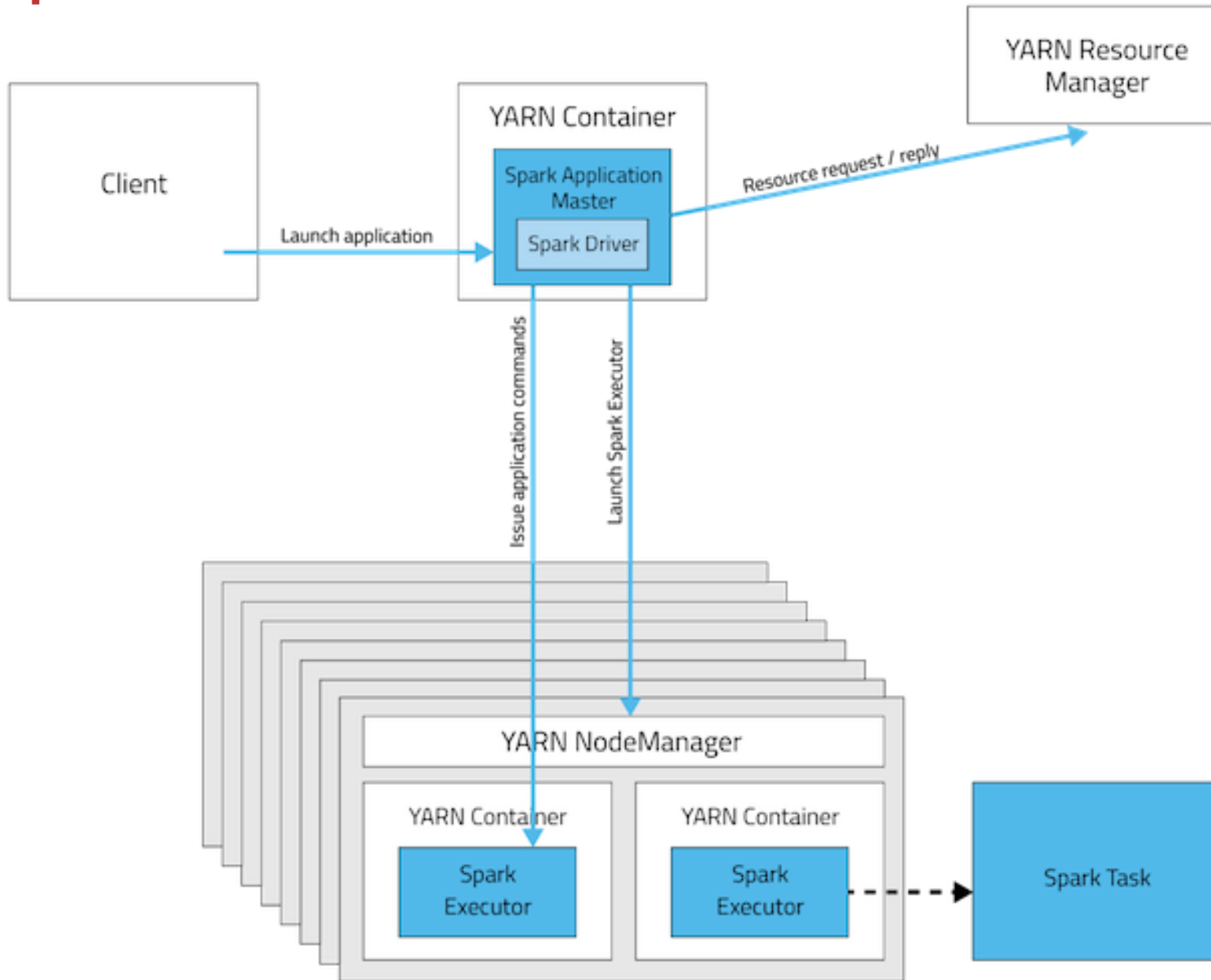
## Where and how it runs

- › on the cluster – full parallelism
  - mode client
  - mode cluster
- › locally – parallel run on multiple cores
- › determined by parameters `--master`, `--deploy-mode`

# Spark on YARN client mode



# Spark on YARN cluster mode



# Mode client and mode cluster

- › default mode = client
- › **client** – good for an interactive work and debugging (output to the local console)
- › **cluster** – good for production

# Spark configuration – requirement of sources

- › `--name`
- › `--driver-memory`
- › `--num-executors`
- › `--executor-cores`
- › `--executor-memory`

## Example

- › `pyspark --master yarn --deploy-mode client  
--driver-memory 1G  
--num-executors 3 --executor-cores 2  
--executor-memory 3G`



# Source allocation plan – example

## Generally:

- › `--num-cores <= 5`
- › `--executor-memory <= 64 GB`

## Let's say... cluster of 6 nodes, each 16 cores and 64 GB RAM

- › 1 core and 1GB per node → OS  
6 \* 15 cores and 63 GB remaining
- › 1 core for Spark Driver: 6 \* 15 – 1 = 89 cores remaining
- › 89 / 5 ~ 17 executors; each node approx. 3 executors
- › 63 GB / 3 ~ 21 GB memory for an executor
- › due to memory overhead → set 19 GB for an executor

# Thank for your attention

PROFINIT

Profinit, s.r.o.  
Tychonova 2, 160 00 Praha 6



Telefon  
+ 420 224 316 016



Web  
[www.profinit.eu](http://www.profinit.eu)



LinkedIn  
[linkedin.com/company/profinit](https://linkedin.com/company/profinit)



Twitter  
[twitter.com/Profinit\\_EU](https://twitter.com/Profinit_EU)