



# Large Scale Data Processing

## Lecture 4 – Spark

dr hab. inż. Tomasz Kajdanowicz, Piotr Bielak, Roman  
Bartusiak, Krzysztof Rajda

November 18, 2021



# Overview

Spark - batch



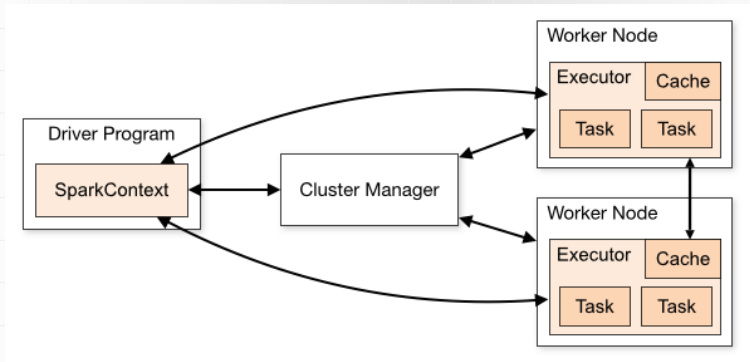
# General

## Spark - batch

- ▶ University of California
- ▶ UC Berkeley AMPLab
- ▶ Matei Zaharia PhD Thesis
- ▶ Huge community
- ▶ JVM

# Architecture

## Spark - batch



# Architecture

## Spark - batch

- ▶ one Driver - many Workers
- ▶ Each application in separate JVM
- ▶ Driver needs to be accesible from workers

# Architecture

## Spark - batch

- ▶ Application - our main()
- ▶ Driver - executes our main(), schedules DAG
- ▶ Executor - each worker spawns executors in order to run tasks of our application

### Tune executors per worker

- ▶ too small executors - unnecessary overhead
- ▶ too big executors - IO issues, failure recovery issues
- ▶ keep balance
  - ▶ 5 cores per executor?
  - ▶ leave core for IO (HDFS, Lustre, ...)
  - ▶ leave resources for application manager and other overheads



# DAG

## Spark - batch

- ▶ Directed Acyclic Graph
- ▶ Represents computations
- ▶ We are not doing computations in our code, we are creating DAGs and executing them

# Architecture

## Spark - batch

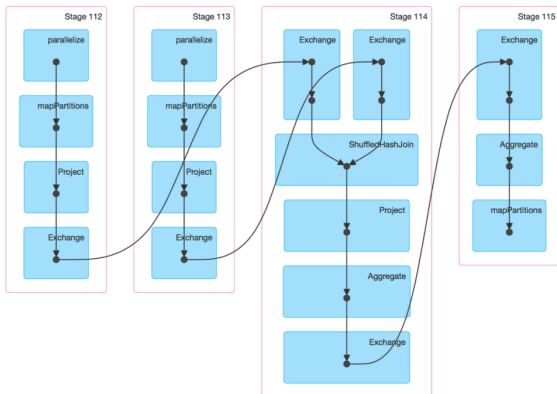
### Details for Job 8

Status: SUCCEEDED

Completed Stages: 4

▶ Event Timeline

▼ DAG Visualization



# Architecture

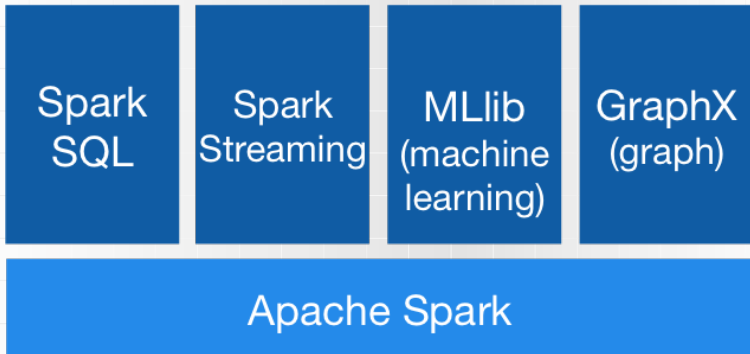
## Spark - batch

### Application:

- ▶ Job
  - ▶ Stage (eg. map)
    - ▶ Task
    - ▶ Task
    - ▶ ...
  - ▶ Stage (eg. reduce)
  - ▶ ...
- ▶ Job
- ▶ ...

# Spark components

Spark - batch



# Spark - batch

## Spark - batch

- ▶ Gather data for some time
- ▶ get a set that is limited and bounded
- ▶ process whole set at once (does not mean that processing will happen by loading everything to the memory etc.)

# Runtimes

## Spark - batch

- ▶ Standalone
- ▶ YARN
- ▶ Mesos
- ▶ Kubernetes

# Runtimes - Standalone

Spark - batch

Used in WCSS (utilizing pdsdsh)

Simply:

- ▶ Put up master
- ▶ Take master address
- ▶ Put up nodes using master address

# Runtimes - YARN

## Spark - batch

- ▶ Comes from Hadoop
- ▶ Primarily only for Hadoop scheduling
- ▶ MapReduce V2



# Runtimes - MESOS

## Spark - batch

- ▶ UC Berkeley
- ▶ Used by Twitter, AirBnB...
- ▶ Full abstraction over resources

# Runtimes - Myriad

Spark - batch

- ▶ Mesos + YARN on same infrastructure
- ▶ YARN running in Mesos

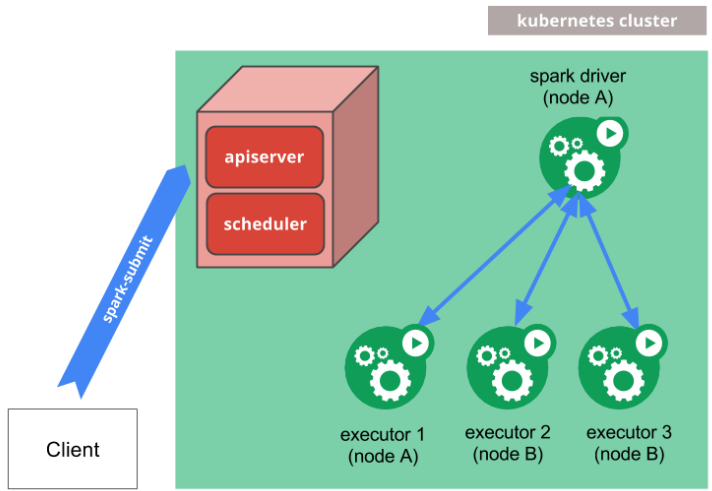
# Runtimes - Kubernetes

## Spark - batch

- ▶ Spark  $\geq$  2.3
- ▶ Kubernetes  $\geq$  1.6

# Runtimes - Kubernetes

## Spark - batch



# Core - primitives

Spark - batch

RDD - resilient distributed dataset

- ▶ HDFS
- ▶ Hadoop API
- ▶ Directly from collections

# Core - primitives

Spark - batch

## Accumulators

- ▶ Shared variables between executors
- ▶ Only add - efficient

# Core - primitives

## Spark - batch

### Broadcast variables

- ▶ Efficient way to distributed read-only data between executors
- ▶ Spark optimizes communication in order to minimize overhead
- ▶ Reduces overhead when data reused between stages

# Core - data partitioning

## Spark - batch

- ▶ Each RDD is divided into partitions
- ▶ Each partition is processed by single executor
- ▶ You should have at least equal number of partitions as the number of CPUs in a cluster (taking into account data set size)



# Core - data partitioning

## Spark - batch

- ▶ Too big partitions - memory issues
- ▶ Too many partitions in comparison to data set size - performance issues
- ▶  $2-3 * \text{numCores}$  of partitions (depends on data set size)
- ▶ For big data sets - increase the number of partitions

# Core - shuffling

## Spark - batch

- ▶ repartitioning
- ▶ expensive
- ▶ disk I/O
- ▶ network I/O
- ▶ serialization/deserialization

# Core - data transformations

## Spark - batch

### Narrow

- ▶ Does not require data shuffling
- ▶ map, filter ...
- ▶ Spark groups narrow transformations - pipelining

### Wide

- ▶ Requires data shuffling
- ▶ groupByKey, ...

# Core - data transformations

Spark - batch

Remember about load balancing!

- ▶ Narrow operations will not cause shuffling
- ▶ Without shuffling data can get skewed
- ▶ Skewed data -> performance problems
- ▶ repartition manually

# Core - reduceByKey, combineByKey, ... vs groupByKey

Spark - batch

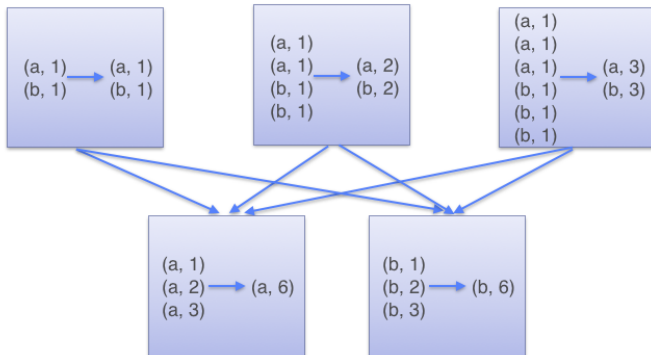
reduceByKey, combineByKey, foldByKey decrease data size that needs to be

- ▶ saved to disk
- ▶ sent over network
- ▶ serialized
- ▶ deserialized

# Core - reduceByKey

Spark - batch

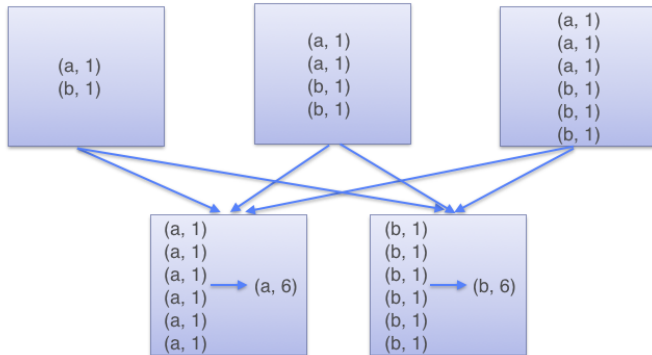
## ReduceByKey



# Core - groupByKey

Spark - batch

## GroupByKey



# Core - persistence

## Spark - batch

- ▶ operations are lazy
- ▶ we need to persist operations results in order to reuse it



```
1 rdd.persist()  
2 \\ or  
3 rdd.cache()
```

- ▶ can increase performance up to 10x



# Core - persistence

## Spark - batch

- ▶ supports Kryo serialization
- ▶ multiple storage levels
- ▶ data can be compressed
- ▶ off-heap memory support

# Core - must remember

## Spark - batch

- ▶ all functions passed to transformations/actions are serialized
- ▶ there is no magic communication mechanism
- ▶ code is executed on variables copies
- ▶ no sync of those var changes

# Core - data actions

## Spark - batch

Materializes results as requires exact rows with data

- ▶ take
- ▶ count
- ▶ first
- ▶ foreach
- ▶ ...

# SparkSQL

## Spark - batch

- ▶ is a module in Apache Spark that integrates relational processing with Spark's functional programming API.
- ▶ lets Spark programmers leverage the benefits of relational processing (e.g., declarative queries and optimized storage), and lets SQL users call complex analytics libraries in Spark (e.g., machine learning)

# SparkSQL

## Spark - batch

- ▶ utilize Spark CORE
- ▶ Represents structured and semistructured data
- ▶ Use
  - ▶ SQL/HiveQL
  - ▶ DataSet API

# SparkSQL - SQL

## Spark - batch

```
1 // Register the DataFrame as a SQL
2 //temporary view
3 df.createOrReplaceTempView("people")
4
5 val sqlDF = spark.sql("SELECT * FROM people")
6 sqlDF.show()
7 // +----+-----+
8 // | age|   name|
9 // +----+-----+
10 // |null|Michael|
11 // | 30|   Andy|
12 // | 19|  Justin|
13 // +----+-----+
14
```

# SparkSQL - DataSet API

## Spark - batch

```
1 case class Person(name: String, age: Long)
2
3 // Encoders are created for case classes
4 val caseClassDS = Seq(Person("Andy", 32))
5                       .toDS()
6 caseClassDS.show()
7 // +-----+-----+
8 // |name|age|
9 // +-----+-----+
10 // |Andy| 32|
11 // +-----+-----+
12
```

# SparkSQL - DataSet API

Spark - batch

```
1 // Encoders for most common types are automatically
2 // provided by importing spark.implicit._
3 val primitiveDS = Seq(1, 2, 3).toDS()
4 primitiveDS.map(_ + 1).collect() // Returns: Array
5   (2, 3, 4)
```



# SparkSQL - DataSet API

## Spark - batch

```
1 // DataFrames can be converted to a Dataset by
2 // providing a class. Mapping will be done by name
3 val path = "examples/src/main/resources/people.json"
4 val peopleDS = spark.read.json(path).as[Person]
5 peopleDS.show()
6 // +----+-----+
7 // | age|   name|
8 // +----+-----+
9 // |null|Michael|
10 // | 30|   Andy|
11 // | 19|  Justin|
12 // +----+-----+
13
```

# SparkSQL - DataSet API

Spark - batch

```
1 val teenagers = peopleDS.where('age >= 10)
2   .where('age <= 19)
3   .select('name).as[String]
4 teenagers.show
5 //      +-----+
6 //      | name |
7 //      +-----+
8 //      |Justin|
9 //      +-----+
10
```

# SparkSQL - DataSet API

Spark - batch

```
1 val symbol = 'someSymbol  
2 // symbol: Symbol = 'someSymbol '  
3
```

# MLLib

## Spark - batch

- ▶ API mix
- ▶ Features:
  - ▶ data loading
  - ▶ data processing
  - ▶ ml methods
  - ▶ ...

# MLLib - DataFrames API

Spark - batch

- ▶ pipelines
- ▶ friendly
- ▶ optimizations
- ▶ uniform

# MLLib

## Spark - batch

- ▶ classification
  - ▶ binary
  - ▶ multi-class
  - ▶ multi-label
- ▶ regression
- ▶ clustering
- ▶ collaborative filtering
- ▶ frequent-pattern mining

# MLLib

## Spark - batch

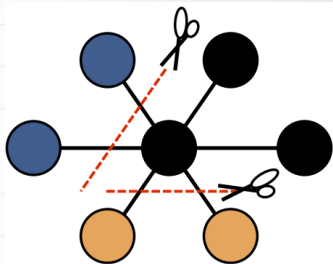
- ▶ hyper-parameters search
- ▶ cross validation
- ▶ train-test split

# GraphX

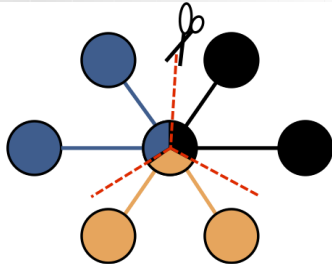
Spark - batch

- ▶ graph representations on spark
- ▶ based on RDD API
- ▶ a little of graphs algorithms





Edge Cut



Vertex Cut

# GraphFrames

Spark - batch

- ▶ based on DataFrame API
- ▶ should be faster than GraphX
- ▶ smaller API
- ▶ in some places, use GraphX under the hood

# spark-shell

Spark - batch

- ▶ shell for Spark
- ▶ created context
- ▶ spark API imported

# Zeppelin

## Spark - batch

- ▶ notebooks for Spark
- ▶ create, or connect to remote context
- ▶ Helium for visualization
- ▶ collaboration
- ▶ scheduler
- ▶ custom dependencies

# Spark notebooks

Spark - batch

- ▶ more like iPython notebooks
- ▶ more built-in visualizations

# Spark not for everything

Spark - batch

- ▶ Graph structured data
- ▶ edge list with data on 'from' node

# Spark not for everything

## Spark - batch

- ▶ Use spark to find  $n$  level neighbours of node
- ▶ find node by name to get identifier  $id$
- ▶ find neighbours of node  $id$
- ▶ collect identifiers
- ▶ repeat, until get  $n$  levels
- ▶ collect all data
- ▶ query for data objects for each node

# Spark not for everything

## Spark - batch

- ▶ Directly call DB using recursive approach
- ▶ Do calls in parallel using Cats-effects
- ▶ Batch the queries in order to optimize the parallelism
- ▶ utilize stream based processing



# Next week

Spark - batch

▶ Spark streaming



# Large Scale Data Processing

## Lecture 4 – Spark

dr hab. inż. Tomasz Kajdanowicz, Piotr Bielak, Roman  
Bartusiak, Krzysztof Rajda

November 18, 2021