

THE DATA SCIENCE LAB

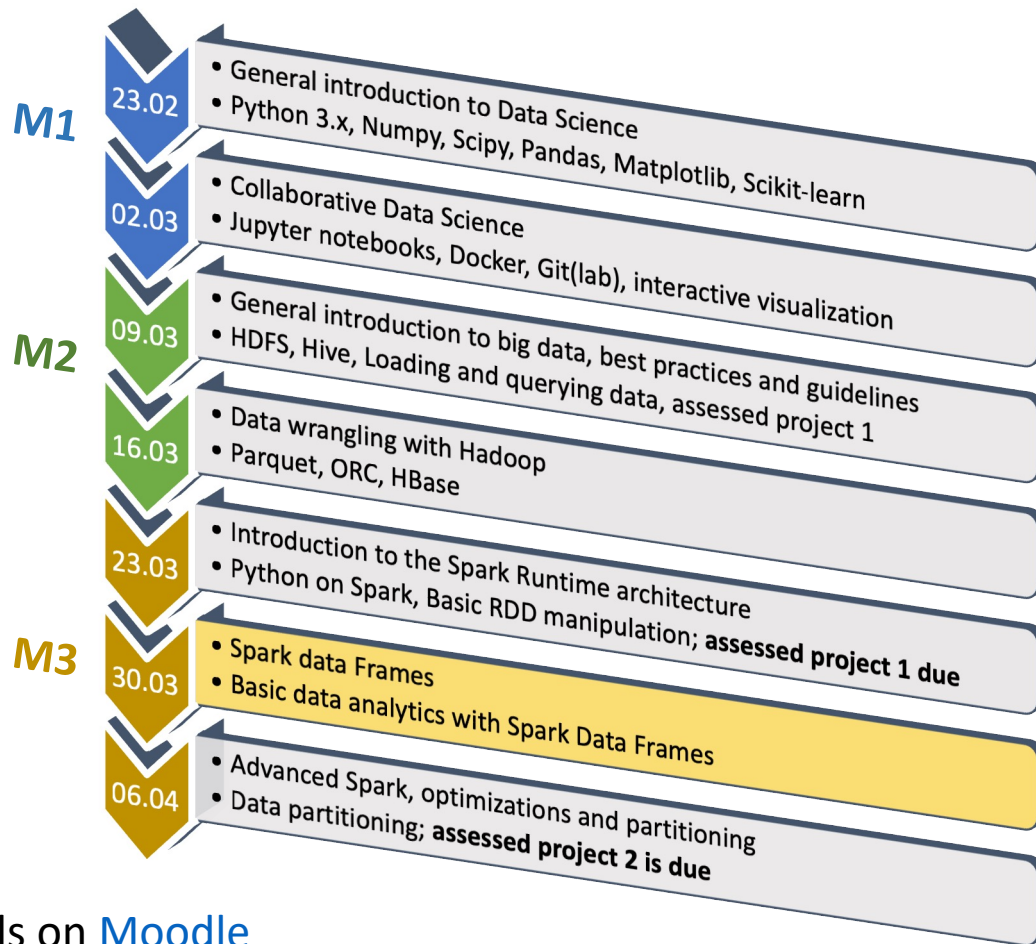
Data Wrangling with Hadoop

COM 490 – Spring 2022

Week 6

THIS CLASS WILL BE RECORDED

Agenda Spring 2022



*Details on [Moodle](#)

Week 5 Recap

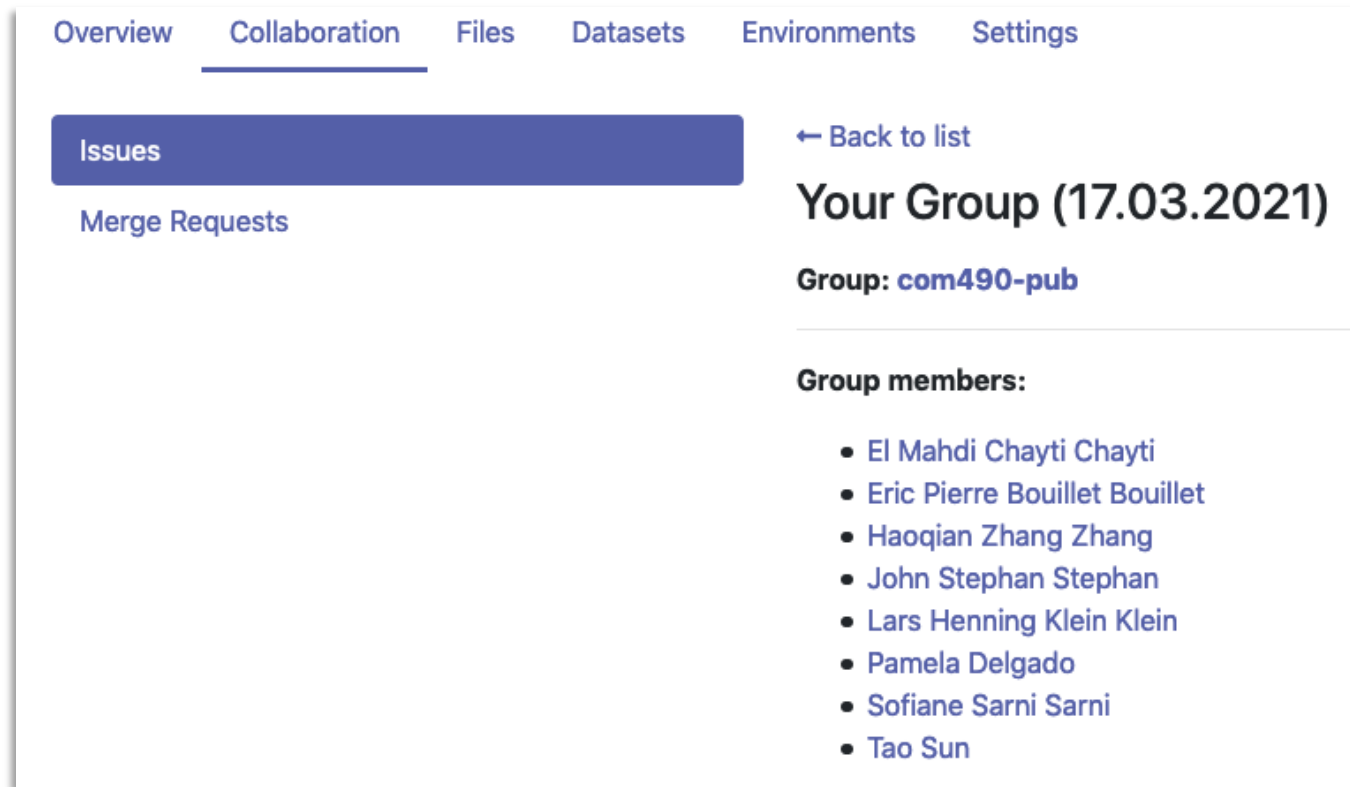
- What is Spark?
- RDDs
 - Immutable
 - Resilient
 - Lineage
- Operations on RDDs
 - Transformations
 - Actions
 - Lazy execution
- Exercises to get started using Gutenberg corpus

Week 5 – Questions?

Week 6 – Homework 2 : Checkpoint

Have you seen this in your homework project ?

If not, contact us !



The screenshot displays the DVC Collaboration web interface. At the top, there are navigation tabs: Overview, Collaboration (which is selected and underlined), Files, Datasets, Environments, and Settings. On the left side, under the Collaboration tab, there are two buttons: 'Issues' (highlighted in a dark blue box) and 'Merge Requests'. The main content area on the right shows a '← Back to list' link, followed by the title 'Your Group (17.03.2021)'. Below this, it states 'Group: com490-pub'. A section titled 'Group members:' follows, containing a bulleted list of eight names: El Mahdi Chayti Chayti, Eric Pierre Bouillet Bouillet, Haoqian Zhang Zhang, John Stephan Stephan, Lars Henning Klein Klein, Pamela Delgado, Sofiane Sarni Sarni, and Tao Sun.

Today's Agenda

- Introduction to Data Frames
- DataFrame demo
- DataFrames and PySpark under the hood
- Exercises week 6
 - Guttenberg corpus part II

Introduction to Spark DataFrames

RDD Revisited

- Resilient (lineage)
- Distributed (partitioned)
- Unstructured (mostly), rows or key-row pairs
- Type safe
 - Scala's compiler optimization
- Use of lambda functions
- Fine grained control – tell spark how to transform a data
 - low level – more responsibility to the programmer:
 - Decide transformations and actions
 - Which part of the data
 - In what order

Spark DataFrames

- Distributed Collections of Data

- Organized into rows of named columns

Structured

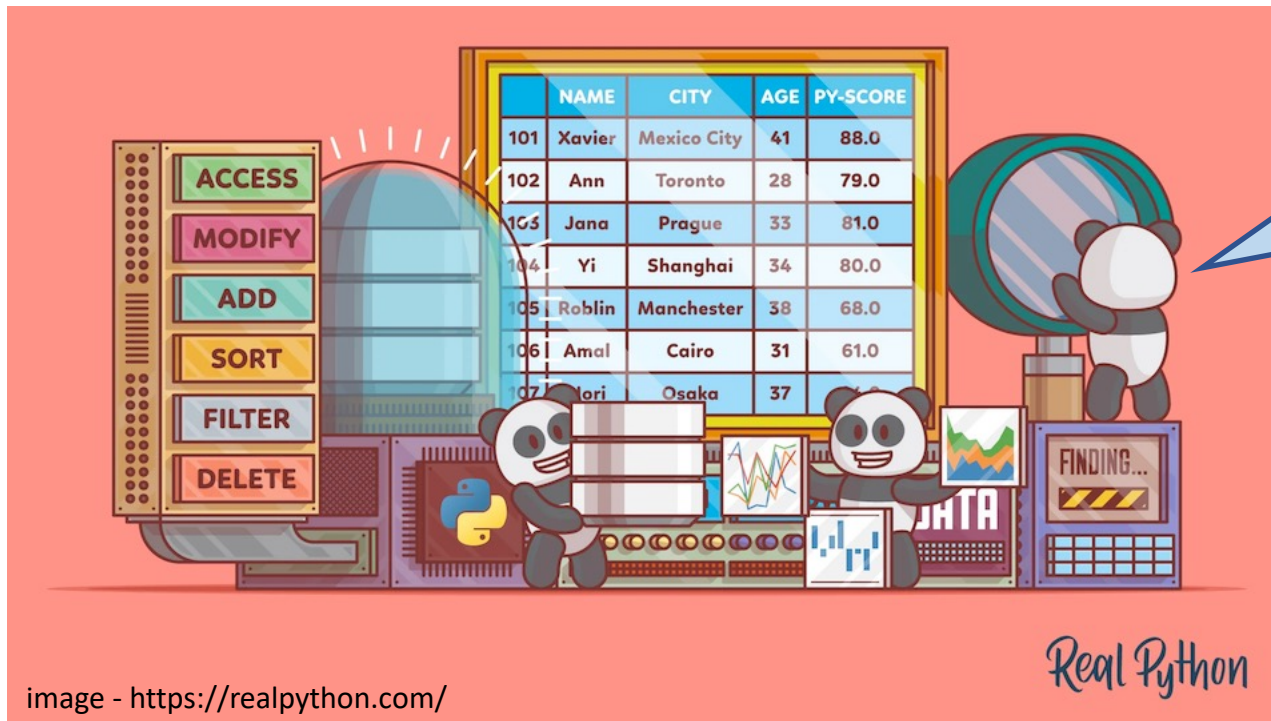
- Very much like relational database Tables

- Optimized for relational-type of queries on tables (logical plan optimization)

Col 1	Col 2	Col 3
1	a	10:00
2	b	11:00
3	c	12:00
4	d	13:00
5	e	14:00

Origin of Data Frames

Spark Data Frame API Inspired by R and Python's Pandas



Want to join the
Big Data party!

What are Spark Data Frames

- Inspired by R and Python Pandas

SOURCE

- Local File Systems
- Distributed File Systems (HDFS)
- Cloud Storage (S3)
- External data bases
- Spark RDD

DATA FORMAT – out of the box

- TEXT
 - JSON
 - CSV
 - Parquet
 - ORC
 - Hive Table
- + **Other with plugins** (Avro, ElasticSearch, Cassandra, ...)

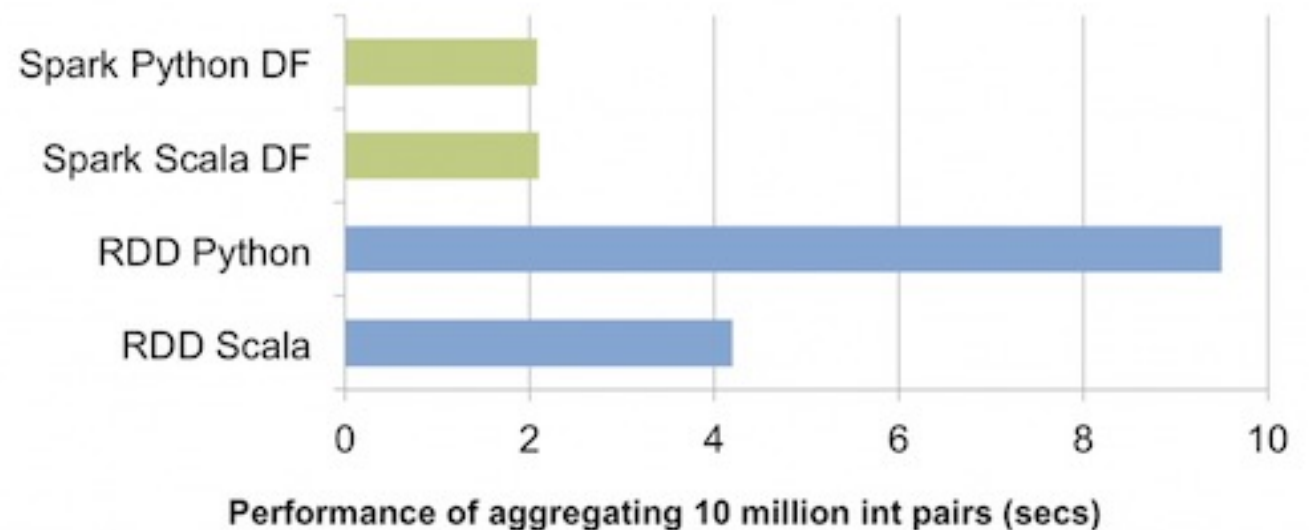
- Parallelism & query optimizer, unlike R and Python

Why Spark Data Frames – RDD vs DataFrame

Resilient Distributed Dataset (RDD)	Data Frame
Structured & unstructured data	Structured data (table, named columns)
Schema must be declared manually	Auto discovery of file schema
Lambda functions (map, reduce)	Declarative, almost as SQL queries
Lower level language	Higher level language
No built-in other than generic compiler optimization. Must be done manually	Execution optimization
Type safety at compile time	Type safety at run-time (e.g. trying to access a non-existing column)

Spark DataFrame performance

- In DataFrame data is optimally managed off-heap (off JVM)
 - No need for Java/Scala (de)serialization when accessing object, unlike RDD
 - Avoid garbage collection, unlike RDD
- Aggregation (group by) is harder and not as efficient with RDD. In comparison exploration analysis is quick and easier on large DataFrame



[databrkick blog 2015](#)

(*) to be taken with a grain of salt

Which one should I use? RDD vs DataFrame

- Use RDD for operations that requires low level functionalities and control on unstructured data
- Use DataFrame for high level (SQL like) operations on structured data

DataFrame under the hood

PySpark

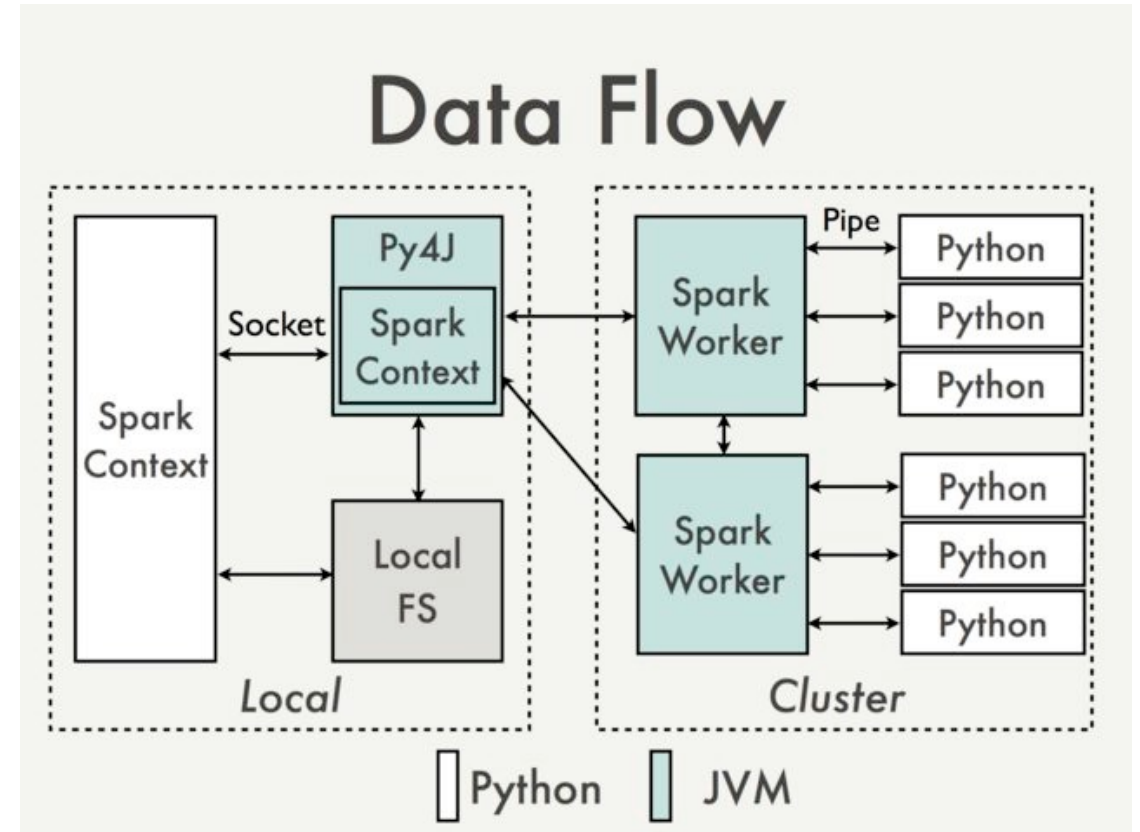
- PySpark - Python front end API for Spark



- Interface RDDs with Python
- Py4J – python library to dynamically access JVM objects
- Compatible with
 - PySparkSQL – SQL query library for DataFrame
 - MLib – Machine learning library
 - GraphFrames – Graph processing based on DataFrames (Graphx is on RDDs)

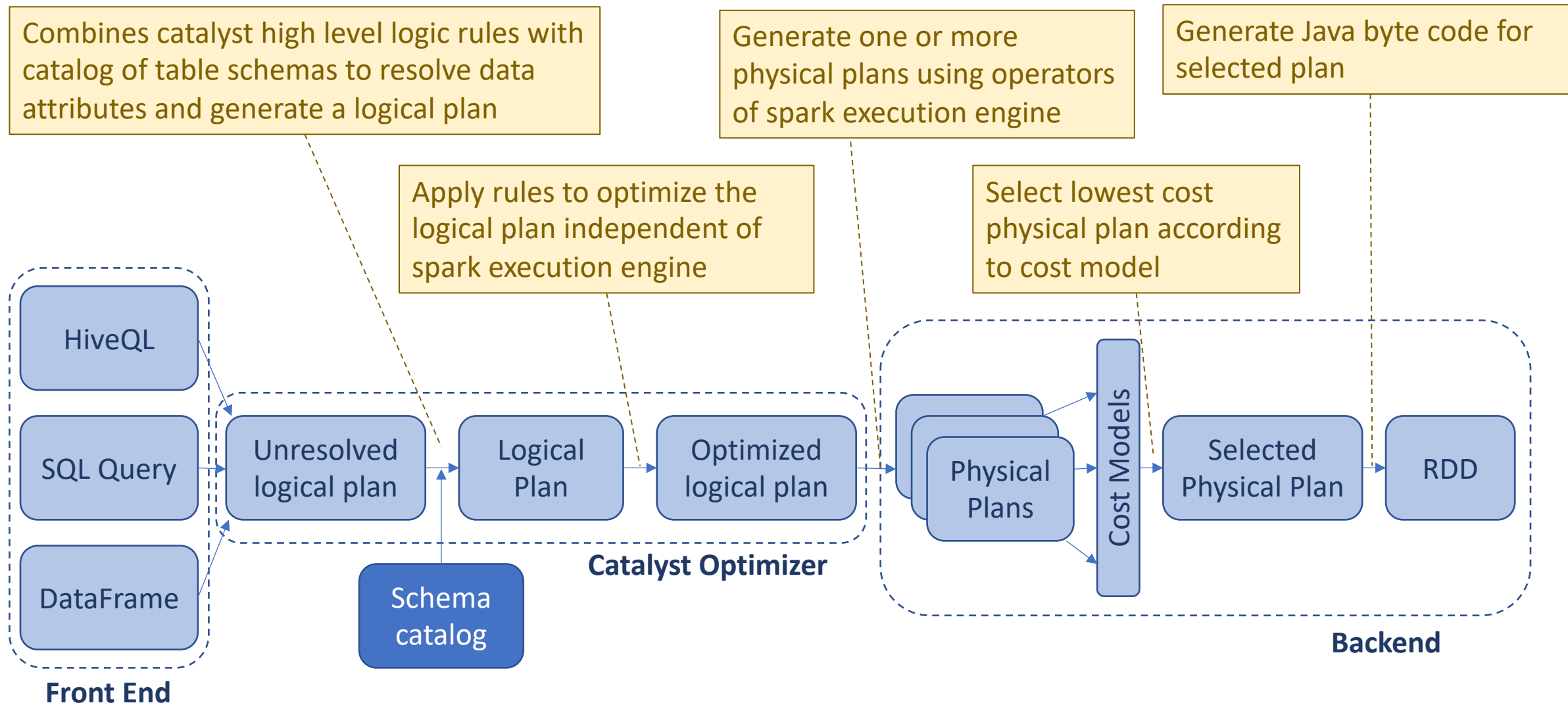
PySpark

- Spark workers pull data from source into JVM
- Data is actually processed into python subprocesses
- (de)serialization and streaming at every step



[PySpark internal](#)

Catalyst Optimizer



A parting word on Spark DataSets

- DataSet = extensions of DataFrame with convenience of RDD.
 - Strong type safety
 - RDD with Spark SQL optimized execution engine
 - Operate on serialized data (no deserialization overhead)
- Available only on Scala and Java
 - Since 1.6 DataFrame on Scala and Java are alias for DataSet[row]

Useful References

- Spark docs <http://spark.apache.org/docs/latest>
- DataFrames and code generation
<https://medium.com/virtuslab/spark-sql-under-the-hood-part-i-26077f85ebf0>
- Python Spark DataFrames starter documentation
https://spark.apache.org/docs/latest/api/python/getting_started/quickstart_df.html
- Spark MLlib guide <https://spark.apache.org/docs/latest/ml-guide.html>

Start your engines

<https://dslab2022-renku.epfl.ch/projects/com490/lab-course>