

Intro to PySpark Workshop

Garren Staubli
Sr. Data Engineer
@gstaubli



#PySparkWorkshop Resources: garrens.com/pyspark124



Do I know what I'm talking about?

Working with Spark since 2015

- Batch analytics in Spark + Hive, Pig and Hadoop MapReduce
- Real-time big data reporting using Spark/Impala/CDH
- Spark Structured Streaming + ML apps for real-time decision making

 **50+** answers on
stackoverflow
for Spark



Main Points

- Apache Spark
- Sample App Walkthrough
- Interactive Azure Jupyter Notebook
- Python-specific Spark advice
- Resources to continue learning



About Apache Spark

Apache Spark™ is a fast and general engine for large-scale data processing.

Lightning-fast cluster computing

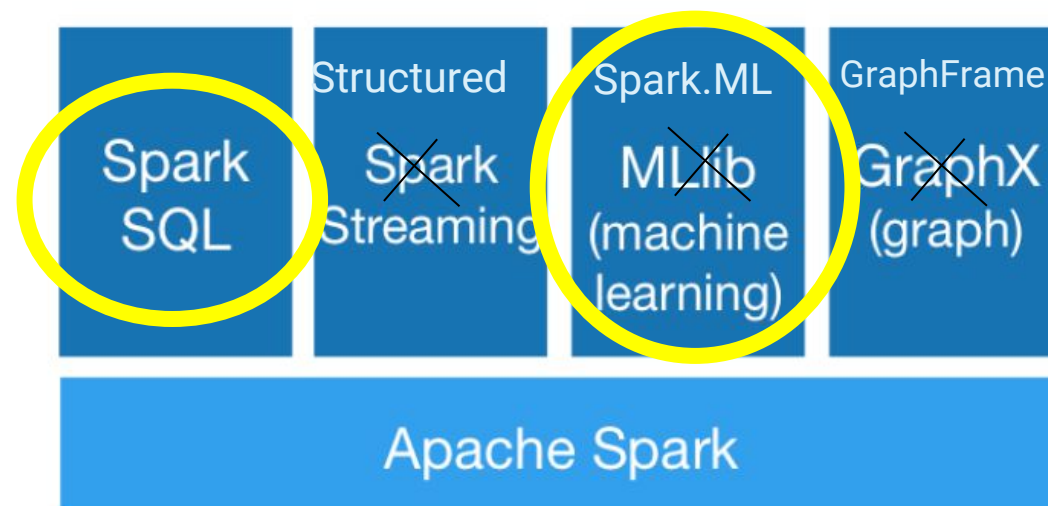
Lazily Evaluated

- Transforms vs Actions

Immutable

Transformations <i>(lazy)</i>	Actions
orderBy	show
filter	count
groupBy	take
select	collect
drop	save
join	

Transformations contribute to a query plan,
but nothing is executed until an action is called

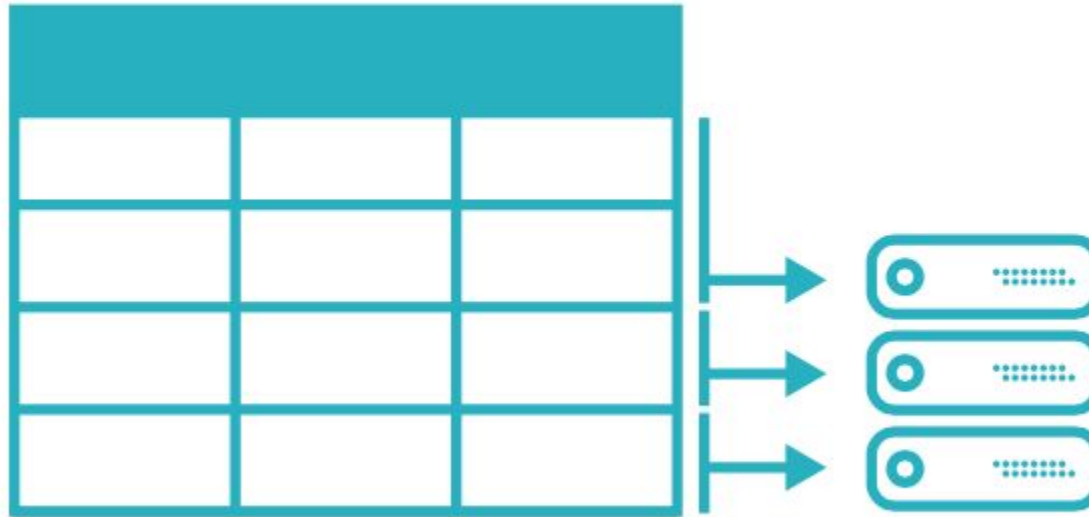


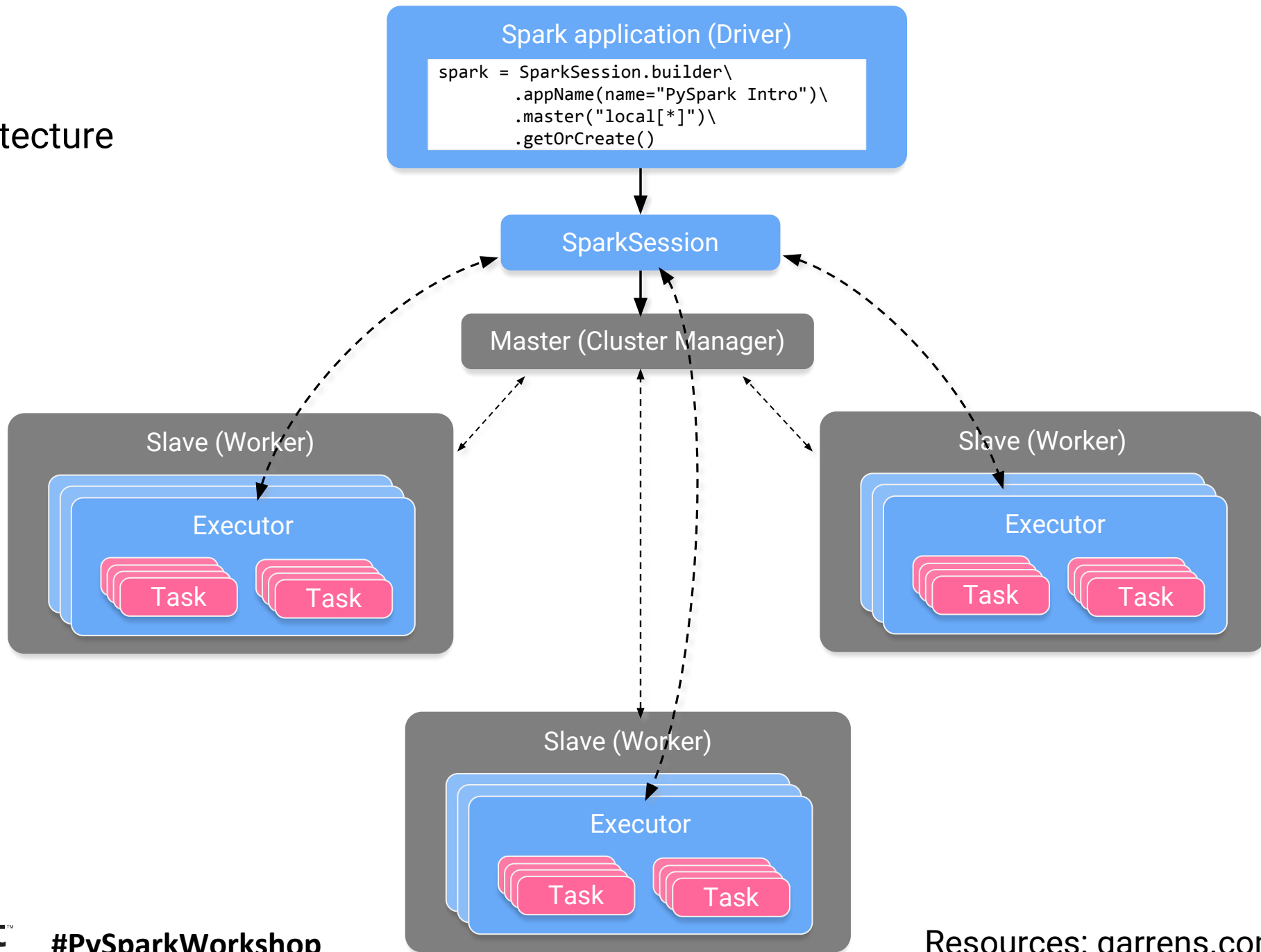
About Apache Spark

Spreadsheet on a single machine



Table or DataFrame partitioned across servers in a data center





About Apache Spark | Spark SQL



is **not** about SQL
is about **more** than
SQL

About Apache Spark | 2 Kinds of Actions

2 kinds of Actions



VS



distributed

driver

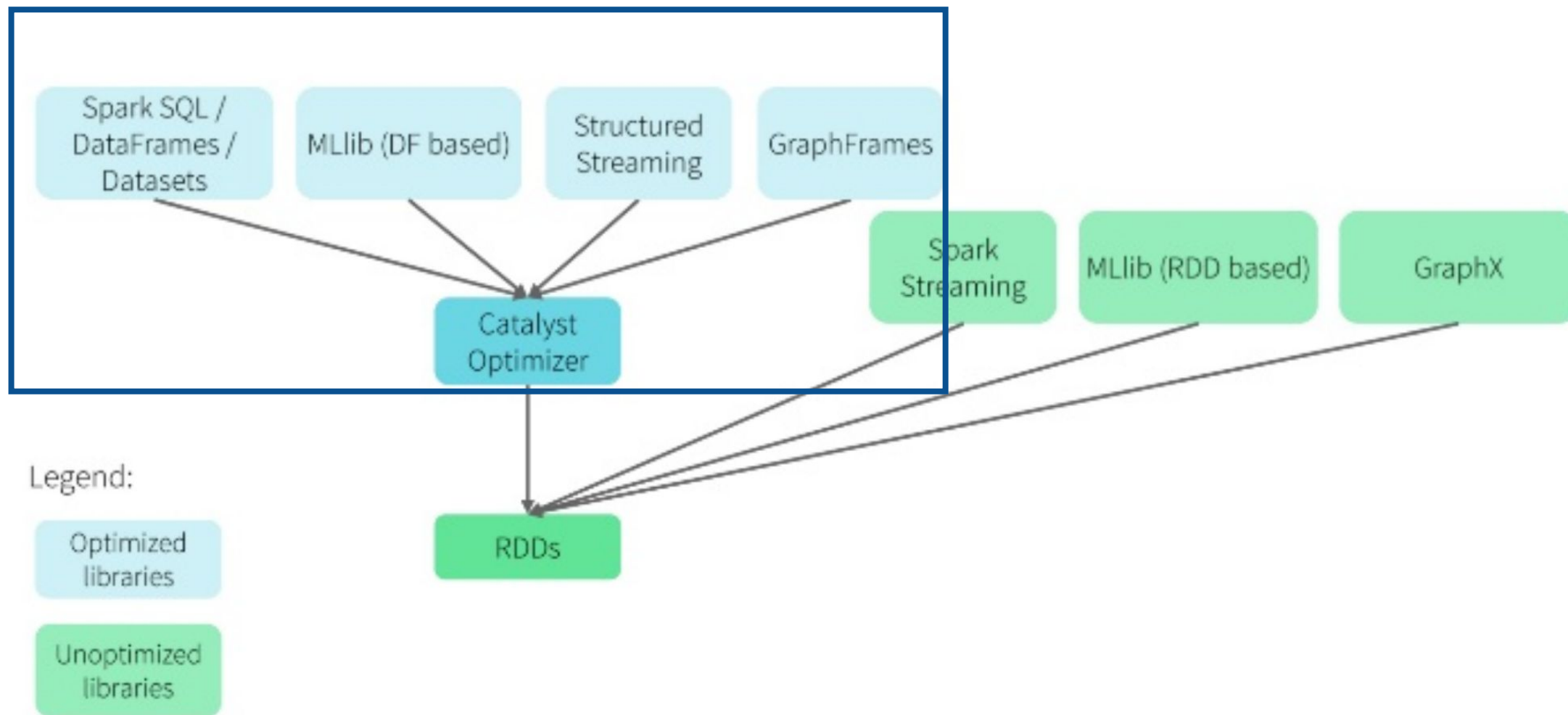
occurs across the cluster

result must fit in driver JVM

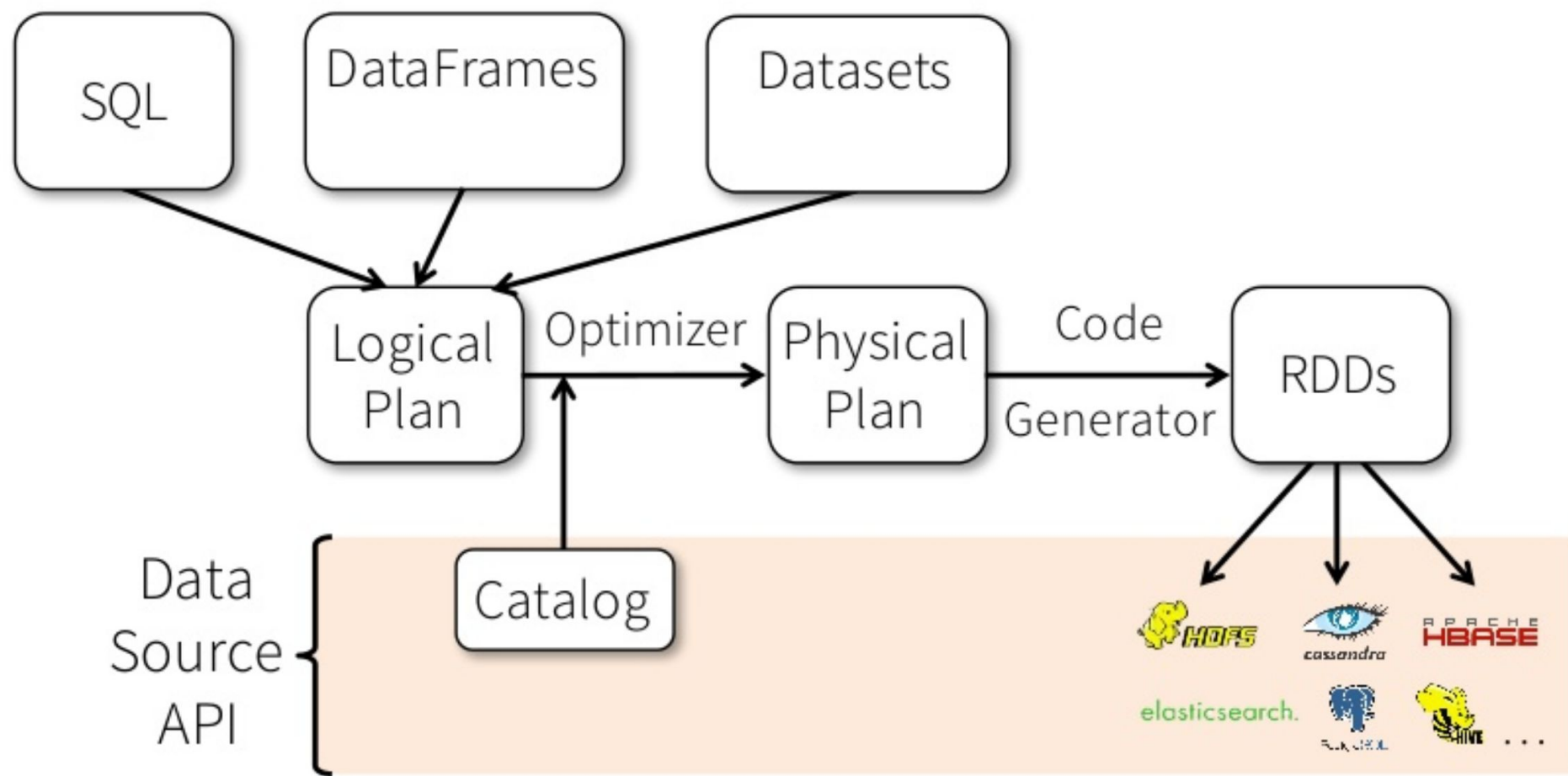
saveAsTextFile, (HDFS, S3, SQL, NoSQL, etc.)

collect, count, reduce, take, show..

About Apache Spark | Modern vs Legacy

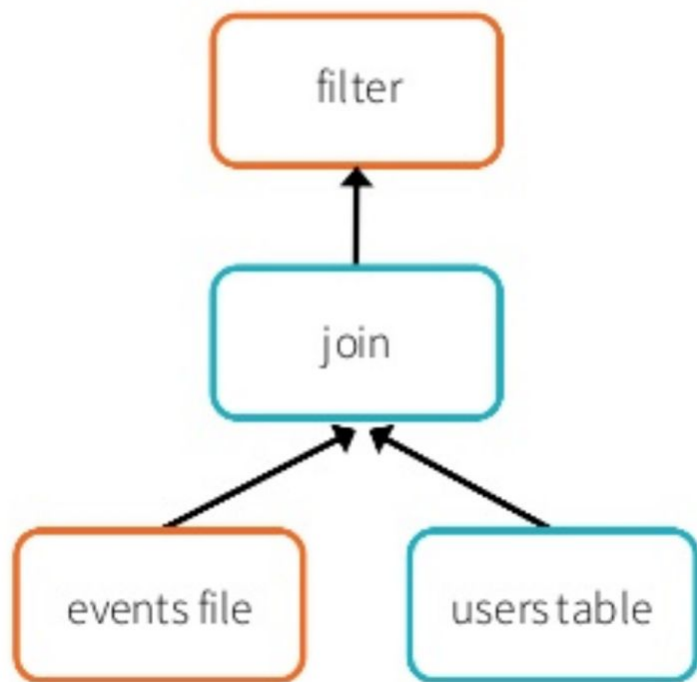


About Apache Spark | Modern Optimization

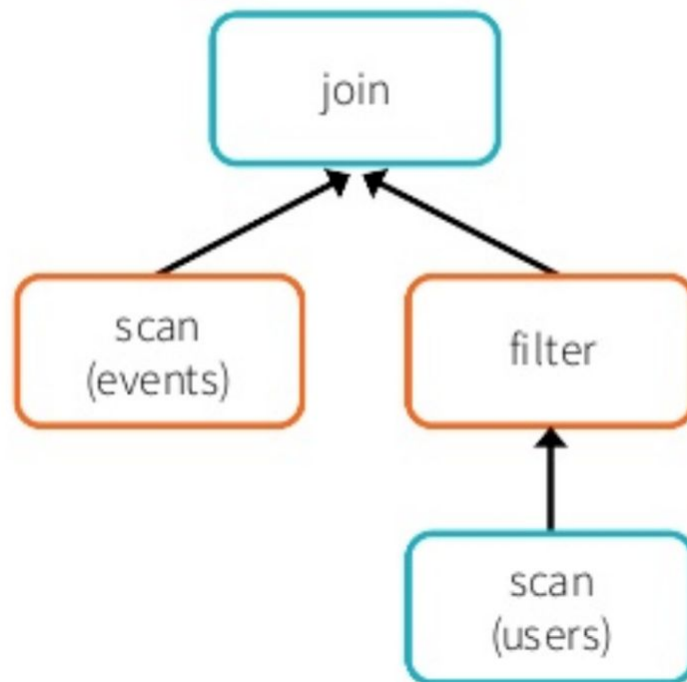


About Apache Spark | Planning

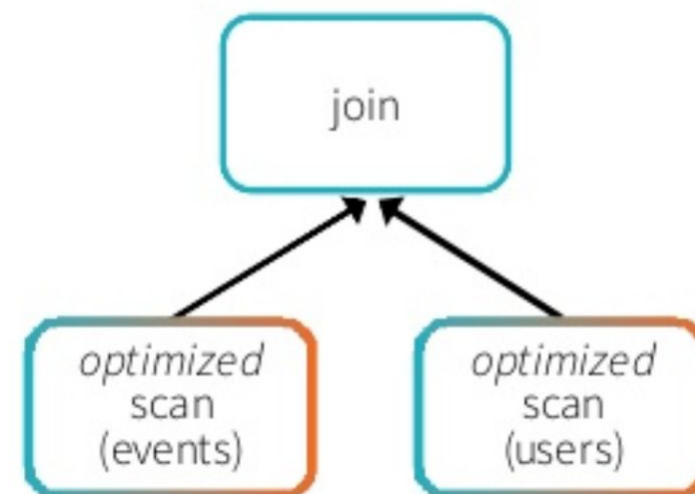
Logical Plan



Physical Plan



Physical Plan
with Predicate Pushdown
and Column Pruning



Walkthrough | Create Spark Session

```
from pyspark.sql import SparkSession
spark = SparkSession.builder\
    .appName(name="PySpark Intro")\
    .master("local[*]")\
    .getOrCreate()
```

Deploy modes:
Local, standalone, YARN,
Mesos and Kubernetes

SparkSession - in-memory
SparkContext

[Spark UI](#)

Version

v2.2.1

Master

local[*]

AppName


PySpark Intro

Walkthrough | Read CSV into DataFrame

```
green_trips = spark.read\  
    .option("header", "true")\  
    .option("inferSchema", "true")\  
    .csv("green_tripdata_2017-06.csv")
```

Forces eager evaluation;
default is **false**

Walkthrough | Behind the Scenes: UI

 2.2.0

Jobs

Stages

Storage

Environment

Executors

SQL

PySparkShell application UI

Spark Jobs (?)

User: garrenstaubli
Total Uptime: 2.5 min
Scheduling Mode: FIFO

Completed Jobs

Event Timeline

Shows when jobs started and ended and when executors joined or left. Drag to scroll. Click Enable Zooming and use mouse wheel to zoom in/out.

Job Id ▾	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
1	csv at NativeMethodAccessorImpl.java:0	2018/01/23 22:05:21	4 s	1/1	8/8
0	csv at NativeMethodAccessorImpl.java:0	2018/01/23 22:05:21	0.7 s	1/1	1/1

Walkthrough | Behind the Scenes: UI



Jobs

Stages

Storage

Environment

Executors

SQL

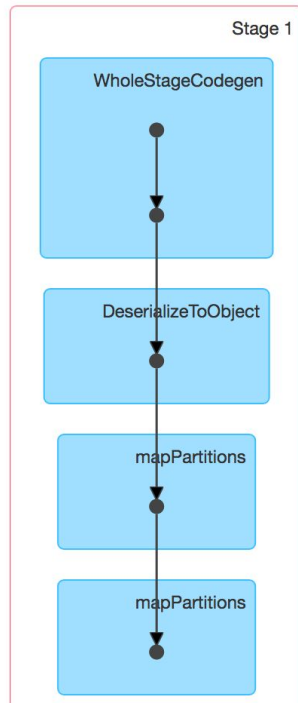
Details for Job 1

Status: SUCCEEDED

Completed Stages: 1

▶ Event Timeline

▼ DAG Visualization



Walkthrough | DataFrame Schema

green_trips.printSchema()

Eagerly evaluated (inferSchema = true)

```
root
|-- VendorID: integer (nullable = true)
|-- lpep_pickup_datetime: timestamp (nullable = true)
|-- lpep_dropoff_datetime: timestamp (nullable = true)
|-- store_and_fwd_flag: string (nullable = true)
|-- RatecodeID: integer (nullable = true)
|-- PULocationID: integer (nullable = true)
|-- DOLocationID: integer (nullable = true)
|-- passenger_count: integer (nullable = true)
|-- trip_distance: double (nullable = true)
|-- fare_amount: double (nullable = true)
|-- extra: double (nullable = true)
|-- mta_tax: double (nullable = true)
|-- tip_amount: double (nullable = true)
|-- tolls_amount: double (nullable = true)
|-- ehail_fee: string (nullable = true)
|-- improvement_surcharge: double (nullable = true)
|-- total_amount: double (nullable = true)
|-- payment_type: integer (nullable = true)
|-- trip_type: integer (nullable = true)
```

Lazily evaluated (inferSchema = false)

```
root
|-- VendorID: string (nullable = true)
|-- lpep_pickup_datetime: string (nullable = true)
|-- lpep_dropoff_datetime: string (nullable = true)
|-- store_and_fwd_flag: string (nullable = true)
|-- RatecodeID: string (nullable = true)
|-- PULocationID: string (nullable = true)
|-- DOLocationID: string (nullable = true)
|-- passenger_count: string (nullable = true)
|-- trip_distance: string (nullable = true)
|-- fare_amount: string (nullable = true)
|-- extra: string (nullable = true)
|-- mta_tax: string (nullable = true)
|-- tip_amount: string (nullable = true)
|-- tolls_amount: string (nullable = true)
|-- ehail_fee: string (nullable = true)
|-- improvement_surcharge: string (nullable = true)
|-- total_amount: string (nullable = true)
|-- payment_type: string (nullable = true)
|-- trip_type: string (nullable = true)
```




You guessed it...
We're hiring!

BlueprintTM

CONSULTING SERVICES

- 2015 #1
- 2016 #1

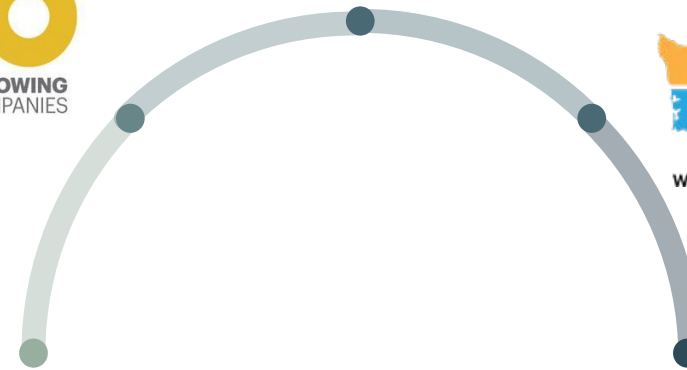


- 2014
- 2015
- 2016



- 2016

- 2014
- 2015 #1
- 2016 #1
- 2017 #4



- 2015 #373
- 2016 #166
- 2017 #161



#PySparkWorkshop

Resources: garrens.com/pyspark124