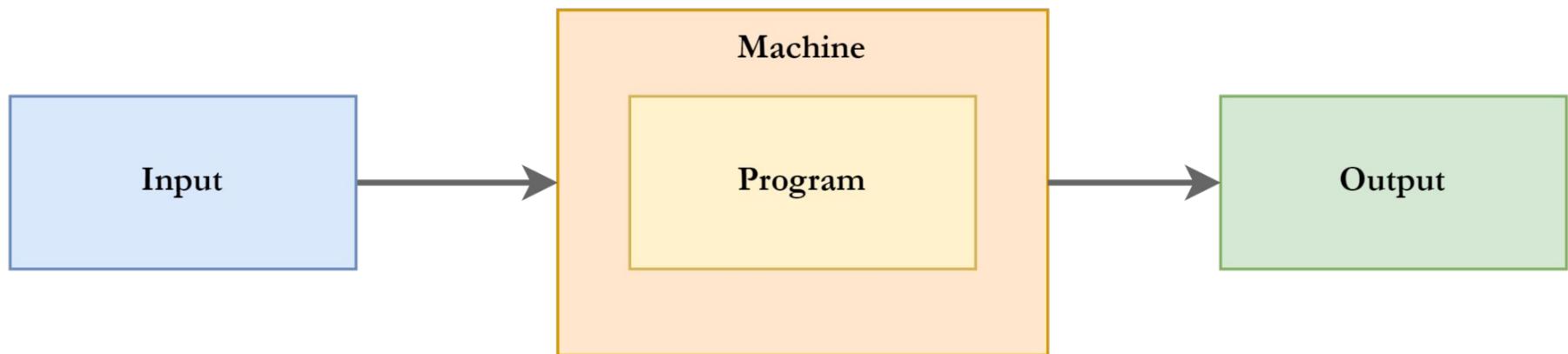


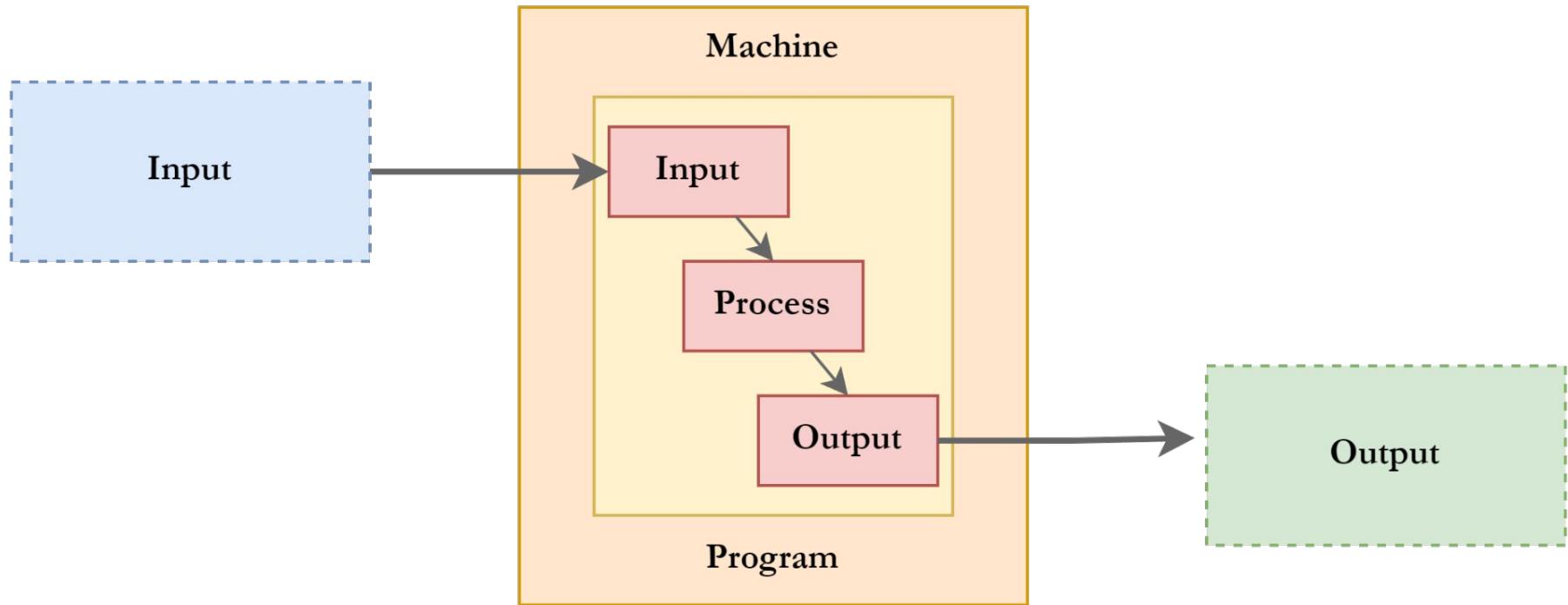
# Data Processing with PySpark

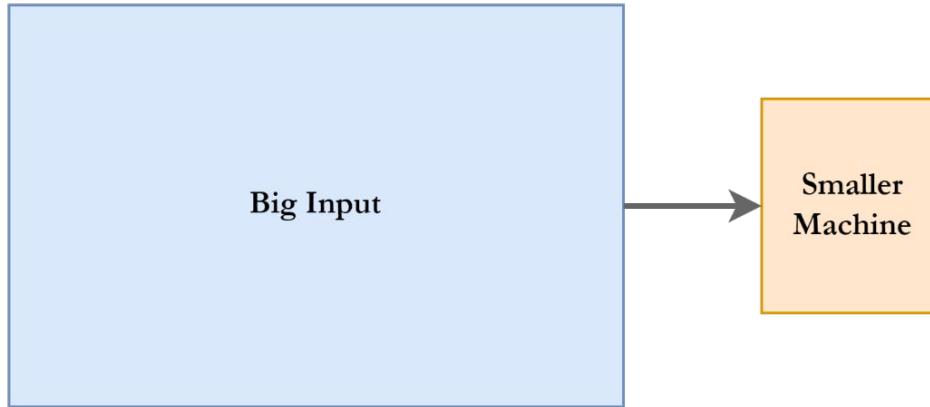
# Introduction

# Data Processing Model



Enough CPU, Memory and Disk

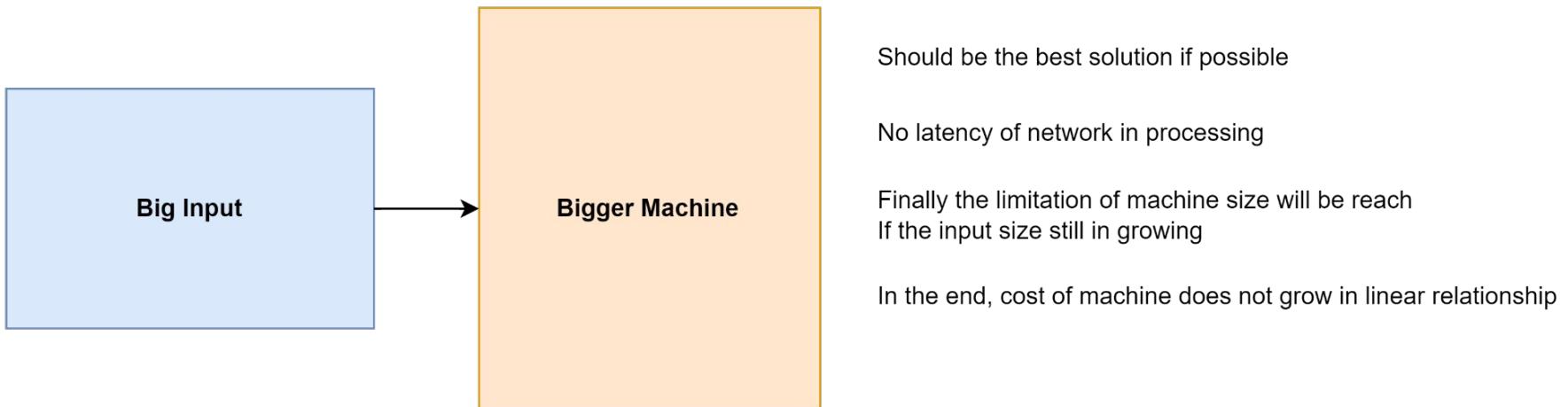




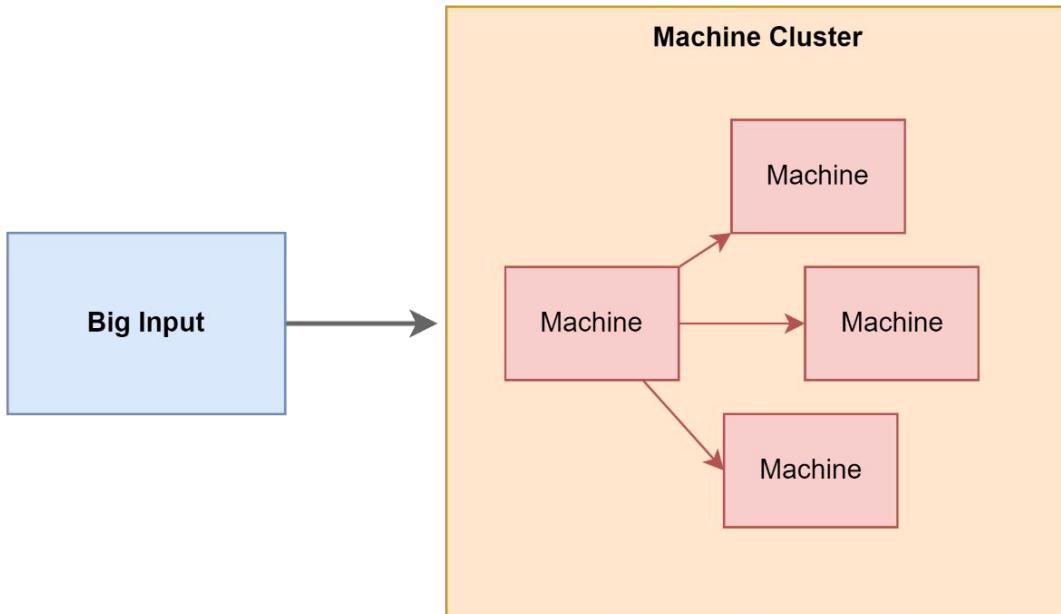
Problem cannot divided to produce partial outputs

Partial outputs must be reduced to the final solution

# Vertical Scaling



# Horizontal Scaling



It is possible to buy an affordable machine size in market

There is a challenge in how to make a program running in each machine to collaborate for processing together

Network cost in communication between machines, but it is worth to lessen the processing time

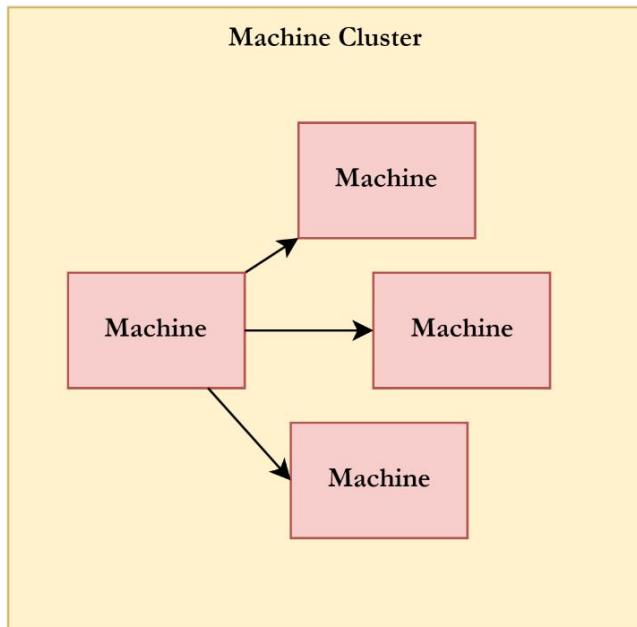
Today, available cloud solution (public or private) make this solution easier to apply

## Condition to choose processing solutions

Dataset type	Size range	Fits in RAM?	Fits on local disk?
Small dataset	Less than 2–4 GB	Yes	Yes
Medium dataset	Less than 2 TB	No	Yes
Large dataset	Greater than 2 TB	No	No

# Hardest Parts and Challenges

Even if nowaday we have many programming languages which suport network programming.. There are remaining hardest problem need to design and resolve.



Network Fallacies

Resiliency

Require experience and skills to achieve them

Efficient

It is not a good decision to make a program from scratch

Maintainable



Apache Spark is a unified computing engine and a set of libraries for parallel data processing on computer clusters.

Structured  
Streaming

Advanced  
Analytics

Libraries &  
Ecosystem

Datasets

Structured APIs  
DataFrames

SQL

RDDs

Distributed Variables

- Unified
- Compute Engine
- Libraries

# Unified

Supports data processing from simple to complex tasks

Supports scalability in data processing tasks:

- Single Core, Single Machine
- Multiple Core, Single Machine
- Multiple Core, Multiple Machine

Support multiple programming languages

- Python
- R
- Scala
- Java
- SQL
- C# (Microsoft)
- Go

# Libraries

With core libraries provided, there are additional libraries extending Spark, which are supported by community

- SparkML
- GraphFrame

# Compute Engine

Focus on computation processing and caching temporary solution for accelerate computation

No concept in long term storage in Spark, but it can integrate with other solutions

Other storage solutions: S3, Azure Data Lake, RDBMS, NoSQL

# What we will focus in this course

Structured  
Streaming

Advanced  
Analytics

Libraries &  
Ecosystem

Structured APIs

Datasets

DataFrames

SQL

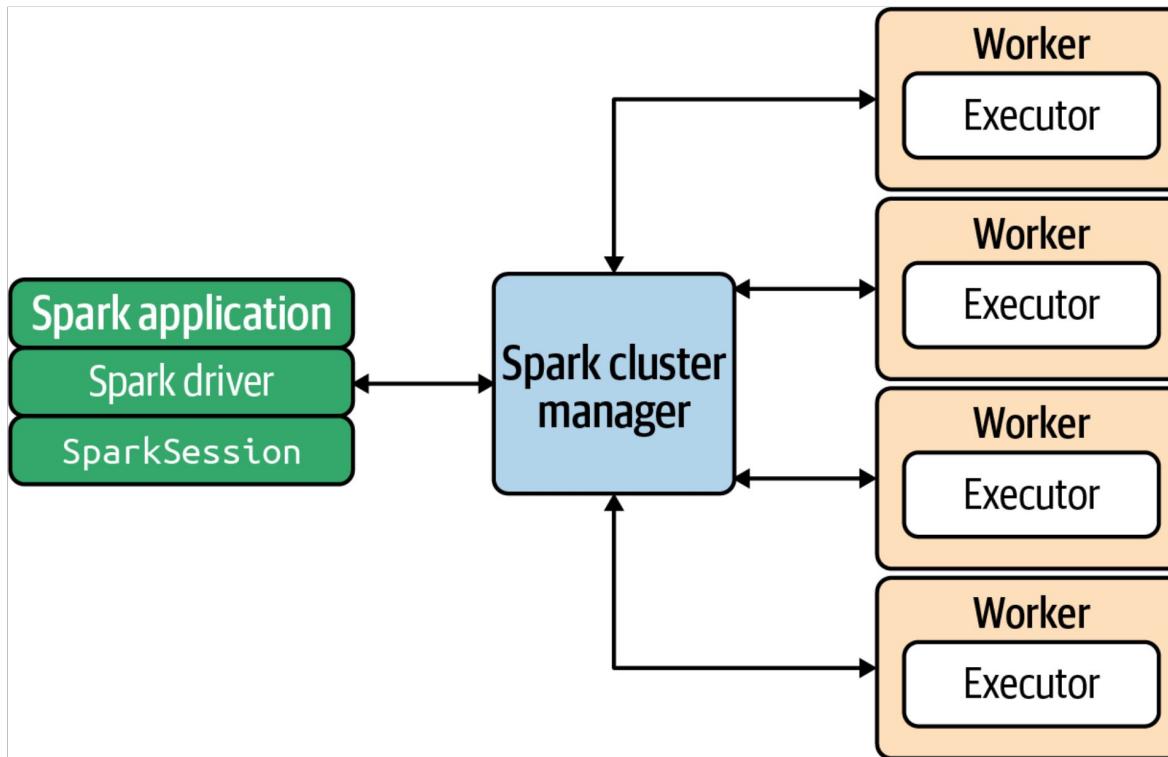
Low-level APIs

RDDs

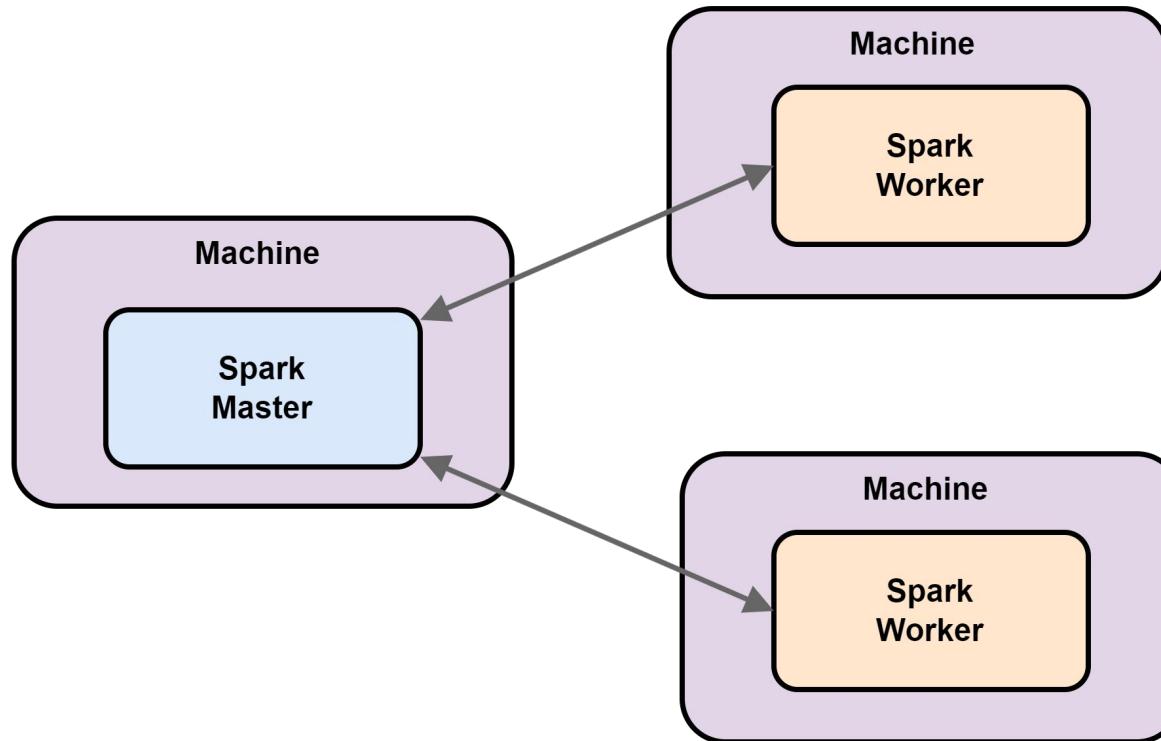
Distributed Variables

# Spark Cluster

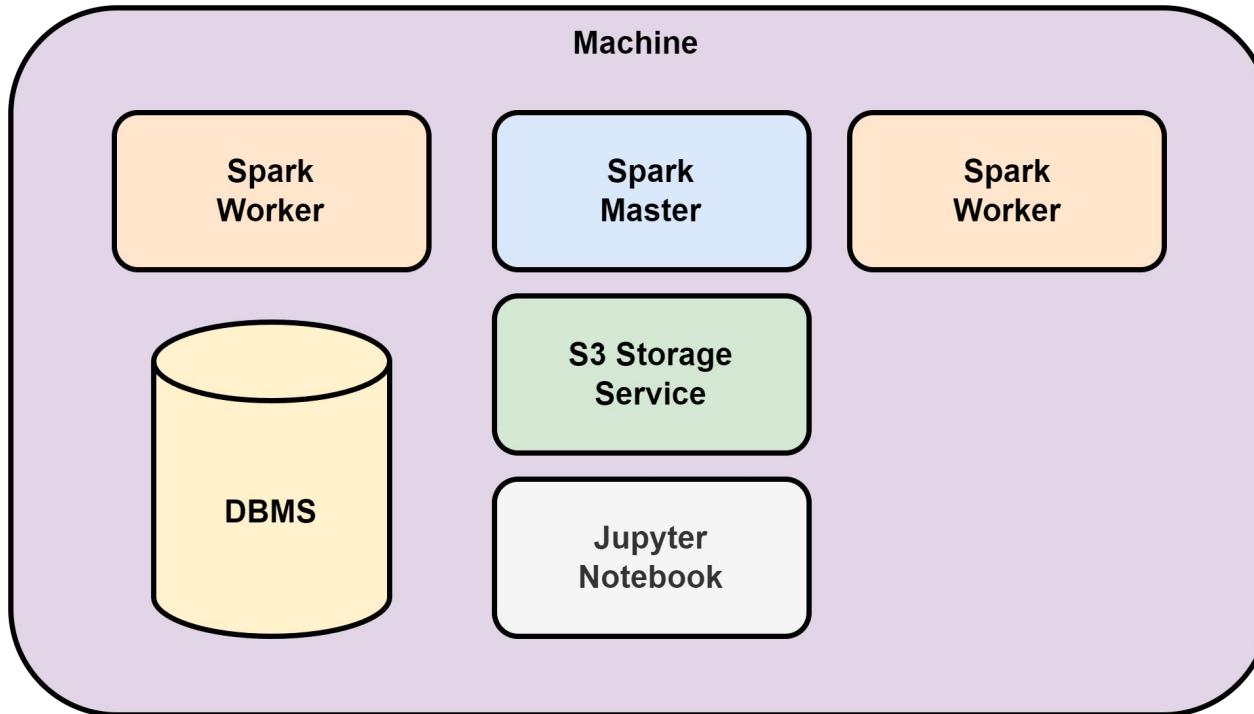
# Spark Cluster Architecture



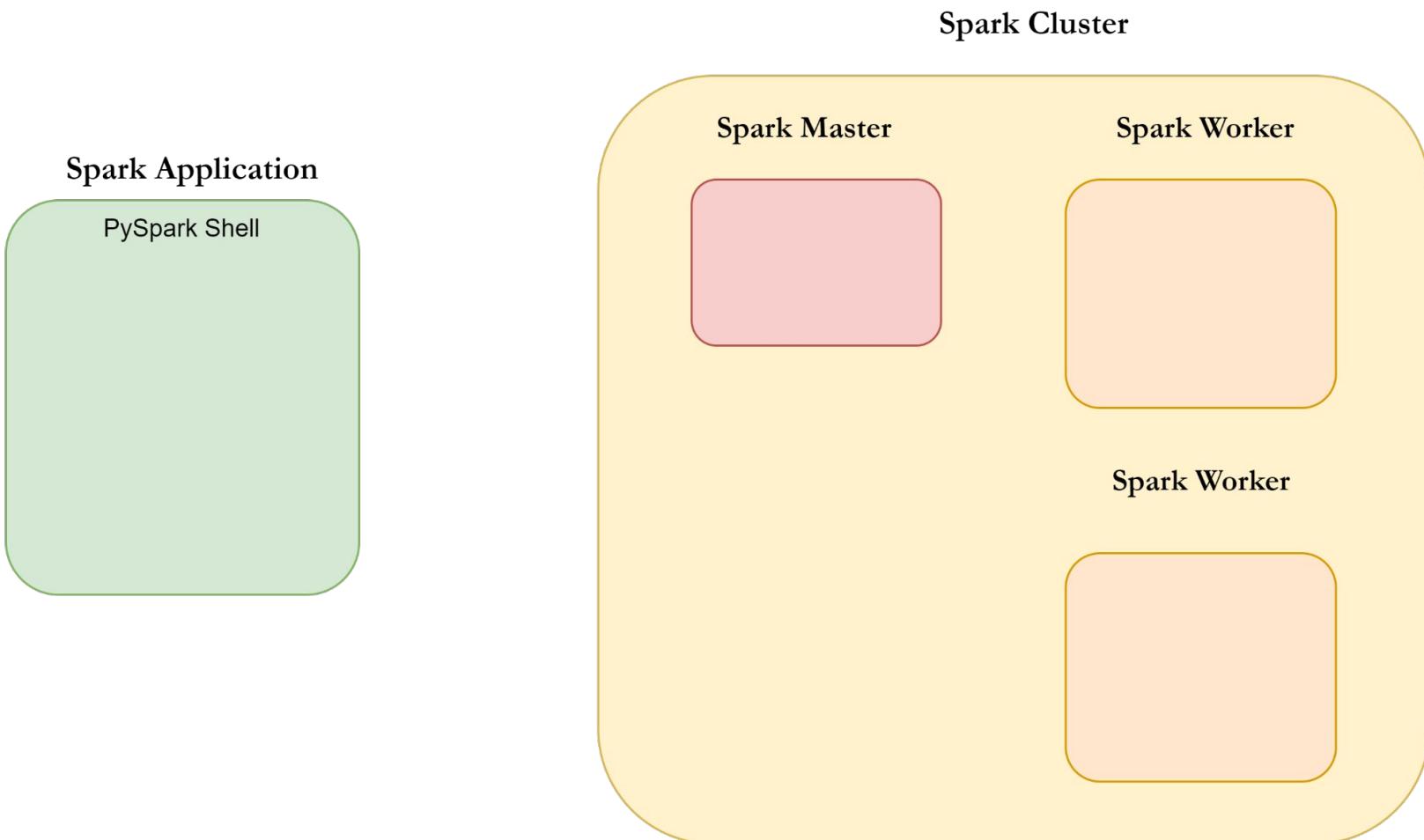
# Real Spark Cluster

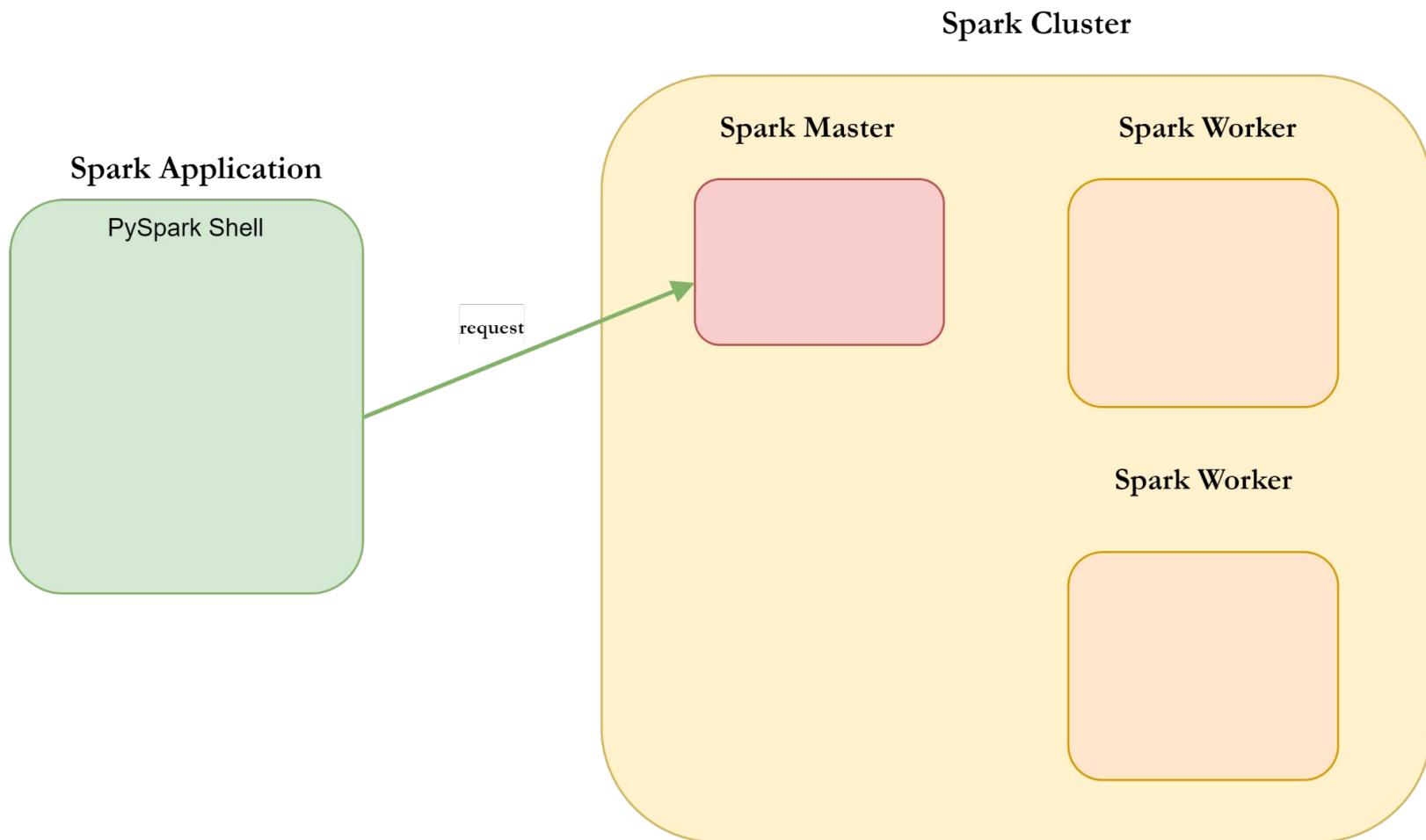


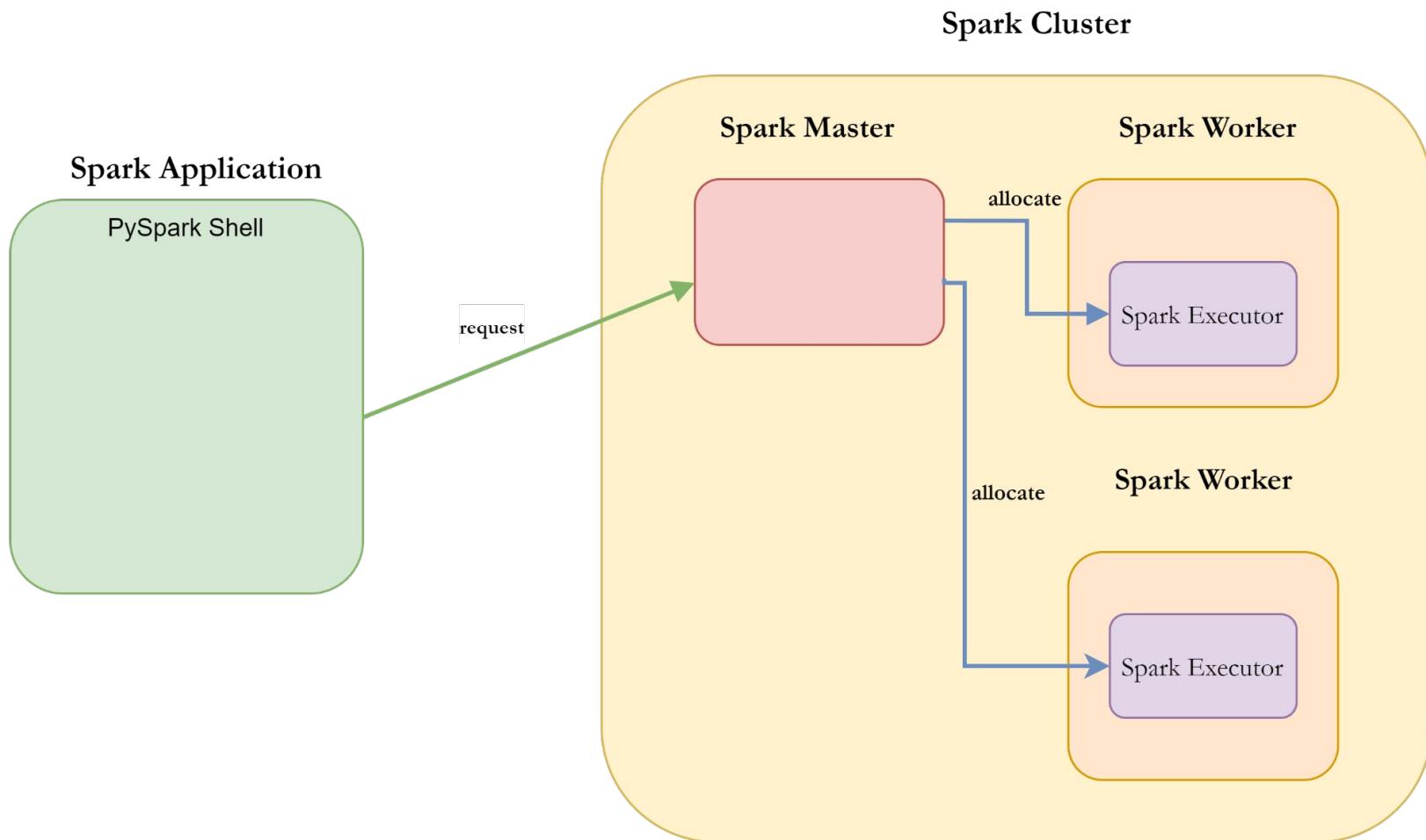
# Poorman's Spark Cluster



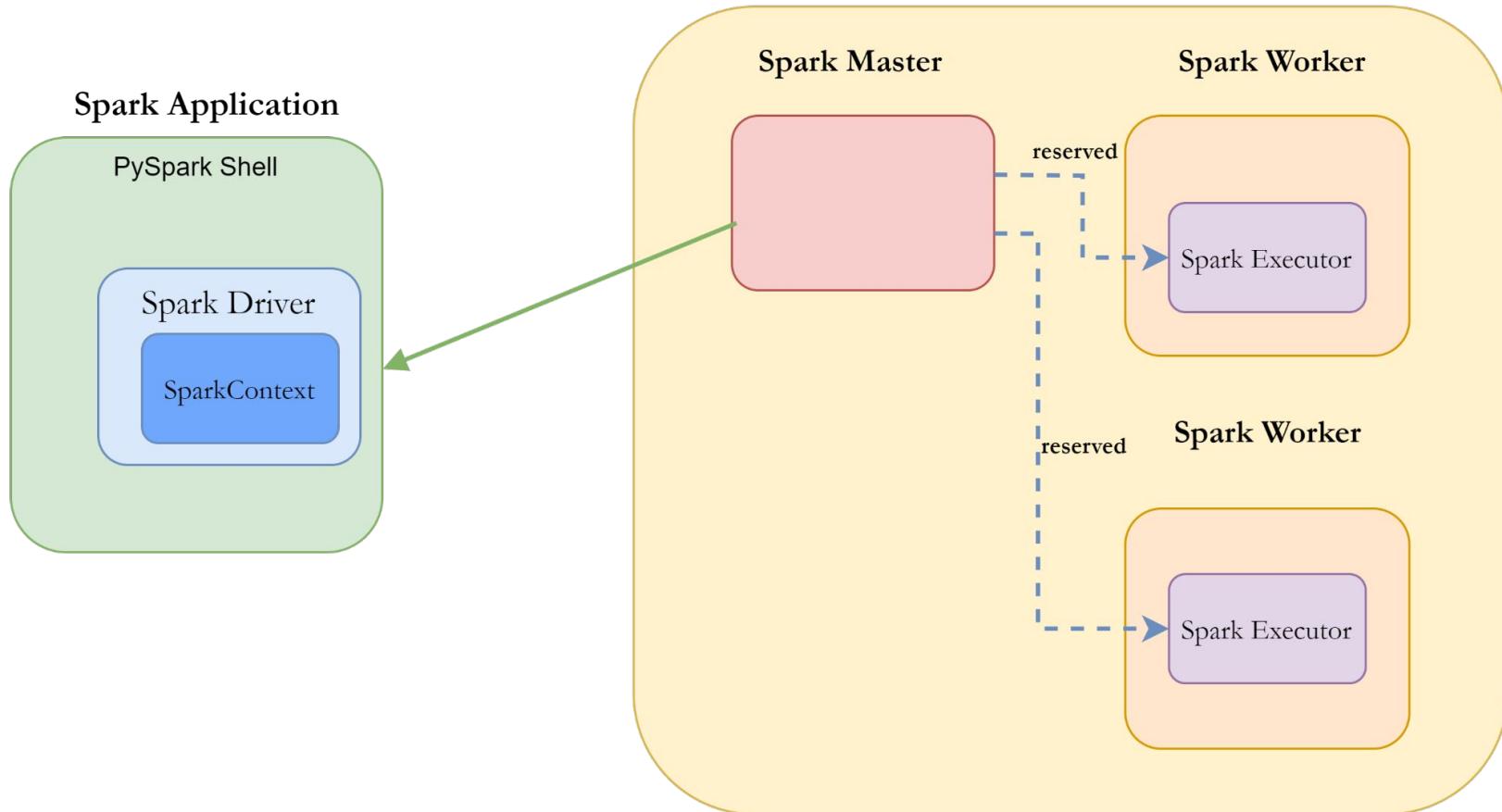
# PySparkShell

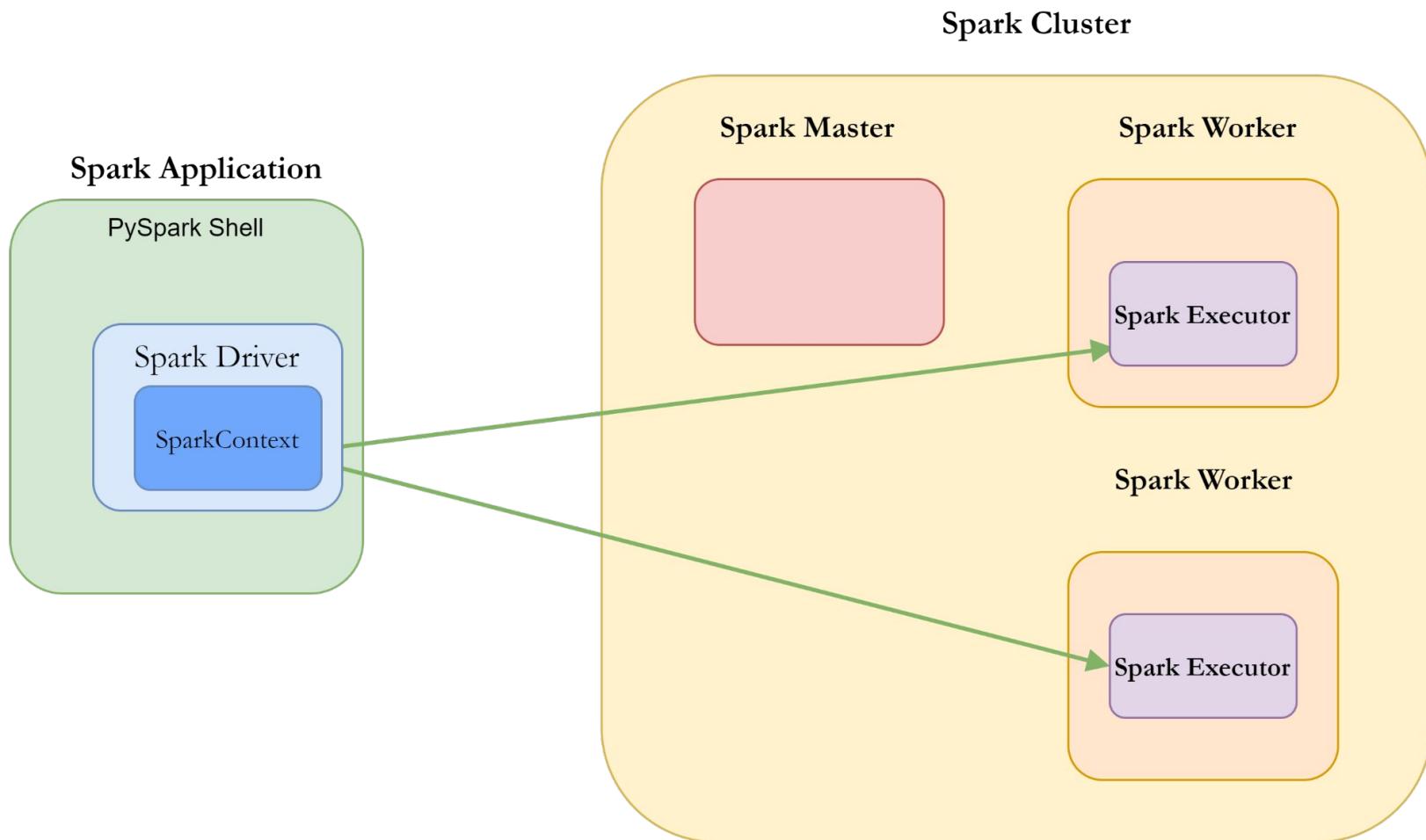




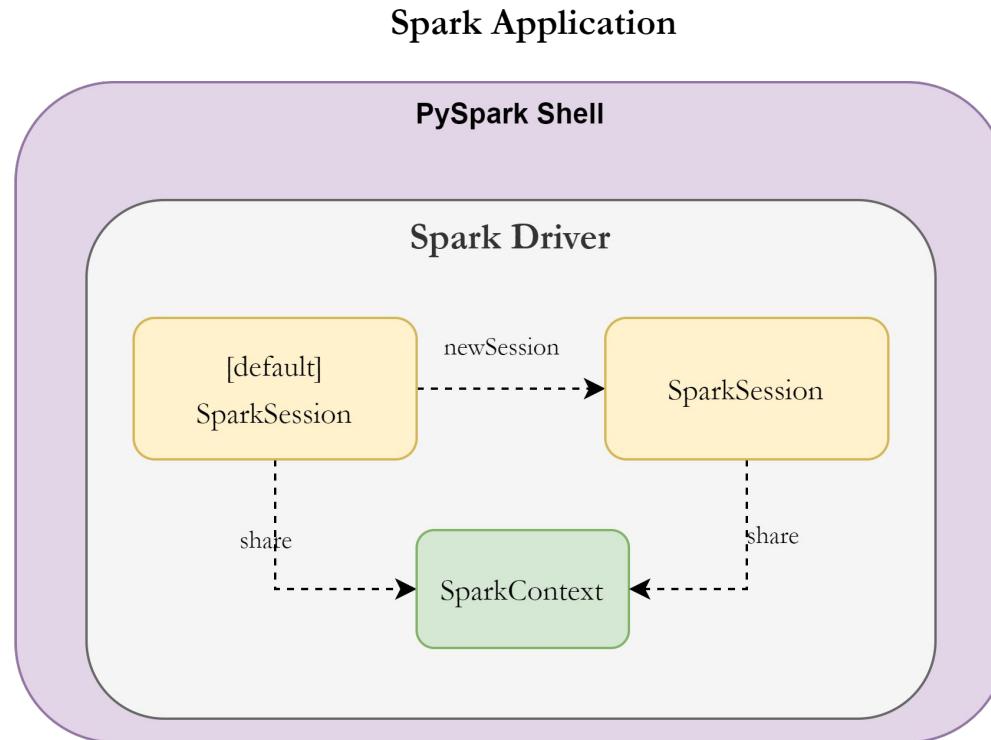


# Spark Cluster





# What is SparkSession?



# Where is SparkSession and SparkContext?



```
graph LR; sparkContext[spark.sparkContext] <--> SparkContext[SparkContext]; SparkContext --- RDD[RDD]
```

The diagram illustrates the relationships between `spark.sparkContext` and `SparkContext`. A double-headed arrow connects them, and a single-headed arrow points from `SparkContext` to `RDD`.

# Spark App

## **How to develop Spark Application**

1. Create a SparkSession instance by your program
2. How to run a Spark application
  - 2.1 Run it directly from a machine which can access all machines in Spark cluster.
  - 2.2 Use spark-submit script in Spark package (the application cannot override spark.master configuration)

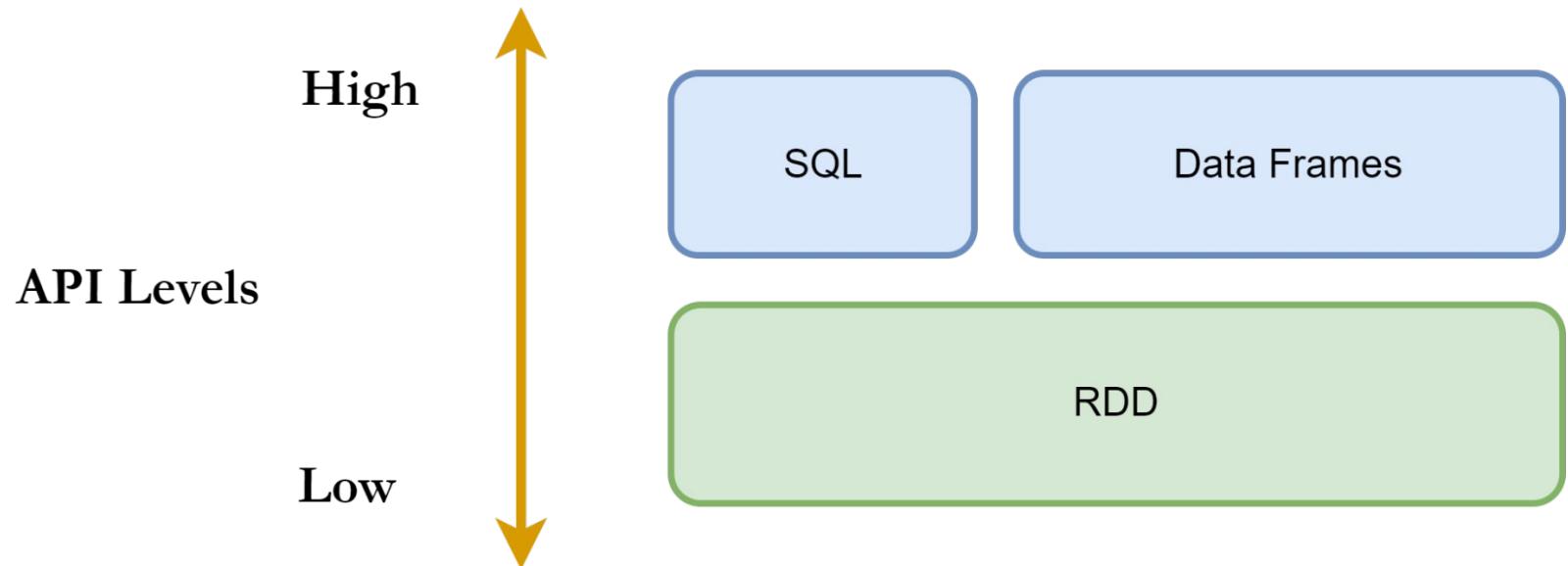
```
from pyspark.sql import SparkSession

spark = (
    SparkSession
        .builder
        .appName("spark_app_name")
        .master("spark_master_url")
        .config("config_name", "config_value")
        .getOrCreate()
)

spark # SparkSession instance
spark.sparkContext # SparkContext instance
```

# Spark APIs

# PySpark APIs Evaluation



# How are data treated in each APIs

## RDD

- Processing data like **list of object**
- object type can be primitive types like int, bool, string, or dict, list
- Data are **distributed** and **partitioned** in cluster

## DataFrame

- Processing data like **tables and rows with schema** (e. g. Excel or R DataFrame),
- Data are **distributed** and **partitioned** in cluster (as RDD is an underlying layer)

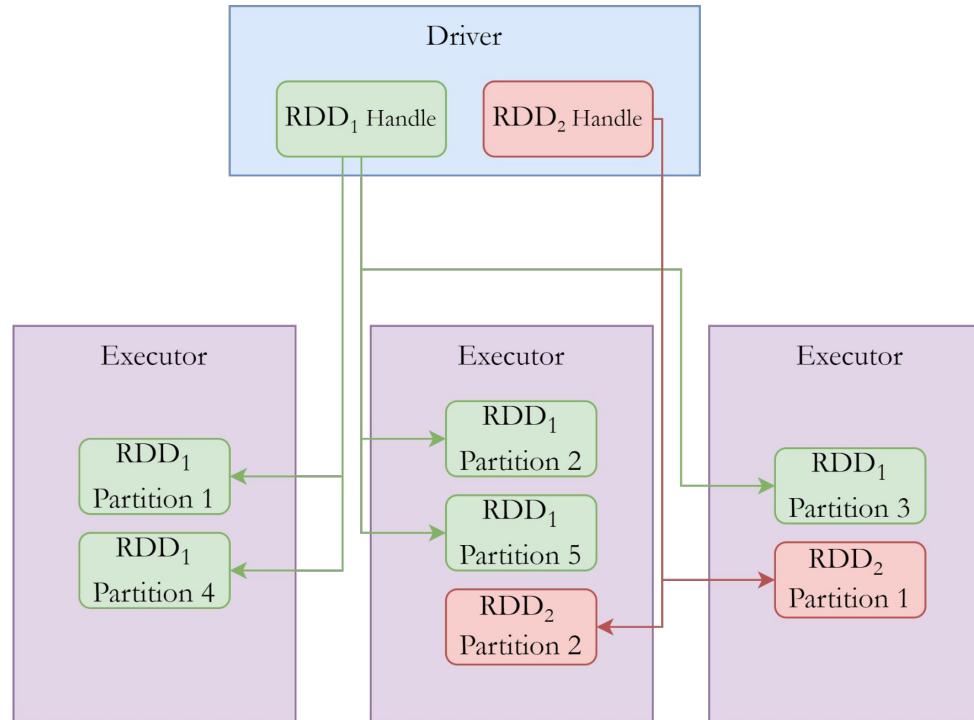
## SparkSQL

- Processing data like **tables and rows with schema** (e. g. Excel or R DataFrame)
- Data are **distributed** and **partitioned** in cluster (as RDD is an underlying layer)

# RDD

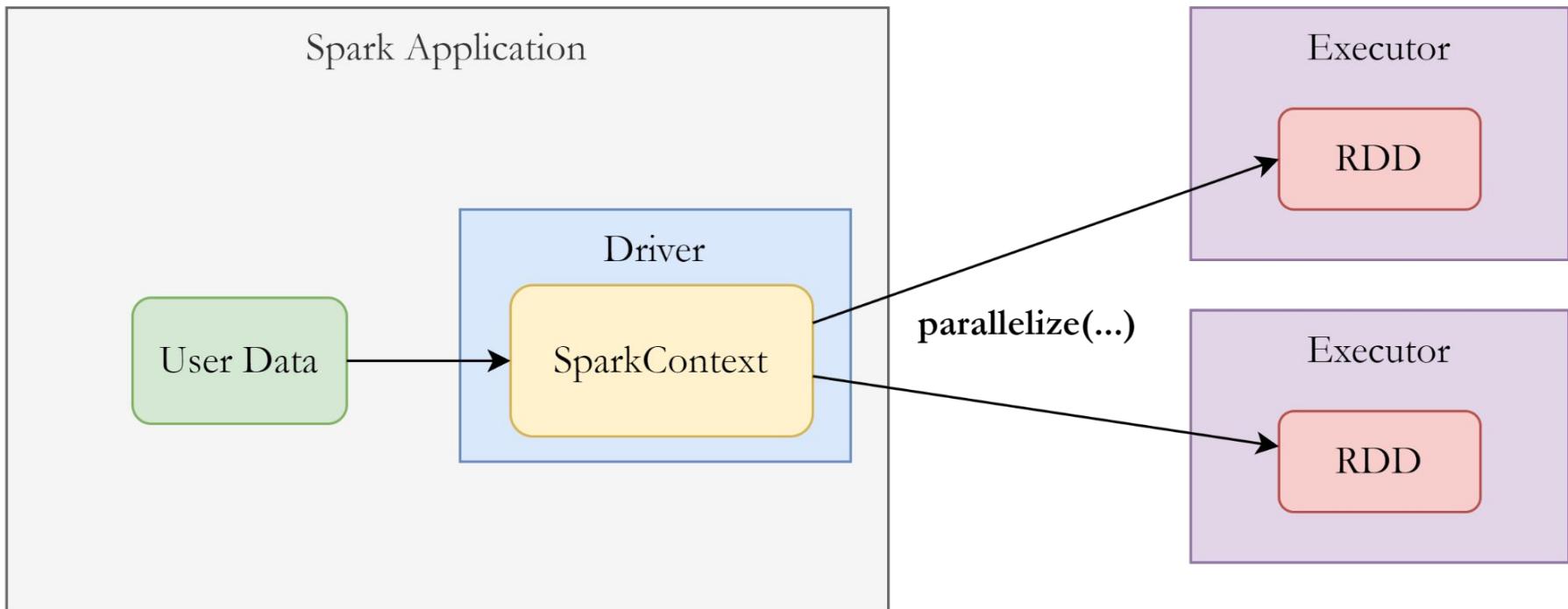
# RDD: Resilient Distributed Dataset

## Distributed

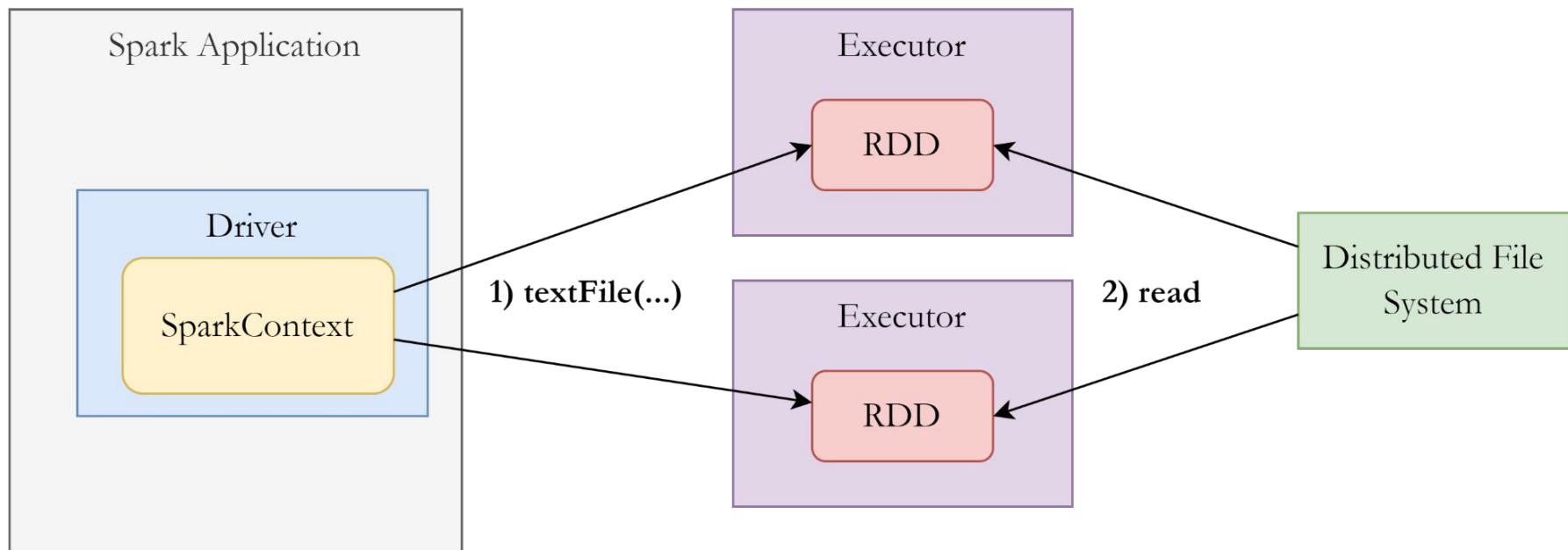


# How to create a RDD

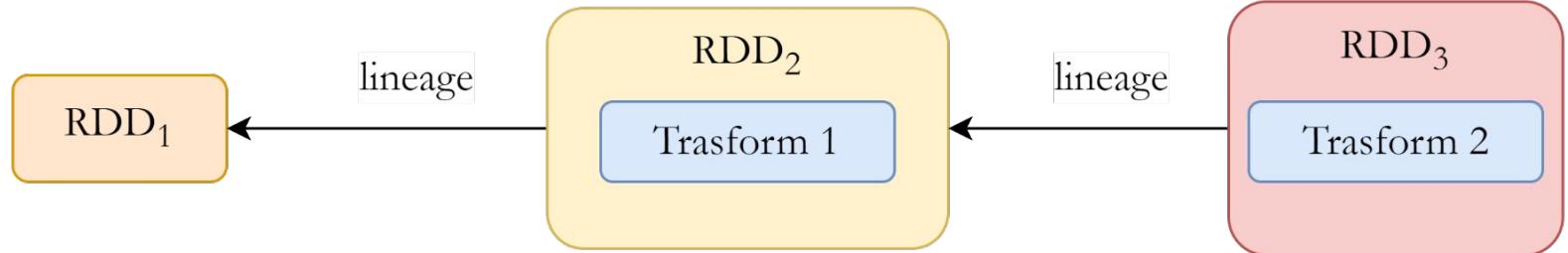
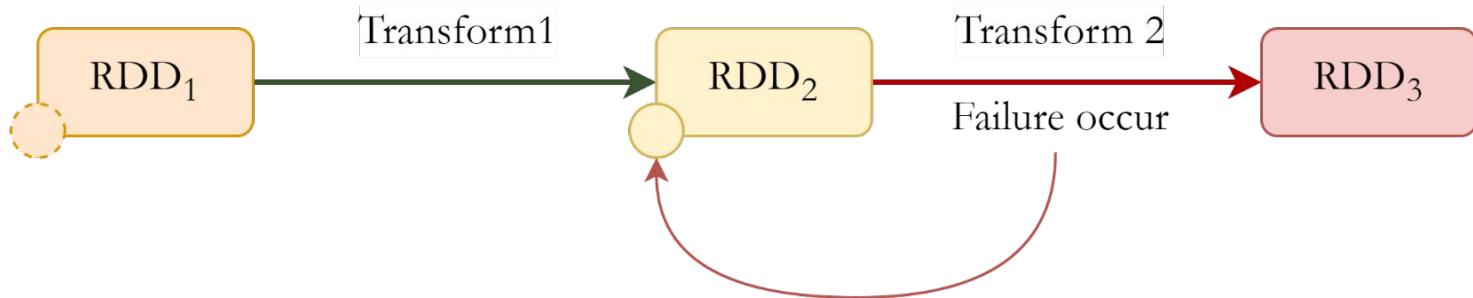
## Manually distributed

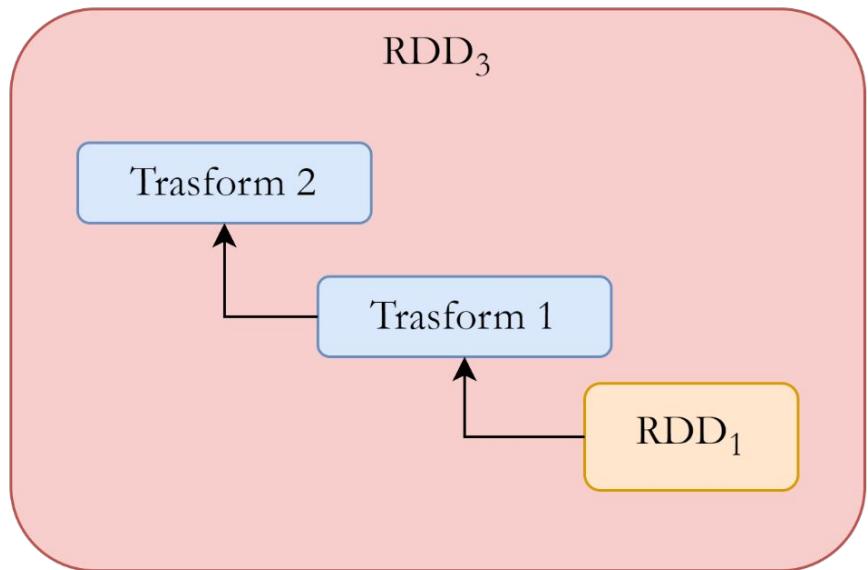


## Spark provided function



## Resilient





Prefer memory-only processing,  
but it can use disk if necessary

Prefer recomputation to cache temporary result,  
but it can cache for accelerating computation time.

# RDD Operation (Lazy Evaluation)

## How to call RDD operation

In Spark Application

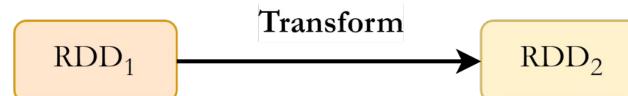
[rdd\_handle].[operation\_name](arg1, arg2, ...)

## Transform

Generate another RDD (create a data processing lineage)

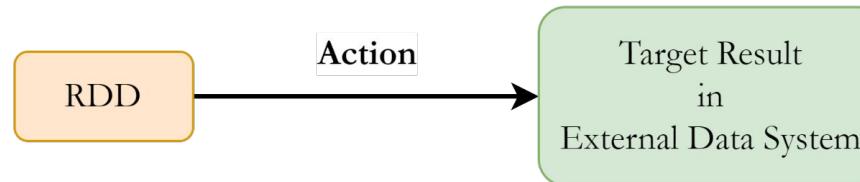
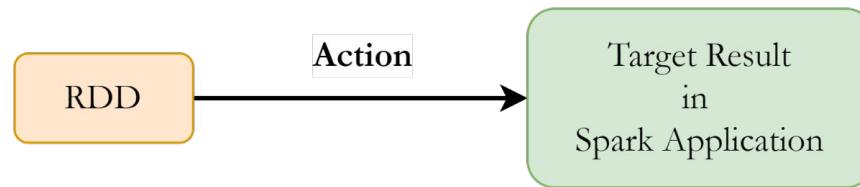
No operation happens in Spark cluster

Allow Spark to optimize operations when execute all combined operations in lineage

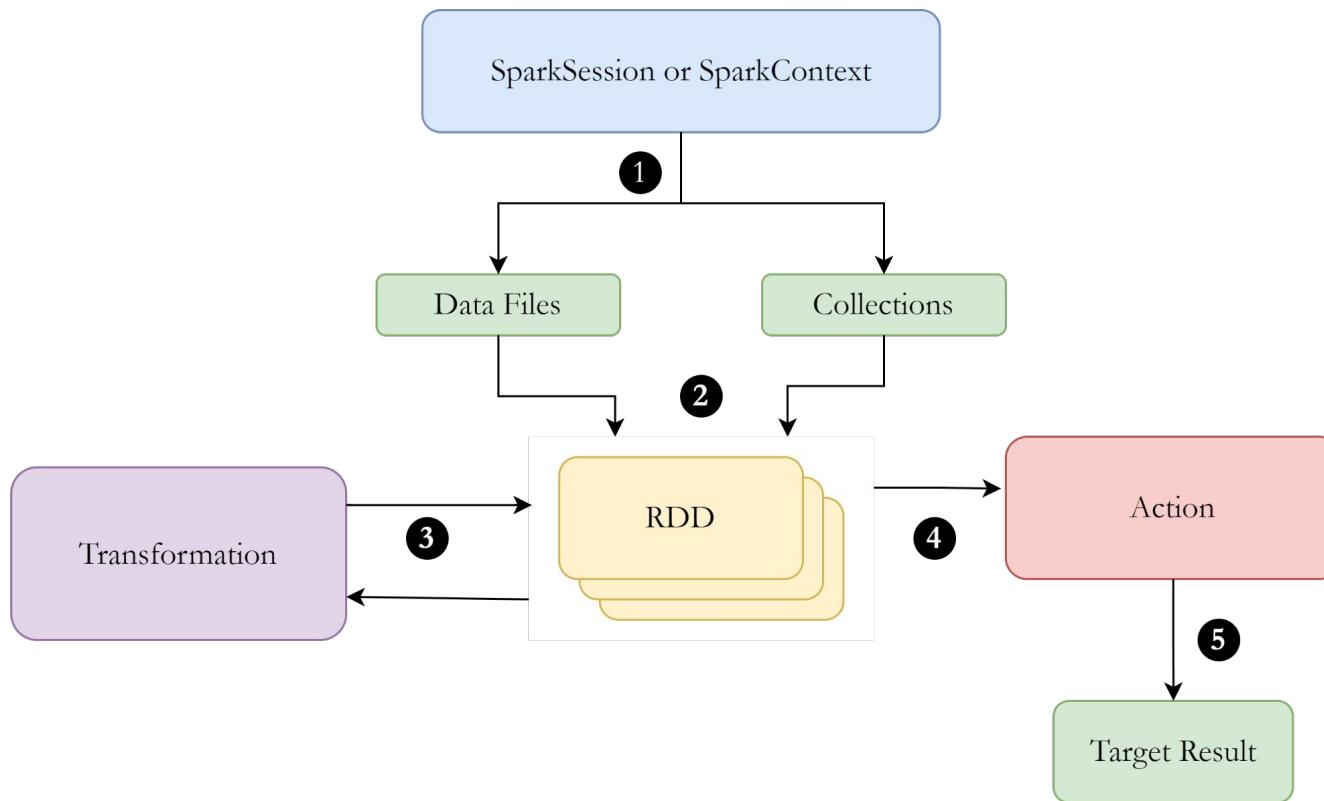


## Action

- Trigger the Spark cluster to processing as to the processing lineage in RDD which is called
- Some actions send result back the driver  
(Beware: memory limit on the driver machine)
- Some actions have the executors write the result to specified destination  
(Beware: the destination should support intensive write if the data is big)

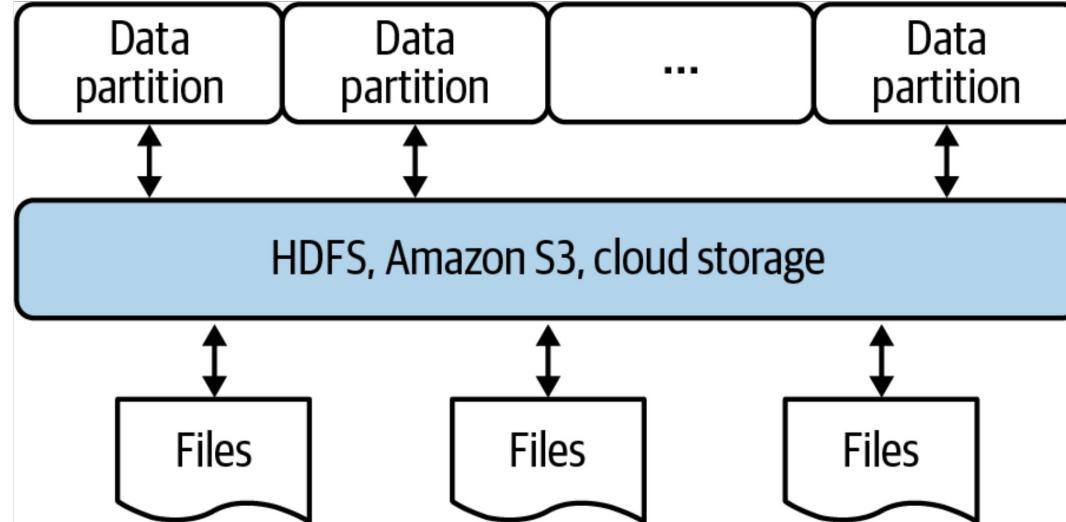


## Working with RDD summary

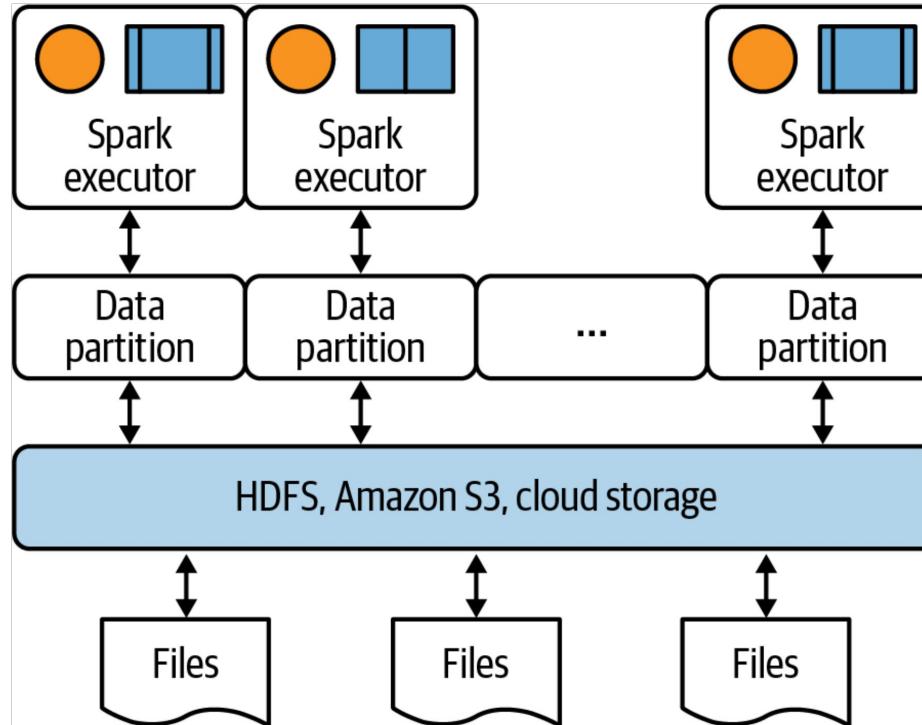


# Partition and Shuffle

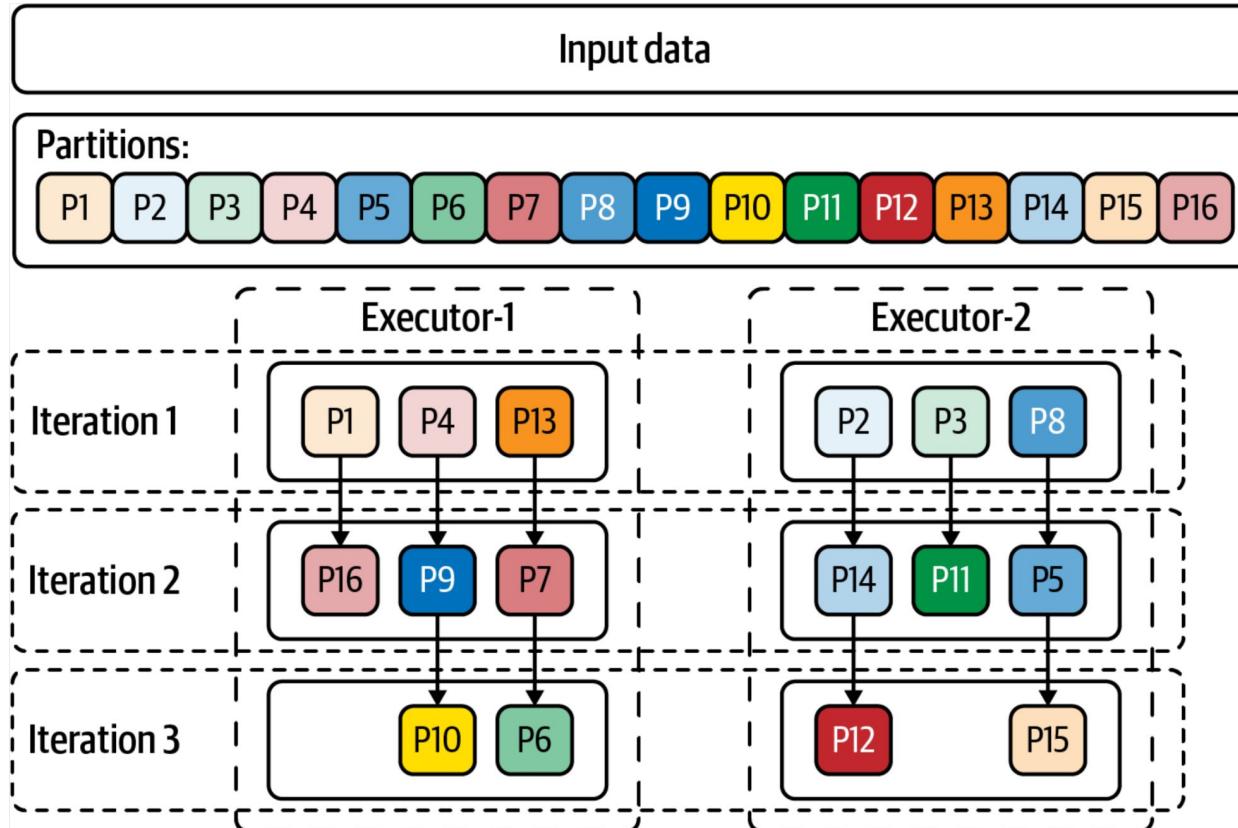
Source: [Data Algorithms in Spark](#)



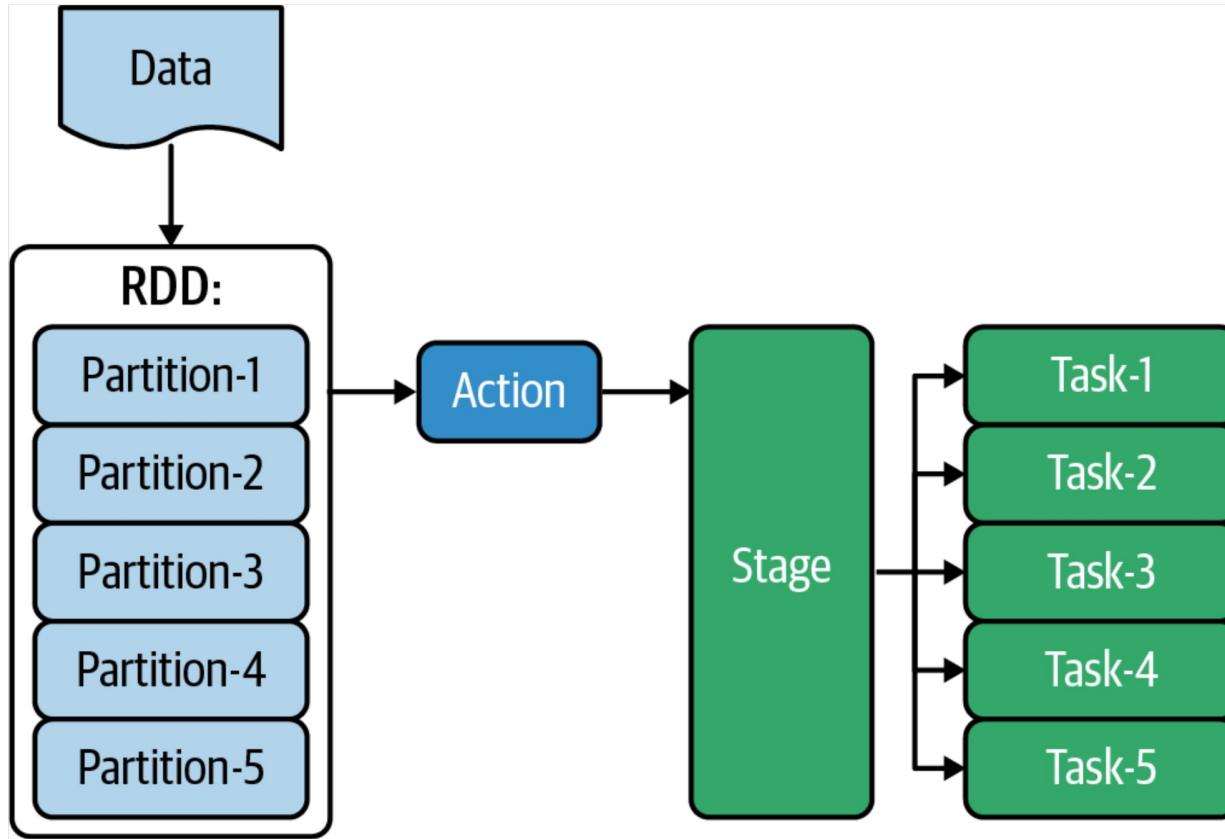
Source: [Data Algorithms in Spark](#)



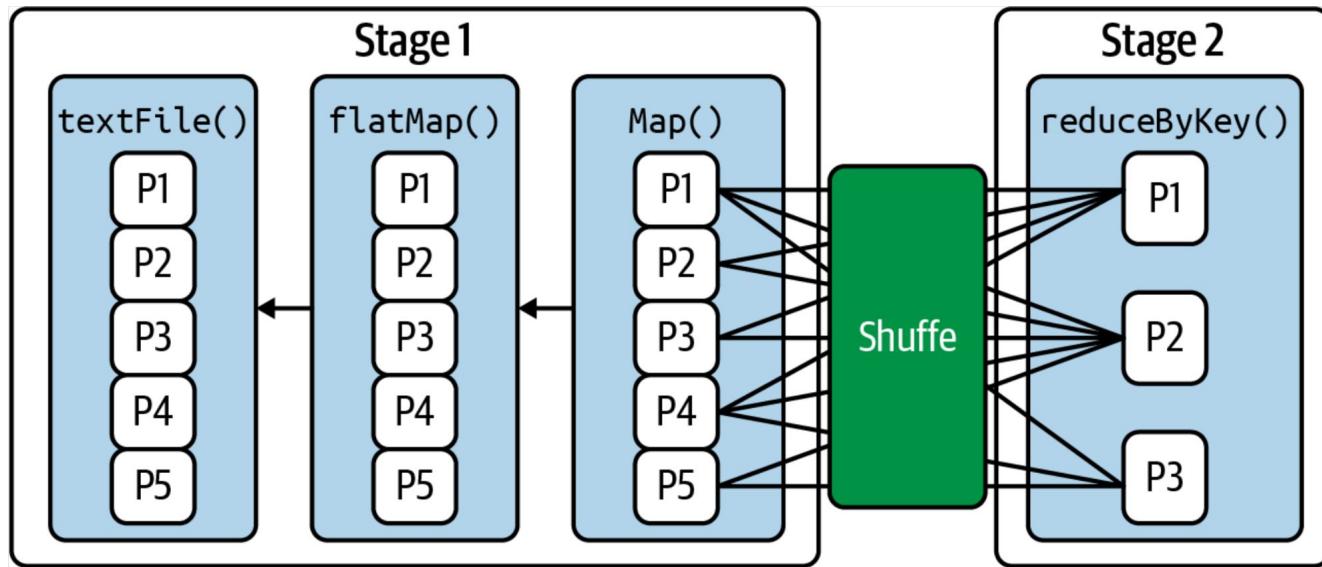
Source: [Data Algorithms in Spark](#)



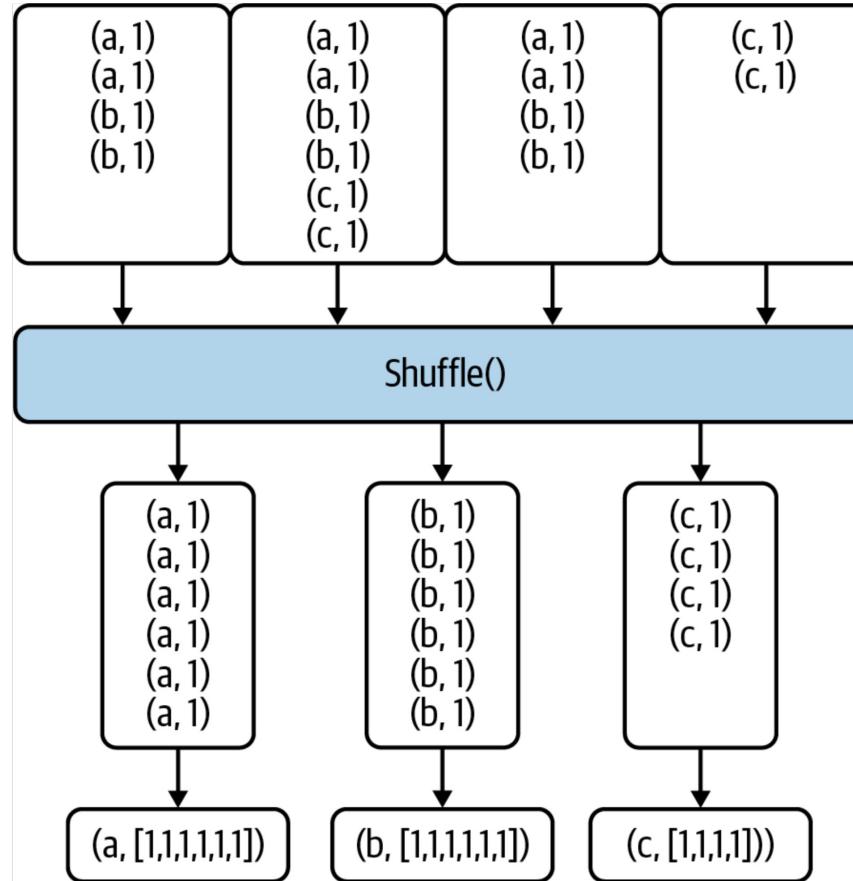
Source: [Data Algorithms in Spark](#)



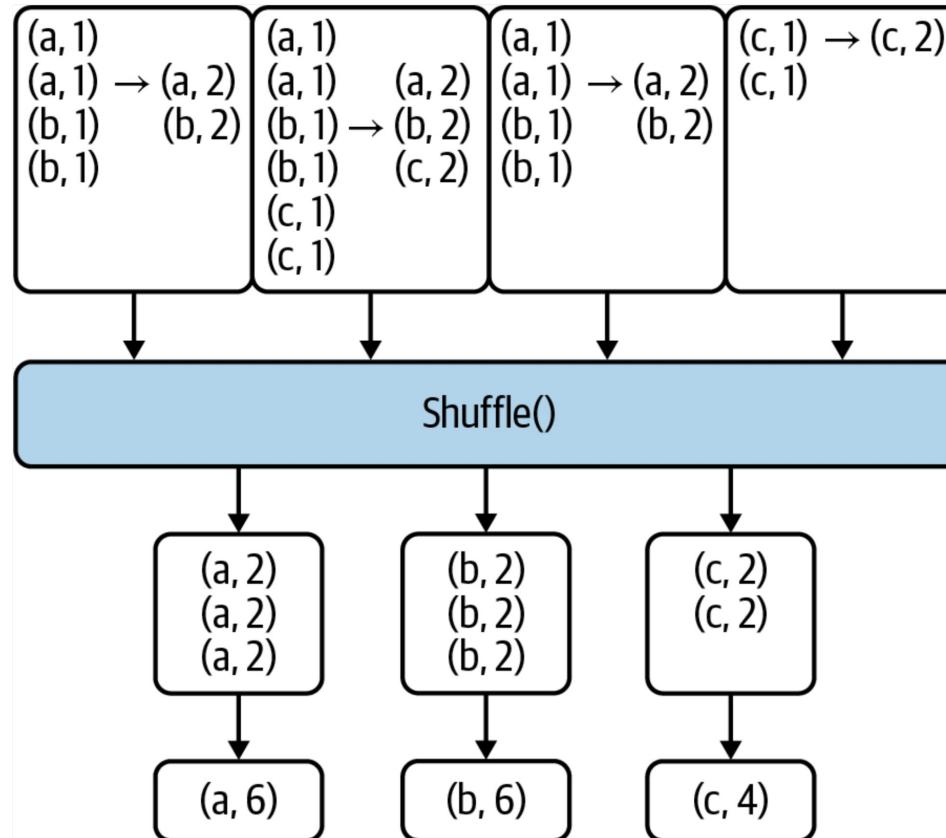
Source: [Data Algorithms in Spark](#)



Source: [Data Algorithms in Spark](#)

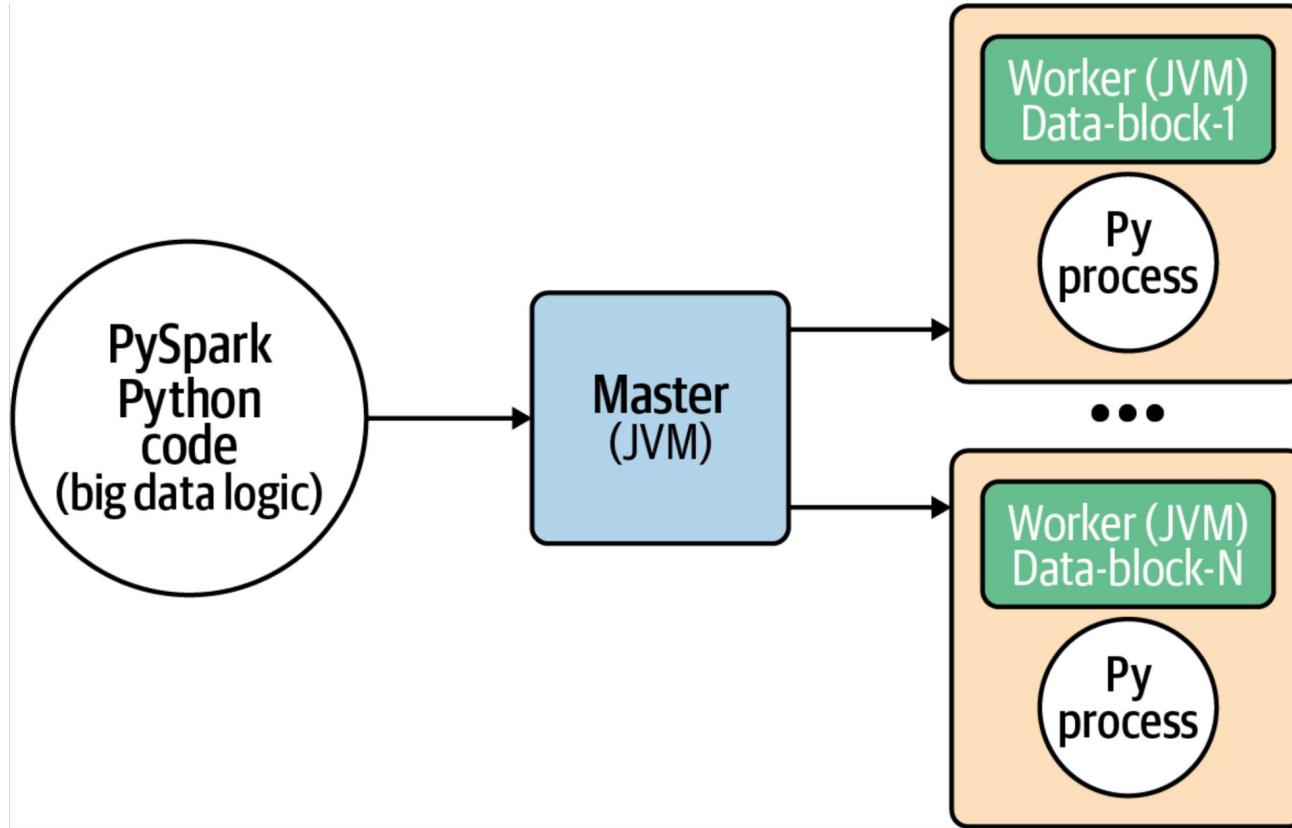


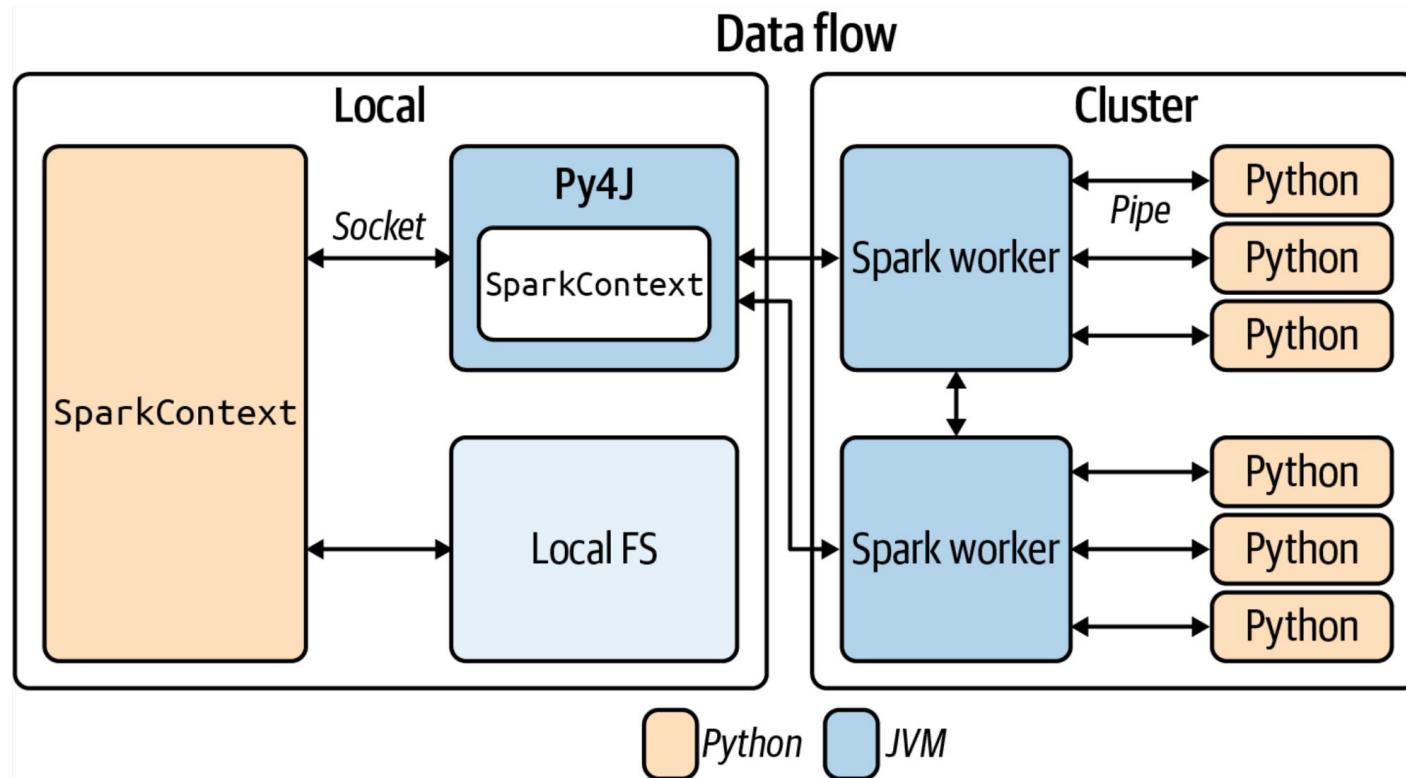
Source: [Data Algorithms in Spark](#)



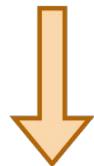
# DataFrame

Source: [Data Algorithms in Spark](#)





`rdd.[transform_or_action] ([python function])`



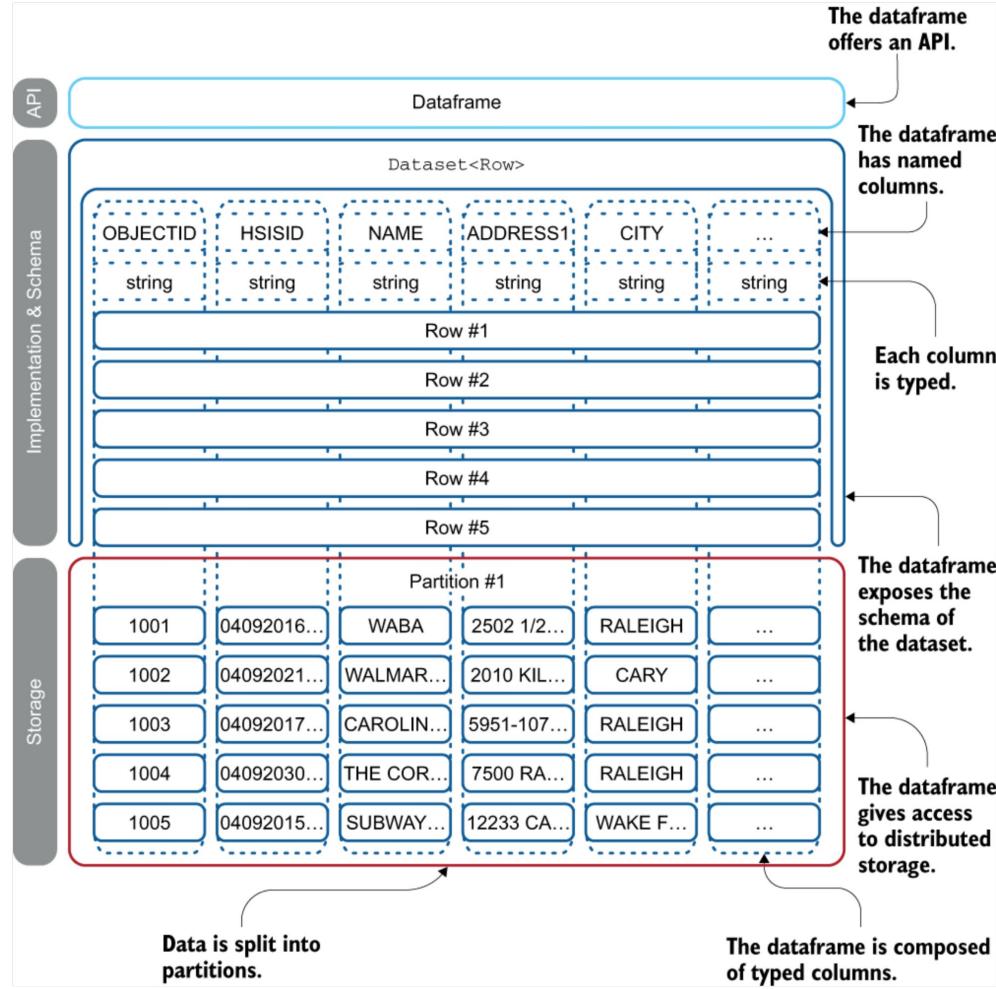
Intra-Process Serialization Cost

`dataframe.[operation](arg1, arg2, ...)`

No intra-process serialization cost

Almost popular tools have concepts of DataFrame



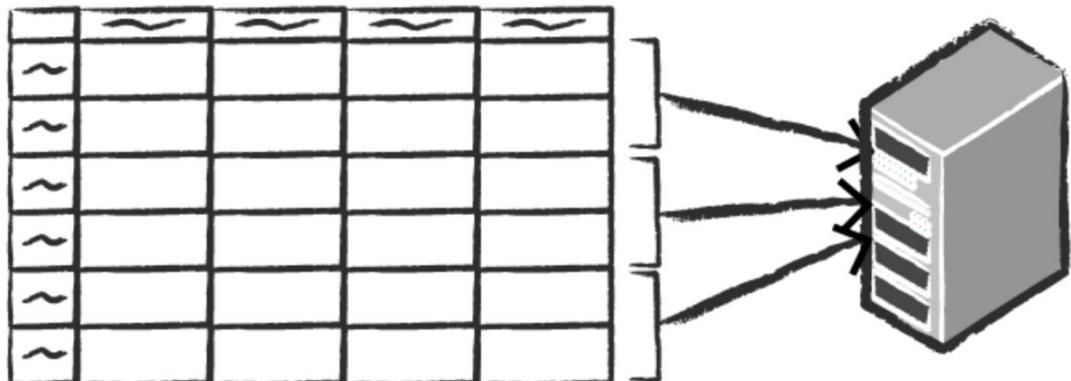


Source: [Spark The Definitive Guide](#)

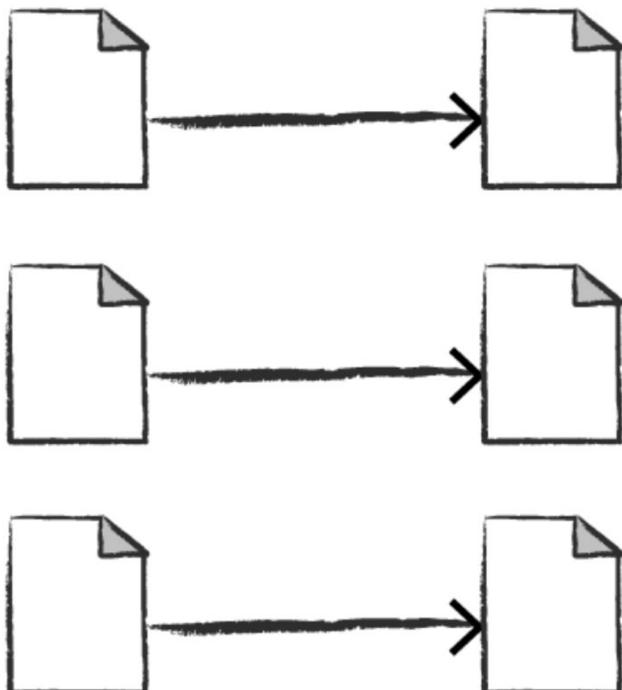
Spreadsheet on  
a single machine



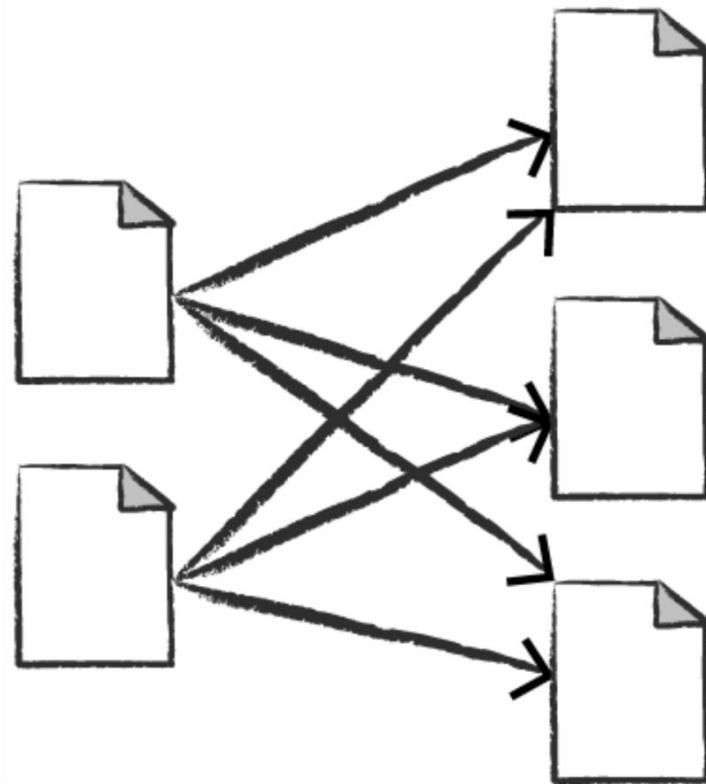
Table or Data Frame  
partitioned across servers  
in a data center

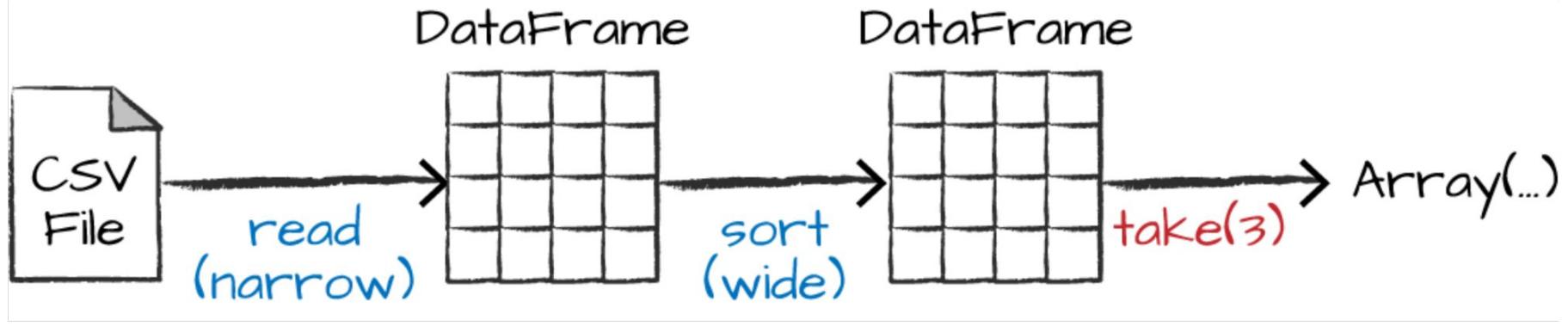


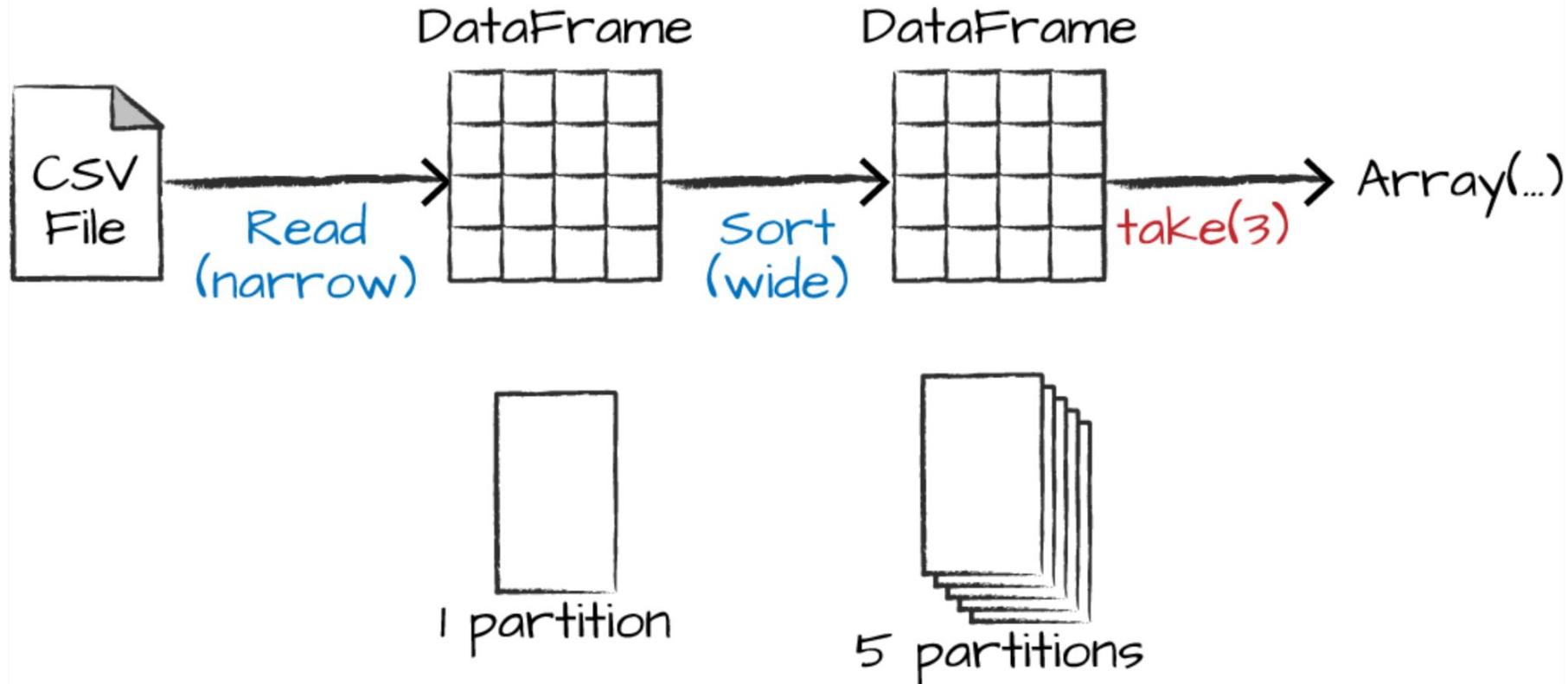
Narrow transformations  
1 to 1

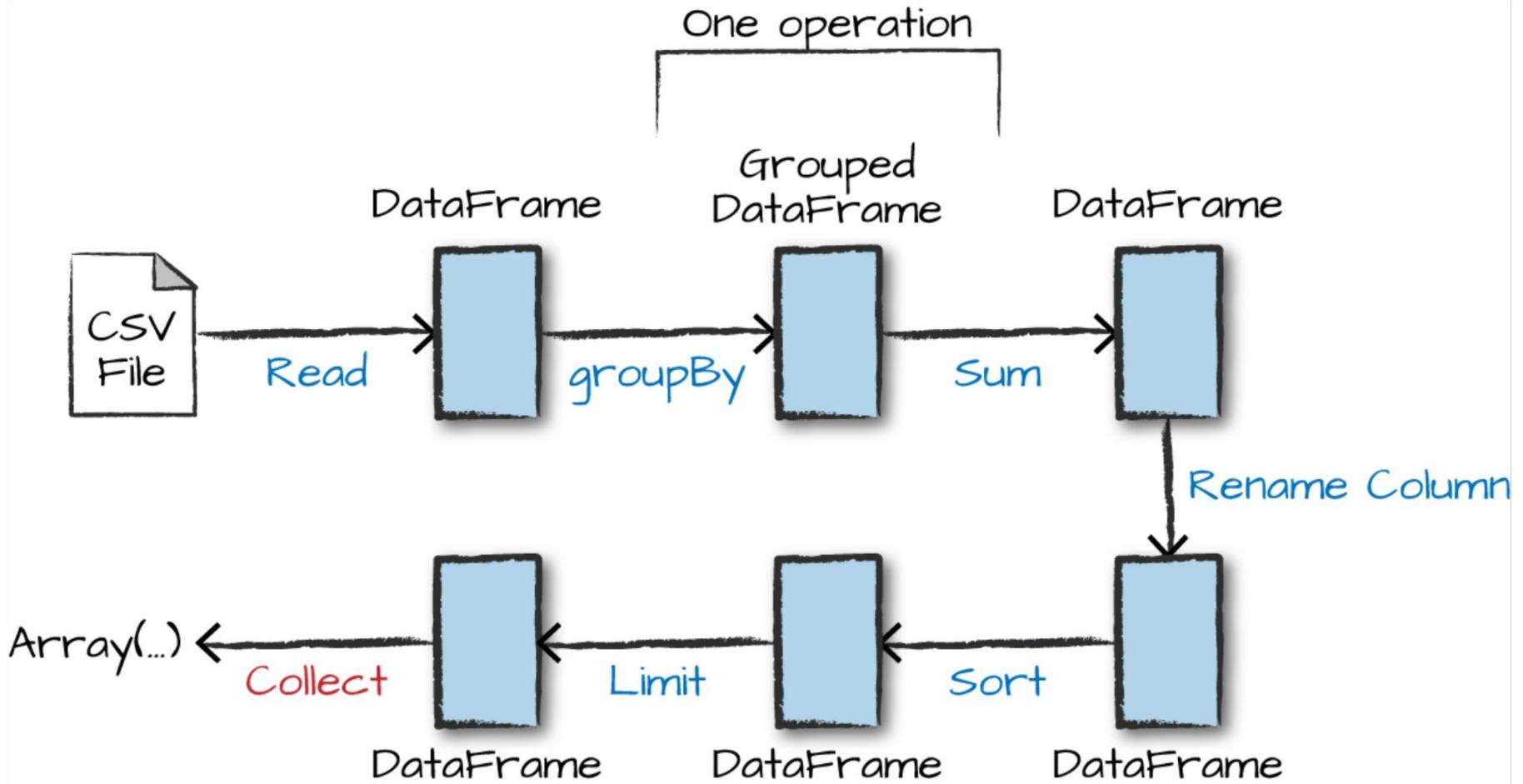


Wide transformations  
(shuffles) 1 to N



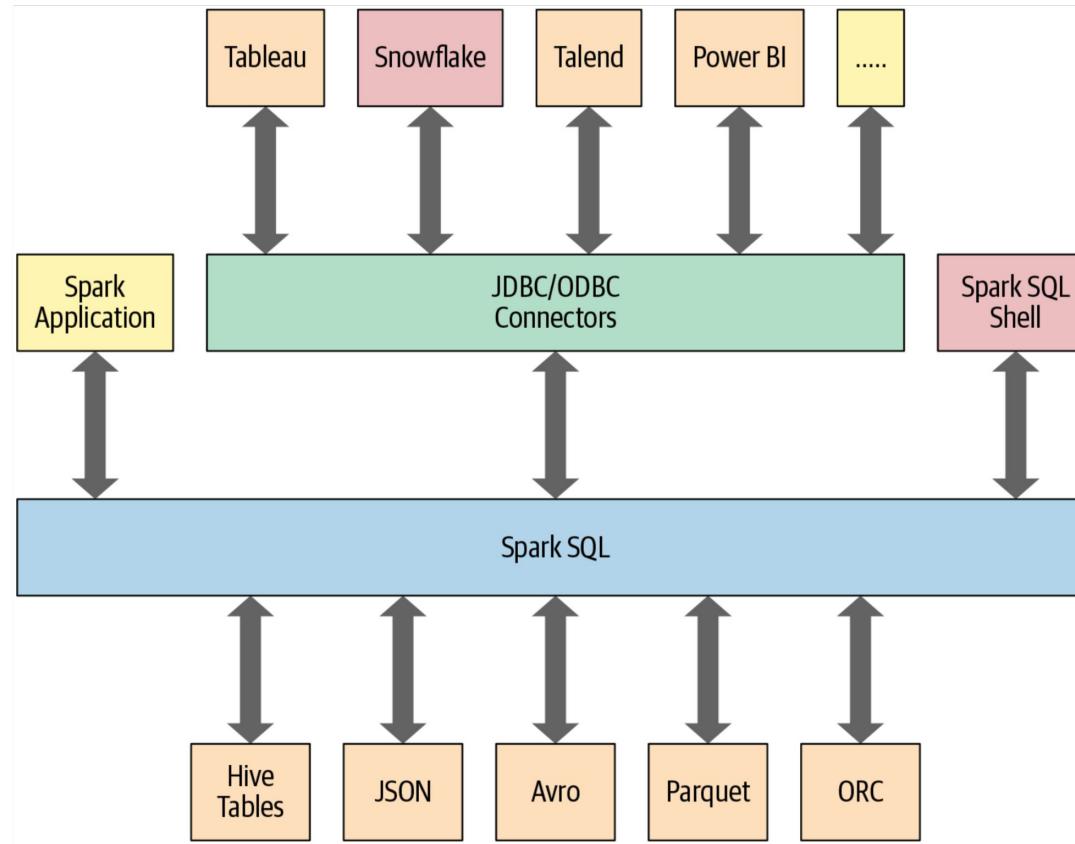






# SparkSQL

Source: [Learning Spark](#)



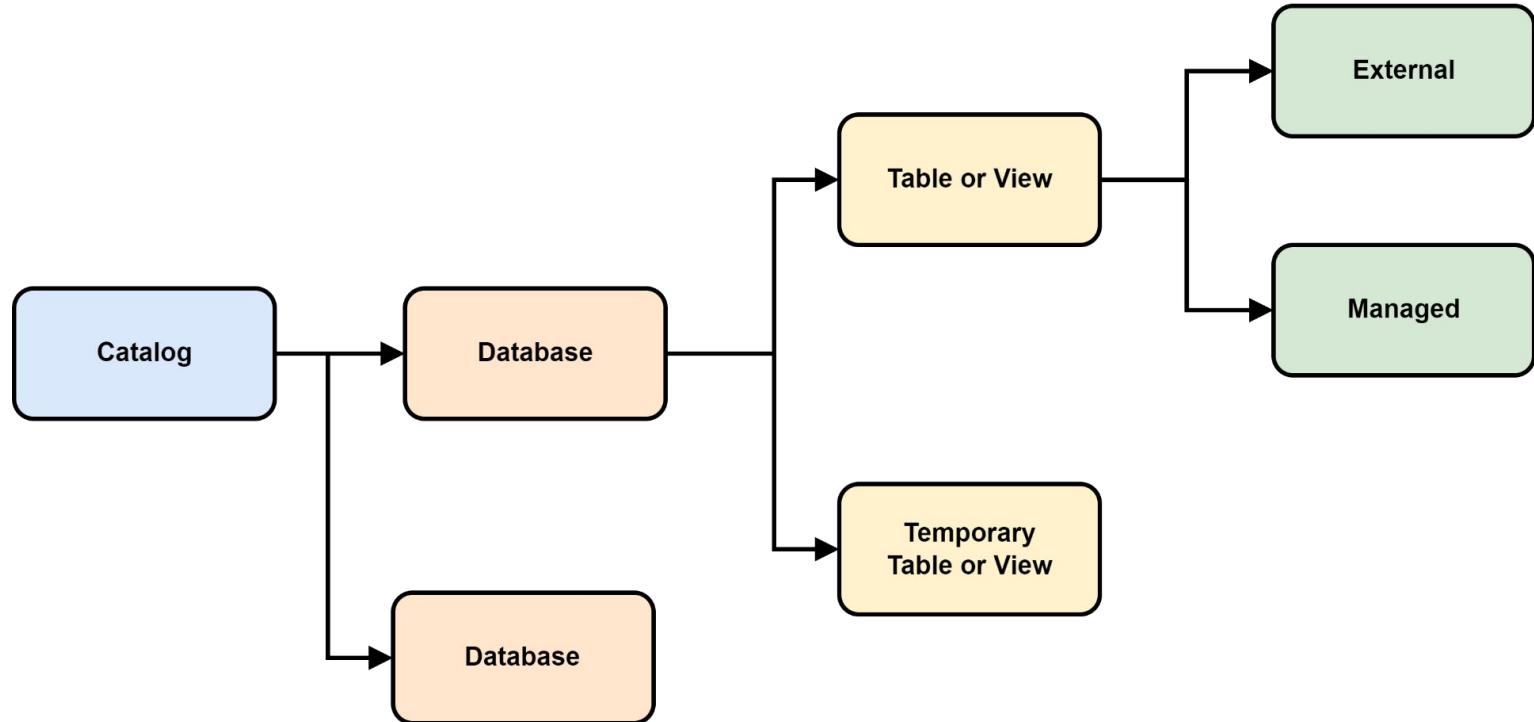
**spark.sql([SQL query])** → **DataFrame\***

How spark knows about databases and tables to process SQL ?

\* SQL command like INSERT does not return a dataframe

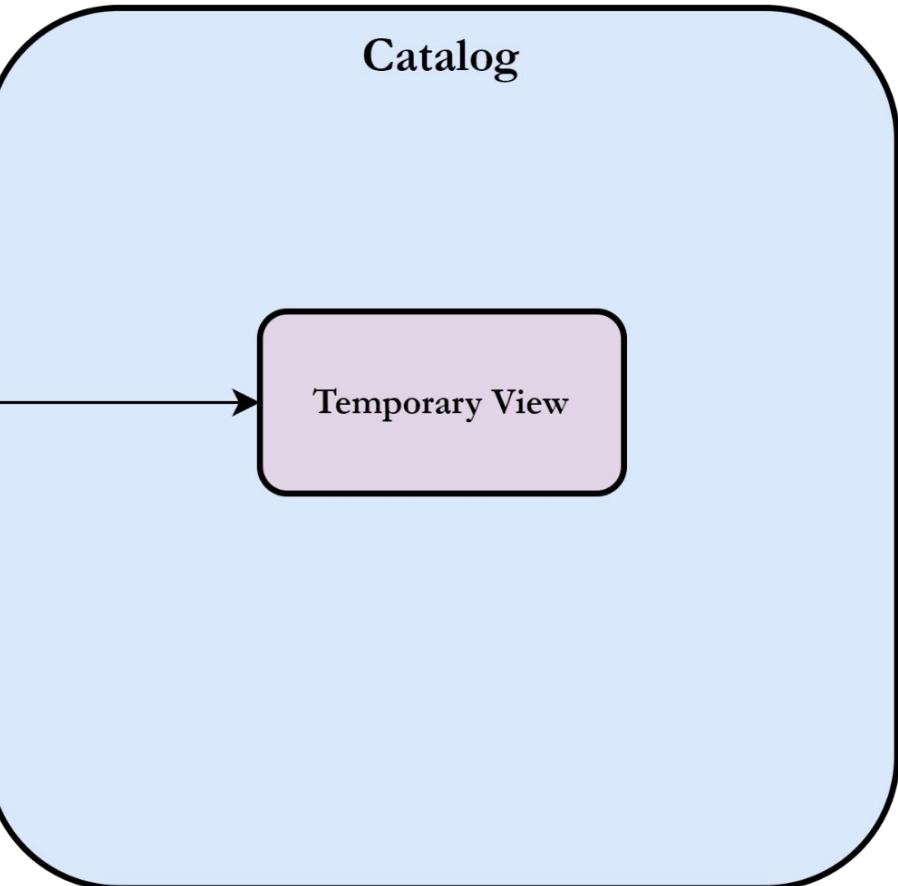
# Spark Catalog

Catalog is just a metadata how to use a stored data: access, read, write



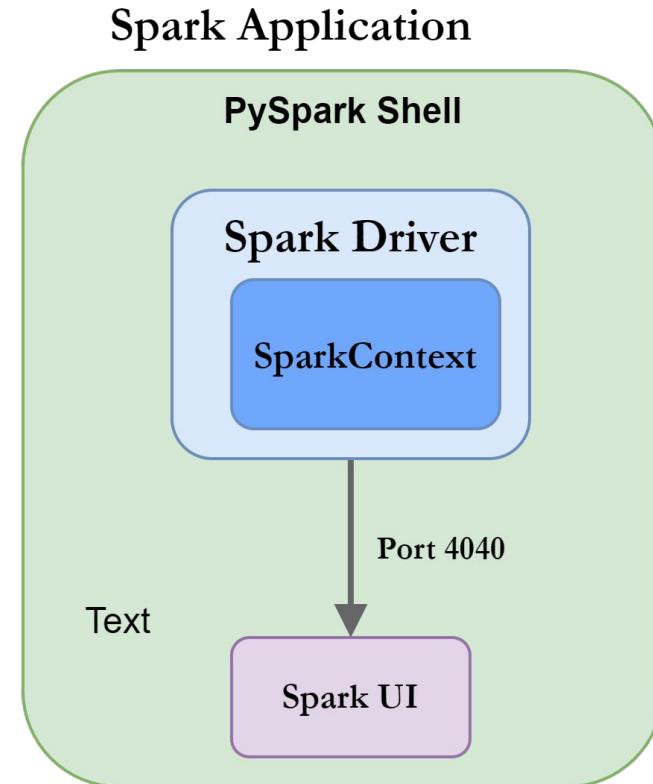
Catalog

`createOrReplaceTempView()`



# SparkUI

## Monitoring and Tracking How an **Application** interact with **Spark**



# Jobs

APACHE Spark 3.0.0-preview2

Jobs Stages Storage Environment Executors SQL

## Details for Job 0

Status: SUCCEEDED  
Associated SQL Query: 0  
Completed Stages: 1

▶ Event Timeline  
▼ DAG Visualization

Stage 0

```
graph TD; A[BatchScan] --> B[WholeStageCodegen]; B --> C["mapPartitionsInternal"];
```

Completed Stages (1)

# Stages

APACHE Spark 3.0.0-preview2 Jobs Stages Storage Environment Executors SQL

## Details for Stage 0 (Attempt 0)

Total Time Across All Tasks: 0.2 s  
Locality Level Summary: Process local: 1  
Associated Job Ids: 0

▼ DAG Visualization

Stage 0

```
graph TD; A[BatchScan<br/>DataSourceRDD [0]<br/>count at NativeMethodAccessorImpl.java:0] --> B[MapPartitionsRDD [1]<br/>count at NativeMethodAccessorImpl.java:0]; B --> C[WholeStageCodegen<br/>MapPartitionsRDD [2]<br/>count at NativeMethodAccessorImpl.java:0]; C --> D[mapPartitionsInternal<br/>MapPartitionsRDD [3]<br/>count at NativeMethodAccessorImpl.java:0]
```

▶ Show Additional Metrics

# Environment

Environment		livy-session-0 application UI
Name	Value	
Java Home	/home/gitpod/sdkman/candidates/java/11.0.19.fx-zulu	
Java Version	11.0.19 (Azul Systems, Inc.)	
Scala Version	version 2.12.15	
Spark Properties		
Name	Value	
spark.app.id	app-20230714040120-0000	
spark.app.initial.file.urls	spark://localhost:35007/files/py4j-0.10.9.5-src.zip.spark://localhost:35007/files/pyspark.zip	
spark.app.initial.jar.urls	spark://localhost:35007/jars/netty-codec-dns-4.1.86.Final.jar.spark://localhost:35007/jars/netty-transport-native-epoll-4.1.86.Final-linux-x86_64.jar.spark://localhost:35007/jars/netty-codec-memcache-4.1.86.Final.jar.spark://localhost:35007/jars/netty-transport-sctp-4.1.86.Final.jar.spark://localhost:35007/jars/netty-handler-4.1.86.Final.jar.spark://localhost:35007/jars/netty-resolver-dns-4.1.86.Final.jar.spark://localhost:35007/jars/ivy-repl_2.12-0.8.0-incubating-SNAPSHOT.jar.spark://localhost:35007/jars/kryo-shaded-4.0.2.jar.spark://localhost:35007/jars/netty-codec-http-4.1.86.Final.jar.spark://localhost:35007/jars/netty-handler-ssl-ocsp-4.1.86.Final.jar.spark://localhost:35007/jars/netty-resolver-dns-classes-macos-4.1.86.Final.jar.spark://localhost:35007/jars/netty-codec-mqtt-4.1.86.Final.jar.spark://localhost:35007/jars/netty-resolver-4.1.86.Final.jar.spark://localhost:35007/jars/netty-codec-redis-4.1.86.Final.jar.spark://localhost:35007/jars/netty-codec-haproxy-4.1.86.Final.jar.spark://localhost:35007/jars/netty-transport-udt-4.1.86.Final.jar.spark://localhost:35007/jars/netty-transport-native-kqueue-4.1.86.Final.jar.spark://localhost:35007/jars/netty-buffer-4.1.86.Final.jar.spark://localhost:35007/jars/netty-codec-xml-4.1.86.Final.jar.spark://localhost:35007/jars/netty-common-4.1.86.Final.jar.spark://localhost:35007/jars/commons-codec-1.9.jar.spark://localhost:35007/jars/ivy-thriftserver-session-0.8.0-incubating-SNAPSHOT.jar.spark://localhost:35007/jars/netty-transport-native-unix-common-4.1.86.Final.jar.spark://localhost:35007/jars/ivy-api-0.8.0-incubating-SNAPSHOT.jar.spark://localhost:35007/jars/netty-handler-proxy-4.1.86.Final.jar.spark://localhost:35007/jars/netty-resolver-dns-native-macos-4.1.86.Final-osx-aarch_64.jar.spark://localhost:35007/jars/netty-transport-4.1.86.Final.jar.spark://localhost:35007/jars/netty-codec-4.1.86.Final.jar.spark://localhost:35007/jars/netty-codec-socks-4.1.86.Final.jar.spark://localhost:35007/jars/netty-transport-native-kqueue-4.1.86.Final-osx-aarch_64.jar.spark://localhost:35007/jars/ivy-co-2.12-0.8.0-incubating-SNAPSHOT.jar.spark://localhost:35007/jars/netty-resolver-dns-native-macos-4.1.86.Final-osx-x86_64.jar.spark://localhost:35007/jars/minlog-1.3.0.jar.spark://localhost:35007/jars/objenesis-2.5.1.jar.spark://localhost:35007/jars/netty-transport-classes-epoll-4.1.86.Final.jar.spark://localhost:35007/jars/ivy-co-2.12-0.8.0-incubating-SNAPSHOT.jar.spark://localhost:35007/jars/netty-rsc-0.8.0-incubating-SNAPSHOT.jar.spark://localhost:35007/jars/netty-transport-rtx-4.1.86.Final.jar.spark://localhost:35007/jars/netty-codec-stomp-4.1.86.Final.jar.spark://localhost:35007/jars/netty-codec-smtp-4.1.86.Final.jar	
spark.app.name	livy-session-0	
spark.app.startTime	1689307279327	