



# SPARK BASICS SIMPLY EXPLAINED

## PART 1: HOW SPARK BOOSTS YOUR BIG DATA PROCESSING!

Data with  
Nikk the Greek





- *What are the features of Spark?*
- *How does it work lightning fast?*

Data with  
Nikk the Greek



# *The key features of* **APACHE Spark™**

# Simple. Fast. Scalable. Unified.

## Key features



### Batch/streaming data

Unify the processing of your data in batches and real-time streaming, using your preferred language: Python, SQL, Scala, Java or R.



### SQL analytics

Execute fast, distributed ANSI SQL queries for dashboarding and ad-hoc reporting. Runs faster than most data warehouses.



### Data science at scale

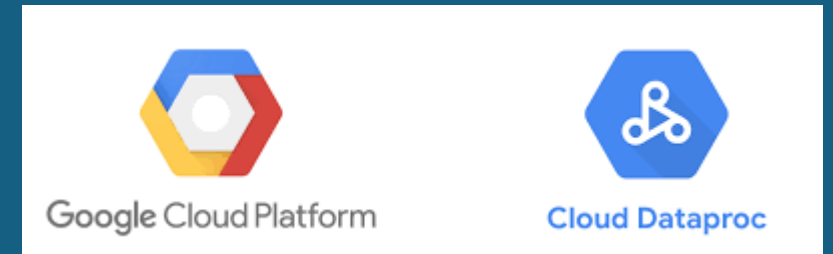
Perform Exploratory Data Analysis (EDA) on petabyte-scale data without having to resort to downsampling



### Machine learning

Train machine learning algorithms on a laptop and use the same code to scale to fault-tolerant clusters of thousands of machines.

# Cloud providers support Spark



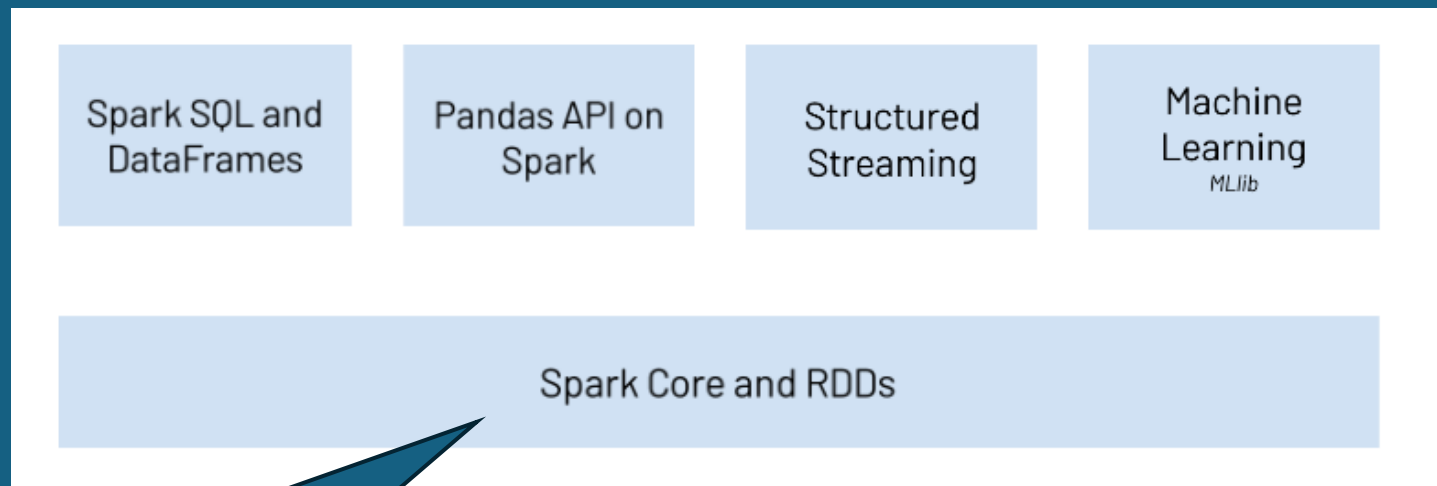
Spark SQL and  
DataFrames

Pandas API on  
Spark

Structured  
Streaming

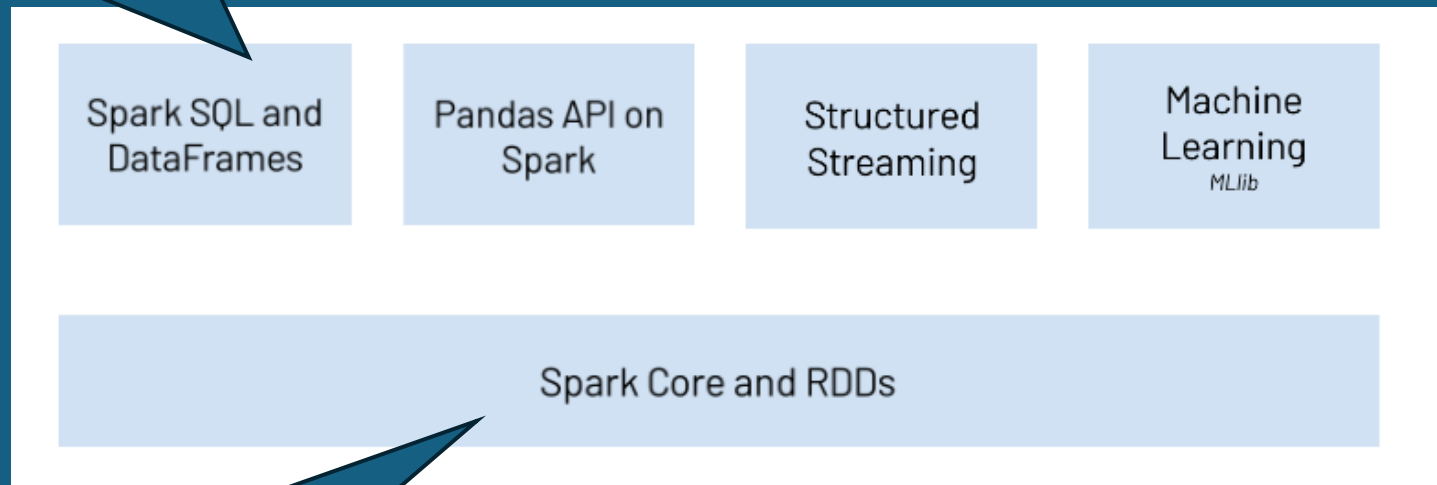
Machine  
Learning  
MLlib

Spark Core and RDDs



- Underlying execution engine
- It provides RDDs (Resilient Distributed Datasets) and in-memory computing capabilities with fault tolerance
- DataFrames are built on top of RDDs with multiple optimizations

- working with (semi-)structured data using DataFrames
- Leverage SQL warehouse capabilities with a hive metastore
- Uses lazy evaluation to leverage auto optimization using a Catalyst Optimizer

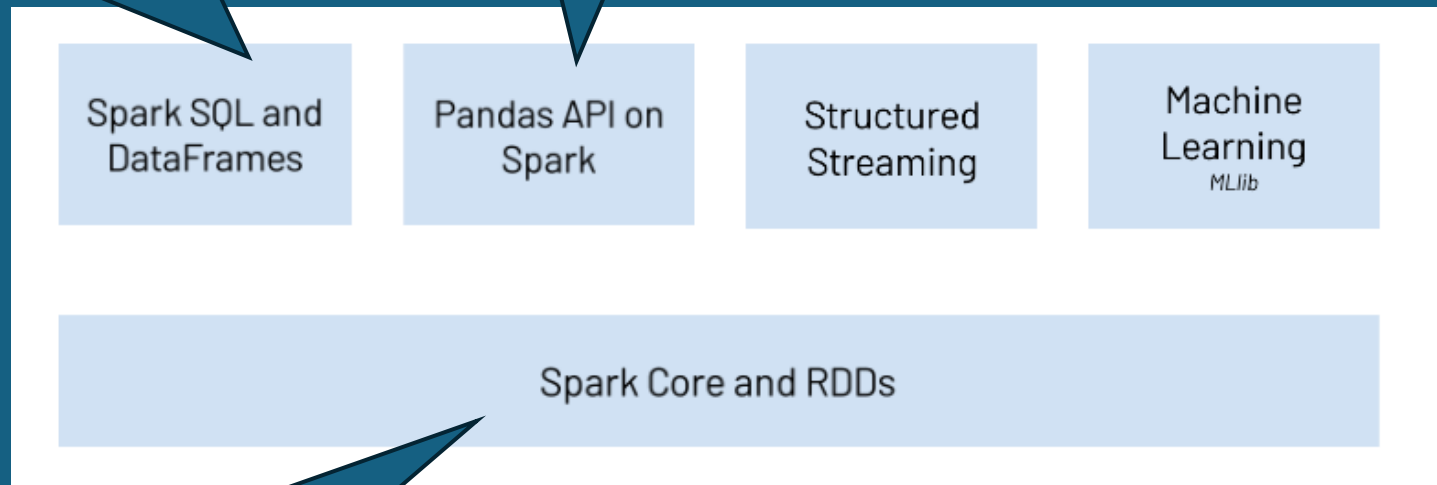


- Underlying execution engine
- It provides RDDs (Resilient Distributed Datasets) and in-memory computing capabilities with fault tolerance
- DataFrames are built on top of RDDs with multiple optimizations



- working with (semi-)structured data using DataFrames
- Leverage SQL warehouse capabilities with a hive metastore
- Uses lazy evaluation to leverage auto optimization using a Catalyst Optimizer

- Use your Pandas code directly on Spark without any modification

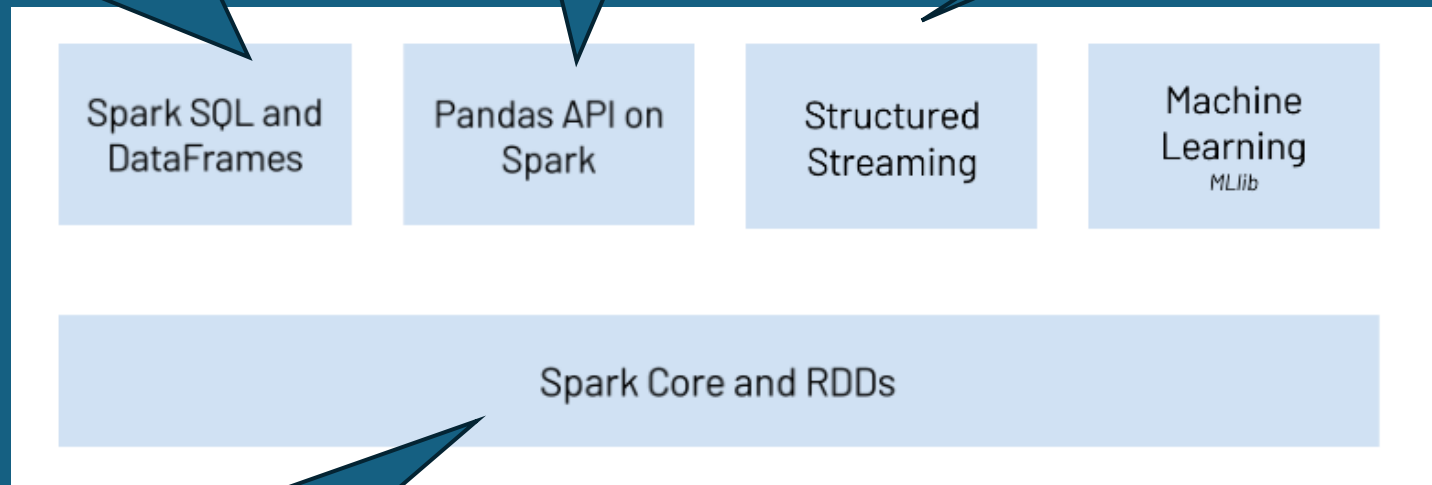


- Underlying execution engine
- It provides RDDs (Resilient Distributed Datasets) and in-memory computing capabilities with fault tolerance
- DataFrames are built on top of RDDs with multiple optimizations

- working with (semi-)structured data using DataFrames
- Leverage SQL warehouse capabilities with a hive metastore
- Uses lazy evaluation to leverage auto optimization using a Catalyst Optimizer

- Use your Pandas code directly on Spark without any modification

- stream processing engine built on the Spark SQL engine.

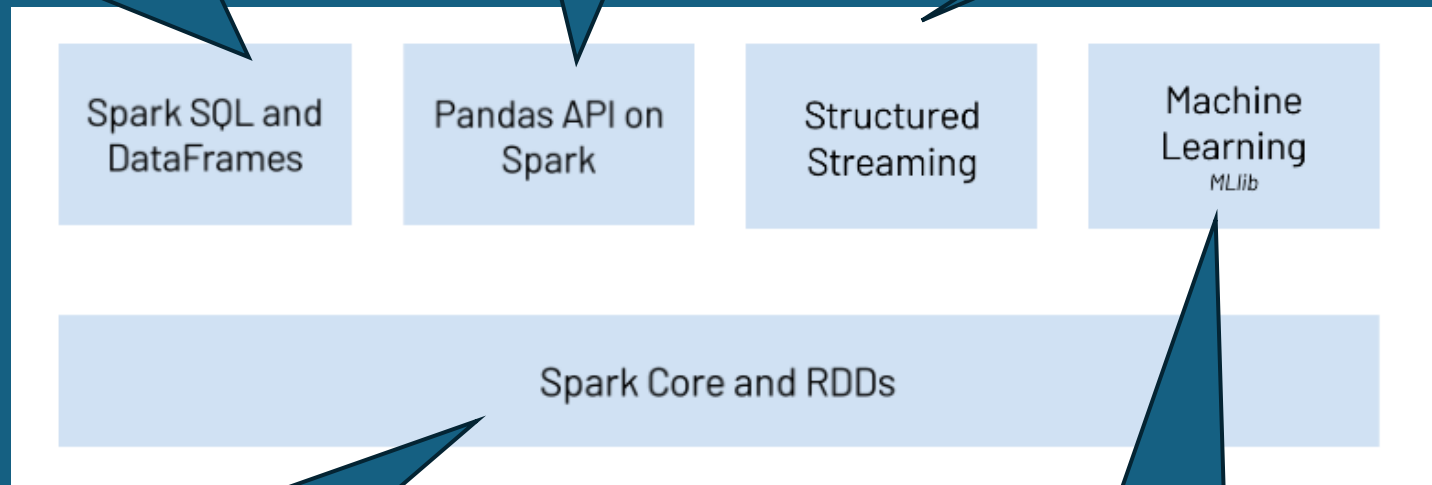


- Underlying execution engine
- It provides RDDs (Resilient Distributed Datasets) and in-memory computing capabilities with fault tolerance
- DataFrames are built on top of RDDs with multiple optimizations

- working with (semi-)structured data using DataFrames
- Leverage SQL warehouse capabilities with a hive metastore
- Uses lazy evaluation to leverage auto optimization using a Catalyst Optimizer

- Use your Pandas code directly on Spark without any modification

- stream processing engine built on the Spark SQL engine.



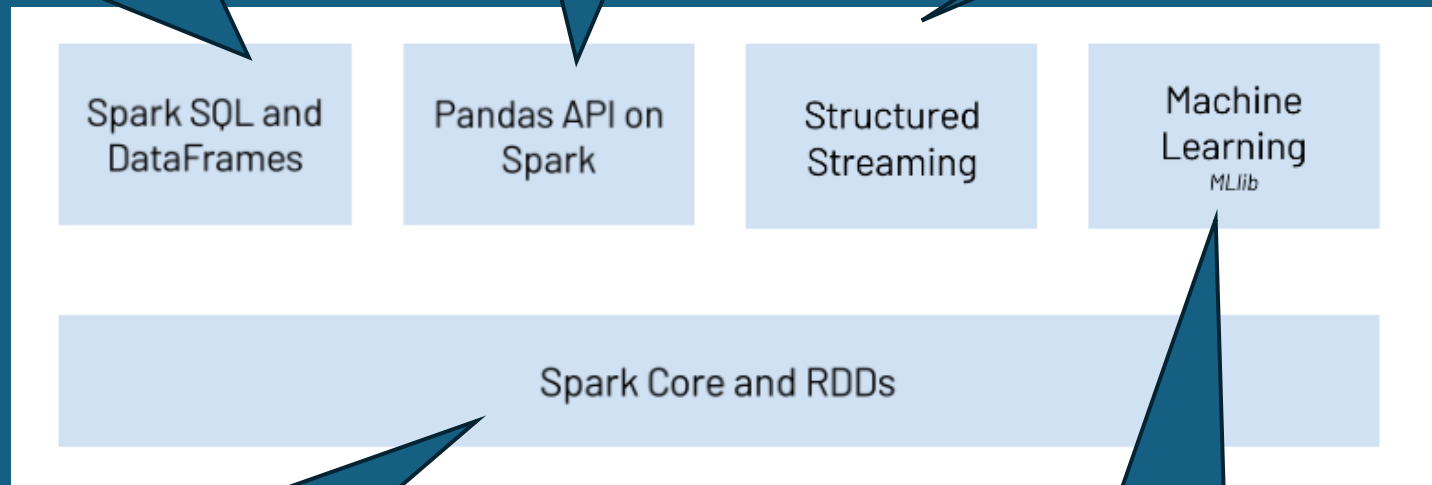
- Underlying execution engine
- It provides RDDs (Resilient Distributed Datasets) and in-memory computing capabilities with fault tolerance
- DataFrames are built on top of RDDs with multiple optimizations

- A scalable machine learning library built on top of Spark

- working with (semi-)structured data using DataFrames
- Leverage SQL warehouse capabilities with a hive metastore
- Uses lazy evaluation to leverage auto optimization using a Catalyst Optimizer

- Use your Pandas code directly on Spark without any modification

- stream processing engine built on the Spark SQL engine.



GraphX: Graph computation engine on Spark

- Underlying execution engine
- It provides RDDs (Resilient Distributed Datasets) and in-memory computing capabilities with fault tolerance
- DataFrames are built on top of RDDs with multiple optimizations

- A scalable machine learning library built on top of Spark

# Why is lightning fast?



*In short: The way Spark **distributes** and **optimizes** data processing in memory instead on storage like Hadoop*



*In a 100TB sorting Benchmark Spark 3x faster with 10x less resources compared to Hadoop*



*Spark can be up to 100x faster than Hadoop*

# Count a small dataset

Load data



col1	col2
1	A
2	B
3	A
4	B
5	B
6	B
7	A
8	A
9	B
10	B
11	A
12	A

Filter col1 < 10



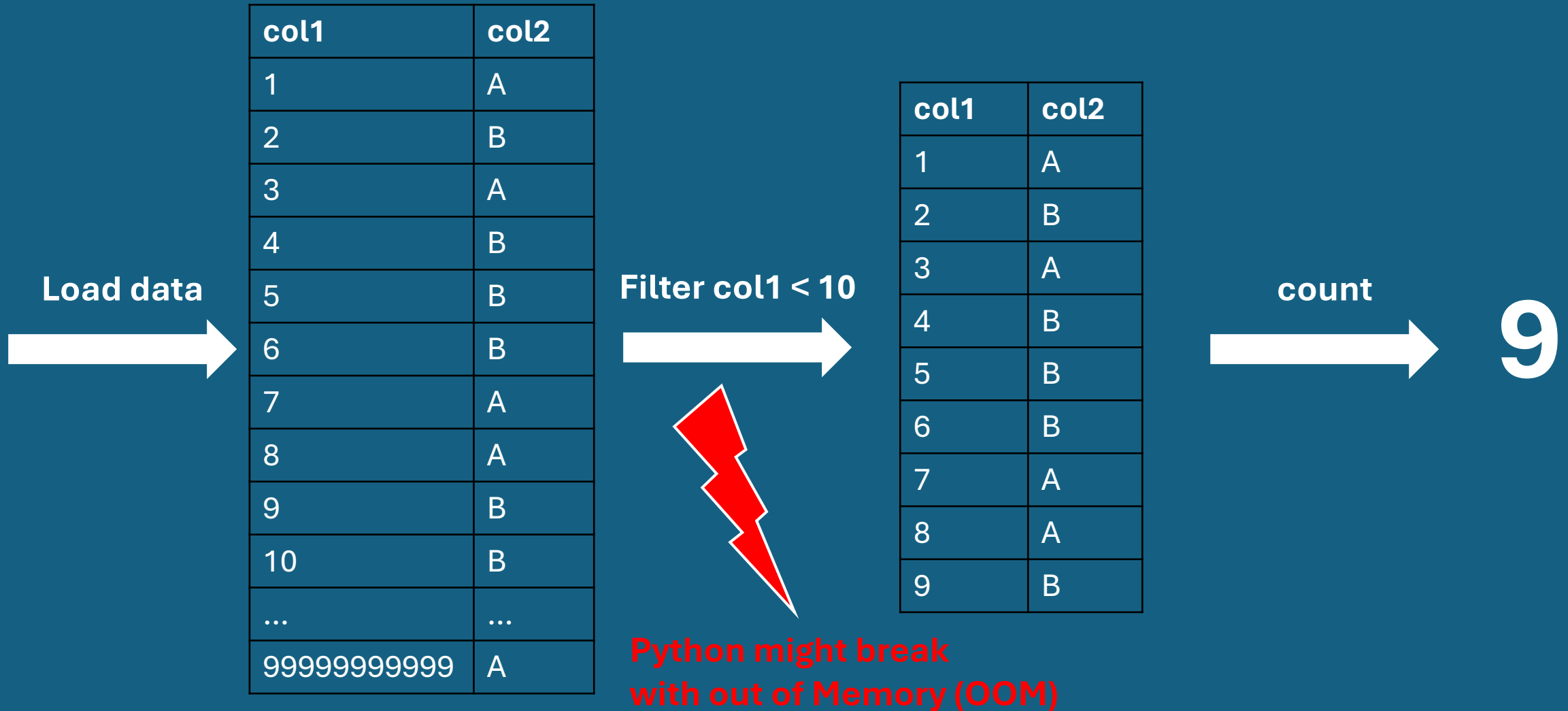
col1	col2
1	A
2	B
3	A
4	B
5	B
6	B
7	A
8	A
9	B

count



9

# Count a big dataset – Python Out of Memory



# Spark parallizes by splitting data into partitions

**Spark Partition**

col1	col2
1	A
2	B
3	A

col1	col2
4	B
5	B
6	B

col1	col2
7	A
8	A
9	B

col1	col2
10	B
...	...
999...	A

Load data



Filter col1 < 10



col1	col2
1	A
2	B
3	A

col1	col2
4	B
5	B
6	B

col1	col2
7	A
8	A
9	B

count



9

Adjust partition sizes  
based on your  
ressources available



# Spark Architecture

	col1	col2		
1	col1	col2		
2	4	7	col1	col2
3	5	8	A	B
6	9	11	B	A
	12		A	

...

**Driver**

**Executor 1**

**Core1**

**Core2**

**Executor 2**

**Core1**

**Core2**

# Spark Architecture

**Driver**

- One core executes one partition at the same time. This is called a **task**!
- Speed can be easily „scaled“ by the number of total cores

**Executor 1**

col1	col2
7	A
8	A
9	B

**Core1**

**Core2**

**Executor 2**

col1	col2
4	B
5	B
6	B

**Core1**

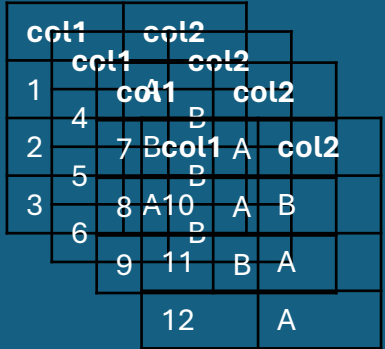
**Core2**

col1	col2
10	B
11	A
12	A

col1	col2
1	A
2	B
3	A

**Task**

# Spark Architecture



...

**Driver**

- One core executes one partition at the same time. This is called a **task**!
- Speed can be easily „scaled“ by the number of total cores

**Executor 1**

**Core1**

**Core2**

**Core3**

**Core4**

**Executor 2**

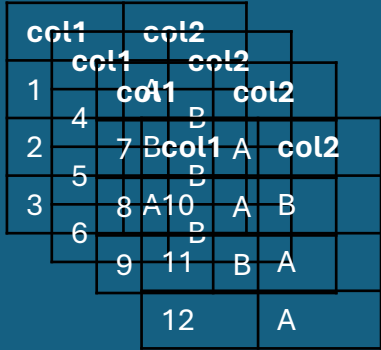
**Core1**

**Core2**

**Core3**

**Core4**

# Spark Architecture



...

**Driver**

- One core executes one partition at the same time. This is called a **task**!
- Speed can be easily „scaled“ by the number of total cores

**Executor 1**

**Core1**

**Core2**

**Executor 2**

**Core1**

**Core2**

**Executor 3**

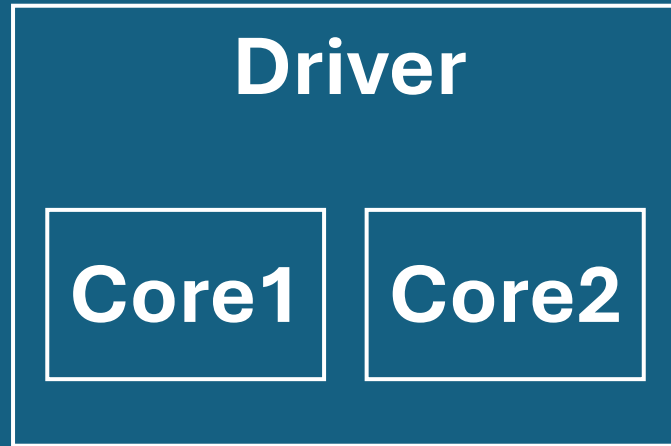
**Core1**

**Core2**

# Single Node Cluster

	col1	col2			
1					
2	4				
3	5	7	B		
	6	8	A	10	A
		9	11	B	A
			12		A

...



# Python vs. PySpark – Lazy Evaluation

```
# Load parquet data from a defined path
pdf = pd.read_parquet(path)

# Filter data
pdf = pdf[pdf["id"] > 10000]

# Count
pdf.shape[0]
```

- Python loads all data in the first line into Memory
- After this the filtering and count is performed
- For small datasets Python might be faster

```
# Load parquet data from a defined path
sdf = spark.read.parquet(path)

# Filter data Transformation
sdf = sdf.filter(f.col("id") > 10000)

# Count Action
sdf.count()
```

- Spark is loading, filtering and counting data „lazily“ only with the count (called action)
- Based on all identified steps the „Catalyst Optimizer“ finds the most efficient execution plan

# Summary

Supports multiple languages like Python, SQL, Scala, Java

Spark is well supported by cloud providers as quasi standard for big data

A Pandas API on Spark is available

ML, streaming and batch processing is supported

Spark Architecture is build by a driver and one or multiple executors with cores to execute tasks

High scalability and flexibility due to parallisation based on partitions executed within tasks on the available cores

Lazy Evaluation until an action like count is executed for optimizing the query execution plan using the Catalyst Optimizer