



Databricks PySpark

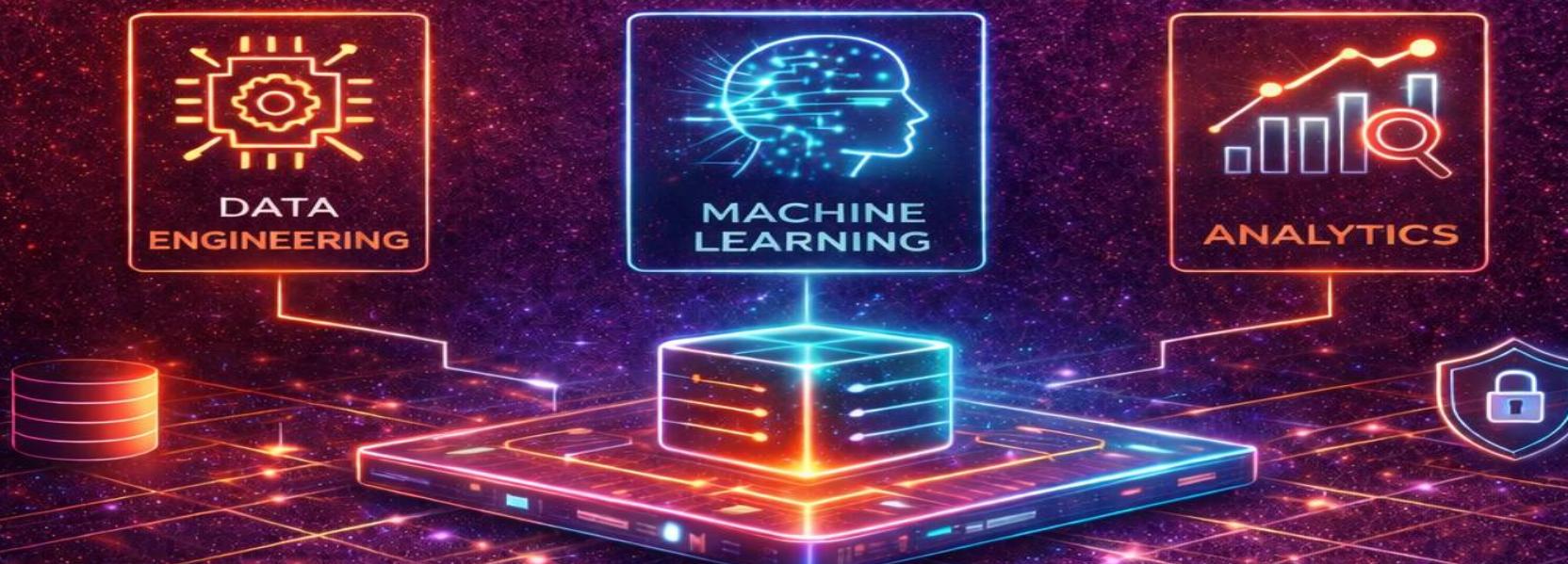
# Introduction to Databricks & PySpark



# Introduction to Databricks & PySpark



Is a **unified data analytics platform**  
built on Apache Spark



- ✓ Ingest, process, analyze, and govern data at scale
- ✓ Process data
- Analyze data

# Databricks Features



## Databricks – Core Platform Features

### Unified Analytics Platform



- Single platform for DE,
- DS, ML, & Analytics
- Eliminates silos

### Apache Spark Engine



- Optimized runtime
- Batch & Streaming

### Lakehouse Architecture



- Data Lake + Data Warehouse
- Structured & unstructured data
- Built on Delta Lake

### Delta Lake



**DELTA LAKE**

- ACID transactions
- Schema enforcement
- Time travel



# Databricks Features



## Databricks – Enterprise & AI Capabilities

### Data Engineering



- ETL / ELT pipelines
- Job scheduling
- Cloud storage integration

### Machine Learning & AI



- MLflow
- Model lifecycle
- Python, Scala, R

### Databricks SQL & BI



- Serverless SQL
- Dashboards
- Power BI & Tableau

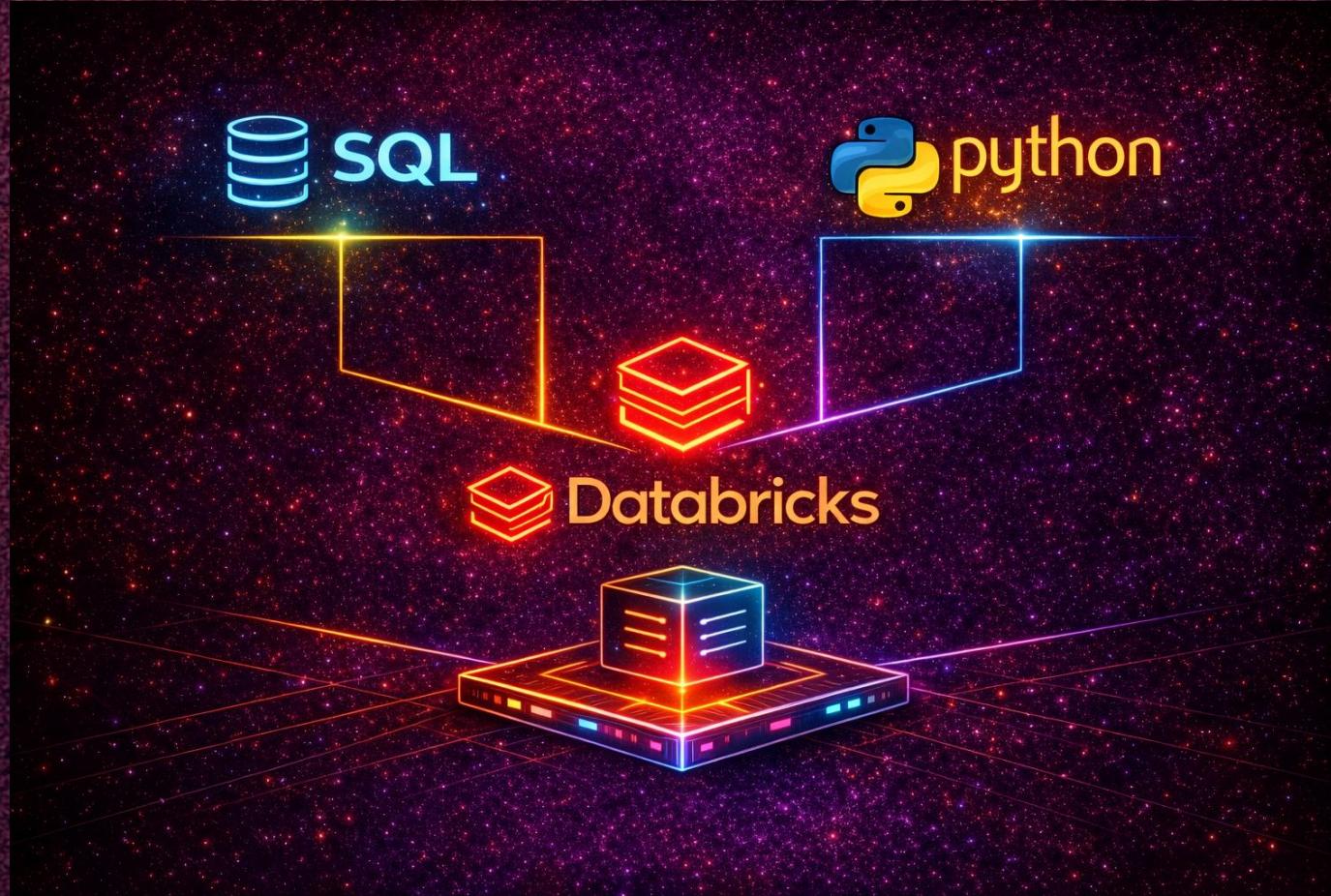
### Governance & Security



- Unity Catalog
- Fine-grained access
- Lineage & auditing



# Pre-Requisites



*Let's begin a great journey!*





# Introduction to Apache Spark

*Fast • Distributed • In-Memory Analytics Engine*

# Why Apache Spark?

- *In-memory processing*
- *Faster than MapReduce*
- *Scales to big data*
- *Batch & Streaming support*



# MapReduce Challenges



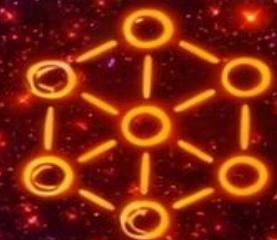
*Slow  
processing*



*Disk I/O  
bottleneck*



*Complex  
programming*



*Difficult to  
scale*

# The Evolution of APACHE SPARK

2009	2010	2013	2014	2015	2016
BIRTH OF SPARK	OPEN SOURCE RELEASE	APACHE INCUBATION	APACHE SPARK 1.0	DATAFRAMES & SPARK SQL	SPARK 2.0 (MAJOR MILESTONE)
UC Berkeley   AMPLab	Spark released as open source	Donated to Apache Software Foundation	Top-level Apache project	DataFrame API introduced	Unified APIs: RDD: DataFrame-Dataset
Faster than Hadoop MapReduce	Early academic & Startup adoption:	- Community-driven development begins	- Spark Core - Streaming, MLlib. - GraphX - Early Spark SQL	- Catalyst Optimizer added - SQL becomes first-class API	- Dataset API introduced - Structured Streaming introduced
Written in Scala	- RDD (Resilient Distributed Dataset)				




RDDs → DataFrames → Datasets      STRUCTURED STREAMING

# The Evolution of APACHE SPARK

2017–2018	2020	2020	2021–2022	2023	2024–2025
STREAMING & PERFORMANCE	SPARK 3.0	SPARK 3.0	LAKEHOUSE ERA	SPARK 3.4+	MODERN SPARK
<ul style="list-style-type: none"><li>Structured Streaming becomes stable</li><li>Event-time processing</li><li>Windowing &amp; watermarking</li><li>Cloud adoption increases</li></ul>	<ul style="list-style-type: none"><li>Adaptive Query Execution (AQE)</li><li>ANSI SQL compliance</li><li>Dynamic partition pruning</li></ul>	<ul style="list-style-type: none"><li>ANSI SQL compliance</li><li>Dynamic partition pruning</li><li>Faster joins &amp; queries</li></ul>	<ul style="list-style-type: none"><li>Deep cloud storage integration</li><li>ADLS • S3 • GCS</li><li>Delta Lake, Iceberg, Hudi</li><li>Foundation of Lakehouse architecture</li></ul>	<ul style="list-style-type: none"><li>Faster Python execution</li><li>Pandas API on Spark</li><li>Better observability</li></ul>	<ul style="list-style-type: none"><li>SQL, Streaming &amp; ML engine</li><li>Cloud-native analytics</li><li>AI feature pipelines</li><li>Industry standard for big data</li></ul>
					

2017–2018 STREAMING & PERFORMANCE      2020      SPARK 3.4+      MODERN SPARK

# Apache

# spark

## Core Components

**CDU**

Cloud And Data Universe

# Apache Spark

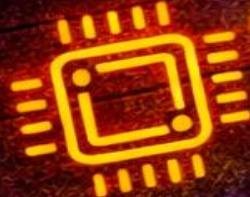
*Core Components*



Cloud And Data Universe

# Apache Spark

## Core Components



### SPARK CORE

- The foundation of Spark
- Distributed task scheduling
- Memory management
- Fault tolerance
- RDD abstraction
- Interaction with storage systems  
(HDFS, S3, ADLS, etc.)



### SPARK SQL

- SQL queries
- DataFrame & Dataset APIs
- Catalyst Optimizer
- ANSI SQL support
- Integration with Hive Metastore
- Analytics, BI workloads, ETL transformations

**CDU**

Cloud And Data Universe

# Apache Spark

## Core Components

### STRUCTURED STREAMING



- Modern streaming engine
- Event-time processing
- Windowing & watermarking
- Exactly-once guarantees
- Same API as batch (DataFrames)



### MLLIB

- Classification
- Regression
- Clustering
- Recommendation systems
- Feature engineering



### GRAPHX

- Graph processing framework
- Graph algorithms
- PageRank
- Connected components
- Social network analysis

**CDU**

Cloud And Data Universe

# Apache Spark

Language APIs



PYTHON PySpark



SCALA

Spark (native)



JAVA

Spark Java API



SparkR

CDU

Cloud And Data Universe

# Spark Architecture



# Spark Architecture



# What is a Spark Application?

A **Spark Application** is a single, complete program that runs on Apache Spark to process data.

In simple words: **your code + Spark engine + cluster resources** working to complete a task.

## — Formal Definition (Interview-Ready) —

- A **Spark application** is a **user program** that runs on Apache Spark,
- Consisting of **one driver program** and **one or more executor processes**,
- Used to perform **distributed data processing**.



# Driver Program

(Brain of Spark)



The **Driver** is where your Spark application starts.

## Responsibilities:

- ✓ Creates **SparkSession / SparkContext**
- ✓ Converts your code into a logical plan
- ✓ Splits jobs into **stages & tasks**
- ✓ Sends **tasks to executors**
- ✓ Collects **results**
- 📌 **Only ONE driver per Spark application**

# SparkContext / SparkSession



## Entry point to Spark

### Talks to:

- ✓ Cluster Manager
- ✓ Executors

### Tracks:

- ✓ RDDs
- ✓ DataFrames
- ✓ Jobs

In modern Spark, you mainly use **SparkSession**.

# Cluster Manager

## The Middle Layer in Spark

### Role of Cluster Manager

The Cluster Manager acts as the mediator between the **Driver** and the **Executors**.

- ✗ It does **NOT** process data.
- ✓ It manages resources.

### What the Cluster Manager Actually Does

#### 1 Driver requests resources

- ✓ CPU cores
- ✓ Memory



#### 2 Cluster Manager allocates executors

- ✓ Decides where executors run
- ✓ Decides how many executors are created



#### 3 Executors are launched on worker nodes

- ✓ Cluster Manager starts them
- ✓ Executors then register back with the Driver

After this, Driver talks directly to Executors  
Cluster Manager steps aside.

# Executors



## Executors

- ✓ JVM processes running on worker nodes
  - ✓ Execute tasks
  - ✓ Store data in **memory or disk**
  - ✓ Send results back to the driver
- 📌 Executors live only as long as the application runs



Cloud And Data Universe

# ★ Spark Tasks

## What is a Task?

- ✓ A task is the smallest unit of work in Spark
- ✓ Runs on a single data partition
- ✓ Executed by an executor



## Key Characteristics

- ✓ One task = one partition
- ✓ Tasks run in parallel across executors
- ✓ Tasks are **stateless**
- ✓ Tasks cannot be shared across executors

## How Tasks Are Created

- ✓ Triggered by a Spark action (`count()`, `collect()`, `save()`)
- ✓ A job is divided into **stages**
- ✓ Each stage contains multiple tasks



Cloud And Data Universe



# What is a Job in Spark?

A Spark Job is a unit of execution that is triggered when an action is called on a Spark DataFrame, Dataset, or RDD.

In simple words:

A job is the work Spark decides to execute when you ask for a result.

## Formal Definition (Interview-Ready)

- ✓ A Spark job is a computation triggered by an **action**, which Spark executes by dividing the work into multiple stages and tasks.



## How a Spark Job is Created

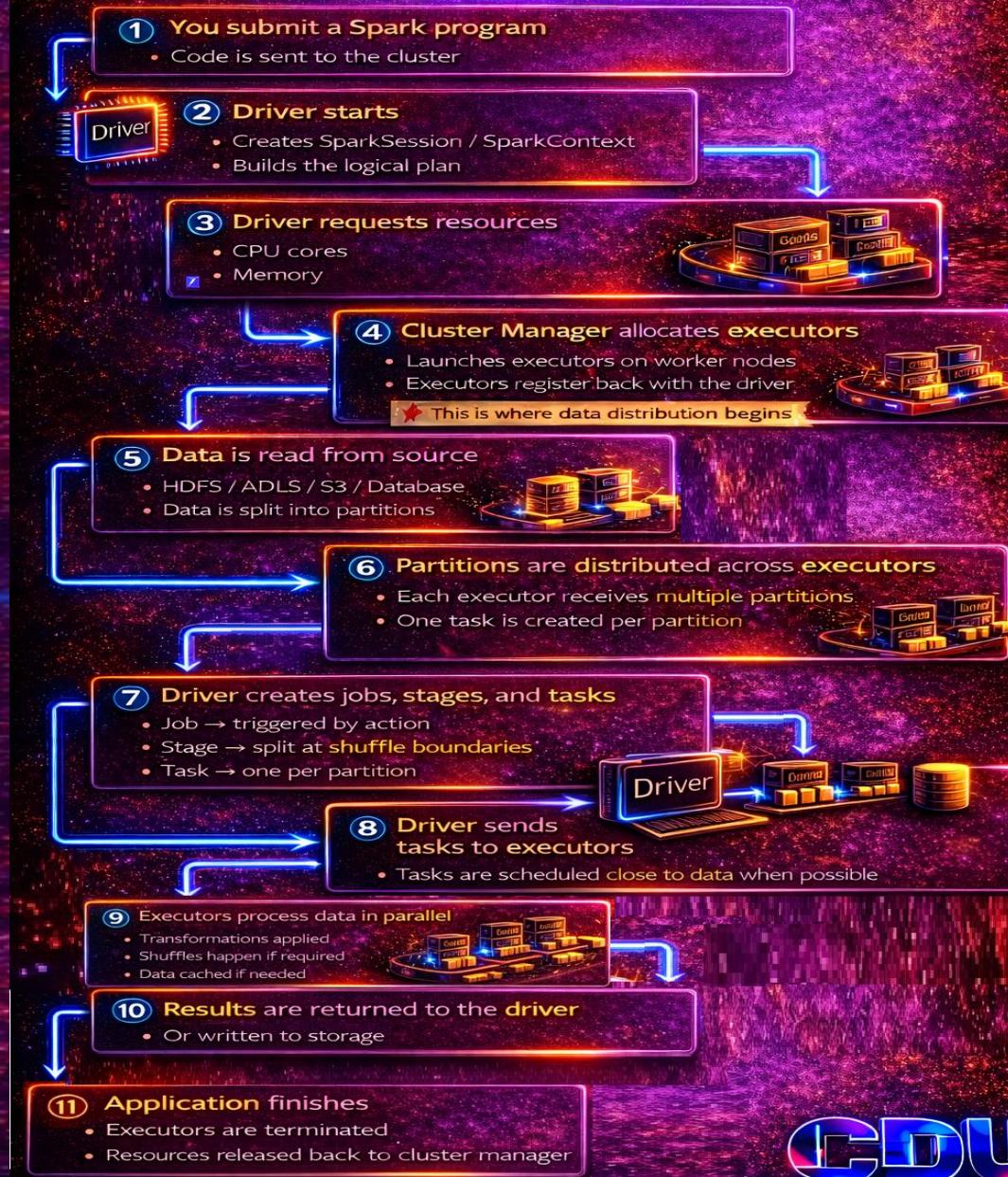
- 1 You write Spark transformations
- 2 Spark builds a logical plan (lazy execution)
- 3 You call an action  

- 4 Spark creates a job
- 5 Job is split into stages
- 6 Stages are split into tasks
- 7 Tasks run on executors



Cloud And Data Universe

# ❖ Spark Application Execution Flow



**Time for  
Practical Execution**





Databricks

# Assistant

Your AI Copilot



Cloud And Data Universe



# What Databricks Assistant does

*Understands Your Workspace*





# Data & analytics assistance

- ✓ Generate Databricks SQL queries
- ✓ Explain query execution plans
- ✓ Help with window functions, joins, aggregations
- ✓ Suggest performance improvements  
(AQE, caching, partitioning)

## Example:

- 💬 Optimize this Pyspark join"
- 💬 Convert this SQL query to  
PySpark DataFrame API"





# Code Help (Core Strength)

## You can ask it to:

- ✓ Write PySpark / Spark SQL / Python code
- ✓ Convert SQL ↔ PySpark
- ✓ Optimize slow Spark jobs
- ✓ Explain existing notebook code
- ✓ Fix errors and exceptions

## Example:

- “ Optimize this PySpark join”
- “ Convert this SQL query to PySpark DataFrame API”





# Platform & governance help

- ✓ Explain Unity Catalog
- ✓ Help with permissions & access control
- ✓ Answer questions on  
Delta Lake, time travel,  
schema enforcement





# ML & AI support

- ✓ Help with MLflow
- ✓ Feature engineering in Spark
- ✓ Model training & tracking examples.
- ✓ Explain notebooks used.  
in ML pipelines



# *Creating RDD from list*



# Understanding RDD Partitions



## ⚡ Parallelism

- Allows many tasks to be executed in parallel

## ⚙️ Optimization

- Number of partitions impacts performance

## ✓ Fault Tolerance

- Individual partitions can be recomputed if lost

# Managing PARTITIONS



## ⚡ Parallel Processing

- Optimizes resource usage

## 📦 Load Balancing

- Handles skewed data distribution

## \_PARTITION/Coalesce

- Scale partitions up or down using: `repartition()` / `coalesce()`

# Unity Catalog



— Cloud And Data Universe —

**We're at a transition point in Databricks.**

As DBFS is no longer the recommended approach,  
we need to adopt

# Unity Catalog

as the industry standard  
to manage files and tables.



## X Why DBFS Is No Longer the Best Option



- DBFS is workspace-scoped, not account-scoped



- Limited governance & access control



- Not designed for enterprise-grade data security



- Difficult to manage permissions across teams



- Not aligned with modern lakehouse governance



DBFS = legacy convenience, not enterprise standard

# ✅ Why Unity Catalog Is the Industry Standard



✓ Centralized data & file governance



✓ Fine-grained access control (catalog → schema → table)



✓ Works across all Databricks workspaces



✓ Integrated with cloud storage (ADLS / S3 / GCS)



✓ Built for security, compliance, and scalability



✓ **Unity Catalog = governance-first architecture**



# How We Should Manage Files & Tables Now

## ◆ Storage Layer



- ✓ Use external locations (ADLS / S3 / GCS)
  - Controlled via Unity Catalog

## ◆ Data Objects



- ✓ Managed Tables
- ✓ External Tables
- ✓ Volumes (for files)

## ◆ Access Control



- ✓ Centralized permissions
- ✓ Identity-based access
- ✓ Audit-ready

**“ Unity Catalog is not optional anymore  
— it’s the foundation of modern Databricks.**



**Over the next few videos you will learn basics  
on metastore, its objects and setting up Unity Catalog.**

**This is important before we move ahead  
in Databricks & PySpark.**

# Introduction to DataFrames

ID	Name	Salary

## DataFrame



# What are DataFrames in Spark?

A DataFrame is a **distributed collection** of data organized into named columns, similar to table in a database or a DataFrame in Pandas, built to work at **big-data scale**.

ID	Name	Salary
1	John	50000
2	Jane	60000
3	Mike	70000
4	Sarah	80000

DataFrame

👉 Internally, **DataFrames** are optimized and distributed across a **Spark cluster**.



# Key Characteristics of DataFrames

## 1 Schema-Based

- ✓ Data has **column names** and **data types**
- ✓ Enables validation and optimization

RDD

ID	Name	Salary
1	Samy	1000
2	John	2000
3	Mike	3000
4	Alice	4000

## 2 Distributed

- ✓ Data is split into **partitions**
- ✓ Processed in **parallel** across executors



## 3 Immutable

- ✓ Any transformation creates a new DataFrame

ID	Samy
1	1000
2	2000
3	3000
4	4000

ID	Salary
1	1000
2	2000
3	3000
4	4000

## 4 Optimized (Very Important)

- ✓ Catalyst Optimizer → query optimization
  - ✓ Tungsten Engine → memory & CPU optimization
- This makes DataFrames much faster than RDDs.



# DataFrame vs RDD

## (Interview Favorite)

Feature	DataFrame	RDD
Schema	Yes	No
Optimization	Automatic	Manual
Performance	Faster	Slower
API level	High-level	Low-level
SQL support	Yes	No



# Different Ways to Create a DataFrame in Spark

Broadly, DataFrames can be created from:

1



Files

(csv, json, parquet, delta)

2

Tables

ID	Name	Salary
1	John	50000
2	Jane	60000

3 RDDs



4

Databases  
(using jdbc)



Databases (using jdbc)



# RDD vs DataFrame vs Dataset

## — Which to Use When?

Scenario	Best Choice
Unstructured data	RDD
Structured data	DataFrame
SQL / analytics	DataFrame
Maximum performance	DataFrame
Type safety	Dataset
Python users	DataFrame
Fine-grained control	RDD



# Conclusion on DataFrames

**DataFrames** are the most powerful and production-ready abstraction in Spark — optimized, scalable, and built for real-world data workloads.



## Optimized

Leverage Catalyst & Tungsten for peak performance



## Scalable

Handle massive datasets with ease



## Production-Ready

Industry-proven and widely adopted



## Built for Real Workloads

Ideal for ETL, analytics, and reports



RDDs give control, **DataFrames** give power, and **optimization** makes the **difference**.



Cloud And Data Universe



# Create DataFrame from RDD



# Create DataFrame from CSV



**DataFrame**

<b>id</b>	<b>name</b>	<b>age</b>
1	Alice	30
2	Bob	25



— Cloud And Data Universe