

Spark Overview and Features

A deep dive into Apache Spark — the distributed data processing engine powering modern analytics, machine learning, and real-time data pipelines.



What is Apache Spark?

Apache Spark is a **distributed data processing engine**. Spark distributes data and computation across multiple machines.

It is designed for:



Large-scale data processing



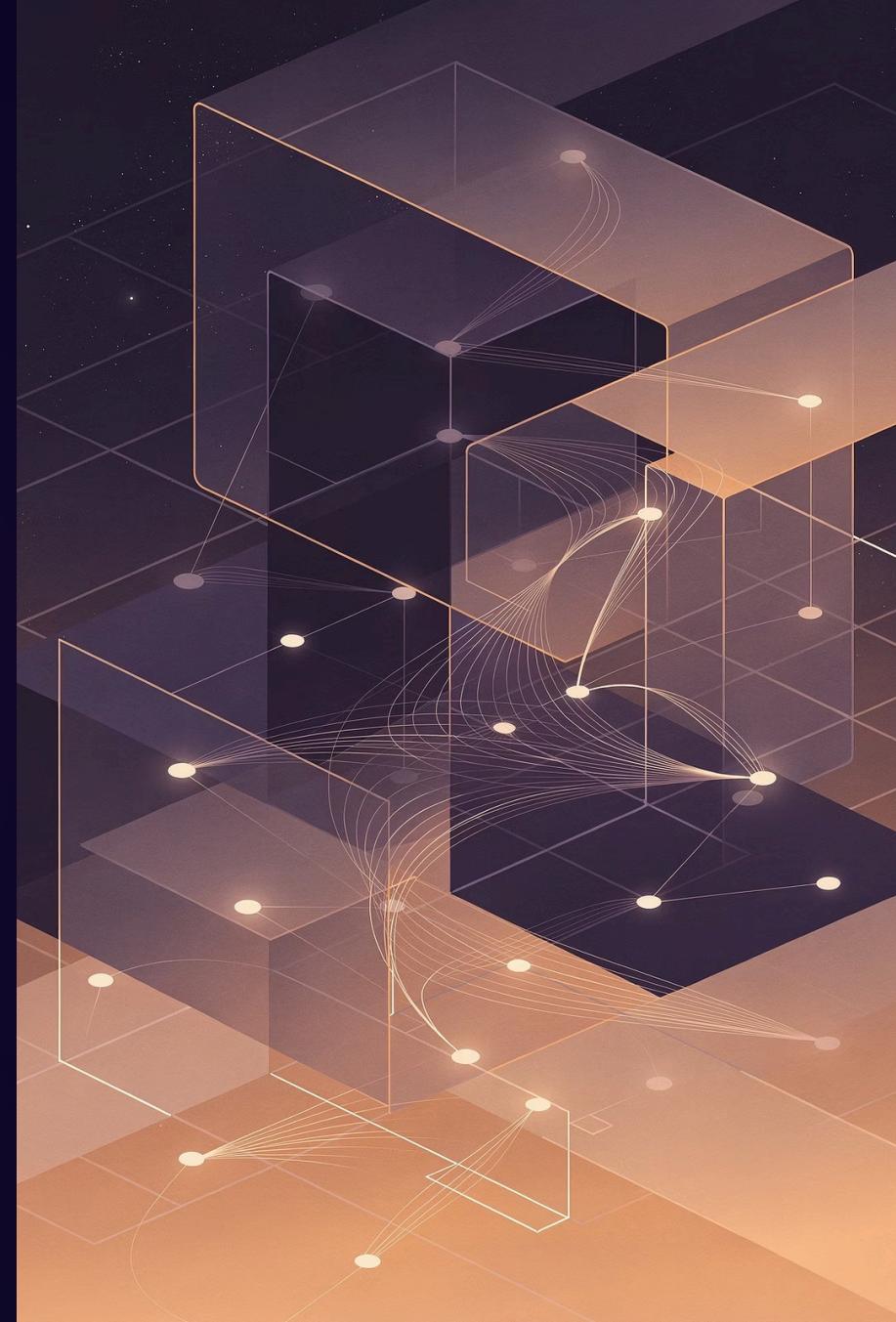
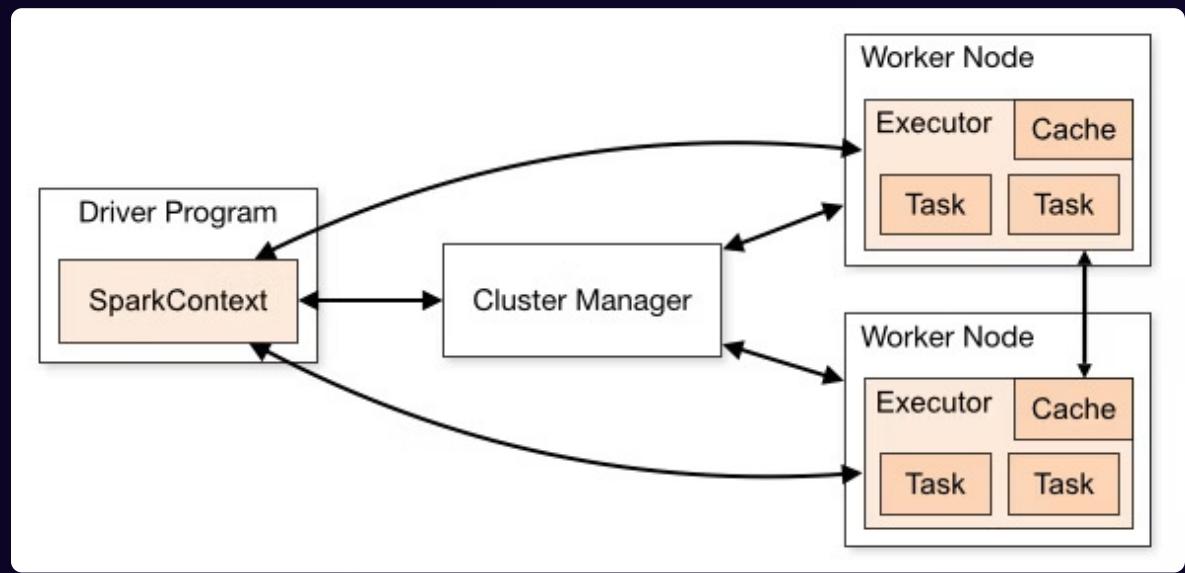
Parallel computation



Batch and real-time analytics



SQL and Machine Learning workloads



Why Spark Was Created

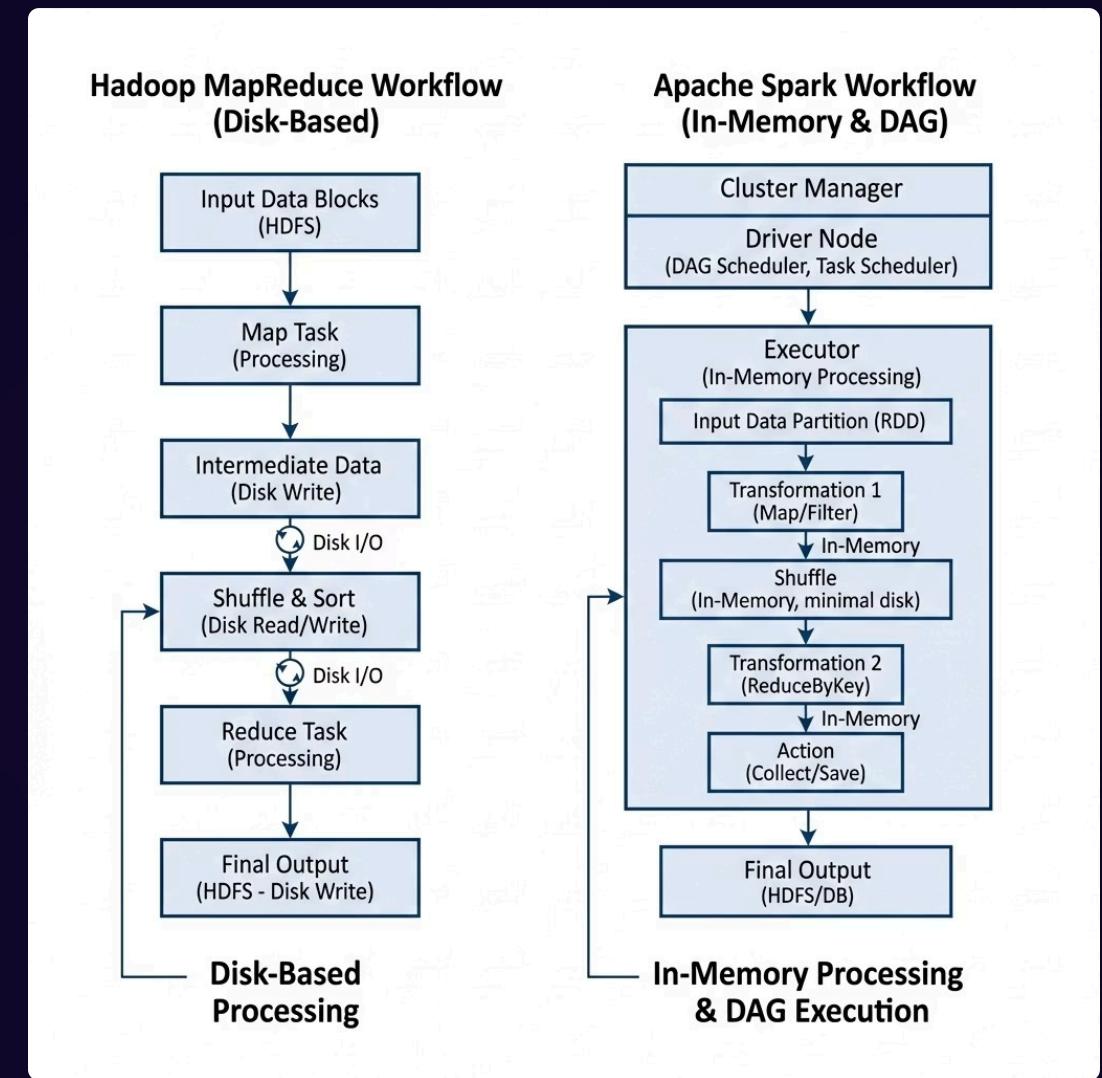
Before Spark, big data processing relied on **Hadoop MapReduce**. Spark was built to improve speed and simplify distributed computing.

Problems with MapReduce

- Disk-heavy processing
- High latency
- Complex programming model
- Separate engines for SQL, ML, Streaming

Spark's Solution

Spark was designed from the ground up to address every one of these pain points — delivering in-memory speed, a unified engine, and a developer-friendly API.



Core Features of Spark

In-Memory Processing

- Stores intermediate results in memory
- Reduces disk I/O
- Improves performance

Lazy Evaluation

- Builds execution plan first
- Executes only when action is triggered

More Spark Features



Fault Tolerance

Uses lineage to recompute lost partitions



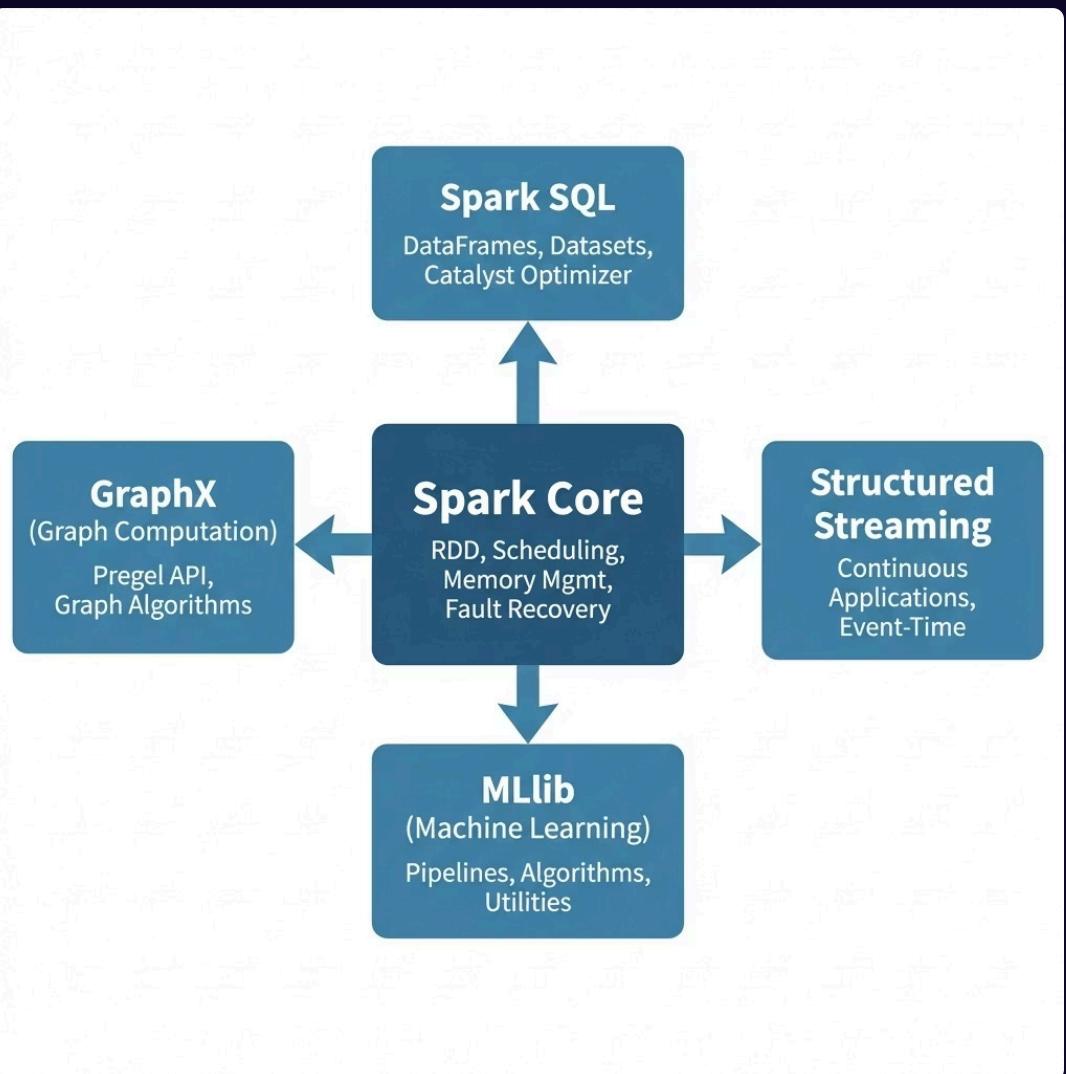
Scalability

Works on laptop, clusters, and in the cloud



Unified Engine

SQL · Streaming · ML · Graph processing





Importance of Spark in Data Processing

Modern data systems require a unified platform that handles every workload. **Spark provides all of this in one engine.**



Batch processing



Real-time analytics



Interactive queries



Machine learning pipelines

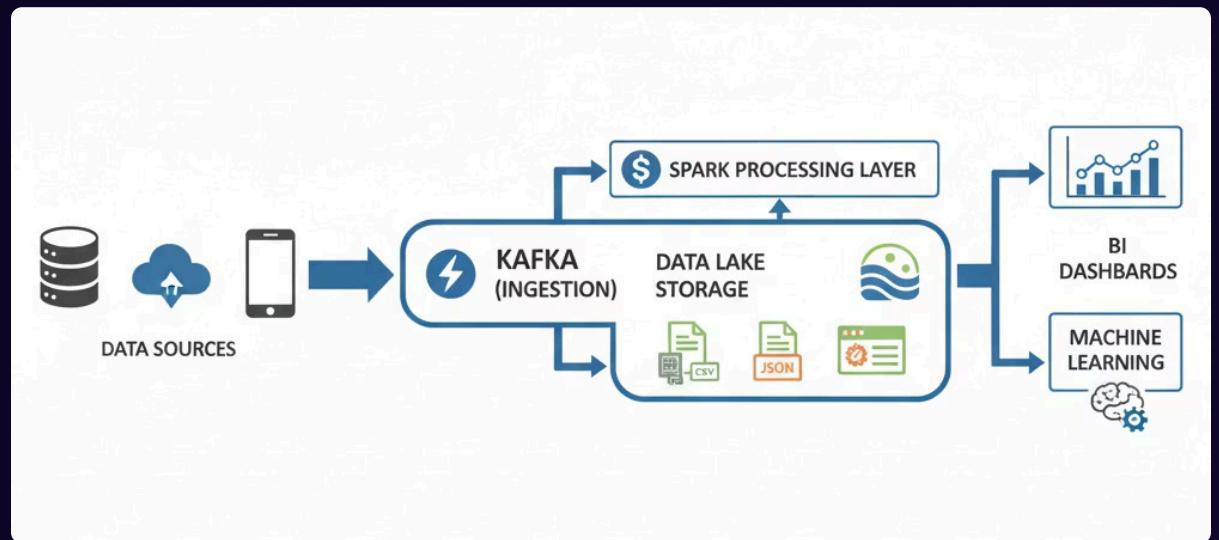
Where Spark is Used

Use Cases

- ETL pipelines
- Data lakes
- Data warehouses
- Streaming applications
- ML feature engineering

Spark integrates with:

- Kafka
- Hive
- Cloud storage
- Kubernetes



Hadoop vs Spark Stack

Hadoop Stack

01

HDFS

Storage

02

YARN

Resource Manager

03

MapReduce

Processing

Characteristics:

Disk-based

Batch-oriented

Slower for iterative tasks

Spark Stack

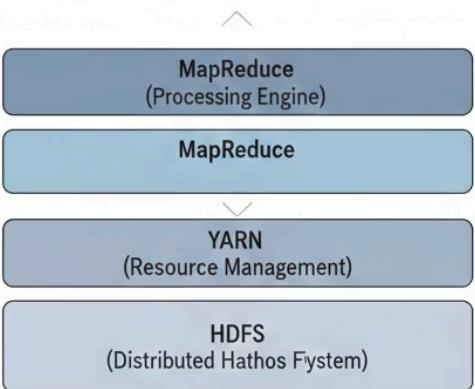
Spark Components

- Spark Core
- Spark SQL
- Structured Streaming
- MLlib
- GraphX

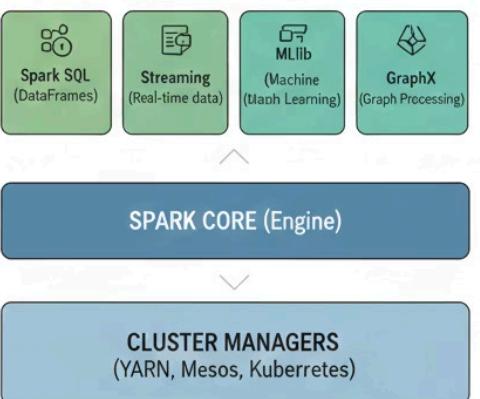
Runs on:

- HDFS
- Cloud Storage
- YARN
- Kubernetes

TRADITIONAL HADOOP STACK



MODERN SPARK STACK



Key Differences

Feature	Hadoop MapReduce	Spark
Processing	Disk-based	Memory-first
Speed	Slower	Faster
Iterative Jobs	Inefficient	Efficient
Streaming	Separate system	Built-in
Ease of Use	Complex	Developer-friendly

When to Use Hadoop vs Spark

Use Hadoop When:

- You need distributed storage (HDFS)
- Pure batch workloads
- Legacy ecosystem compatibility

Use Spark When:

- You need fast analytics
- You need streaming + batch
- You build ML pipelines
- You need interactive SQL

Summary

Distributed Analytics Engine

Spark is a distributed analytics engine.

Performance Improvement

It improves performance over Hadoop MapReduce.

Unified Workloads

It supports batch, streaming, SQL, and ML.

Modern Data Architectures

It is widely used in modern data architectures.

