

3. Benefits of Apache Spark & Why It's So Popular

Apache Spark is one of the most widely used big data processing frameworks because of its **speed, versatility, and ease of use**. It has gained massive popularity due to its ability to handle **batch processing, real-time analytics, machine learning, and graph processing** within a single ecosystem.

◆ Key Benefits of Apache Spark

1. Speed (100x Faster than Hadoop)

- Spark processes data **in-memory**, eliminating the need for frequent disk I/O.
 - Uses **Directed Acyclic Graph (DAG) scheduling** for optimized task execution.
 - Can achieve up to **100x speedup** over Hadoop's MapReduce.
-

2. Unified Framework for Multiple Workloads

- Supports **batch processing (Spark Core)**, **real-time streaming (Spark Streaming)**, **SQL-based querying (Spark SQL)**, **machine learning (MLlib)**, and **graph processing (GraphX)** within one framework.
 - Reduces the need for multiple tools, making it **simpler to develop and maintain**.
-

3. Ease of Use

- Provides high-level APIs in **Python (PySpark)**, **Scala**, **Java**, and **R**.
 - Allows developers to use **SQL queries (Spark SQL)** instead of complex programming.
 - Features **DataFrames and Datasets**, making it intuitive for data analysts and scientists.
-

4. Real-Time Processing Capabilities

- Unlike Hadoop MapReduce, which only supports batch processing, Spark **processes streaming data in near real-time** using **Spark Streaming**.
 - Works well with **Kafka, Flume, and Kinesis** for real-time event processing.
-

5. Works with Hadoop & Without Hadoop

- Can run **on top of Hadoop (HDFS, YARN)** for data storage and resource management.
 - Can also work **independently on cloud services (AWS, Azure, GCP) or Kubernetes**.
 - Supports **local mode**, allowing development on a laptop before deploying to a cluster.
-

6. Scalability & Fault Tolerance

- Can scale from a **single machine to thousands of nodes** in a cluster.
 - Provides **automatic recovery** from node failures using **RDD fault tolerance**.
-

7. Integration with Other Big Data Tools

- Works with **HDFS, Apache Hive, HBase, Cassandra, and Amazon S3** for storage.
 - Supports **machine learning libraries like TensorFlow, PyTorch, and MLflow**.
 - Can be integrated with **BI tools (Tableau, Power BI)** for visualization.
-

Why Is Spark So Famous?

- ✓ **Faster than Hadoop** due to in-memory computing.
- ✓ **Supports real-time data processing**, unlike Hadoop.
- ✓ **Easy to use** with SQL, Python, and DataFrames.
- ✓ **Unified framework** for batch, streaming, ML, and graph analytics.
- ✓ **Works on-prem and in the cloud**, making it flexible.

]4.What is Databricks?

Databricks is a **cloud-based data analytics and AI platform** built on **Apache Spark**. It provides a **managed and optimized version of Spark**, along with additional tools for data engineering, data science, machine learning, and real-time analytics.

Founded by the **creators of Apache Spark**, Databricks simplifies **big data processing** by offering a **fully managed, scalable, and collaborative environment** on cloud platforms like **AWS, Azure, and Google Cloud**.

◆ Key Offerings of Databricks Over Apache Spark

Feature	Apache Spark	Databricks
Setup & Management	Requires manual cluster setup & tuning	Fully managed & auto-optimized
Performance Optimization	Needs manual tuning for speed	Uses Photon Engine for better performance
Data Storage	Works with HDFS, S3, etc.	Supports Delta Lake for ACID transactions
Ease of Use	Requires scripting & CLI	Provides Notebooks with UI-based workflows
Collaboration	No built-in team collaboration	Supports real-time collaboration & versioning
Streaming Support	Spark Streaming	Enhanced Streaming with auto-scaling
Security & Compliance	Needs manual setup	Built-in role-based access control (RBAC) , encryption
Machine Learning	MLlib, TensorFlow, PyTorch	Databricks MLflow, AutoML, Model Serving
Cloud Support	Can be deployed manually	Fully managed on AWS, Azure, GCP

◆ Unique Features of Databricks

1. Managed Apache Spark

- No need to manually set up and tune clusters.
- **Auto-scaling and performance tuning** improve speed and efficiency.

2. Delta Lake (Better Than Parquet or HDFS)

- **Adds ACID transactions** to Spark's data lake.
- Enables **versioning, schema enforcement, and time travel**.

3. Photon Engine (Performance Boost)

- Optimized execution engine that runs **faster than open-source Spark**.
- Improves **SQL and ML workloads** with better CPU efficiency.

4. Collaborative Notebooks

- Built-in **Jupyter-like Notebooks** for Python, SQL, Scala, and R.
- Supports **team collaboration, version control, and automated workflows**.

5. Unified Data & AI Platform

- Combines **big data, machine learning, and analytics** in one platform.
- Integrates with **BI tools like Tableau, Power BI, and Looker**.

6. Security & Compliance

- Built-in **access control, encryption, and governance**.
 - Compliant with **GDPR, HIPAA, and SOC 2** standards.
-

Why Choose Databricks Over Apache Spark?

- ✓ **Fully managed** → No need to manually configure Spark.
- ✓ **Faster performance** → Photon Engine and Delta Lake.
- ✓ **Easier collaboration** → Shared notebooks & team workflows.
- ✓ **Advanced machine learning** → MLflow, AutoML, and Model Serving.
- ✓ **Better governance** → Security, access control, and compliance.

Requirement for spark setup on window

- 1.JDK 8 or JDK 11
- 2.Python 3.6 or higher
- 3.Hadoop Winutils
- 4.Spark binaries
- 5.Environment variables (JAVA_HOME, SPARK_HOME, HADOOP_HOME, Path)
- 6.Python IDE (PyCharm, VS Code, or Jupyter Notebook)

1. What is a Transformation Function?

A **Transformation** is a **lazy operation** that applies a function to an RDD or DataFrame and returns a new RDD/DataFrame **without executing immediately**. Spark optimizes transformations by building a **Directed Acyclic Graph (DAG)** and executes them only when an action is triggered.

✅ Key Points:

- Transformations are **lazy (delayed execution)**.
- They create a new RDD/DataFrame from an existing one.
- Execution happens only when an **action** is called.
- Types: **Narrow Transformations** (no shuffle) & **Wide Transformations** (shuffle involved).

💡 Example:

```
python  
CopyEdit  
rdd2 = rdd.map(lambda x: x * 2) # No execution yet (Lazy)
```

2. What is an Action Function?

An **Action** triggers the **execution of the DAG** and returns a result to the driver or writes data to an external storage. Unlike transformations, actions **force computation** and return values.

✅ Key Points:

- Actions **trigger execution** of transformations.
- They **return values** or **write results** to external storage.
- Without actions, Spark **does not execute any transformations**.

💡 Example:

```
python  
CopyEdit  
result = rdd2.collect() # Triggers execution and returns data
```

3. List of Important Spark Functions

🔥 Transformations (Lazy Operations - Create New DataFrames/RDDs)

Function	Type	Description	Example
<code>map()</code>	Narrow	Applies function to each element	<code>rdd.map(lambda x: x*2)</code>
<code>filter()</code>	Narrow	Filters elements based on condition	<code>rdd.filter(lambda x: x > 10)</code>
<code>flatMap()</code>	Narrow	Splits elements into multiple outputs	<code>rdd.flatMap(lambda x: x.split(" "))</code>
<code>distinct()</code>	Wide	Removes duplicate elements	<code>rdd.distinct()</code>
<code>groupByKey()</code>	Wide	Groups elements by key (causes shuffle)	<code>rdd.groupByKey()</code>
<code>reduceByKey()</code>	Wide	Aggregates values by key	<code>rdd.reduceByKey(lambda a, b: a + b)</code>
<code>sortByKey()</code>	Wide	Sorts RDD based on key	<code>rdd.sortByKey()</code>
<code>mapValues()</code>	Narrow	Applies function only to values (key unchanged)	<code>rdd.mapValues(lambda x: x+1)</code>
<code>union()</code>	Wide	Combines two RDDs	<code>rdd1.union(rdd2)</code>
<code>intersection()</code>	Wide	Returns common elements	<code>rdd1.intersection(rdd2)</code>
<code>subtract()</code>	Wide	Returns elements in <code>rdd1</code> but not in <code>rdd2</code>	<code>rdd1.subtract(rdd2)</code>
<code>cartesian()</code>	Wide	Computes Cartesian product of two RDDs	<code>rdd1.cartesian(rdd2)</code>
<code>repartition(n)</code>	Wide	Increases partitions (shuffle involved)	<code>rdd.repartition(4)</code>
<code>coalesce(n)</code>	Narrow	Reduces partitions (avoids shuffle)	<code>rdd.coalesce(2)</code>

⚡ Actions (Trigger Execution - Return Values or Data)

Function	Description	Example
----------	-------------	---------

<code>collect()</code>	Returns all elements to driver	<code>rdd.collect()</code>
<code>count()</code>	Returns the number of elements	<code>rdd.count()</code>
<code>first()</code>	Returns the first element	<code>rdd.first()</code>
<code>take(n)</code>	Returns first <code>n</code> elements	<code>rdd.take(5)</code>
<code>reduce()</code>	Aggregates elements using function	<code>rdd.reduce(lambda a, b: a + b)</code>
<code>foreach()</code>	Applies function to each element (no return)	<code>rdd.foreach(print)</code>
<code>show()</code>	Displays DataFrame content	<code>df.show()</code>
<code>countByKey()</code>	Counts elements per key	<code>rdd.countByKey()</code>

Optimization Functions (Caching, Broadcasting, and Partitioning)

Function	Description	Example
<code>cache()</code>	Caches RDD/DataFrame in memory	<code>rdd.cache()</code>
<code>persist(level)</code>	Stores RDD in memory/disk with levels	<code>rdd.persist(StorageLevel.MEMORY_AND_DISK)</code>
<code>unpersist()</code>	Removes cached data from memory	<code>rdd.unpersist()</code>
<code>checkpoint()</code>	Saves RDD to disk to prevent recomputation	<code>rdd.checkpoint()</code>
<code>broadcast()</code>	Broadcasts a small variable to executors	<code>sc.broadcast(my_dict)</code>
<code>accumulator()</code>	Shared variable for aggregation	<code>acc = sc.accumulator(0)</code>

DataFrame-Specific Functions

Function	Description	Example
<code>select()</code>	Selects specific columns	<code>df.select("name", "age")</code>
<code>filter()</code> / <code>where()</code>	Filters rows based on condition	<code>df.filter(df.age > 25)</code>
<code>groupBy()</code>	Groups by column values	<code>df.groupBy("department").count()</code>
<code>agg()</code>	Aggregates data	<code>df.groupBy("dept").agg(avg("salary"))</code>
<code>orderBy()</code>	Sorts rows	<code>df.orderBy(df.salary.desc())</code>
<code>withColumn()</code>	Adds/modifies a column	<code>df.withColumn("new_col", df["age"] + 5)</code>
<code>drop()</code>	Drops a column	<code>df.drop("age")</code>
<code>join()</code>	Joins two DataFrames	<code>df1.join(df2, "id", "inner")</code>

Spark SQL Functions

Function	Description	Example
<code>createOrReplaceTempView()</code>	Creates a temporary SQL table	<code>df.createOrReplaceTempView("employees")</code>
<code>sql()</code>	Runs an SQL query	<code>spark.sql("SELECT * FROM employees WHERE age > 30")</code>

Window Functions (Advanced Aggregations)

Function	Description	Example
<code>row_number()</code>	Assigns a unique row number	<code>row_number().over(Window.partitionBy("dept").orderBy("salary"))</code>

<code>rank()</code>	Assigns rank (handles ties)	<code>rank().over(Window.partitionBy("dept").orderBy("salary"))</code>
<code>dense_rank()</code>	Similar to <code>rank()</code> , but without gaps	<code>dense_rank().over(Window.partitionBy("dept").orderBy("salary"))</code>
<code>lead()</code>	Gets next row's value	<code>lead("salary", 1).over(Window.partitionBy("dept").orderBy("salary"))</code>
<code>lag()</code>	Gets previous row's value	<code>lag("salary", 1).over(Window.partitionBy("dept").orderBy("salary"))</code>