# Class 1: Introduction to Big Data and Hadoop

## 1. Introduction to Big Data

Definition of Big Data

Big Data refers to datasets that are too large, complex, and dynamic to be processed by traditional data processing tools. It's characterized by the three Vs: Volume, Velocity, and Variety.

**Volume:** Think of Big Data as the massive scale of data, often exceeding petabytes. For instance, consider the colossal amount of global social media interactions happening every second.

**Velocity:** This aspect pertains to the breakneck speed at which data is generated and ingested. A prime example is the constant stream of Twitter posts or the rapid influx of sensor data from Internet of Things (IoT) devices.

**Variety: Big Data encompasses diverse data types, including structured, semi-structured, and unstructured data. For instance, think of customer reviews (unstructured) and financial spreadsheets (structured).**

**Why Big Data Matters**

**Understanding the significance of Big Data in today's world is crucial:**

Improved Decision-Making: **Big Data empowers better decision-making. Think about optimizing supply chains based on real-time demand data or making precise financial forecasts.**

Predictive Analytics: **It enables predictive modeling. Consider a telecom company predicting customer churn to proactively retain customers.**

Personalized Recommendations: **Companies like Amazon and Netflix use Big Data to provide tailored recommendations, enhancing user engagement and boosting sales.**

**Challenges in Traditional Data Processing**

**Traditional data processing methods face limitations:**

Legacy Systems: **Dealing with legacy systems is challenging as they struggle to handle the onslaught of Big Data. Imagine a company hampered by outdated infrastructure as data continues to grow.**

Slower Processing Times: **Large datasets can lead to frustratingly slow processing times, impacting business operations. Think of a retail company grappling with delayed sales reporting.**

Scalability Issues: **Scaling traditional databases is no easy task. Picture a company that encountered difficulties when trying to expand its database infrastructure.**

High Infrastructure Costs: **Maintaining vast data centers can be cost-prohibitive. Visualize a company that saved significantly on infrastructure costs by embracing Big Data technologies.**

**Real-World Example: Netflix's Recommendation System**

**Netflix harnesses Big Data for its recommendation system:**

Data Collection: **Netflix collects and processes diverse user data, including viewing history and user ratings.**

Personalization Algorithms: **It employs intricate algorithms for personalized recommendations, utilizing collaborative filtering and content-based filtering.**

Impact on Business: **Netflix's recommendation system significantly boosts user retention and engagement.**

2. Introduction to Hadoop

**What is Hadoop?**

**Hadoop is an open-source framework designed for processing Big Data. It comprises core components like HDFS, MapReduce, and YARN.**

Addressing Big Data Challenges: **Hadoop is adept at addressing challenges such as scalability, fault tolerance, and cost-effectiveness. It accomplishes this by distributing data and tasks across a cluster.**

**Real-World Example: Facebook's Data Analysis**

**Facebook utilizes Hadoop for data analysis:**

Data Sources: **Facebook collects a diverse array of data, including user interactions and posts.**

Hadoop at Facebook: **The company employs Hadoop clusters in its data processing pipelines and analytics applications.**

Impact on Facebook: **Hadoop-based data analysis has yielded valuable insights, enhancing the Facebook platform.**

## 3. Hadoop Use Cases

**Real-World Applications of Hadoop**

**Discover how various industries leverage Hadoop:**

Finance: **Hadoop aids in fraud detection and risk analysis, benefitting financial institutions.**

Healthcare: **It's instrumental in patient analytics and drug discovery, improving patient care.**

Retail: **Hadoop optimizes supply chain management and inventory forecasting, enhancing efficiency.**

**Benefiting from Scalability and Cost-Efficiency**

**Hadoop's scalability and cost-effectiveness make it a valuable asset for organizations:**

Scalability: **Hadoop scales horizontally by adding more nodes, facilitating data growth.**

Cost-Efficiency: **It utilizes clusters of commodity hardware, reducing infrastructure costs.**

**Real-World Example: Airbnb's Data Analysis**

**Airbnb employs Hadoop for data analysis:**

Data Sources: **They collect diverse data, including host and guest information, with the challenge of handling this massive dataset.**

Analytics Applications: **Airbnb leverages Hadoop for various analytics tasks, improving user experience through data-driven insights.**

## 4. Hadoop Ecosystem

**Overview of the Hadoop Ecosystem**

**Explore the broader Hadoop ecosystem, including core components like HBase, Hive, Pig, and Spark:**

Core Components: **Understand how each component fits into the Hadoop ecosystem and their roles.**

**Real-Time Processing with Spark**

**Dive deeper into Apache Spark as a pivotal Hadoop ecosystem component:**

Real-Time Processing: **Spark's prowess in real-time data processing and streaming capabilities.**

**Real-World Example: Twitter's Real-Time Analytics with Spark**

**See how Twitter uses Apache Spark for real-time analytics:**

Data Streams: **Twitter's processing of real-time data streams and handling immense tweet volumes.**

Spark's Efficiency: **Spark's speed and efficiency advantages in specific use cases.**

Linux and Hadoop commands are related in the sense that Hadoop, which is primarily designed to run on Unix-like operating systems including Linux, provides its own set of command-line utilities and tools for managing and interacting with Hadoop clusters. These Hadoop commands are typically run within a Linux terminal or shell. Here's how they are connected:

1. **Linux as the Operating System**: Hadoop is often deployed on Linux-based operating systems because of their stability, scalability, and compatibility with Hadoop's design principles. Many Hadoop distributions are optimized for Linux environments.
2. **Terminal or Shell**: In Linux, users interact with the operating system through a terminal or shell. This is where you can execute both Linux commands and Hadoop commands.
3. **Hadoop Commands**: Hadoop provides its own set of command-line utilities for tasks such as managing the Hadoop Distributed File System (HDFS), running

MapReduce jobs, and interacting with Hadoop services. Some commonly used Hadoop commands include `hadoop fs` for HDFS operations, `hadoop jar` for running MapReduce jobs, and various `hdfs` and `yarn` commands for cluster management.

4. **Linux Commands for General Tasks**: In addition to Hadoop-specific commands, Linux provides a wide range of general-purpose commands for tasks like file manipulation, process management, user administration, and networking. Users may use Linux commands to navigate the file system, manage permissions, and perform other tasks related to Hadoop cluster administration.

Here are a few examples to illustrate the connection:

- To navigate to the Hadoop configuration directory, you might use a Linux command like `cd /etc/hadoop`.
- To list files in HDFS, you would use a Hadoop command like `hadoop fs -ls /user`.
- To create a directory in HDFS, you would use `hadoop fs -mkdir /user/new_directory`.
- To stop a Hadoop service like the ResourceManager, you might use `yarn rmadmin -shutdown`.

So, Linux is the underlying operating system on which Hadoop is typically installed, and it provides the environment in which you run both general Linux commands and Hadoop-specific commands to manage and work with Hadoop clusters and data. Familiarity with Linux commands is essential for Hadoop cluster administrators and users.

These commands will help you navigate the file system, manage files and directories, and perform common tasks:

1. `pwd` (Print Working Directory):
   - Shows the current directory (folder) you are in.
2. `ls` (List):
   - Lists files and directories in the current directory.
   - Common options:
     - `ls -l`: Long format, shows detailed file information.
     - `ls -a`: Lists hidden files (those starting with a dot).
3. `cd` (Change Directory):
   - Allows you to change the current directory.
   - Examples:
     - `cd /path/to/directory`: Change to an absolute path.
     - `cd ..`: Move up one directory.

- **cd ~**: Change to your home directory.

4. **mkdir** (Make Directory):
   - Creates a new directory.
   - Example: **mkdir new_directory**.

5. **touch**:
   - Creates an empty file.
   - Example: **touch new_file.txt**.

6. **cp** (Copy):
   - Copies files or directories from one location to another.
   - Example: **cp file.txt /path/to/destination**.

7. **mv** (Move/Rename):
   - Moves files or directories to a new location or renames them.
   - Example: **mv old_file.txt new_file.txt** (renaming).

8. **rm** (Remove):
   - Deletes files or directories.
   - Be cautious with this command, as deleted files are usually not recoverable.
   - Example: **rm file.txt** (removing a file).

9. **rmdir**:
   - Deletes an empty directory.
   - Example: **rmdir empty_directory**.

10. **rm -r**:
    - Deletes a directory and its contents recursively.
    - Example: **rm -r directory_to_delete**.

11. **cat** (Concatenate):
    - Displays the content of a text file.
    - Example: **cat file.txt**.

12. **more** and **less**:
    - Allows you to view the contents of a file page by page.
    - Example: **less long_text_file.txt**.

13. **head** and **tail**:
    - Display the beginning or end of a file, respectively.
    - Example: **head -n 10 file.txt** (displays the first 10 lines).

14. **grep**:
    - Searches for a specific pattern or text within files.
    - Example: **grep "keyword" file.txt**.

15. **find**:
    - Searches for files and directories in a specified location.
    - Example: **find /path/to/search -name "file_pattern"**.

16. **chmod** (Change Mode):
    - Changes file permissions.
    - Example: **chmod 644 file.txt** (gives read and write permissions to the owner and read-only to others).
17. **chown** (Change Owner):
    - Changes the owner of a file or directory.
    - Example: **chown new_owner:new_group file.txt**.

These are some of the fundamental Linux commands that will help you get started with basic file and directory operations. Remember to exercise caution when using commands like **rm** and **chmod** to avoid unintentional data loss or security issues.

Here are some basic Hadoop commands that are frequently used in Hadoop's distributed file system (HDFS) and MapReduce operations:

**HDFS (Hadoop Distributed File System) Commands:**

1. **hadoop fs -ls**:
    - Lists files and directories in HDFS.
    - Example: **hadoop fs -ls /user**.
2. **hadoop fs -mkdir**:
    - Creates a new directory in HDFS.
    - Example: **hadoop fs -mkdir /user/new_directory**.
3. **hadoop fs -copyFromLocal**:
    - Copies a file or directory from the local file system to HDFS.
    - Example: **hadoop fs -copyFromLocal localfile.txt /user/hadoop/hdfsfile.txt**.
4. **hadoop fs -copyToLocal**:
    - Copies a file or directory from HDFS to the local file system.
    - Example: **hadoop fs -copyToLocal /user/hadoop/hdfsfile.txt localfile.txt**.
5. **hadoop fs -mv**:
    - Moves a file or directory within HDFS.
    - Example: **hadoop fs -mv /user/oldfile.txt /user/newfile.txt**.
6. **hadoop fs -rm**:

- Deletes a file or directory in HDFS.
- Example: `hadoop fs -rm /user/filetodelete.txt`.

7. `hadoop fs -cat`:
   - Displays the content of a file in HDFS.
   - Example: `hadoop fs -cat /user/hadoop/hdfsfile.txt`.

8. `hadoop fs -tail`:
   - Displays the last part of a file in HDFS.
   - Example: `hadoop fs -tail /user/hadoop/hdfsfile.txt`.

9. `hadoop fs -get`:
   - Copies files from HDFS to the local file system.
   - Example: `hadoop fs -get /user/hadoop/hdfsfile.txt localfile.txt`.

10. `hadoop fs -put`:
    - Copies files from the local file system to HDFS.
    - Example: `hadoop fs -put localfile.txt /user/hadoop/hdfsfile.txt`.

## MapReduce Commands:

11. `hadoop jar`:
    - Submits a MapReduce job to the Hadoop cluster using a JAR file.
    - Example: `hadoop jar mymapreduce.jar input_dir output_dir`.

12. `hadoop job -list`:
    - Lists all currently running MapReduce jobs.

13. `hadoop job -kill`:
    - Terminates a running MapReduce job.
    - Example: `hadoop job -kill job_id`.

14. `hadoop fs -getmerge`:
    - Merges all files in an HDFS directory into a single local file.
    - Example: `hadoop fs -getmerge /user/hadoop/output_dir localfile.txt`.

15. `hadoop fs -du`:
    - Shows the disk usage of files and directories in HDFS.
    - Example: `hadoop fs -du /user/hadoop`.

16. `hadoop fs -count`:
    - Counts the number of directories, files, and bytes used in HDFS.
    - Example: `hadoop fs -count /user/hadoop`.

These are some of the basic Hadoop commands used for managing files and running MapReduce jobs in Hadoop. Depending on your specific use case and the Hadoop ecosystem tools you're using (like Hive, Pig, or Spark), you may also use additional commands and utilities.

# Class 2: Hadoop Distributed File System (HDFS)

**Introduction to HDFS:**

- **What is HDFS?**
  - HDFS stands for Hadoop Distributed File System.
  - It's a distributed, scalable, and fault-tolerant file system designed to store vast amounts of data across multiple commodity hardware nodes.
- **Key Features:**
  - **Scalability:** HDFS can handle petabytes of data by distributing it across a cluster of machines.
  - **Fault Tolerance:** Data is replicated across nodes to ensure reliability.
  - **High Throughput:** Optimized for batch processing.
  - **Streaming Data Access:** Suitable for large files.
  - **Economical:** Uses commodity hardware.

**HDFS Architecture:**

- **Components:**
  - **NameNode:** Master server that manages the file system namespace and regulates access to files.
  - **DataNode:** Slave nodes that store data and perform read/write operations.
  - **Block:** Data is split into blocks (typically 128MB or 256MB) and distributed across DataNodes.
  - **Replication:** Blocks are replicated to multiple DataNodes for fault tolerance.
- **How It Works:**
  - The NameNode stores metadata (file structure, permissions) and coordinates file access.
  - DataNodes store the actual data blocks and send heartbeats to the NameNode.
  - Clients communicate with the NameNode to locate data blocks and directly access DataNodes for data.

**HDFS Commands:**

- HDFS provides a set of command-line utilities for interacting with the file system.
- Common HDFS commands include `ls`, `mkdir`, `put`, `get`, `rm`, `mv`, and `cp`.

- Example: To list files in a directory: **hadoop fs -ls /path/to/directory**

## Data Replication:

- HDFS replicates data blocks to ensure fault tolerance.
- The default replication factor is usually 3 (data stored on three DataNodes).
- Replication provides data durability even if some nodes fail.

## Hands-on: Setting up HDFS, Basic Commands:

- In the hands-on session, you will set up HDFS on your Hadoop cluster (single-node).
- You'll learn how to use HDFS commands to create directories, upload files, and manage the file system.

```
Useful hive queries:


--Creating a Database:
CREATE DATABASE database_name;

--Listing Databases:
SHOW DATABASES;

--Switching to a Database:
USE database_name;

--Creating a Table:
CREATE TABLE table_name (
    column1 datatype,
    column2 datatype,
    ...
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE;


--Listing Tables in a Database:
SHOW TABLES;

--Display metadata/columns
DESCRIBE table_name;

--loading Data into a Table:
LOAD DATA LOCAL INPATH '/home/maria_dev/customers_data.csv'  INTO TABLE
table_name;
```

```sql
--simple SELECT Query:
SELECT * FROM table_name;

--Aggregating Data (e.g., COUNT, SUM, AVG):
SELECT COUNT(*), SUM(column_name), AVG(column_name) FROM table_name GROUP BY
column_name;


--Filtering Data with WHERE Clause:
SELECT * FROM table_name WHERE column_name = 'value';

--Sorting Data:
SELECT * FROM table_name ORDER BY column_name;


--Joining Tables:
SELECT a.column1, b.column2 FROM table1 a JOIN table2 b ON a.key = b.key;


--Deleting Data from a Table:
DELETE FROM table_name WHERE condition;


--Dropping a Table:
DROP TABLE table_name;


--Dropping a Database:
DROP DATABASE database_name;


--Running Custom Scripts (e.g., Python UDFs):
ADD FILE script_name.py;
SELECT my_udf(column_name) FROM table_name;

-- Creating Views:
CREATE VIEW view_name AS SELECT * FROM table_name WHERE condition;




DROP DATABASE IF EXISTS customer_orders CASCADE;

--Create new database
CREATE DATABASE IF NOT EXISTS customer_orders;


--use database
USE customer_orders;

--Creates orders table
CREATE TABLE orders (
    Customer_id INT,
    customer_unique_id INT,
    customer_zip_code_prefix STRING,
```

```
    customer_city STRING,
    customer_state STRING
    -- Add more columns as needed
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE;


--Load data -(file is in your local machine Linux )
LOAD DATA LOCAL INPATH '/home/maria_dev/customers_dataset.csv' INTO TABLE
orders;

--DESCRIBE orders
DESCRIBE orders;

--total number of records in the orders table
SELECT COUNT(*) FROM orders;
```

Apache Hive:

Apache Hive was developed by Facebook to address the growing need for querying large datasets stored in the Hadoop Distributed File System (HDFS). Here is a brief history of its development and evolution:

Apache Hive was created by Facebook in 2007 to handle large data volumes using a SQL-like interface on top of Hadoop's MapReduce framework. Hive simplifies querying and analyzing big data, making it accessible for users familiar with SQL.

Apache Hive is a data warehousing solution built on top of the Hadoop ecosystem, primarily used for querying and managing large datasets stored in Hadoop's HDFS (Hadoop Distributed File System). It provides a SQL-like interface to query data, which makes it accessible for users familiar with SQL. Hive abstracts the complexity of Hadoop's MapReduce programming model, allowing users to write queries in a language called HiveQL (Hive Query Language).

Here's a chart summarizing the differences between Hive and MapReduce:

| Feature | Hive | MapReduce |
|---|---|---|
| Abstraction Level | High-level | Low-level |
| User Interface | SQL-like language (HiveQL) | Java (or other languages via streaming) |
| Ease of Use | User-friendly, suitable for SQL users | Requires programming knowledge |
| Primary Function | Data warehousing, querying, and analysis | General-purpose data processing |

| Feature | Hive | MapReduce |
|---|---|---|
| Execution | Translates HiveQL into MapReduce, Tez, or Spark jobs | Direct execution of map and reduce functions |
| Use Case | Data analysts querying large datasets | Developers performing custom processing |
| Optimization | Built-in query optimizer | Manual optimization by users |
| Data Handling | Supports partitioning, bucketing, and various storage formats (ORC, Parquet, Avro) | Processes data in HDFS; manual data management |
| Integration | Part of the Hadoop ecosystem, integrates with other tools like Tez and Spark | Core part of Hadoop's data processing framework |

## 1. Data Warehouse

**Description**: A centralized repository that stores structured data from multiple sources, optimized for query and analysis. It typically contains historical data and supports business intelligence activities. **Examples**:

- **Amazon Redshift**: A fully managed data warehouse service in the cloud.
- **Google BigQuery**: A serverless, highly scalable, and cost-effective multi-cloud data warehouse.
- **Snowflake**: A cloud-based data warehousing platform that allows storage and analysis of data.

## 2. Data Lake

**Description**: A large repository that stores raw data in its native format until it is needed. It can handle structured, semi-structured, and unstructured data and is often used for big data analytics. **Examples**:

- **Amazon S3 (Simple Storage Service)**: Often used as a data lake due to its scalability and durability.
- **Azure Data Lake Storage**: A scalable and secure data lake solution that integrates with Azure services.
- **Apache Hadoop**: An open-source framework that allows for the distributed storage and processing of large data sets across clusters of computers.

## 3. Data Mart

**Description**: A subset of a data warehouse focused on a specific business area or department. It is designed to meet the needs of a particular group of users. **Examples**:

- **Sales Data Mart**: Contains sales data to support the sales department's analysis and reporting needs.
- **Finance Data Mart**: Stores financial data to assist the finance team with financial reporting and analysis.
- **Marketing Data Mart**: Holds marketing data to help the marketing department with campaign analysis and customer segmentation.

## 4. Database

**Description**: An organized collection of structured data, typically managed by a database management system (DBMS). Databases can be relational (SQL) or non-relational (NoSQL). **Examples**:

- **MySQL**: An open-source relational database management system.
- **MongoDB**: A NoSQL database known for its flexibility and scalability, often used for handling large volumes of unstructured data.
- **PostgreSQL**: An open-source relational database known for its advanced features and SQL compliance.

## 5. Data Lakehouse

**Description**: A data management architecture that combines the benefits of data lakes and data warehouses. It provides the ability to store raw data like a data lake while offering the management and optimization features of a data warehouse. **Examples**:

- **Databricks Lakehouse Platform**: Combines the capabilities of data lakes and data warehouses for unified data management.
- **Snowflake (with Unstructured Data Support)**: Allows for the storage and querying of both structured and unstructured data, effectively serving as a lakehouse.
- **Google BigQuery Omni**: A multi-cloud analytics solution that integrates data lakes and warehouses for comprehensive data management and analysis.

Here's a comparison chart of the important data storage concepts:

| Concept | Description | Example(s) |
|---|---|---|
| Data Warehouse | Centralized repository optimized for query and analysis, storing structured data from multiple sources. | Amazon Redshift, Google BigQuery, Snowflake |
| Data Lake | Large repository storing raw data in its native format, suitable for big data analytics. | Amazon S3, Azure Data Lake Storage, Apache Hadoop |
| Data Mart | Subset of a data warehouse focused on specific business areas or departments. | Sales Data Mart, Finance Data Mart, Marketing Data Mart |
| Database | Organized collection of structured data managed by a DBMS, relational or non-relational. | MySQL, MongoDB, PostgreSQL |
| Data Lakehouse | Architecture combining data lake and data warehouse features, offering unified data management. | Databricks Lakehouse Platform, Snowflake (with Unstructured Data Support), Google BigQuery Omni |

# Why Hive

Hive is used for several reasons:

1. **Ease of Use**: Provides a SQL-like interface, making it easier for users familiar with SQL to perform data analysis and querying tasks without deep knowledge of Hadoop's complexities.
2. **Scalability**: Can handle large datasets efficiently due to its integration with Hadoop's scalable storage and processing capabilities.
3. **Extensibility**: Supports user-defined functions (UDFs) to perform custom operations and complex data transformations.
4. **Integration**: Can integrate with various tools and frameworks in the Hadoop ecosystem, including Pig, HBase, and others.
5. **Data Warehousing**: Facilitates data summarization, ad-hoc queries, and the analysis of large datasets, making it suitable for data warehousing applications.

# Hive Architecture

Hive architecture consists of several key components:

1. **Metastore**: Stores metadata about tables, partitions, columns, and data types. It uses an RDBMS for this purpose and is critical for query compilation and execution.

2. **Driver**: Manages the lifecycle of HiveQL queries. It handles query compilation, optimization, and execution.
3. **Compiler**: Translates HiveQL queries into a directed acyclic graph of MapReduce or Tez tasks.
4. **Optimizer**: Applies various optimization techniques like predicate pushdown, partition pruning, and bucketing to enhance query performance.
5. **Executor**: Executes the compiled query plans using Hadoop's MapReduce or Tez.
6. **CLI, UI, and Thrift Server**: Interfaces for interacting with Hive. The CLI (Command Line Interface) and UI (User Interface) allow direct interaction, while the Thrift Server enables programmatic access via JDBC/ODBC.

# Hive Features

Hive offers a variety of features that make it powerful and flexible:

1. **SQL-like Language (HiveQL)**: Simplifies the querying process for users by providing a language similar to SQL.
2. **Partitioning**: Improves query performance by organizing data into partitions, allowing for faster access and management.
3. **Bucketing**: Further decomposes data into more manageable chunks within partitions, optimizing query execution.
4. **Indexing**: Provides faster query results by indexing specific columns.
5. **User-Defined Functions (UDFs)**: Allows users to write custom functions for specific operations, extending Hive's capabilities.
6. **Storage-based Metadata**: Stores metadata in a relational database, speeding up query compilation and execution.
7. **Extensibility**: Integrates with various Hadoop tools and supports different storage formats like ORC, Parquet, and Avro.
8. **Security**: Implements security features such as authentication, authorization, and integration with Kerberos.

**Summary of Discussion:**

```
Class 1: Introduction to Big Data and Hadoop

  1. Introduction to Big Data
     - Definition of Big Data
     - Why Big Data Matters
     - Challenges in Traditional Data Processing
     - Real-World Example: Netflix's Recommendation System

  2. Introduction to Hadoop
     - What is Hadoop?
     - Addressing Big Data Challenges
     - Real-World Example: Facebook's Data Analysis

  3. Hands-On: Setting Up Hadoop
     - Setting Up a Hadoop Cluster
       - Environment Setup
       - Hadoop Configuration
       - Basic Hadoop Commands
```

```
4. Hadoop Use Cases
   - Real-World Applications of Hadoop
     - Finance
     - Healthcare
     - Retail
   - Benefiting from Scalability and Cost-Efficiency
   - Real-World Example: Airbnb's Data Analysis


5. Hadoop Ecosystem
   - Overview of the Hadoop Ecosystem
     - Core Components
   - Real-Time Processing with Spark
     - Real-Time Processing
   - Real-World Example: Twitter's Real-Time Analytics with Spark
```

Hadoop Architecture

Hadoop, at its core, is an open-source framework for distributed storage and processing of large datasets. It was designed to tackle the challenges posed by Big Data. In this section, we will delve into the architecture of Hadoop, understanding its components and the principles that make it a powerful tool for handling massive datasets.

Components of a Hadoop Cluster

A Hadoop cluster is composed of several key components, each with a specific role in the data processing pipeline. Let's explore these components:

1. NameNode:
   - The NameNode is like the brain of the Hadoop Distributed File System (HDFS). It manages the metadata and namespace of the file system.
   - It keeps track of the structure and organization of files, directories, and blocks. However, it doesn't store the actual data of these files.
2. DataNode:
   - DataNodes are the workhorses of the HDFS. They store the actual data in the form of blocks.

- These nodes send regular updates to the NameNode about the health and availability of data blocks they store.

3. ResourceManager:
   - The ResourceManager is responsible for managing and allocating cluster resources.
   - It schedules applications based on available resources and ensures fair allocation among competing applications.

4. NodeManager:
   - NodeManagers are responsible for monitoring resource usage on individual cluster nodes.
   - They report this information back to the ResourceManager, enabling it to make informed decisions about resource allocation.
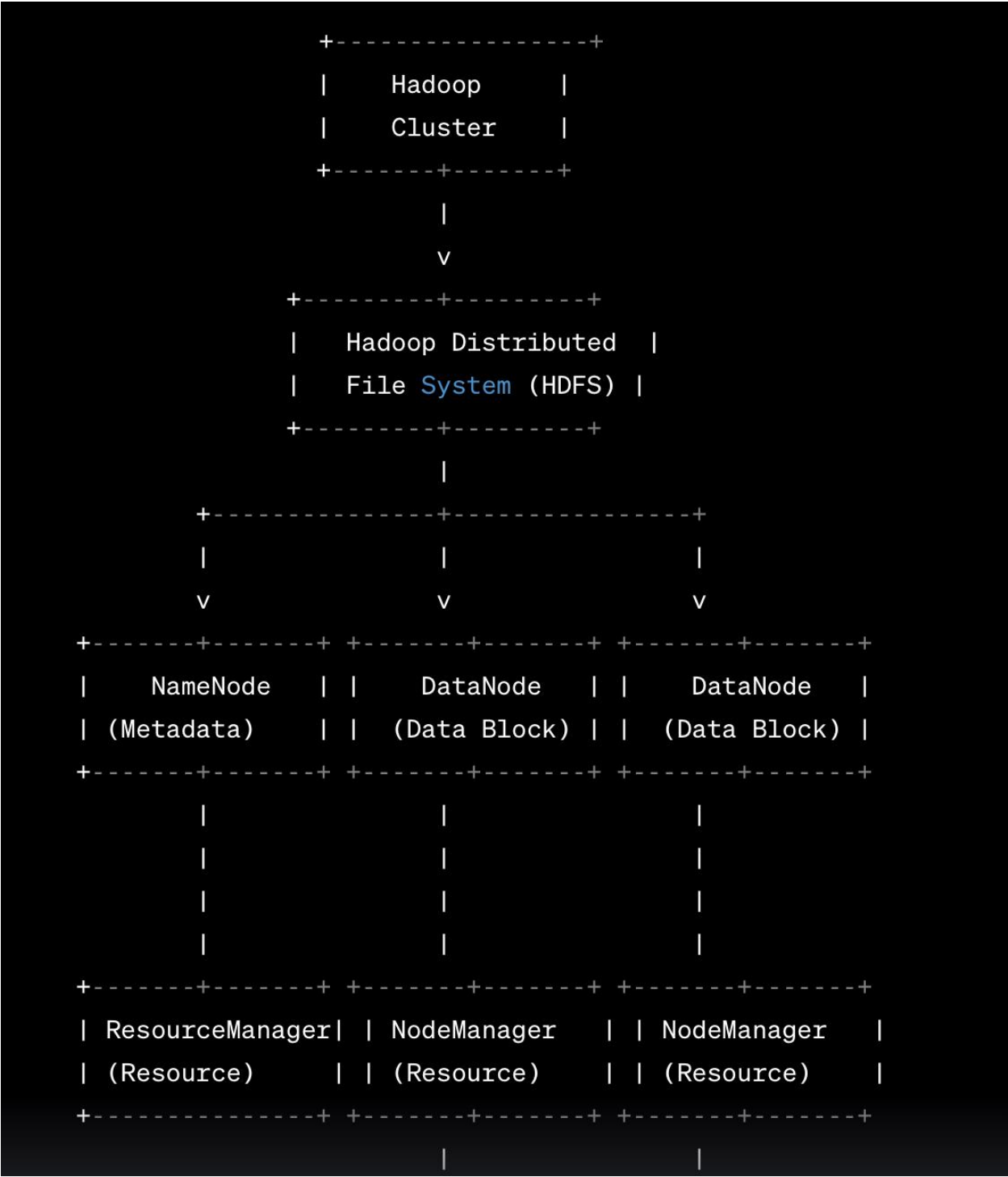
Single-Node Cluster vs. Multi-Node Cluster

Understanding the concept of a single-node cluster is crucial to grasp Hadoop's scalability.

- Single-Node Cluster:
  - In a single-node cluster, all Hadoop services run on a single machine. This setup is ideal for learning and development but doesn't leverage the distributed power of Hadoop.
  - It's like having a mini version of Hadoop on your local machine. You can experiment, write code, and test algorithms without the need for a complex multi-node setup.
  - However, it lacks the fault tolerance and data processing capabilities of a multi-node cluster.
- Multi-Node Cluster:
  - In a multi-node cluster, Hadoop's components are distributed across multiple machines, forming a robust and scalable ecosystem.
  - This architecture allows Hadoop to handle vast amounts of data by dividing it into smaller blocks and distributing them across the cluster.
  - It provides fault tolerance, as even if one node fails, data can be retrieved from replicas stored on other nodes

```
                    +-----------------+
                    |     Hadoop      |
                    |     Cluster     |
                    +-------+---------+
                            |
                            v
                  +---------+---------+
                  |  Hadoop Distributed  |
                  |  File System (HDFS) |
                  +---------+---------+
                            |
              +-------------+----------------+
              |             |                |
              v             v                v
      +-------+-------+ +-------+-------+ +-------+-------+
      |   NameNode    | |   DataNode    | |   DataNode    |
      | (Metadata)    | | (Data Block)  | | (Data Block)  |
      +-------+-------+ +-------+-------+ +-------+-------+
              |                 |                |
              |                 |                |
              |                 |                |
              |                 |                |
      +-------+-------+ +-------+-------+ +-------+-------+
      | ResourceManager| | NodeManager   | | NodeManager   |
      | (Resource)     | | (Resource)    | | (Resource)    |
      +---------------+ +-------+-------+ +-------+-------+
                              |                |
```

```
              |            |

              |            |

              |            |

      +----+-------+  +----+-------+

      | Application|  | Application|

      | Master     |  | Master     |

      | (Job)      |  | (Job)      |

      +-----------+  +-----------+
```

What is Spark:

Apache Spark, is an open-source, fast, and general-purpose distributed data processing framework. It was developed in response to the limitations of the Hadoop MapReduce model, designed to address various shortcomings and enhance the performance of big data processing.

Apache Spark is a powerful framework that offers speed, ease of use, and versatility for processing large datasets. It addresses the shortcomings of traditional data processing tools, making it a popular choice for big data analytics, machine learning, and real-time data processing.

Spark's key characteristics and features:

1. **Speed:** Spark is known for its high processing speed. It achieves this through in-memory data processing, reducing the need to read data from disk, which is a significant bottleneck in traditional data processing frameworks like Hadoop.
2. **Ease of Use:** Spark provides high-level APIs in multiple programming languages, making it accessible to a wide range of users. You can work with Spark using languages like Python, Java, Scala, or R.
3. **Versatility:** Spark is not limited to batch processing. It can handle various workloads, including batch processing, real-time data streaming, machine learning, and graph processing, all within a single platform.
4. **Distributed Processing:** Spark distributes data across a cluster of machines and processes it in parallel. It manages task distribution and recovery in case of node failures, ensuring fault tolerance.
5. **Resilient Distributed Datasets (RDDs):** RDDs are Spark's fundamental data structure. They are a distributed, fault-tolerant collection of data that can be processed in parallel. RDDs are at the core of Spark's processing capabilities.

6. **Library Ecosystem:** Spark comes with a rich ecosystem of libraries and tools, including Spark SQL for structured data processing, MLlib for machine learning, GraphX for graph processing, and Spark Streaming for real-time data processing.
7. **Master/Worker Architecture:** Spark operates on a cluster of computers, where one machine (the master) manages the distribution of tasks to multiple worker nodes. This architecture allows for scalable and parallel processing.
8. **Actions and Transformations:** Spark provides two types of operations on RDDs:
   - Transformations: These create new RDDs from existing ones through operations like map, filter, and reduce.
   - Actions: These return values or store data, like count, collect, or save.

Why Spark:

Apache Spark has gained popularity and become a popular choice for big data processing and analytics due to several compelling reasons:

Spark's combination of speed, ease of use, versatility, and a thriving open-source community has made it the preferred choice for many organizations and data professionals when dealing with big data processing, analytics, and machine learning. Its ability to process data in-memory and support multiple workloads sets it apart as a powerful and flexible framework.

1. **Speed:** Spark is known for its exceptional processing speed. It processes data in-memory, reducing the need to read from and write to disk, which can be a major bottleneck in traditional data processing models. The ability to cache data in memory allows Spark to achieve performance improvements of up to 100 times faster than Hadoop MapReduce for certain applications.
2. **Ease of Use:** Spark provides high-level APIs in multiple programming languages, including Python, Java, Scala, and R. This makes it accessible to a broad audience of data engineers, data scientists, and developers. Additionally, Spark offers interactive shells for quickly prototyping and testing code.
3. **Versatility:** Spark is a versatile framework capable of handling various workloads. It supports batch processing, real-time data streaming,

machine learning, graph processing, and SQL-based data analysis, all within a single platform. This versatility reduces the need for separate tools and simplifies the technology stack.

4. **Distributed Data Processing:** Spark operates in a distributed cluster environment, enabling it to efficiently process large datasets across multiple machines. It automatically handles data distribution and fault tolerance, ensuring reliable and efficient processing.

5. **Resilient Distributed Datasets (RDDs):** RDDs are a fundamental data structure in Spark. They are distributed collections of data that can be processed in parallel. RDDs are resilient, meaning they can recover from node failures and provide a consistent and fault-tolerant data processing model.

6. **Rich Ecosystem:** Spark has a rich ecosystem of libraries and tools that extend its capabilities. These include Spark SQL for structured data processing, MLlib for machine learning, GraphX for graph processing, and Spark Streaming for real-time data processing.

7. **Community and Support:** Spark has a vibrant and active open-source community, ensuring continuous development, improvements, and support. This community has contributed to its extensive documentation, resources, and a wide range of third-party libraries and connectors.

8. **Compatibility:** Spark is designed to work seamlessly with existing Hadoop data. It can read data from HDFS, Hive, HBase, and more, making it easy to integrate with existing big data environments.

9. **Scalability:** Spark is highly scalable, capable of handling a wide range of data sizes, from small datasets to petabytes of information. As data needs grow, Spark clusters can be expanded to accommodate larger workloads.

Spark core components:

Spark Core API:
the Spark Core API provides the core capabilities and abstractions that are essential for distributed data processing, making Spark suitable for a wide range of data processing tasks, including batch processing, real-time processing, machine learning, and more. (Supports all languages R, SQL, Python, Scala and Java)

**1.Spark SQL DataFrames:**
Spark Data Frames are a higher-level abstraction built on top of Spark's core Resilient Distributed Dataset (RDD) API. They provide a more structured and efficient way to work with structured and semi-structured data. DataFrames are conceptually equivalent to tables in a relational database or data frames in R or Python.

key features and concepts related to Spark SQL DataFrames:

- Structured Data: DataFrames represent structured data with named columns and data types. This structure is similar to a table in a relational database.
- Schema Inference: Spark SQL can automatically infer the schema of a DataFrame by examining the data. Alternatively, you can define the schema manually.
- API: DataFrames offer a rich set of high-level API functions for performing various operations on the data, such as filtering, aggregation, and transformation. You can work with DataFrames using SQL-like expressions, which makes it more accessible to those familiar with SQL.
- Integration with Spark Ecosystem: DataFrames seamlessly integrate with other components of the Spark ecosystem, such as Spark Streaming, Spark Machine Learning (MLlib), and Spark GraphX.
- Support for Multiple Data Sources: DataFrames can be used to read and write data from and to various data sources, including Parquet, Avro, ORC, JSON, and Hive tables. This makes it easy to work with data in different formats and storage systems.
- Optimizations: DataFrames leverage the Catalyst optimizer, which performs query optimization, and Tungsten, a physical query execution engine. These optimizations lead to better performance and query execution.
- Interoperability: DataFrames can seamlessly interoperate with existing RDDs, allowing you to leverage the power of DataFrames while working with more complex operations on RDDs when necessary.
- Support for SQL: You can register a DataFrame as a temporary table, allowing you to run SQL queries on it. This enables a smooth transition for SQL users into the Spark ecosystem.
- Immutable: Like RDDs, DataFrames are immutable, meaning that any operation on a DataFrame results in the creation of a new DataFrame, rather than modifying the existing one.
- Python, Scala, and Java APIs: DataFrames are available in multiple languages, making Spark accessible to a wide range of developers.
- Example of using DataFrames in Python:

python

```python
# Create a DataFrame from a list of dictionaries
data = [{"name": "Alice", "age": 30}, {"name": "Bob", "age": 25}]
df = spark.createDataFrame(data)

# Show the content of the DataFrame
df.show()

# Select and filter data
df.select("name", "age").filter(df["age"] > 26).show()
```

In this example, we created a DataFrame, displayed its content, and performed a selection and filter operation using the DataFrame API.

Spark SQL DataFrames offer a powerful and user-friendly way to work with structured data in Spark, making it easier for data engineers and data scientists to perform data processing and analysis tasks.

**2.Streaming:**
Spark Streaming is a component of Apache Spark that enables real-time data processing and analysis. It allows you to process and analyze data as it arrives, providing near-real-time insights and making it suitable for various use cases, such as monitoring, recommendation systems, fraud detection, and more.

key aspects and concepts related to Spark Streaming:

- Micro-Batch Processing: Spark Streaming divides the real-time data into small batches, processes each batch, and then produces results. This micro-batch approach is designed for efficiency and fault tolerance.
- DStreams (Discretized Streams): DStreams are the fundamental data structure in Spark Streaming. They are a sequence of data arriving over time, and you can perform various operations on DStreams, such as mapping, filtering, and windowing.
- Data Sources: Spark Streaming can consume data from various sources, including Apache Kafka, Flume, HDFS, Twitter, and custom data sources. It can ingest data from these sources in real-time, process it, and store or display the results.
- Transformations: You can apply high-level operations to DStreams, which are similar to those in Spark's core API. Common operations include map, reduceByKey, updateStateByKey, and join. These operations allow you to perform computations on the data streams.
- Windowed Operations: Spark Streaming supports windowed operations that enable you to process data over a sliding time window. This is particularly useful for time-based analyses.
- Output Operations: You can define output operations to save the processed data to external systems or display it. For example, you can store results in a database, write to a file, or send alerts.
- Integration with Spark Ecosystem: Spark Streaming integrates seamlessly with the broader Spark ecosystem. You can leverage machine learning libraries, SQL processing, and graph processing on the processed real-time data.
- Exactly-Once Semantics: Spark Streaming provides exactly-once processing semantics to ensure that each data record is processed exactly once, even in the presence of failures.
- Fault Tolerance: Spark Streaming is designed with fault tolerance in mind. It can recover from node or application failures and continue processing data without losing any.

- Event-Time Processing: Spark Streaming can process data based on event time, allowing you to deal with out-of-order data, late-arriving data, and data from multiple sources.

## 3.Apache Spark MLlib:

Spark MLlib (Machine Learning Library) is a scalable machine learning library built on top of the Apache Spark platform. It provides a wide range of machine learning and data mining algorithms and tools to support various data analysis and modeling tasks. Here are some key aspects and concepts related to Spark MLlib:

1. **Scalable:** Spark MLlib is designed for distributed and scalable machine learning. It can handle large datasets and distribute computations across a cluster of machines.
2. **ML Algorithms:** MLlib includes a variety of machine learning algorithms, including classification, regression, clustering, recommendation, and more. Some of the popular algorithms include linear regression, decision trees, k-means clustering, and collaborative filtering.
3. **APIs:** MLlib provides two main APIs for machine learning: a high-level API called the DataFrame-based API and a lower-level API using RDDs (Resilient Distributed Datasets). The DataFrame-based API is more user-friendly and commonly used.
4. **Data Preparation:** MLlib offers tools for data preprocessing, feature extraction, and feature selection. This helps in cleaning and transforming raw data into a suitable format for training machine learning models.
5. **Pipeline API:** The Pipeline API allows you to build workflows for feature extraction, transformation, and model training. It simplifies the process of assembling complex data processing pipelines.
6. **Model Selection:** MLlib includes tools for hyperparameter tuning, cross-validation, and model selection. You can evaluate and choose the best-performing models for your specific use case.
7. **Integration:** MLlib can be seamlessly integrated with other Spark components, such as Spark SQL, Spark Streaming, and Spark GraphX. This integration allows you to combine machine learning with other data processing tasks.
8. **Library Extensibility:** While MLlib provides a broad set of machine learning algorithms, you can also extend it by incorporating custom algorithms and libraries.
9. **Support for Data Types:** MLlib can work with a variety of data types, including numerical, categorical, and text data. It provides methods to handle different data types efficiently.

10. **Model Persistence:** Trained machine learning models can be saved and loaded, allowing you to reuse models for making predictions on new data.

## 4. GraphX:

GraphX is a component of Apache Spark that provides a distributed graph processing framework. It's designed for graph-based computations and analytics on large-scale graph data. GraphX allows you to create and manipulate graph structures, perform graph algorithms, and apply graph analytics to your data.

Here's an overview of GraphX:

1. **Graph Abstraction**: GraphX introduces the Resilient Distributed Property Graph, which extends the basic graph data structure with attributes associated with each vertex and edge.
2. **Graph Algorithms**: It provides a collection of graph algorithms, such as PageRank, connected components, and shortest paths, that you can apply to your graph data.
3. **Graph Analytics**: You can perform various graph analytics tasks, like community detection, subgraph isomorphism, and more.
4. **Distributed Processing**: GraphX leverages the distributed computing capabilities of Apache Spark, making it suitable for processing large-scale graphs across a cluster of machines.
5. **GraphFrames**: GraphX is often used in combination with DataFrames, allowing you to work with structured data and graph data seamlessly in a single Spark application.

Here are details for each step in the lifecycle of a Spark application:

1. User Code: Write code in a supported language (Scala, Java, Python, or R) to define data transformations and processing logic.
2. Spark Driver: Spark application execution starts with the Spark driver, running the main() function and managing task scheduling. The driver is responsible for orchestrating the entire application and interacts with the Cluster Manager.
3. Cluster Manager: Choose a cluster manager (e.g., Mesos, YARN, standalone) to allocate cluster resources. The cluster manager is in charge of resource allocation and job coordination, ensuring that the application has access to the necessary CPU and memory resources.

4. Cluster Nodes: Spark applications run on a cluster of nodes, each equipped with multiple CPU cores and memory. These nodes form the distributed computing environment where the application is executed.
5. Data Distribution: The dataset is divided into partitions and distributed across cluster nodes. Data partitions are stored on different nodes to enable parallel processing.
6. Execution Plan: The Spark driver generates an execution plan that defines how the application should be executed. This plan is essentially a directed acyclic graph (DAG) of transformations and actions.
7. Task Execution: Tasks are executed in parallel on worker nodes, processing data partitions. Each task performs the operations defined in the execution plan, which can include filtering, mapping, aggregation, and more.
8. Data Processing: Tasks perform data transformations and actions, such as filtering, mapping, aggregating, or joining, as specified in the user's code.
9. Fault Tolerance: Spark ensures fault tolerance using lineage information. Lineage tracks the sequence of transformations applied to the data, allowing lost data partitions to be recomputed if needed, ensuring the integrity of the application's results.
10. Result Aggregation: The results of the tasks are collected and aggregated. This stage combines the individual outputs from each task into the final result set.
11. Application Completion: The Spark application is considered complete when all tasks are executed and the final results are obtained. The driver oversees the completion of the application.
12. Resource Cleanup: The Cluster Manager releases allocated resources, and any cached data is cleared from memory. This stage helps ensure that resources are efficiently managed.
13. Output Handling: The application's results can be saved to external storage systems, displayed to the user, or passed to another application for further processing. This step depends on the specific use case and requirements of the application.

RDDs (Resilient Distributed Datasets):

1. **Resilient**: RDD stands for "Resilient Distributed Dataset." The "resilient" part means that RDDs are fault-tolerant. If a node in the cluster fails, Spark can recover lost data partitions because of the lineage information it maintains.
2. **Distributed**: RDDs are distributed collections of data across a cluster of machines. Data is divided into partitions, and each partition is processed in parallel by different worker nodes.
3. **Dataset**: RDDs are similar to datasets, representing a collection of data that can be operated upon. However, unlike traditional datasets, RDDs can be processed in a distributed and fault-tolerant manner.
4. **Transformations**: Transformations in Spark are operations that create new RDDs from existing ones. For example, map, filter, and reduceByKey are common

transformations. These operations are lazily evaluated, meaning they don't execute immediately but build up a logical execution plan.
5. **Actions**: Actions in Spark are operations that trigger computation and return results to the driver program. Examples of actions include count, collect, and saveAsTextFile. Actions lead to the execution of the logical plan created by transformations.
6. **Parallel Processing**: RDDs are designed for parallel processing. Operations on RDDs are automatically parallelized across the cluster, allowing for efficient utilization of the available resources.
7. **Caching**: You can persist (cache) an RDD in memory for reuse. This is particularly useful for iterative algorithms and interactive data analysis, as it avoids recomputing the same data multiple times.
8. **Lineage**: RDDs store information about how they were derived from other datasets. This lineage information enables Spark to recompute lost data partitions in case of node failures, ensuring fault tolerance.
9. **Immutability**: RDDs are immutable, which means once created, an RDD cannot be changed. Instead, transformations create new RDDs based on the original ones.

you can create an RDD (Resilient Distributed Dataset) from existing data in various ways, depending on your specific needs and the data source. Here are some common methods for creating RDDs:

```
from pyspark import SparkContext

sc = SparkContext("local", "RDD Example")
data = [1, 2, 3, 4, 5]
rdd = sc.parallelize(data)
```

**Reading from External Storage**
You can create RDDs by reading data from external storage systems like HDFS, local file systems, databases, and more. Spark supports various data sources, and you can use the appropriate methods for reading data. For example, in Python:

```
from pyspark import SparkContext

sc = SparkContext("local", "RDD Example")
rdd = sc.textFile("hdfs://path/to/your/file.txt")
```

**Using Transformation Operations**

You can create RDDs through transformations on existing RDDs. Transformations like map, filter, and union can be used to create new RDDs from one or more existing RDDs. For example:

new_rdd = rdd.map(lambda x: x * 2)

1. **Using External Data Sources**: Spark provides connectors for various data sources such as HBase, Cassandra, and more. You can create RDDs by connecting to these external data sources and extracting data into RDDs.

**Using Pair RDD Operations**

If your data has key-value pairs, you can create Pair RDDs using methods like mapToPair, groupByKey, and reduceByKey. For example:

pair_rdd = rdd.map(lambda x: (x, x * 2))

**Using DataFrames and Datasets**

You can convert DataFrames or Datasets into RDDs using the rdd method. This is useful if you want to work with RDD operations on structured data.

rdd = dataframe.rdd

Spark session:

A SparkSession is a unified entry point in Apache Spark for working with structured data. It was introduced in Spark 2.0 to simplify the configuration, use, and interaction with Spark components, such as DataFrame and Dataset APIs. SparkSession effectively replaces the earlier SQLContext, HiveContext, and StreamingContext used for different Spark tasks. Here's a closer look at SparkSession:

1. **Unified Entry Point**: SparkSession provides a single entry point for various Spark features, including SQL, DataFrames, Datasets, and Streaming. It eliminates the need to create multiple context objects for different Spark functionality.
2. **Configuration**: When you create a SparkSession, you can configure it with various options using `config` settings. These configurations include settings related to Spark's runtime behavior, the number of CPU cores, memory allocation, and more.

Creating Spark session:

```python
from pyspark.sql import SparkSession

# Create a SparkSession with a specific application name
spark = SparkSession.builder.appName("MyApp").getOrCreate()

# Create a DataFrame from a list of dictionaries
data = [{"name": "Alice", "age": 30}, {"name": "Bob", "age": 25}]
df = spark.createDataFrame(data)

# Show the content of the DataFrame
df.show()

# Select and filter data using SQL-like expressions
result = df.select("name", "age").filter(df["age"] > 26)

# Show the filtered result
result.show()
```

In this example:

1. We import the `SparkSession` class from the `pyspark.sql` module.
2. We create a SparkSession named "MyApp" using `SparkSession.builder.appName("MyApp").getOrCreate()`. The `appName` method sets a name for the Spark application.
3. We create a DataFrame named `df` from a list of dictionaries. This DataFrame represents structured data with two columns: "name" and "age."
4. We use the `show` method to display the content of the DataFrame, showing both the "name" and "age" columns.
5. We perform a selection and filter operation on the DataFrame, similar to using SQL. We select only the "name" and "age" columns where the "age" is greater than 26.
6. We use the `show` method again to display the filtered result.

Creating a RDD

```python
from pyspark.sql import SparkSession

# Create a SparkSession
spark = SparkSession.builder.appName("RDDExample").getOrCreate()

# Create an RDD using the parallelize method
data = [1, 2, 3, 4, 5]
rdd = spark.sparkContext.parallelize(data)

# Perform operations on the RDD
squared_rdd = rdd.map(lambda x: x * x)
filtered_rdd = squared_rdd.filter(lambda x: x % 2 == 0)

# Collect and display the results
result = filtered_rdd.collect()
print("Original Data:", data)
print("Squared Data:", squared_rdd.collect())
print("Filtered Data (even numbers):", result)

# Stop the SparkSession
spark.stop()
```

In this example:

1. We create a SparkSession with the name "RDDExample."
2. We create an RDD named `rdd` using the `parallelize` method and provide it with a list of integers (1, 2, 3, 4, 5).
3. We perform transformations on the RDD, such as squaring each element and filtering for even numbers.
4. We collect the results of the transformations using the `collect` method and display the original data, squared data, and the filtered data (even numbers).
5. Finally, we stop the SparkSession to release resources.

These RDD operations and actions are fundamental to Apache Spark's data processing capabilities, allowing users to perform transformations and extract information from large distributed datasets efficiently and in a fault-tolerant manner.

1. **map** Transformation:
   - Description: The **map** transformation applies a given function to each element of the RDD, producing a new RDD with the results of the function applied to each element.
   - Example: In the provided example, **map** is used to square each element in the RDD, creating a new RDD with the squared values.

2. **filter** Transformation:
   - Description: The **filter** transformation filters the elements of an RDD based on a specified condition, creating a new RDD that contains only the elements satisfying the condition.
   - Example: In the example, **filter** is applied to keep only the even numbers from the squared RDD.

3. **reduce** Transformation:
   - Description: The **reduce** transformation aggregates the elements of an RDD by successively applying a binary operation (function) to combine elements two at a time, reducing the RDD to a single value.
   - Example: The **reduce** operation calculates the sum of all elements in the original RDD.

4. **distinct** Transformation:
   - Description: The **distinct** transformation returns a new RDD with distinct (unique) elements from the original RDD.
   - Example: In the provided example, **distinct** is used to obtain a new RDD with unique elements from the original RDD.

5. **collect** Action:
   - Description: The **collect** action retrieves all the elements of an RDD and returns them as a list in the driver program. This action should be used with caution for large RDDs, as it brings all data to the driver.
   - Example: In the example, the **collect** action is used to obtain the results of the **squared_rdd** and **filtered_rdd** transformations as lists.

6. **count** Action:
   - Description: The **count** action returns the total number of elements in the RDD.
   - Example: In the example, the **count** action is used to count the number of even numbers after filtering.

7. **first** Action:
   - Description: The **first** action retrieves the first element from the RDD.

- Example: In the example, the **first** action is used to retrieve the first element of the original RDD.

8. **take** Action:
   - Description: The **take** action retrieves the first n elements of the RDD as a list. It does not require collecting the entire RDD.
   - Example: In the example, the **take** action is used to obtain the first three elements of the original RDD.

```python
from pyspark.sql import SparkSession

# Create a SparkSession
spark = SparkSession.builder.appName("RDDOperationsExample").getOrCreate()

# Create an RDD using the parallelize method
data = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
rdd = spark.sparkContext.parallelize(data)
```

```python
# RDD Transformations
# 1. Map: Square each element
squared_rdd = rdd.map(lambda x: x * x)

# 2. Filter: Keep even numbers
filtered_rdd = squared_rdd.filter(lambda x: x % 2 == 0)

# 3. Reduce: Calculate the sum of all elements
sum_result = rdd.reduce(lambda x, y: x + y)

# 4. Distinct: Get unique elements
distinct_rdd = rdd.distinct()
```

```python
# RDD Actions
# 1. Collect: Retrieve the results as a list
squared_list = squared_rdd.collect()

# 2. Count: Count the number of elements
count = filtered_rdd.count()

# 3. First: Get the first element
first_element = rdd.first()

# 4. Take: Get the first n elements as a list
first_three_elements = rdd.take(3)
```

```python
# Display the results
print("Original Data:", data)
print("Squared Data:", squared_list)
print("Filtered Data (Even Numbers):", filtered_rdd.collect())
print("Sum of All Elements:", sum_result)
print("Distinct Elements:", distinct_rdd.collect())
print("Number of Elements after Filtering (Count):", count)
print("First Element:", first_element)
print("First Three Elements:", first_three_elements)

# Stop the SparkSession
spark.stop()
```

Creating a Data Frame:

```python
from pyspark.sql import SparkSession

# Create a SparkSession
spark = SparkSession.builder.appName("DataFrameExample").getOrCreate()

# Sample data as a list of dictionaries
data = [
    {"name": "Alice", "age": 30},
    {"name": "Bob", "age": 25},
    {"name": "Charlie", "age": 35},
    {"name": "David", "age": 28},
]

# Create a DataFrame from the sample data
df = spark.createDataFrame(data)

# Show the DataFrame
df.show()

# Stop the SparkSession
spark.stop()
```

1. We import `SparkSession` to create a Spark session.
2. We create a Spark session named "DataFrameExample" using `SparkSession.builder.appName("DataFrameExample").getOrCreate()`.
3. We define sample data as a list of dictionaries, where each dictionary represents a row in the DataFrame. Each dictionary contains two columns, "name" and "age."
4. We use `spark.createDataFrame(data)` to create a DataFrame named "df" from the sample data.
5. We use `df.show()` to display the contents of the DataFrame.
6. Finally, we stop the Spark session using `spark.stop()` to release resources.

```
+-------+---+
|   name|age|
+-------+---+
|  Alice| 30|
|    Bob| 25|
|Charlie| 35|
|  David| 28|
+-------+---+
```

To create a DataFrame from a CSV file in Apache Spark, you can use the read.csv method provided by Spark's SparkSession.

```python
from pyspark.sql import SparkSession

# Create a Spark session
spark = SparkSession.builder.appName("CSVtoDataFrame").getOrCreate()

# Load a CSV file into a DataFrame
df = spark.read.csv("your_file.csv", header=True, inferSchema=True)

# Show the DataFrame
df.show()

# Stop the Spark session
spark.stop()
```

1. We create a Spark session using
   SparkSession.builder.appName("CSVtoDataFrame").getOrCreate().
2. We use the spark.read.csv("your_file.csv", header=True, inferSchema=True) method to read
   the CSV file. Replace "your_file.csv" with the path to your CSV file. The header=True
   option specifies that the first row of the CSV file contains column names, and
   inferSchema=True tries to infer the data types of the columns.

3. We display the contents of the DataFrame using **df.show()**.
4. Finally, we stop the Spark session using **spark.stop()** to release resources.

1. First, you need to install PySpark using pip. Here's the command to do that:

```
pip install pyspark
```

2. You can try importing PySpark in a Python script or a Jupyter notebook. If there are no import errors, the installation was successful.

```
try:
    import pyspark
    print("PySpark is installed correctly.")
except ImportError as e:
    print("PySpark is not installed.")
```

You can list the installed packages to see if PySpark is among them.

```
pip list
```

3. Install python packages:

```
# Import necessary packages for Spark session
from pyspark.sql import SparkSession

# Import functions for DataFrame operations
from pyspark.sql import functions as F


from pyspark.sql.types import StructType, StructField, StringType,
IntegerType, FloatType, DoubleType, TimestampType
```

4. # Run a simple command

```
df = spark.createDataFrame([(1, "one"), (2, "two"), (3, "three")], ["number",
"word"])
df.show()
```

5.Loading Data: Load data from various sources such as CSV, JSON, Parquet, etc.

```python
# Load a CSV file into a DataFrame
df = spark.read.csv("path/to/your/file.csv", header=True, inferSchema=True)
df.show()
```

6.

```python
#DataFrame Operations: Perform operations like filtering, selecting, grouping, and aggregating.

# Select specific columns
df_selected = df.select("user_id", "transaction_amount")

# Filter rows
df_filtered = df.filter(df["user_id"] > 10)

# Group by and aggregate
df_grouped =
df.groupBy("user_id").agg(F.sum("transaction_amount").alias("sum_transaction_
amount"))
df_grouped.show()
```

7.
SQL Queries: Register DataFrame as a temporary table and run SQL queries.

```python
# Register the DataFrame as a temporary table
df.createOrReplaceTempView("my_table")

# Run a SQL query
sql_df = spark.sql("SELECT user_id, COUNT(*) FROM my_table GROUP BY user_id")
sql_df.show()
```

Transformations and actions are the two primary types of operations you can perform on a DataFrame in PySpark. Transformations create a new DataFrame from an existing one, while actions trigger the execution of the transformations and return results.

Let's go through some common transformations and actions using a dataset.

1. Import packages:

```python
from pyspark.sql import SparkSession
```

```
from pyspark.sql.types import StructType, StructField, IntegerType,
DoubleType
from pyspark.sql import functions as F
```

2. Create a Spark session

```
# Create a Spark session
spark = SparkSession.builder \
    .appName("TransformationsAndActions") \
    .getOrCreate()
```

3. Define a schema:

```
# Define schema for the data
schema = StructType([
    StructField("user_id", IntegerType(), True),
    StructField("transaction_amount", DoubleType(), True)
])
```

4. Create a Dataframe:

```
# Create the DataFrame manually (since we don't have the actual CSV file, we
simulate the data)
data = [(5, 545.92), (98, 386.07), (30, 577.02)]
df = spark.createDataFrame(data, schema=schema)
```

6. Show DataFrame:

```
# Show the DataFrame
df.show()
```

Transformations:

1.Select Columns: Select specific columns from the DataFrame.

```
df_selected = df.select("user_id", "transaction_amount")
df_selected.show()
```

2. Filter Rows: Filter rows based on a condition.
```
df_filtered = df.filter(df["transaction_amount"] > 500)
df_filtered.show()
```

3. Add New Column: Add a new column with a calculated value.

```
df_with_new_column = df.withColumn("transaction_with_tax",
df["transaction_amount"] * 1.1)
df_with_new_column.show()
```

4.Group By and Aggregate: Group by a column and perform an aggregation.

```
df_grouped =
df.groupBy("user_id").agg(F.sum("transaction_amount").alias("total_amount"))
df_grouped.show()
```

5.Sort Rows: Sort the DataFrame based on a column.

```
df_sorted = df.orderBy("transaction_amount", ascending=False)
df_sorted.show()
```

Actions:

1. Show Data: Display the content of the DataFrame (used above in transformations as well).

```
df.show()
```

2. Collect: Return all the rows as a list of Row objects.

```
collected_data = df.collect()
print(collected_data)
```

3.Count: Count the number of rows in the DataFrame.

```
row_count = df.count()
print(f"Row count: {row_count}")
```

4.Take: Return the first n rows.
```
first_two_rows = df.take(2)
print(first_two_rows)
```

5.Describe: Compute summary statistics.
```
df.describe().show()
```

6. First: Return the first row.
```
first_row = df.first()
print(first_row)
```

1. #SparkSession: The entry point to programming Spark with the DataFrame and SQL API.
#sum: An aggregation function to calculate the sum of a numerical column.
#col: A function to reference a column in a DataFrame.

```
#import pyspark module
# Import necessary packages for Spark session

from pyspark.sql import SparkSession

# Import functions for DataFrame operations
```

```python
from pyspark.sql import functions as F


# Import types for defining schemas
from pyspark.sql.types import StructType, StructField, StringType,
IntegerType, FloatType, DoubleType, TimestampType


# Create a Spark session
spark = SparkSession.builder \
        .appName("VerifyPySparkInstallation") \
        .getOrCreate()


# Load a CSV file into a DataFrame
df = spark.read.csv("/Users/jeykanesh/Desktop/coding/data.csv", header=True,
inferSchema=True)
df.show()


# Register the DataFrame as a temporary table
df.createOrReplaceTempView("my_table")

# Run a SQL query
sql_df = spark.sql("SELECT user_id, COUNT(*) FROM my_table GROUP BY user_id")
sql_df.show()
```