

Hadoop and Spark

Kazi Aminul Islam

Department of Computer Science

Kennesaw State University

Hadoop

What is Apache Hadoop?

- Apache Hadoop is an open-source software utility that allows users to manage big data sets (from gigabytes to petabytes) by enabling a network of computers (or “nodes”) to solve vast and intricate data problems.
- It is a highly scalable, cost-effective solution that stores and processes **structured, semi-structured and unstructured data** (e.g., Internet clickstream records, web server logs, IoT sensor data, etc.).

Benefits of the Hadoop framework include the following:

- Data protection amid a hardware failure
- Vast scalability from a single server to thousands of machines
- Real-time analytics for historical analyses and decision-making processes

Hadoop use cases

Hadoop is most effective for scenarios that involve the following:

- Processing big data sets in environments where data size exceeds available memory
- Batch processing with tasks that exploit disk read and write operations
- Building data analysis infrastructure with a limited budget
- Completing jobs that are not time-sensitive
- Historical and archive data analysis

Apache Sparks

What is Apache Spark?

- [Apache Spark](#) — which is also open source — is a data processing engine for big data sets.
- Like Hadoop, Spark splits up large tasks across different nodes.
- However, it tends to perform faster than Hadoop and it uses random access memory (RAM) to cache and process data instead of a file system. This enables Spark to handle use cases that Hadoop cannot.

Benefits of the Spark framework include the following:

- A unified engine that supports SQL queries, streaming data, [machine learning \(ML\)](#) and graph processing
- Can be [100x faster than Hadoop for smaller workloads](#) via in-memory processing, disk data storage, etc.
- [APIs](#) designed for ease of use when manipulating semi-structured data and transforming data

Spark use cases

Spark is most effective for scenarios that involve the following:

- Dealing with chains of parallel operations by using iterative algorithms
- Achieving quick results with in-memory computations
- Analyzing stream data analysis in real time
- Graph-parallel processing to model data
- All ML applications

Ecosystem

The Hadoop ecosystem

- Hadoop supports advanced analytics for stored data
 - e.g., predictive analysis, [data mining](#), machine learning (ML), etc.
- It enables big data analytics processing tasks to be split into smaller tasks.
- The small tasks are performed in parallel by using an algorithm (e.g., MapReduce), and are then distributed across a Hadoop cluster (i.e., nodes that perform parallel computations on big data sets).

The Hadoop ecosystem consists of four primary modules

1. **Hadoop Distributed File System (HDFS):** Primary data storage system that manages large data sets running on commodity hardware. It also provides high-throughput data access and high fault tolerance.
2. **Yet Another Resource Negotiator (YARN):** Cluster resource manager that schedules tasks and allocates resources (e.g., CPU and memory) to applications.
3. **Hadoop MapReduce:** Splits big data processing tasks into smaller ones, distributes the small tasks across different nodes, then runs each task.
4. **Hadoop Common (Hadoop Core):** Set of common libraries and utilities that the other three modules depend on

The Spark ecosystem

Apache Spark, the largest open-source project in data processing, is the only processing framework that combines data and [artificial intelligence \(AI\)](#).

This enables users to perform large-scale data transformations and analyses, and then run state-of-the-art machine learning (ML) and AI algorithms.

The Spark ecosystem consists of five primary modules

1. **Spark Core:** Underlying execution engine that **schedules and dispatches tasks and coordinates input and output (I/O) operations.**
2. **Spark SQL:** Gathers information about structured data to enable users to **optimize structured data processing.**
3. **Spark Streaming and Structured Streaming:** Both add stream processing capabilities. Spark Streaming takes data from different streaming sources and **divides it into micro-batches for a continuous stream.** Structured Streaming, built on Spark SQL, reduces latency and simplifies programming.
4. **Machine Learning Library (MLlib):** A set of machine learning algorithms for **scalability plus tools for feature selection and building ML pipelines.** The **primary API** for MLlib is **DataFrames**, which provides uniformity across different programming languages like Java, Scala and **Python.**
5. **GraphX:** User-friendly computation engine that enables interactive building, modification and analysis of scalable, graph-structured data.

Comparing Hadoop and Spark

Faster: Spark

Spark is a Hadoop enhancement to MapReduce. The primary difference between Spark and MapReduce is that Spark processes and retains data in memory for subsequent steps, whereas MapReduce processes data on disk. As a result, for smaller workloads, [Spark's data processing speeds are up to 100x faster than MapReduce](#).

Fault Tolerance: Spark

Furthermore, as opposed to the two-stage execution process in MapReduce, Spark creates a Directed Acyclic Graph (DAG) to schedule tasks and the orchestration of nodes across the Hadoop cluster.

This task-tracking process enables fault tolerance, which reapplies recorded operations to data from a previous state.

Six Key Differences

1. **Performance:** Spark is faster because it uses random access memory (RAM) instead of reading and writing intermediate data to disks. Hadoop stores data on multiple sources and processes it in batches via MapReduce.
2. **Cost:** Hadoop runs at a lower cost since it relies on any disk storage type for data processing. Spark runs at a higher cost because it relies on in-memory computations for real-time data processing, which requires it to use high quantities of RAM to spin up nodes.
3. **Processing:** Though both platforms process data in a distributed environment, Hadoop is ideal for batch processing and linear data processing. Spark is ideal for real-time processing and processing live unstructured data streams.

Six Key Differences

4. **Scalability:** When data volume rapidly grows, Hadoop quickly scales to accommodate the demand via Hadoop Distributed File System (HDFS). In turn, Spark relies on the fault tolerant HDFS for large volumes of data.

5. **Security:** Spark enhances security with authentication via shared secret or event logging, whereas Hadoop uses multiple authentication and access control methods. Though, overall, Hadoop is more secure, Spark can integrate with Hadoop to reach a higher security level.

6. **Machine learning (ML):** Spark is the superior platform in this category because it includes MLlib, which performs iterative in-memory ML computations. It also includes tools that perform regression, classification, persistence, pipeline construction, evaluation, etc.