# LEARNING JOURNAL UNIT 2

Apache Spark and Apache Hadoop are two of the most influential frameworks utilized during big data processing. While both aim at processing and viewing large datasets, these two vary significantly with reference to architecture, processing models, speed, and usability. This week's discovery of Apache Spark has brought into view a clearer understanding of how it has an edge, especially with reference to Hadoop, which has been a precursor of distributed data processing.

## Architectural Differences

Apache Spark is different from Hadoop in how it processes data. Spark uses in-memory computing, which allows it to store data in RAM across the cluster and access it quickly for iterative processes or recurring queries. Hadoop uses a disk-based system called the Hadoop Distributed File System (HDFS) for storage and MapReduce for data processing, where intermediate data is spilled back to disk after each processing stage. This basic architectural difference makes Spark significantly faster in most use cases.

## Programming Model Comparison

A second significant difference lies in their programming models. Wherein Apache Spark enjoys broad support of APIs across many languages like Scala, Python, Java, and R, and further comes with a high-level abstraction known as RDD (Resilient Distributed Dataset) which makes it easy to process distributed data but also allows advanced analytics like machine learning (through MLlib), graph processing (through GraphX), and querying of structured data (through

Spark SQL) under one umbrella. Hadoop is MapReduce-based, however, where users need to specify intricate series of data map and reduce jobs, and hence not as easy to use and manage.

## Fault Tolerance

Fault tolerance is also treated differently. Spark uses lineage information to recompute missing data partitions in the event of failure, while Hadoop calls for data replication between HDFS nodes. Although both procedures ensure data reliability, Spark's procedure is more effective in maintaining recomputation over replication when recomputation takes less time than replication.

## Real-Time Processing Capabilities

Also, Spark provides enhanced real-time processing capabilities. It has a built-in streaming library called Spark Streaming, which is capable of scalable, high-throughput, and fault-tolerant stream processing of live data streams. Hadoop is batch-oriented and lacks direct support for real-time data processing. While real-time tools like Apache Storm or Kafka Streams can be added on top of Hadoop, this increases the complexity of the system architecture.

## Speed and Performance

As for which of the two has faster processing and analysis capabilities, Apache Spark is the winner. With its in-memory computation model, it can process jobs 100 times faster than Hadoop MapReduce for certain applications (Zaharia et al., 2016). The speed advantage is particularly significant for iterative algorithms of machine learning where data access is

repeatedly accessed. In Hadoop, each iteration would be re-reading from disk, which introduces latency and is not as efficient.

## Ease of Use and Flexibility

Besides speed, Spark's support for interactive queries and big data analytics also makes it easier and more flexible to use by data scientists as well as engineers. Spark SQL is a case in point because it enables simple execution of SQL-like queries against big datasets and also easy integration with existing data pipelines. Conversely, Hadoop typically requires the use of Pig or Hive, both of which add added complexity and typically take longer to run.

## Scalability and Future-Readiness

Spark is also effective when data sizes are increasing. Spark has petabyte-scale data sets supported effectively in thousands of nodes with its distributed architecture and optimized execution engine. The ecosystem of Spark also further evolves at a rapid rate, with ongoing development and an active community contributing to it, making it an even more future-proofed option for big data analytics (Karau et al., 2020).

## Conclusion

Briefly, Hadoop and Apache Spark are both major big data processing platforms but different in what they process. Hadoop is mature and suitable for batch processing large static data sets. Spark, with in-memory computing, real-time processing of data, and cross-platform analytics platform, offers improved performance and flexibility in the support of existing data

applications. Based on these advantages, the recommended framework is Apache Spark if

performance and quick analysis of big data are needed.

# References

Karau, H., Warren, R., & Hamstra, R. (2020). *High Performance Spark: Best Practices for Scaling and Optimizing Apache Spark*. O'Reilly Media. https://www.oreilly.com/library/view/high-performance-spark/9781491943199/

Zaharia, M., Xin, R. S., Wendell, P., Das, T., Armbrust, M., Dave, A., Meng, X., Rosen, J., Venkataraman, S., Franklin, M. J., Ghodsi, A., Gonzalez, J., Shenker, S., & Stoica, I. (2016). Apache spark: A unified engine for big data processing. Communications of the ACM, 59(11), 56-65. https://doi.org/10.1145/2934664. https://www.researchgate.net/publication/310613994_Apache_spark_A_unified_engine_for_big_data_processing

**Wordcount:** 653