

MapReduce and Apache Spark

CS5229 - Big Data Analytics Technologies



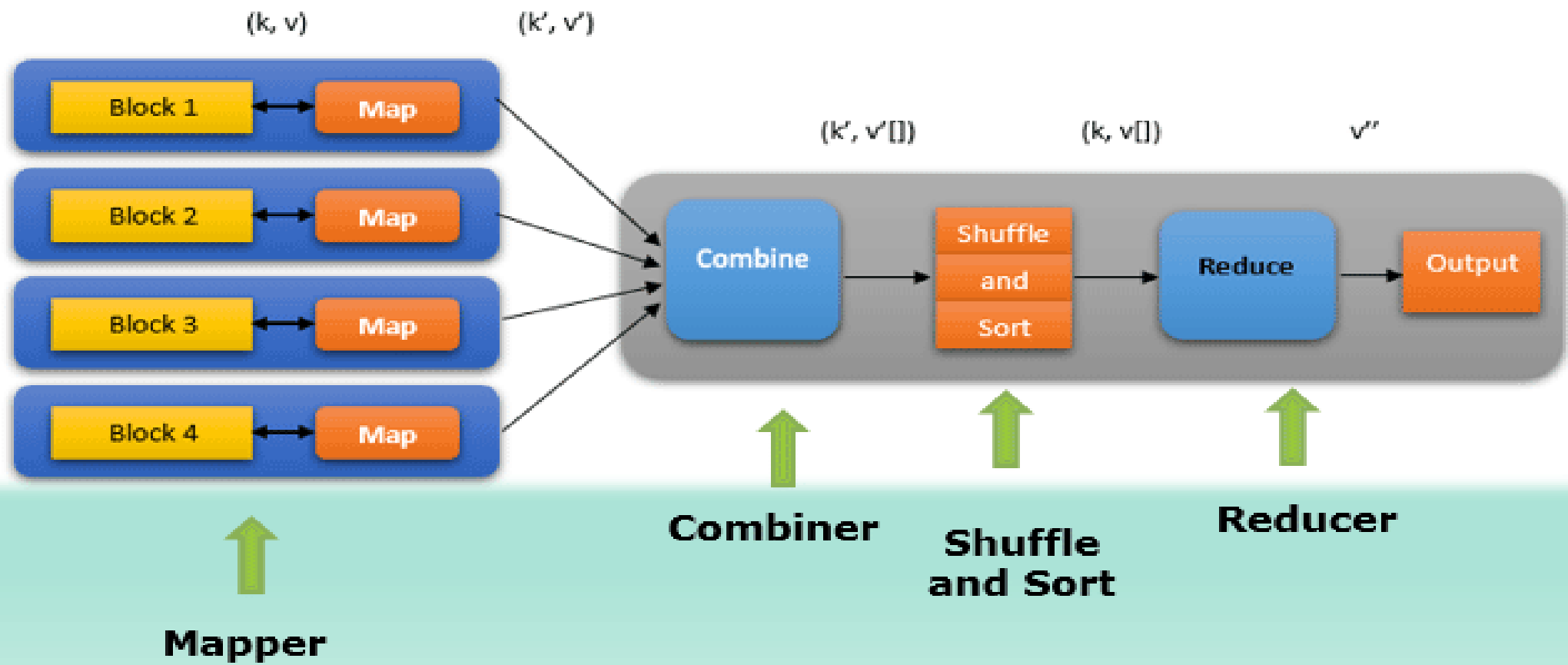
Hasini Weerasooriya

248287L

MSc in Computer Science

University of Moratuwa





MapReduce

- MapReduce is a programming paradigm that enables massive scalability across hundreds or thousands of servers in a Hadoop cluster.
- It performs the processing of large data sets in a distributed and parallel manner.
- MapReduce consists of two distinct tasks – **Map** and **Reduce**

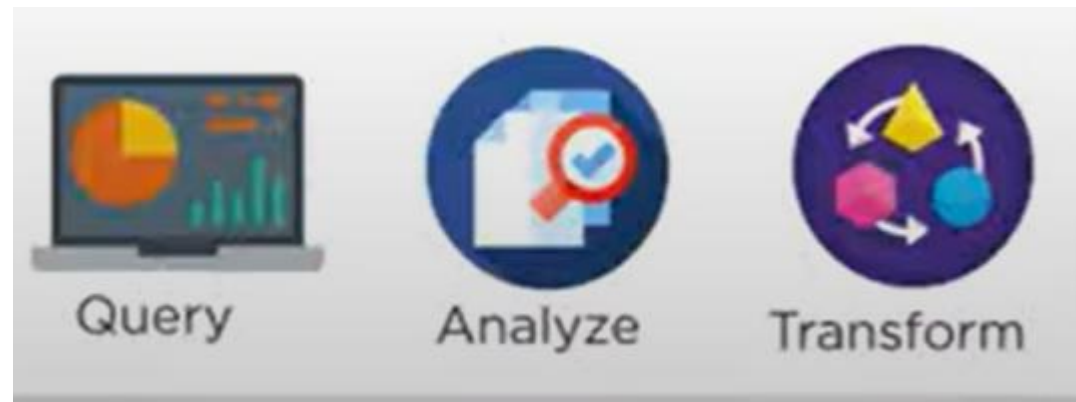
Apache Spark

- Apache Spark is an open-source data processing engine to store and process big data.
- It is built on top of the Hadoop distributed file system.

Support various
Programming Languages



Developers and data scientists incorporate spark into their applications to rapidly query, analyze, and transform data at scale.



DEMO



Amazon EMR

EMR on EC2: Clusters

Clusters (5) Info

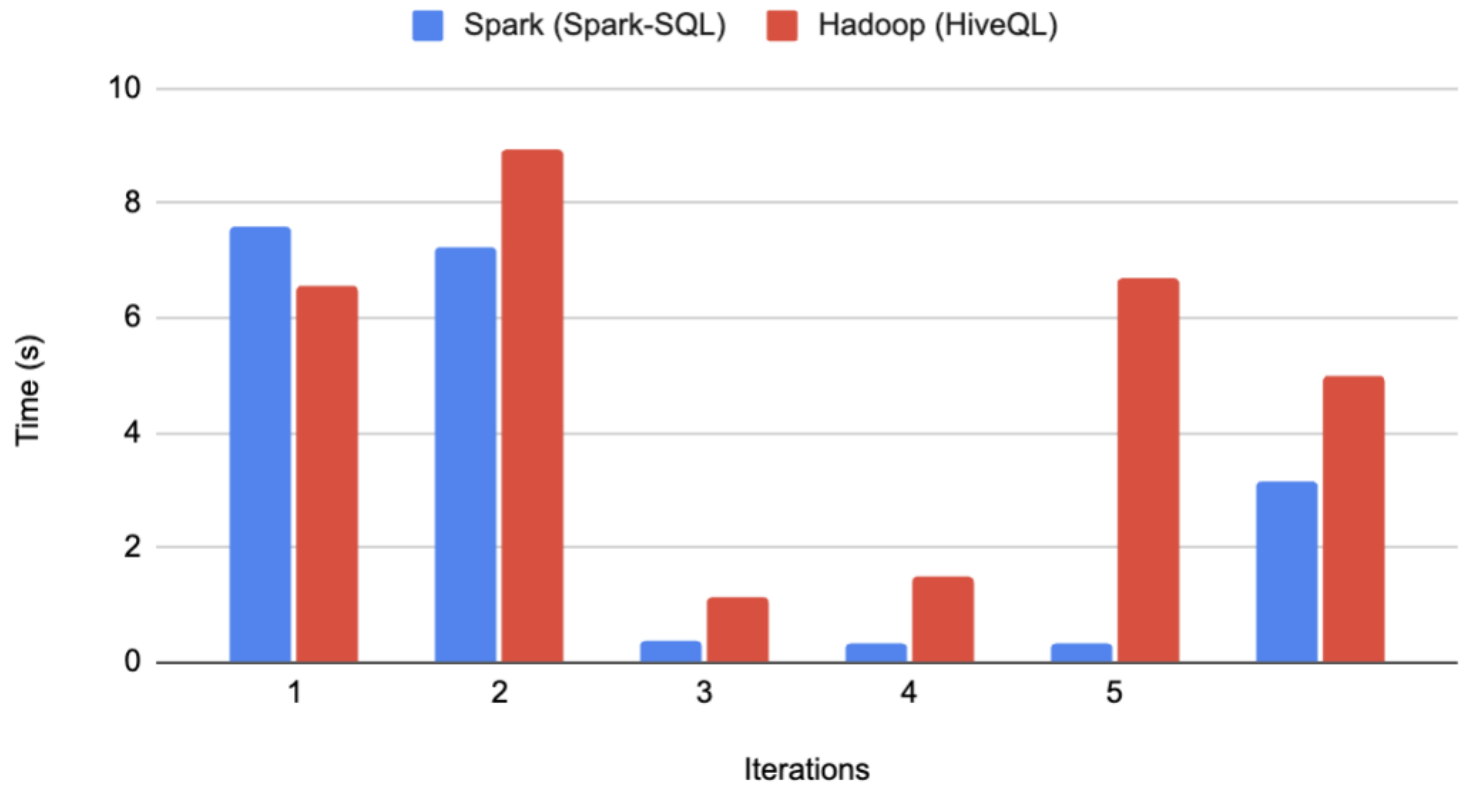
Filter clusters by statusFind clustersFilter clusters by creation date-time

<1>

		Cluster ID	Cluster name	Status	Creation time (UTC+05:30)	Elapsed time	Normalized instance hours
<input type="checkbox"/>			j-161LNYOIV743O	Waiting Ready to run steps	March 04, 2024, 21:16	32 minutes	0
<input type="checkbox"/>			j-17KVIY6C5MKPW	Terminated User request	March 04, 2024, 20:40	53 minutes, 39 seconds	24
<input type="checkbox"/>			j-13RNOUDUCUTHQ	Terminated User request	March 04, 2024, 19:40	58 minutes, 23 seconds	24
<input type="checkbox"/>			j-QEJA2MS9RP2B	Terminated User request	February 13, 2024, 22:54	41 minutes, 33 seconds	24
<input type="checkbox"/>			j-JG68VN4UDV9D	Terminated User request	February 13, 2024, 22:12	26 minutes, 55 seconds	24

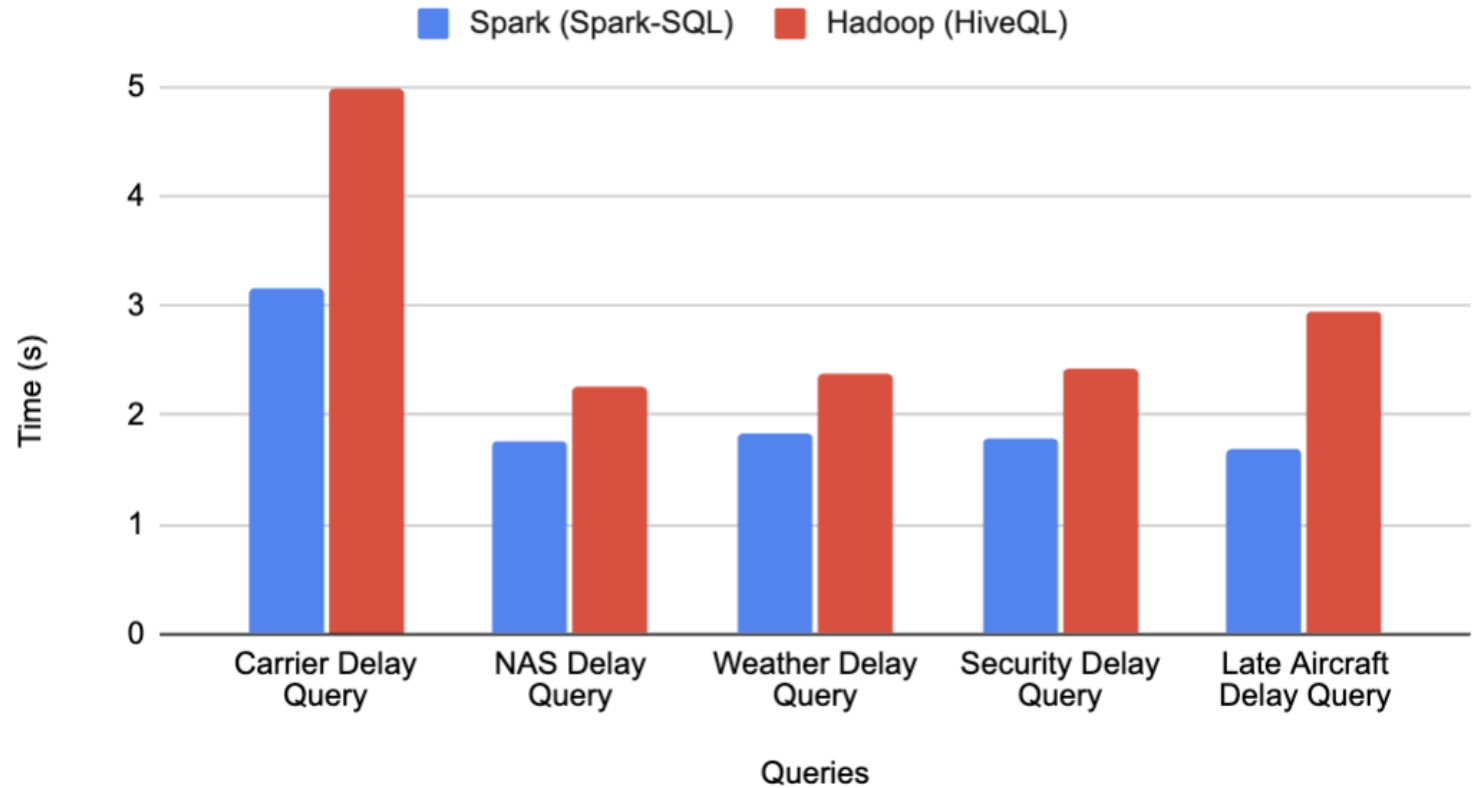
Comparison Results

Running Time vs Iteration



Comparison Results (Cont.)

Running Time vs Queries



MapReduce vs Apache Spark



Processing data using MapReduce in Hadoop is slow

Performs batch processing of data

Compact and Lengthy

Written in Java with more lines of code and takes more time to execute

Doesn't support caching of data

Spark processes data 100 times faster than MapReduce as it is done in-memory

Performs both batch processing and real-time processing of data

Compact and easier than Hadoop

Implemented in Scala with fewer lines of code

Caches the data in-memory & enhances the system performance

Conclusion

- In here we did the comparison using Hadoop and Spark.
- Both used to process big data in different ways.
- Hadoop was created to delegate data processing to several servers instead of running the workload on a single machine.
- Apache Spark is a newer data processing system that overcomes key limitations of Hadoop.



THANK YOU!