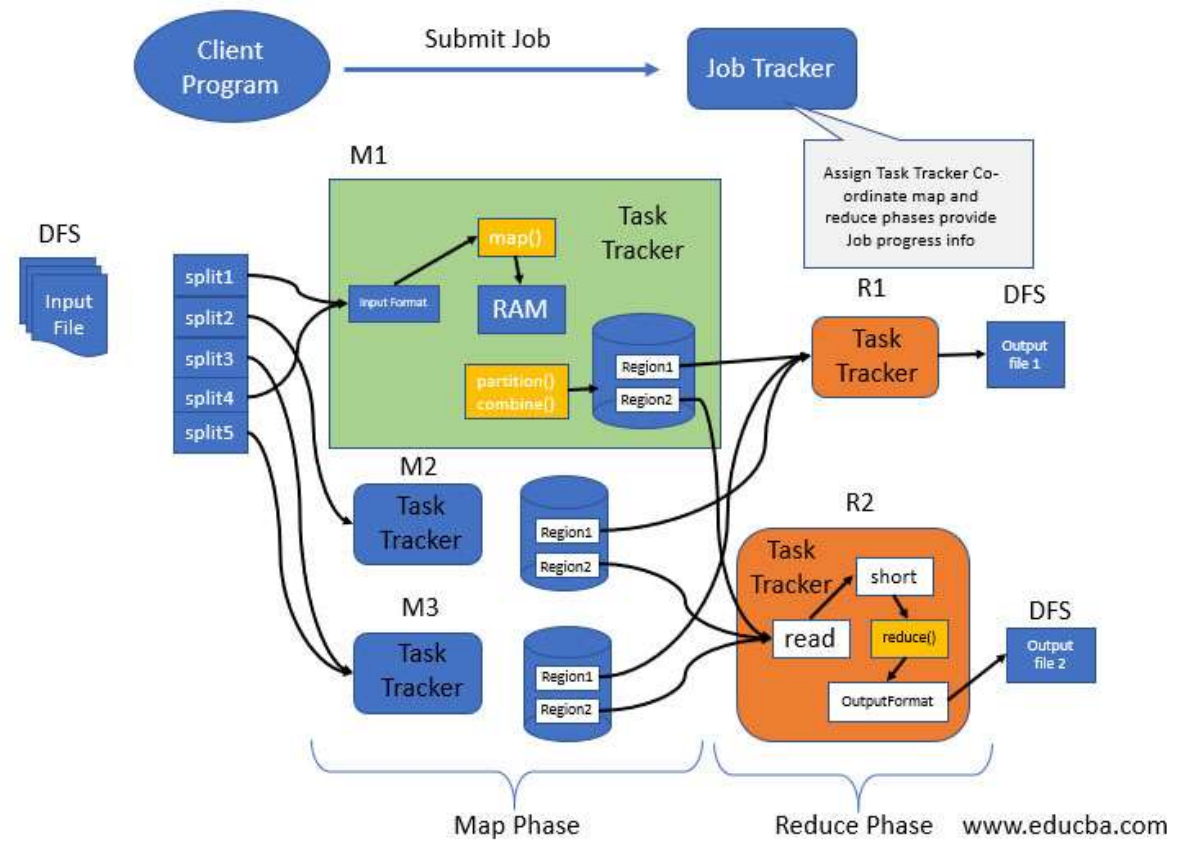# Hadoop MapReduce or Apache Spark?



Presented by Welgama W.I.P.

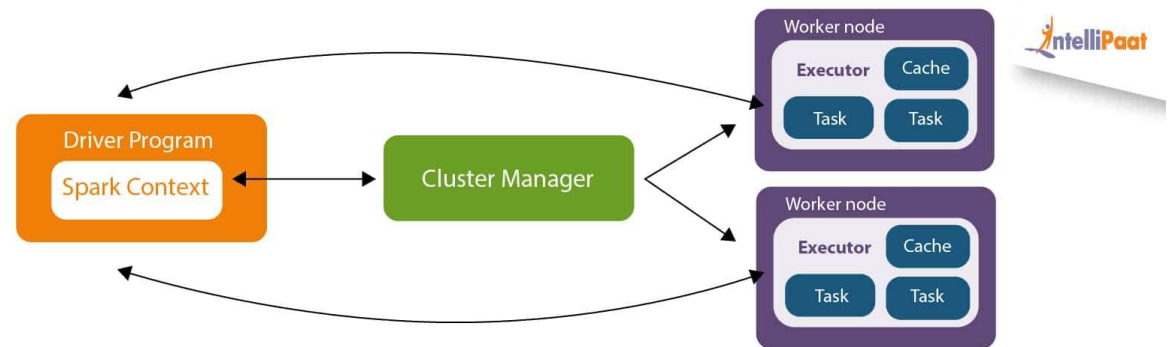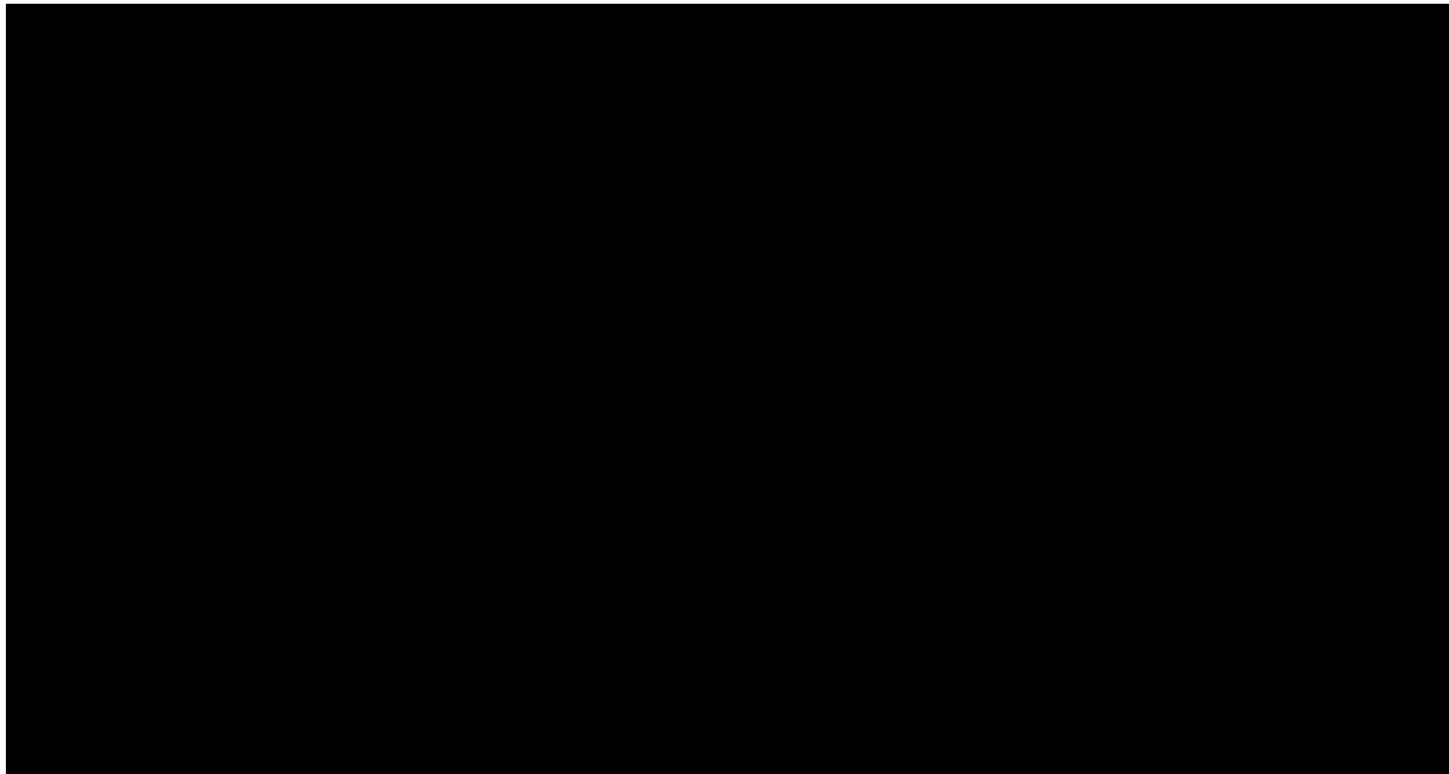239185V

# What is MapReduce?

# What is Spark?

Apache Spark is an open-source, distributed computing system designed for big data processing. It provides an interface for programming in various languages such as Java, Python, and Scala, and supports various data sources including Hadoop Distributed File System (HDFS) and Cassandra. Spark's main abstraction is the Resilient Distributed Dataset (RDD), which enables fault-tolerant and parallel processing of data across a cluster of computers.
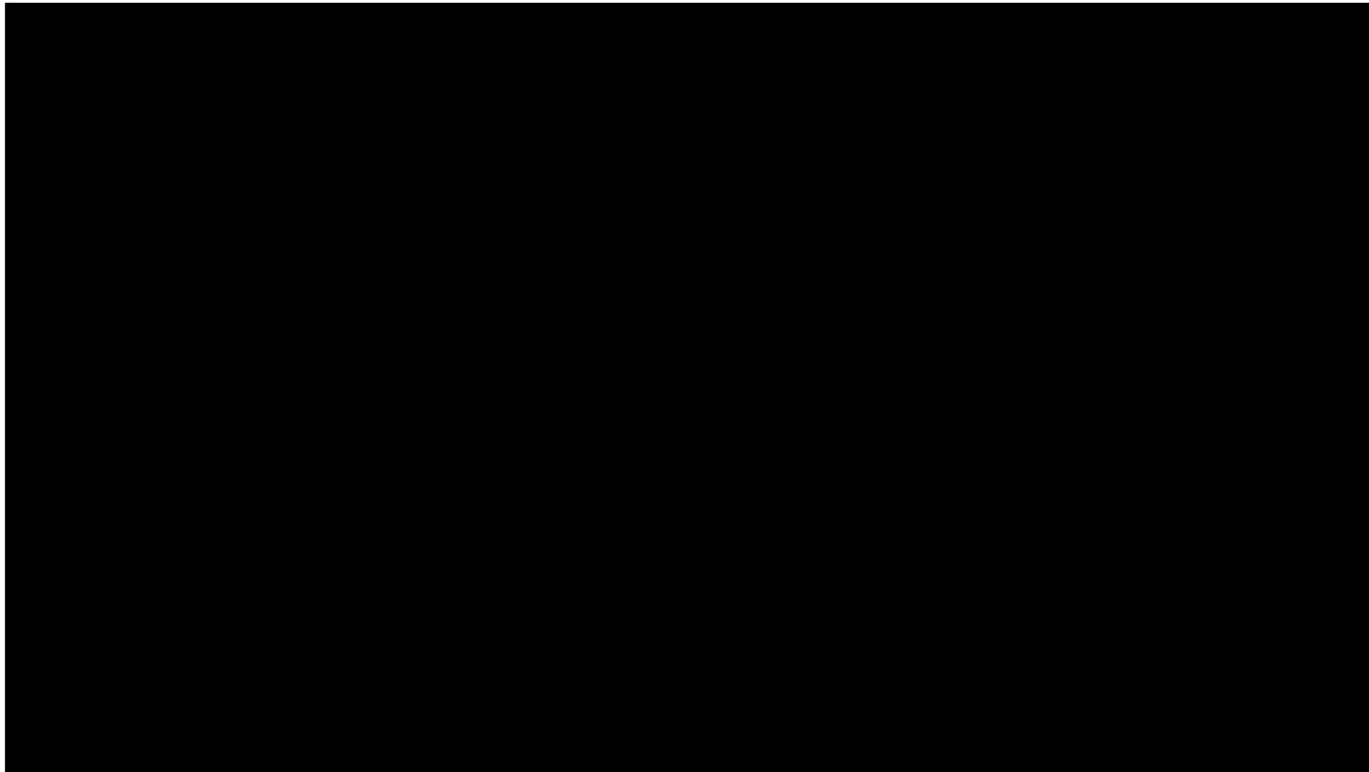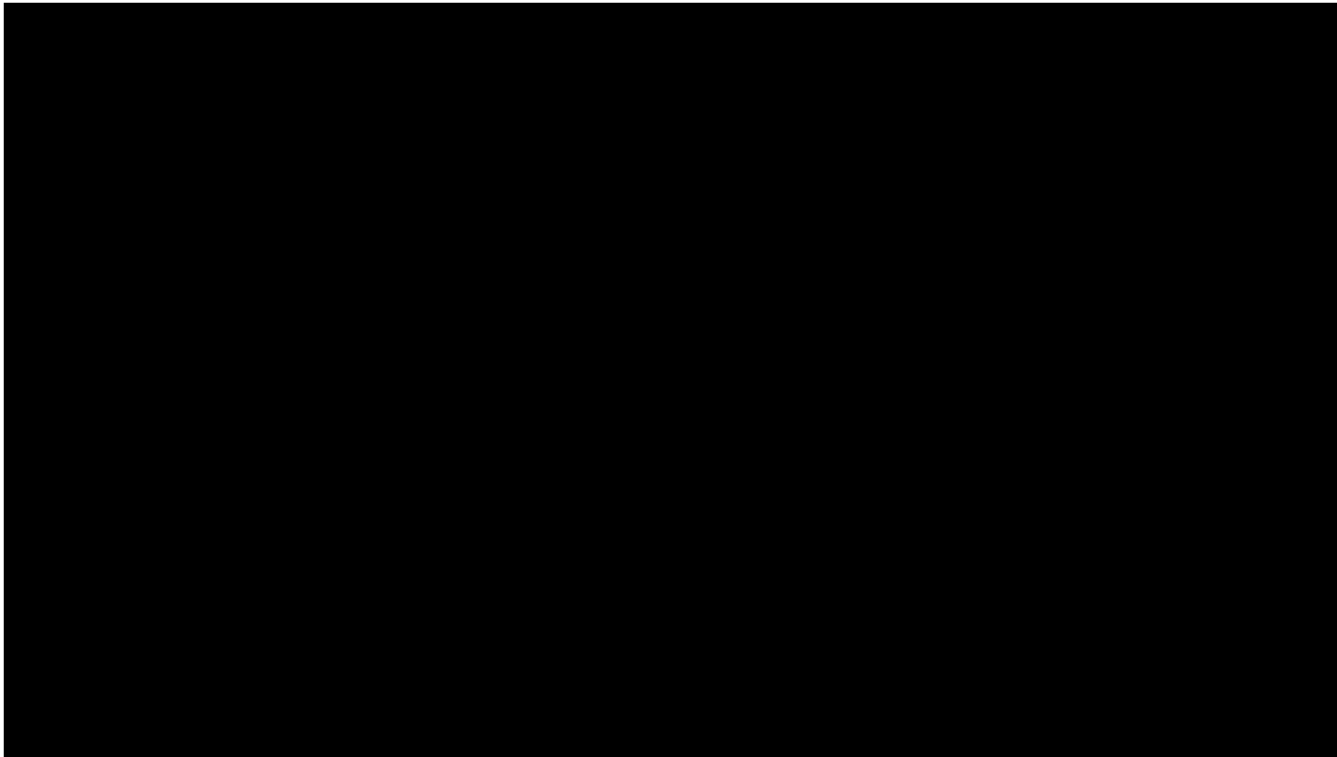
# Demo – Part1
# MapReduce Framework & HiveQL

# Demo – Part2(Using Spark-shell)
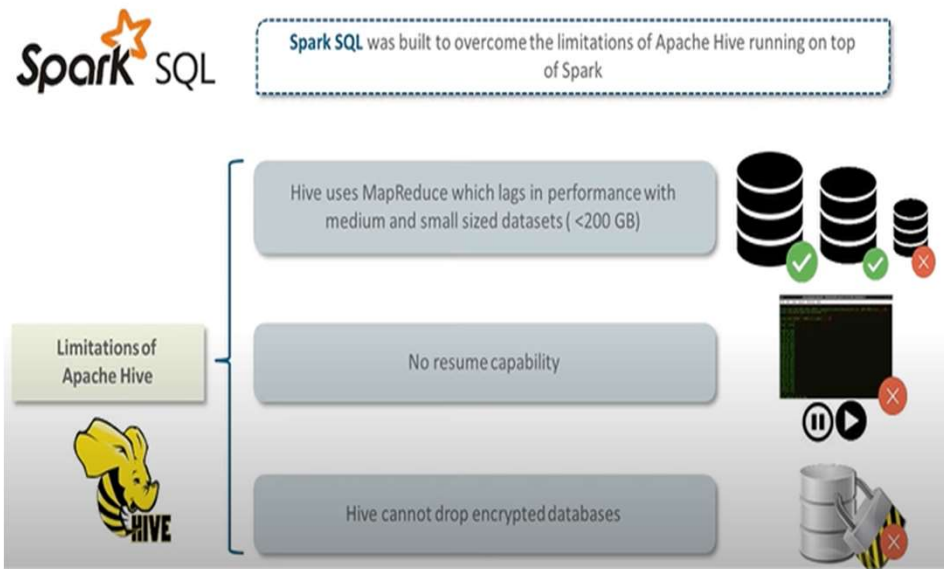## Apache Spark & Spark-SQL
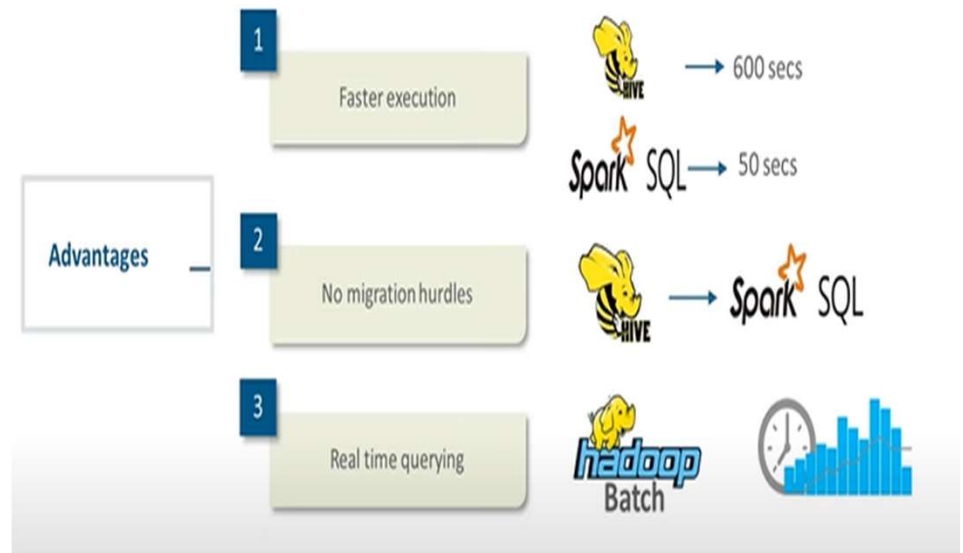
# Demo – Part2(Using Beeline)
## Apache Spark & Spark-SQL
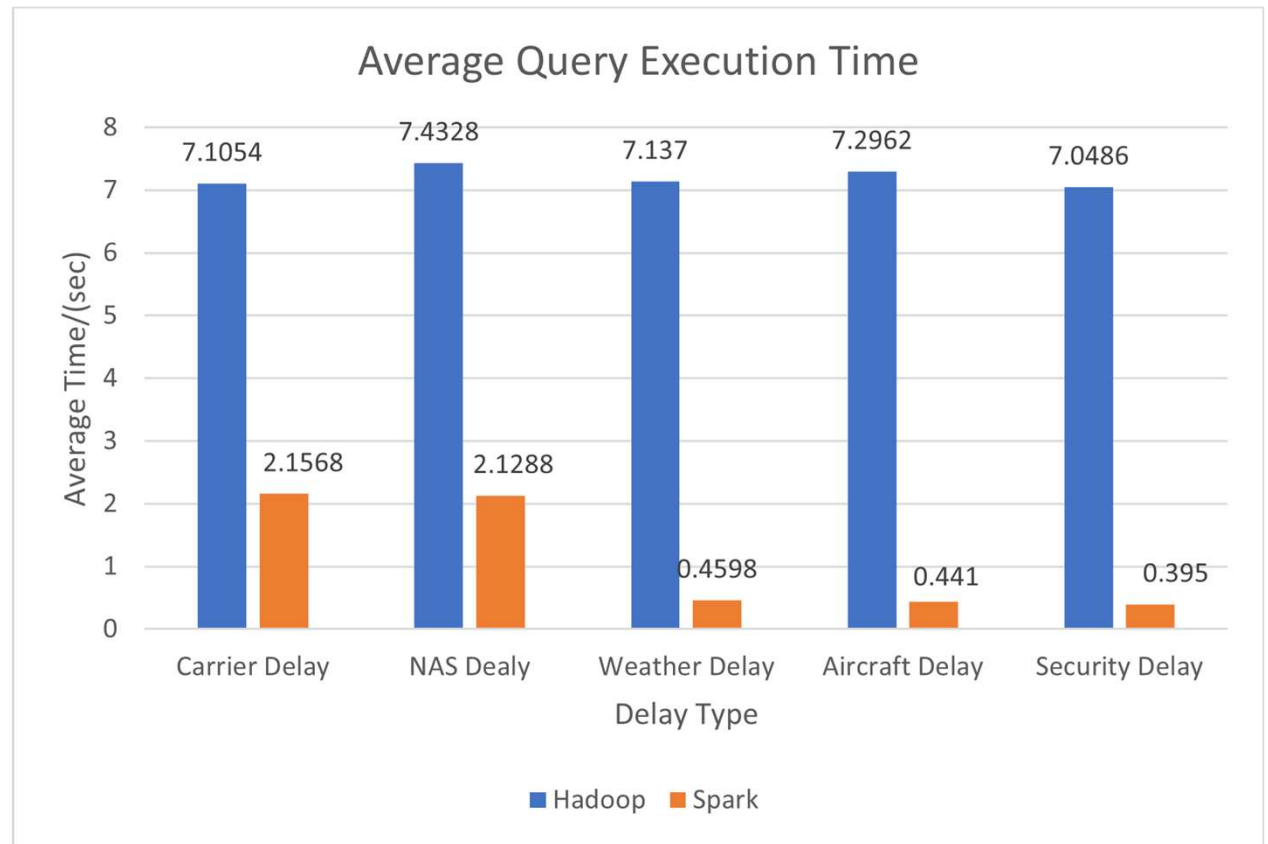
# Comparison

**Limitations of Hive**



**Advantages of Apache Spark**

# Fast Process?

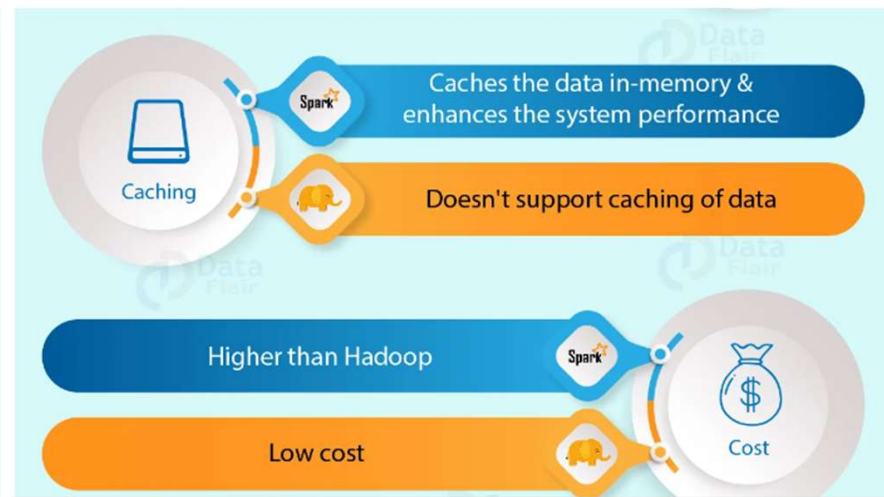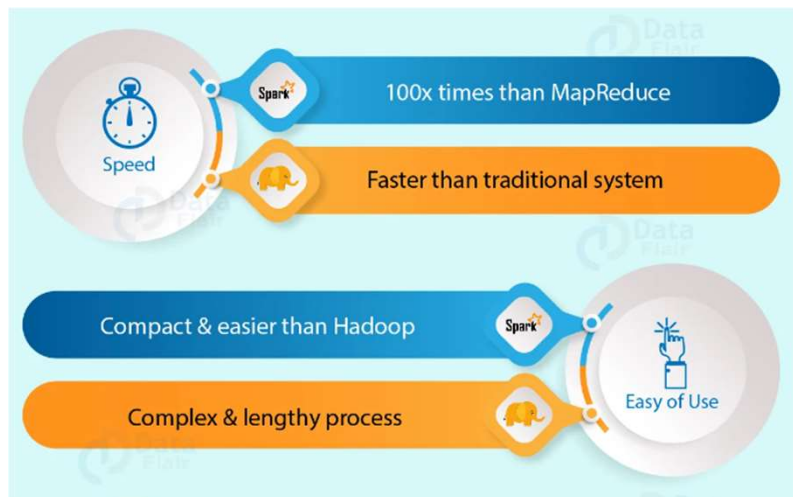# Apache Spark!!



Average Query Execution Time

# Conclusion

THANK YOU!