

Spark Assignmt.

Date:

Page:

- What is Spark and why do we need apache spark
- Spark is an open source, fast, general purpose in memory process for big data process
- Spark is a unified framework which provides all kinds of data process
- It is fast
- Provides high-level programmer
- Supports both real-time and batch processing
- Provides an interactive shell to see & learn about the data.

2) Difference between Hadoop map reduce & Spark.

Spark	Hadoop
→ Process speed is fast	→ Slower when compared to spark
→ In memory process	→ Read/ write from disk
→ Provides high level operator	→ Doesnt provide high level operator
→ Provides an Interactive shell	→ Doesnt provide Interactive Shell.

3) List the areas where Spark & Hadoop map Reduce are good.

- Tasks where map Reduce is good for

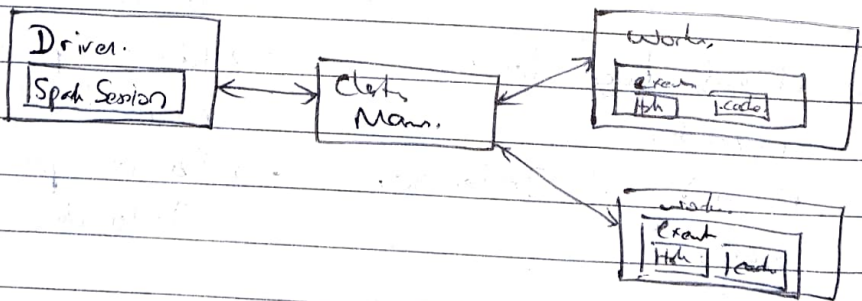
Linear process of huge data sets,
Hadoop map Reduce allow parallel
process of huge amount of data

Task wh Spark is Good for:

- Fast data Process
- Iterative process
- Real time process
- Machine Learning
- Join Datasets

4) How Spark manages the Cluster & Resources

- Spark is Designed to process a large volume of data efficiently and quickly. This distributed system is typically deployed onto a collection of machines which is called as Spark cluster. In Cluster, to efficiently & intelligently manage a collection of Machine, cluster intelligently manages a collection of machine.



5) Write a note on Spark driver & executor

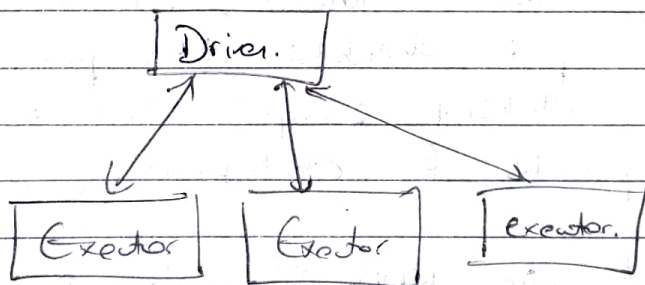
- Spark driver:

Date _____
Page _____

A spark application consists of two parts. The first is application data processing logic exposed using spark APIs and other spark drivers. The Application data Processing logic can be simple or few lines of code to perform.

Spark Executors

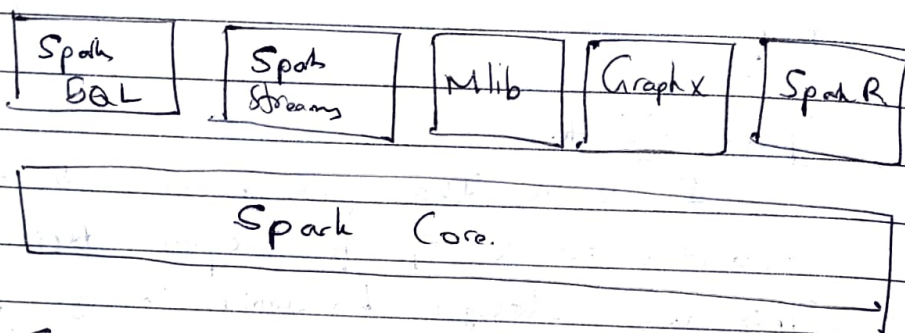
It is a JVM process and is exclusively allocated to a specific spark application. This was designed to avoid sharing spark executors between multiple applications. Lifetime of a spark executor is the lifetime of the application. Spark follows master-slave architecture.



Q) Write a short note on Spark unified stack & explain Spark core.

→ Spark provides a unified data processing engine known as the Spark stack. This stack is built on top of a strong foundation called spark core, which provides all the necessary functions to manage and run distributed applications such as scheduling, coordination & fault tolerance.

In addition it provides powerful generic program abstraction for data process called Resilient Distributed Datasets (RDD's)



Spark Core:

It provides all necessary functions to manage and run Distributed Application such as scheduling coordination & fault tolerance. It is the bedrock of Spark distributed data process engine. The distributed compute Infrastructure is responsible for distribution.

⇒ What are RDD's & Explain Properties of RDD

→ RDD's are immutable fault-tolerant, parallel data structures that let users explicitly persist intermediate results in memory, control their partition to optimize data placement.

→ Immutable RDD's are designed to be immutable which means you can't specifically modify a particular row in the dataset represented by RDD.

→ Fault tolerant:- The ability to process multiple datasets in parallel results requires a cluster of Machines to host & execute the computational logic. Spark automatically takes care of hardware failures.

→ Parallel Data Stru:- It is dividing the TB file into several chunks and executing on each chunk in a parallelized manner.

→ In Memory Computation:- RDD goes beyond speed boundaries by introducing the ability to perform in-memory computation.

Q) What are transformations & Actions?

→ In Spark the core data structures are immutable, so if we want to perform any change, the changes are defined by the transformations.

The actions are the triggers to start the transformations. The Reducetasks, map, flatten are the transformations and collect, show are the actions.

Q) What is ^{Lazy} copy evaluation?

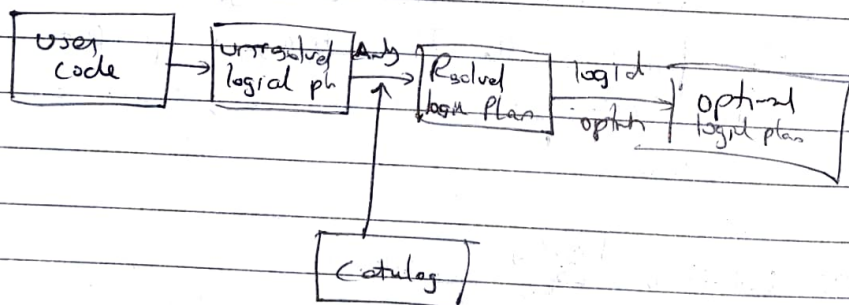
→ Lazy evaluation means that Spark won't execute until the very last moment to execute the graph of computation. In Spark, instead of modifying the data immediately when you express some operation, you build up a plan of transformations that you would like to apply to your source data.

By waiting until the last minute to execute the code, Spark compiles its plan from your raw dataframe transformation to a stream physical plan that will run efficiently.

b) What is logical & Physical plans and explain how structured APIs execute in Spark

→ Logical Planner

The first phase of execution is meant to take user code and convert it into a logical plan. The logical plan only represents a set of abstract transformations that don't refer to the executors or drivers, it's purely to convert the user set of expressions into the most optimized version.



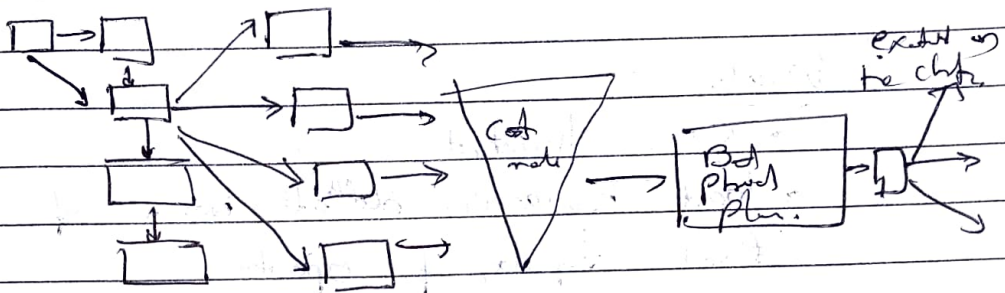
Physical Planner

After creating an optimized (logical plan), Spark then begins physical planning process. Physical plan often called as a spark plan specifies how the logical plan will

execute on cluster by generate different physical execution.

Optimal logical plan

Physical plan.



11) Explain uses of stream processing

- Notification & Alerts
- Real time reports
- Real time decision making
- Update the data to serve in real time
- Online machine learning

12) List out the challenges of stream processing

- Handling a stream out of order states
- Mass large amount of data
- Support high-data throughput
- Attacked load imbalance
- Respect to event at low latency
- Join with external data in other storage system.

13) Explain continuous stream processing model, what are its advantages & Disadvantages

- In continuous processing based system each node in the system is continuously listening to the system is continuously listening to the data streams.

from other nodes and outputs new updates to its child nodes

Advantages

- * offer lowest possible latency

Disadvantages

- * lower throughput
- * load balancing difficult because of fixed topology of operators

14) Explain Micro-Batch Stream processing model, what are its advantages & disadvantages

- It waits to accumulate small batch of input data then processes each batch in parallel using a distributed collection of tasks

Advantages

- High throughput per node
- Dynamic load Balancing

Disadvantages

- Higher latency due to wait to accumulate a micro-batch

15) Write a short note on DStream & structured stream

→ DStream:-

- Spark original DStream API has been used broadly for stream processing since its first release in 2012

→ API is purely based on process time to handle event time operations

Structured Streaming

→ Higher-level stream API built from the ground up on Spark Struct. API
 → Runs on both continuous and micro batch execution models

16) Write short notes on Data sources & Data sink of structured stream in spark

→ Data Source

→ Apache Kafka

→ File on distributed file system

→ A socket source for TCP

Sink

→ Apache Kafka

→ Almost any file format

→ A console sink for test

→ A memory sink for debugging