

Spark Assignment.

Q) what is Spark ? why do we require Apache spark
Apache Spark is an open source fast, and
general purpose in memory processing engine
designed for big data processing. It was

initially developed at the University of California, Berkeley
in 2009 and later became an Apache project in 2010. Spark
is written in Scala and supports high level APIs in Java,
Scala, Python and R. Spark is known for its in memory

fast processing
processing
speed

Computing Capabilities, which significantly speed up data processing.

why do we require Apache spark?

1) unified frameworks for different tasks : spark can handle different tasks like batch processing & real-time processing that would have previously required separate engines. This makes it a versatile tool for various data processing needs.

2) Speed : spark can process data 100 times faster than Hadoop mapreduce when data is processed in memory and 10 times faster when reading / writing from disk.

3) High level operators : spark provides easy to use higher level operators (like map, filter, etc.) to process data efficiently, which are not available in Hadoop mapreduce.

4) Interactive shell : spark offers an interactive shell (for scala and python) that allows users to explore and process data quickly, making it easier to learn and experiment with datasets.

5) supports for multiple storage systems : spark does not come with its own storage system, but it can work with a variety of storage systems like HDFS, S3 & Cassandra. This flexibility allows it to integrate with existing storage architectures.

6) Difference Btw Spark and hadoop Mapreduce

	Spark	Hadoop Map Reduce
Input	persisting in memory (partitioned data is processed in near memory, avoiding disk)	slow due to frequent disk I/O especially for iterative tasks.
Speed	up to 100 times faster than MapReduce for in-memory operations & 10 times faster for disk based operations	Best for linear batch processing of large datasets.
Latency	best suited for iterative tasks, perform real time analysis graph processing & XML	suitable for processing huge datasets where in memory operations are not feasible
Memory	can handle large data but volume may struggle with very handling large datasets when compared to mapreduce due to memory limitation	Requires more robust code in Java leading to more complex development.
APIs	one of provides high level APIs in Scala, Python, Java, C++, making it easier to write complex operations with considerably less code.	fault tolerance via replication & recovery from disk.
Performance	put fault tolerance via RDD lineage (allow you to recompute lost data).	primarily designed for batch processing with no native mutation support.
Scalability	supports real time processing using with spark streaming	requires third party libraries (eg mahout) but they operate under requires more disk space & I/O and its more economical when immediate results are not required.
Notable	has new built-in machine learning library that operates in memory (memory reuse efficient for iterative computation processes & makes better use of cluster memory).	needs external libraries for graph processing, handles failure by writing intermediate results to disk which ensures consistency but at a performance cost.
Graph Processing	comes with GraphX API for graph computation	
Reliability	Resilient via RDD's (current distributed datasets) when new data to be recomputed is lost.	

- 3) List the areas where Spark & Hadoop Mapreduce are good:
- Areas where Apache Spark is good:
 - > **Batch data processing:** Spark's in-memory processing makes it up to 100 times faster for operations that can be done in memory and 10 times faster for disk-based operations compared to Hadoop.
 - > **Iterative processing:** Spark excels in scenarios that require repeated processing of the same data, such as ML & graph processing. Its Resilient Distributed Datasets (RDDs) allow iterative operations in memory without writing intermediate elements to disk.
 - > **Near Realtime processing:** Spark is ideal for use cases requiring near real-time insights, such as live data streams, because of its Spark Streaming component.
 - > **Graph processing:** Spark's GraphX API is designed for graph computation, which requires iterative algorithms. Spark handles these efficiently with in-memory processing.
 - > **Machine learning:** Spark comes with a built-in MLlib library for machine learning, offering fast iterative data processing.
 - > **Joining datasets:** Spark is well suited for joining datasets, especially when combining smaller datasets quickly due to speed.

Areas where Hadoop Mapreduce is good:

- > **Linear processing of large datasets:** Hadoop Mapreduce is ideal for sequential processing of very large datasets that cannot fit in memory. It breaks data into chunks, processes them in parallel, and writes the results to disk.
- > **Handling very large partitions:** Hadoop Mapreduce can process big datasets because it uses disk-based storage, managing better suited for tasks where the data volume is too large to fit in memory.

- PAGE NO. _____ DATE _____
- > **Cost of efficient processing:** Hadoop Mapreduce is an economical solution when speed is not critical, such as batch processing jobs that can run overnight.
 - > **Data durability & fault tolerance:** MapReduce's reliance on writing intermediate data to disk provides stronger fault tolerance for certain scenarios where data loss is unacceptable. If a node fails, data recovery is straightforward due to the persistence of each chunk.
 - > **Joining very large datasets:** For very large datasets that require a lot of shuffling and sorting, Hadoop Mapreduce may outperform Spark as its disk-based approach can handle greater data volumes without running into memory limits.

How spark manages the cluster & resources?

spark is a distributed computing system that processes large datasets across multiple machines, referred to as spark cluster. Resource management and task scheduling in spark are handled through a resource manager such as YARN or Mesos. The main components that manage the cluster & cluster manager: The cluster manager is responsible for coordinating resources across the machines in the cluster. Spark relies on external resource managers like Apache YARN or Apache Mesos, or can use its own built-in standalone resource manager.

Responsibilities of the cluster manager:

- Managing the creation of resources in the cluster, as well as the available memory & CPU cores.
- Organizing the distribution of work by assigning tasks to co-located nodes on free available resources.

Spark driver is the central coordinator of a spark application. It manages the execution of tasks and interacts with cluster manager.

* Responsibilities of the spark driver:

- Communicates with the cluster manager to launch parallel tasks called executors on the worker nodes.
- Distributes the tasks (pieces of the data processing logic) to the workers.

④ Spark SQL: handles structured data processing and appears both batch and interactive queries using SQL.

⑤ Spark streaming: supports real-time data processing allowing the analysis of live data streams.

⑥ spark MLlib: Spark's built-in machine learning library offering various algorithms for both training & inferences in tasks.

⑦ spark graphx: designed for graph processing it supports complex computations typically required for graph analysis.

⑧ spark R: an API for running R code using integrating spark power with the R language for statistical computing.

Spark core is the central component and the backbone of Apache spark. It provides the necessary functionalities for running and managing distributed applications on a cluster. spark core includes fault tolerance, shuffling, task coordination, and data movement among distributed nodes. It also provides the programming abstractions.

And APIs for development to build large scale data processing applications. Key features of spark core:

① distributed computing infrastructure: Manages the distribution of data in computation across multiple machines.

② distributed computing framework: Manages the distribution of threads task scheduling, coordination and fault tolerance across the distributed data processing system.

③ supports distributed processing - the movement of data between machines.

④ Resilient distributed datasets (RDDs) - RDDs are the foundation of distributed in spark core. They are immutable, fault-tolerant and distributed collections of objects that can be processed in parallel over a spark cluster.

⑤ and distributed collections of objects that can be operated in parallel.

• RDDs always will be explicitly persist intermediate results in memory. Ensuring efficient processing much faster than serial processing. Following follows like machine learning, much faster.

Supports immutability operations.

① In-memcached computation: Spark core deals with memory contention when dealing with both performing transactionalized distributed operations. This is particularly useful for transactional computing and

② fault tolerance: spark core provides fault tolerance by using RDD lineage.

③ In any part of the computation fails, spark can recompute lost data.

④ broad scatter becomes transformation applied to the input data.

⑤ rich set of operations - provide user applicability wide variety of operations from data transformation and cumulative reduce or map.

⑥ user defined join, group by ordinary others, these operations also allows developers to easily express complex data processing operations.

⑦ language support: Spark core provides APIs in native languages namely Scala, Java, Python, and R making it accessible to a wide range of developers.

⑧ what are RDDs? Explain the properties of RDD.

Resistant Distribution Dataset (RDD) are the fundamental data

collections of objects that can be processed in parallel over a spark cluster.

RDDs allow us to perform transformations and actions on large distributed efficiently and with fault tolerance. RDDs provide a high level abstraction.

for distributed data processing and support fault tolerance through lineage they are available in-memory data processing. Which significantly speeds up iterative computations like these machine learning algorithms by keeping the properties of RDDs.

• RDDs are immutable meaning they cannot be changed after they are created.

• RDDs are automatically partitioned fault-tolerant by keeping track of the transfer.

• RDDs always will be explicitly persist intermediate results in memory.

• parallel processing allows data to be distributed across multiple machines.

• memory efficiency follows like machine learning.

Log in to your application and run the command `git pull`. This will download the latest changes from the remote repository.

Run `git status` to see what files have changed. You can use the `-n` option to see the number of changes:

```
git status -n
```

The output will look something like this:

```
On branch master
Your branch is up-to-date with 'origin/master'.
Changes to be committed:
  (use "git add <file>" to include in what's committed)

    modified:   README.md
    modified:   index.html
```

To commit the changes, run `git commit -m "Initial commit"`:

```
git commit -m "Initial commit"
```

Finally, push the changes to the remote repository:

```
git push origin master
```

You can also use the `-u` option to set the upstream branch:

```
git push -u origin master
```

After pushing, you can see the changes on GitHub or your application's dashboard.

If you want to make changes to a specific file, you can use the `git add` command:

```
git add README.md
```

Then commit the changes:

```
git commit -m "Add README.md"
```

Finally, push the changes:

```
git push origin master
```

That's it! You've successfully pushed changes to a GitHub repository.

- (a) Explain how the car's of autumn powalling
 (b) What is the relationship between the time constants of certain processes occurring during autumn, allowing prolonged adaptation to changing conditions.
 (c) What are the main activities and movements: stream powalling and
 (d) What would stimulate plants into this adaptation?
 (e) Explain what activity most strongly stimulates this adaptation.
 (f) How could seasonal variation in light be used to explain this phenomenon?
 (g) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (h) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.

(i) Explain how the car's of autumn powalling
 (j) What is the relationship between the time constants of certain processes occurring during autumn, allowing prolonged adaptation to changing conditions.
 (k) What would stimulate plants into this adaptation?
 (l) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (m) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (n) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (o) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (p) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (q) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (r) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (s) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (t) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (u) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (v) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (w) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (x) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (y) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.
 (z) Explain how seasonal variation in light and day length affect the timing of seasonal events such as flowering and fruiting.

Q(1) State management & event time vs processing time

- Out of order and late \rightarrow Scalability
- Complex event processing & Data consistency \leftarrow exactly once processing \rightarrow windowing & fragmentation
- Monitoring & Debugging & Resource management
- Latency sensitivity \rightarrow Data integration, resilience.

Q(2) Explain continuous stream processing methods. What are its advantages & disadvantages?

Continuous Stream processing is a method in which data is processed as soon as it arrives, typically in real time or near real time. This done by continuously ingesting and analyzing a stream of data events with low latency requirements: real time processing, low latency event driven stateful processing.

Advantages of Continuous Stream Processing:

- 1) Real time insights - Low latency, Scalable & distributed
- 2) Event driven architecture, efficient use of resources.
- 3) Stateful processing

Q(3) Explain batch stream processing methods. What are its advantages and disadvantages?

Q(4) Explain micro batch stream processing method, what are its advantages and disadvantages?

Continuous stream processing method that processes data as it arrives in real time, often in the form of continuous, unbounded data streams. This method enables immediate analysis and action on data as soon as it is processed.

Without a complete dataset (as in batch processing), it is essential for applications that require low latency real time division making and continuous data monitoring.

Key features of continuous stream processing:

Event driven processing, low latency, scalable.

Advantages of Continuous Stream Processing:

- Real time processing, low latency, scalability
- Event driven processing, efficiency with stateful operations
- Complexity in system design, data ordering and late data
- State management complexity, higher processing costs
- Fault tolerance, consistency and exactly one processing
- Monitoring and debugging.

Q(4) Explain micro batch stream processing method, what are its advantages and disadvantages?

Micro Batch Stream Processing:

Micro Batch Stream processing is a batched stream data flow and divided into small time based batches for processing. Instead of processing after event by event (as in continuous stream processing), the incoming data is collected for a short interval and then published as a single batch.

Key features of Micro Batch Stream Processing: Time based batching, near real time processing (fast turnaround).

- Advantages of New Batch Stream processing
- 1) Simplified fault tolerance 2) Never Real time processing
 - 3) Simple System Design 4) Efficient Resource Usage
 - 5) Windowing and Aggregation 6) Scalable and Distributed
- Disadvantages of new Batch Stream processing
- 1) Increased latency 2) Limited Real time capabilities
 - 3) Handling late and out of order data 4) Complex Aggregation Across Batches 5) Higher Resource Consumption
 - 6) Resource contention during peaks.

Q15) Write short notes on Dstreaming and Structured Streaming
 Dstream (discretized stream) is the core abstraction in Apache Spark Streaming, used to represent a continuous stream of data in micro batches. Dstreams are built on top of Resilient Distributed Datasets (RDDs) where incoming data is split into small batches and each batch is processed as an RDD. This method allows Spark Streaming to provide fault-tolerant, Scalable and real-time stream processing.

Stream = `ssc.streaming("localhost", 9999)`.

Words = `stream.flatMap(lambda line: line.split(" ")).`

wordcounts = `words.map(lambda word: (word, 1)).reduceByKey(lambda x, y: x+y)`

`ssc.start()`.

Advantages of Dstreams:

- * easy integration with RDDs
- * fault tolerance

Disadvantages of Dstreams:

- * high latency
- * limited real time processing
- * complexity with Advanced use cases.

Structured Streaming is the next generation Stream processing engine introduced in Apache Spark 2.0. It builds on top of Spark SQL's DataFrame and Dataset APIs, providing a declarative approach to stream processing. Unlike Dstream's Structured Stream is designed for both real-time and batch processing offering low latency and continuous processing capabilities with strong exactly once guarantees.

Key features of Structured Streaming:

- * Declarative API
- * Continuous processing
- * Event Time handling
- * Fault tolerance
- * Exactly Once Semantics

cursor = `spark.readStream.format("socket").option("host", "localhost").option("port", 9999)`.

words = `cursor.selectExpr("split(value, ' ')").alias("word")`

wordcount = `words.groupBy("word").count()`.

query = `wordCounts.writeStream.outputMode("complete").format("console").start()`.

Query. awaitTermination().

Advantages of Structured Streaming:

- * Unified API for Batch and Streaming
- * Low latency
- * Event time support
- * Exactly once Semantics

Disadvantages of Structured Streaming:

- * Complexity for low level control
- * Maturity

Q16) Write short notes on Data Sources and Data sinks of Structured Streaming in spark.

Data sources are they entry points through which streaming data is ingested into spark's Structured Streaming environment.

Structured streaming supports a variety of data sources allowing us to process real-time data from different systems and formats.

Common data sources:

- file source & socket source & Kafka source & Rate sources

Data sinks in Structured streaming are the endpoints to which processed streaming data is written in ~~so~~ it is sent to a data sink for storage, display or further analysis. Structured streaming supports a range of data sinks to suit various output requirements.

- * file sink & console sink & Kafka sink & Memory sink
- * foreground sink