Which commonly available Hadoop distribution is offered as a provision-able service through both Amazon Web Services Elastic MapReduce (EMR) and Google Compute Engine?



MapR

- Cloudera
- O Hortonworks
- Hadapt

Which of the following vendors offers an MPP Data Warehouse appliance?



Tableau

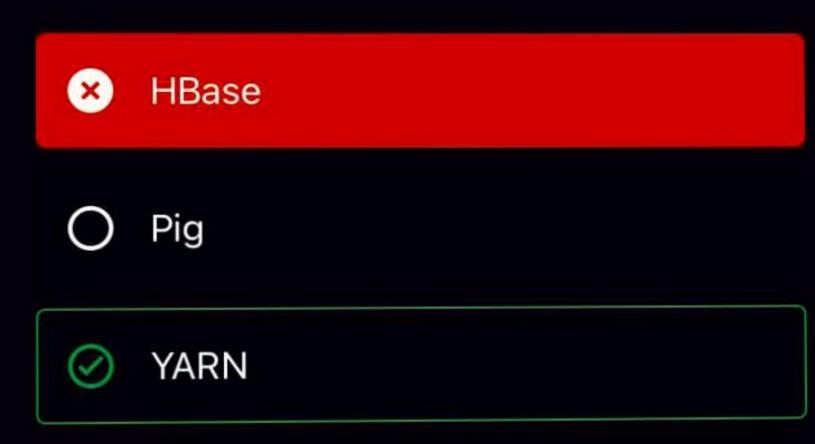
Teradata

O Informatica

Pig was developed by Facebook as a data warehouse application for Hadoop.



What new resource manager was released in HDFS 2.0?



Amabari

Massively Parallel Processing (MPP) data warehouses are always queried using imperative code.



Which of the following is a use case for Sqoop?



Bulk transfers of data of structured data into HDFS

- O Moving large amounts of log files into HDFS
- O Analysis unstructured or semistructured data
- O Allows random read-and-write in MapReduce

Which command shares the same functionality as HDFS dfs commands?



Hadoop fs

- HDFS dfsadmin
- Hadoop dfs
- Hadoop ds

HDFS allows for files to be updated.



True

Hadoop hybrid products allow MapReduce querying of relational data and/or co-locate Hadoop and a relational database engine on the same cluster.



Which of the following is not a typical Hadoop stack member/distribution component?



```
id: integer
name: string
account_idlist: array<string>
```

What is the query you would run to get the list of all accounts along with the customer ids for all customers?

O select id, explode(account_idlist) as account_ids from customers;

> select id, account_ids from customers lateral view

explode(account_idlist) as account_ids;

×

select id, account_ids from customers

lateral view explode(account_idlist) accountsTable as account_ids;

Which of these scheduling policies have the longest wait time?

C Last In First Out Scheduler

Capacity Scheduler

Sair Scheduler

First In First Out

Orders: 100GB

Products: 500MB

Reviews: 1GB

If the most common join operation performed was on the Products and Reviews tables and you wanted to make the join operation between these tables as fast as you possibly could, what design decisions would you take?



Bucket the table on the product id column, have the number of buckets in one table be a multiple of another and join on the product id column.



Bucket the table on the product id column, have an equal number of buckets on each table and sort each bucket on product id, join on the product id column.



Which of these is the difference between a managed and an external table?

Managed table metadata is stored in
 HDFS, and external table metadata is stored in the metastore.

Managed table records are stored in the warehouse directory, and external table records are stored elsewhere in HDFS.

Managed table records are stored in the metastore, and external table records are stored in HDFS.

What is the container that YARN uses when scheduling tasks on the cluster?

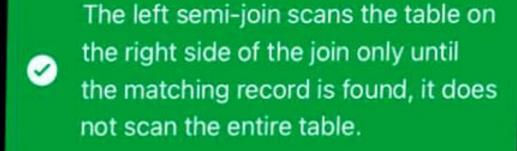
- The memory and CPU required across all nodes in the cluster to complete the task
- The logical unit of all the resources that a process needs, memory and CPU
- O The memory and CPU required on the name node to process a task
- O The storage capacity of a node which is required to process a task

See results

What is the default location for Hive data?

- O /user/hive/warehouse on the local file system
- /user/hive/warehouse on HDFS
- O /warehouse on the local file system
- /warehouse on HDFS

Why is the left semi-join faster than the IN/EXISTS subqueries?



 column in the join operation which is more efficient than using multiple columns like in the IN/EXISTS subqueries.

The left semi-join uses only one

The left semi-join runs entirely in memory which makes it more performant.

What are the advantages of using the Hive Beeline interface over the older Hive command line?

- O Faster and more efficient queries because of additional parallelism
- Support for multiple concurrent clients and more robust security
- O Better formatting and color coding on the terminal
- O Greater fault tolerance due to more replicas of data

Question 1 of 10

Which of the following is NOT a reason that map-only joins are faster than joins which run both map and reduce phases?

Complete MapReduce operations

have additional steps such as shuffle

reduce phases which gets eliminated in map-only joins.

and sort between the map and

8

- O Map-only joins reduce data transfer within a cluster.
- O Improved processing time because one phase is entirely eliminated.
- The memory overhead of a
 MapReduce job is much greater than a Map-only job.

×

Which are possible ways in which join performance can be optimized?

Ø

Structuring joins such that the reduce phase of the underlying MapReduce is eliminated

- O Matching of join records is made really fast by caching data
- O Ensuring that multiple processes work on data in large tables
- O reduce phase is run on more than one machine

How would you specify a window between the first row in a result set and the current row?



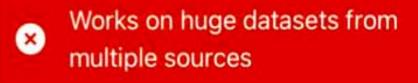
rows between 6 preceding and current row



rows between unbounded preceding and current row

- O nothing, this is the default specification
- rows between unbounded preceding and and unbounded following

Analytical processing jobs include all of these characteristics EXCEPT:



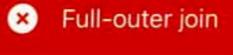
- O Reads data, does not usually edit or update it
- Deals mostly with recent data,
 collected in the last few hours or days
- Tends to involve long running jobs

×

Which of the following is an important reason that Hive stores data in a denormalized form?

- Joins are performant and can be performed to get all data related to an entity in one go.
- Redundancy is not acceptable even though disk space is cheap.
- Read operations are more common in Hive, so disk seeks should be minimized for better performance.

When the table on the right side of a join operation is much smaller than the table on the left side of a join operation, which of the following are map-only joins?



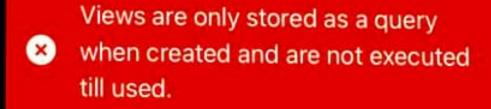
- Right-outer join
- Left-outer join



What is the name of the tool used to automate deployments in the cloud?

- O Cloudera Distribution including Hadoop
- Cloudera CDH
- Cloudera Director
- Cloudera Manager

Which of these is NOT a characteristic of views in Hive?



- O Views can include one or more tables in the query.
- Data in the view is frozen in time and does not change along with the underlying data in the table.
- The structure is frozen in time and not affected by changes in the underlying table.

Which of these is NOT a supported scalar type in Pig?



Why are joins with one join column faster than joins with 2 join columns?



The number of MapReduce jobs run to perform the join is equal to the number of join columns. Fewer columns, fewer jobs; so faster.

- The amount of data processed is increases as the number of join columns increase. Less data; faster joins.
- The amount of data loaded into memory is directly proportional to the number of join columns. Fewer columns, less memory used; so faster.

Question 2 of 10

Which are the names of the HBase processes that run when HBase is started up in pseudo-distributed mode?

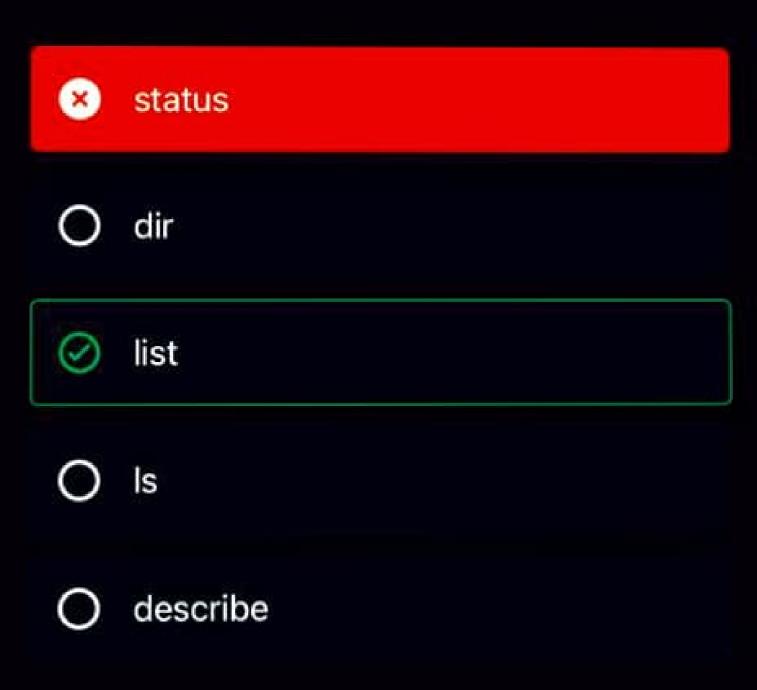
- RegionServer, Master, Zookeeper
- NameNode, DataNode and SecondaryNameNode
- × NodeManager, ResourceManager
- HRegionServer, HMaster, **HQuorumPeer**

What do you call the main program in a Spark application?



- Session
- Context

Which command would you use to see all tables that are present in HBase?



Question 6 of 10



Which of the following are the only 2 operations possible on an RDD?

- Transformation and Action
- Aggregation and Update
- Transformation and Retrieval
- Update and Retrieval
- Aggregation and Action

Which of the following functions is NOT a valid load/store function Pig?

```
**JsonLoader()**

**HBaseStorate()**

**AvroStorage()**

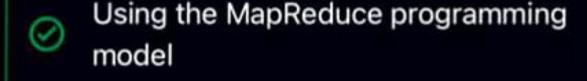
**CSVExcelStorage()**

**RDBMSStorage()**
```

How would you perform grouping and aggregation operations on HBase?



- O Using multiple shell commands piped together
- O Using the Structured Query Language



personal:name = name of each student academic:year = grade in which the student currently studies row key = student id
You want to find the total number of students in each grade. What would be the operation in the reduce phase?

students for each grade key.

Find the maximum number of

×

Find the minimum number of

O Average the number of students for each grade key.

students for each grade key.

Sum the number of students for each grade key.

Question 4 of 10



What does CSD stand for?

- Custom Service Descriptor
- Cloudera Service Description
- Custom Service Definition
- Cloudera Standalone Data

What data structure represents the actions and transformations in a Spark application?

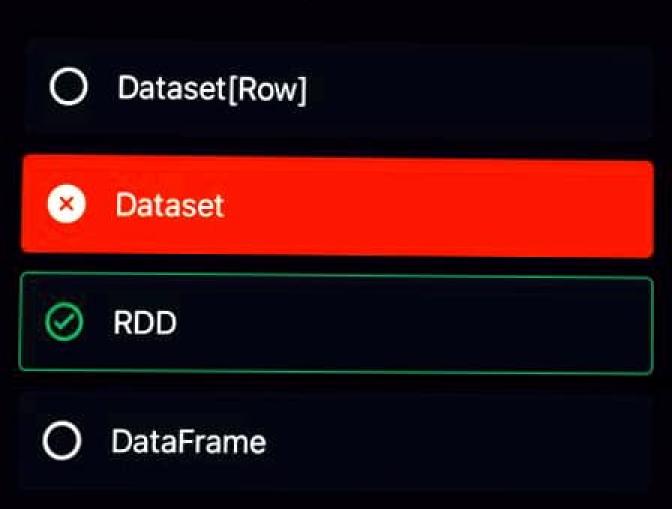


- Maximum Heap
- Queue
- Directed-acyclic Graph
- Undirected Graph

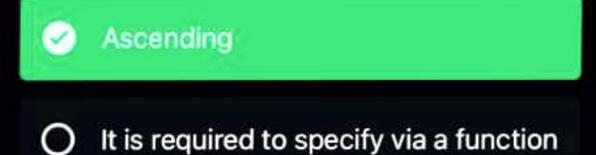
Which is the correct Spark method to use to convert a Python list to an RDD?

- sc.fromDF()
 - Sc.make()
 - O sc.toRDD()
 - sc.parallelize()

What is the name of the original core abstraction in Spark?



Which is the default sort order in Spark?



O Descending



Which format do you use to read from external databases?

- O format('database')
- O format('mysql')
- format('jdbc')
- format('external')

Which of these functions is used to return data back to the driver?



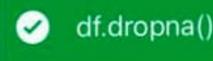
Question 10 of 10



Which is the correct Spark method to use to convert a Python list to an RDD?

- osc.make()
- O sc.toRDD()
- sc.fromDF()
- sc.parallelize()

You're analyzing the data loaded into a dataframe and you find that some records have some fields missing. How will you get rid of these records so you don't process them further?



- O df.filter(df['some_col'] != NaN)
- O df.dropmissing()
 - df['col1', 'col2', 'col3'], keep only
 those columns with no missing information

Which of these are examples of stream processing?



Analyzing new users to see whether the rate of new signups is higher on mobile or the website

- Analyzing reviews to see if the overall sentiment for a product is good or bad
- Analyzing orders on an e-commerce
 site to see what zip codes they're
 concentrated at



Analyzing log messages to see whether all pages of the ecommerce site are up and running

What is the major difference in applying k-means clustering on batch vs streaming data?



Determining how relevant older data is compared with newer data in the stream

- The algorithm needs to scale to handle more data as compared with batch operations.
- O More iterations of the algorithm as newer batches of data arrive
- O The algorithm changes based on how fast the newer data arrives.

What does the explode() function in Spark do?

- O Creates a new batch for each group of elements
- O Creates a new window for each group of elements
- Creates a new row in the dataframe for every element
- O Creates a new column for each element

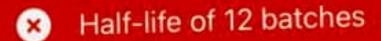
Given a NoSQL database, which component is most properly suited to transfer and adapt the data to stream processing?

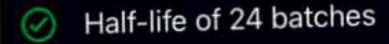


- C Kafka Producer
- Kafka Connect
- C Kafka Broker

You're tracking trending news stories over the last 2 days with your streaming application and you receive a new batch of updates every hour. **News received a day ago is only half as important as the news received in the last hour.

**How would you specify this forgetfulness in your streaming k-means algorithm?





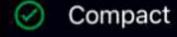
- Half-life of 1000 news articles
- C Half-life of 24 articles



Given a topic that has to store userrelated information (name, address, etc.), which type of clean-up policy would be the most appropriate?



None



O Delete

What is the relationship between a DStream and an RDD?



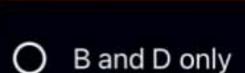
A DStream is made up of a sequence of RDDs, where every RDD contains entities which are received within a batch interval.

- An RDD is a group of DStreams
 where each DStream contains
 entities received within a batch
 interval.
- A DStream is just one RDD that contains all the entities available in a stream.

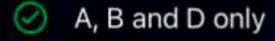
Which of the following are characteristics of RDDs in Spark? A: They are distributed across multiple machines in the cluster B: They can be reconstructed when a node crashes C: They are stored on the hard disks

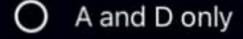
of cluster machines D: They are lazily evaluated

A, B and C only









How does a subscriber in Kafka receive messages from a publisher?

0

A subscriber subscribes to a topic and receives messages from all publishers to that topic

- A subscriber subscribes to a message and receives all messages which are of the same type
- A subscriber subscribes to an individual publisher and receives all messages that the publisher sends out



Which of these operations require real-time processing?



Requests on sale days to an ecommerce site



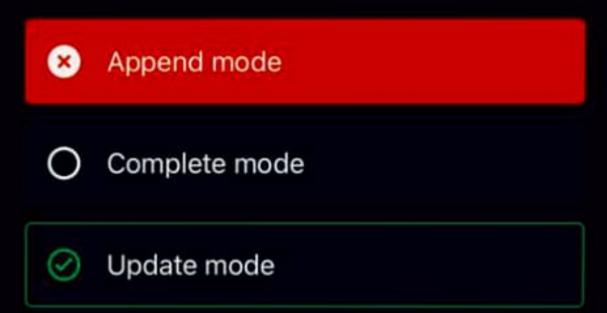
Daily indexing of web pages to serve search operations

- Setting up a heat map of the location
 of the most active Twitter users in the last month
- Running operations to see weekly trends in sales on an e-commerce site

Your e-commerce sites sends every transaction with product id as a stream to a Spark application. The data sent can be assumed to always arrive on time.

You want to keep a global count of how many items of each product were sold and save this to an RDBMS.

What kind of output mode will you use when you write out this information if you want to minimize updates to the RDBMS?



You're tracking trending news stories over the last 24 hours with your streaming application and you receive news updates every hour. What kind of decay factor will you use to perform streaming k-means clustering?



0 < decayFactor < 1

- O decayFactor = 0
- O decayFactor = 1

Question 2 of 10



What is a broker in Kafka?

- A single topic in a Kafka cluster
- A single node in a Kafka cluster
- A single subscriber in a Kafka cluster
- A single publisher in a Kafka cluster

multiplicative operation represented by this pseudocode:

```
multiple(a, b) {
  return a * b
}
```

**What is the pseudocode for the inverse function **of this in the sliding window operation?

```
subtract(a, b) {

× return a - b
}
```

```
sum(a, b) {
    return a + b
}
```

```
divide(a, b) {

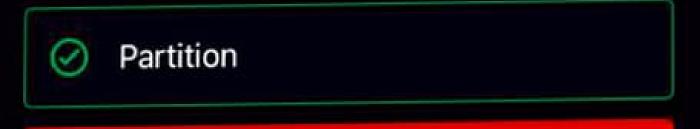
return a/b
}
```

Next question

Question 4 of 10



Deleting records before certain offsets can be done by which of the following?





- O Broker
- Cluster

Given a Kafka cluster composed of five brokers, what is the maximum number of replicas that a topic can have?



- Two Replicas
- One Replica
- Ten Replicas

Which tool should a new custom application be integrated with to produce data to Kafka?



Kafka Connect

- Admin Client
- O REST Proxy



Kafka Producer

How does Kafka allow to implement transactions across microservices?



It can serve as a WAL for distibuted systems

- O It implements a two-phase commit protocol across multiple databases
- It should be used as the only
 database in a system that should serve all user queries directly
- Different databases can coordinate
 update operations by
 communicating via Kafka

What is "stateless stream processing" refer to?



".map()" function in Kafka Streams

O Storing an intermediate state to process each record



Processing records in a stream when each record can be processed independently

All processing implemented using Kafka Streams

How are records ordered in Kafka?



Records are ordered per partition in lexicographical order

- O Records are ordered per topic in the order they were written to Kafka
- Records are ordered per partition in the order they were written to Kafka
- Records are not ordered in Kafka

Question 9 of 10



What approach to analytics allows us to perform ad-hoc queries?

- O Either analytics-on-read or analytics-on-write
- Analytics on-read
- Analytics on-write
- Neither of these approaches

What are the components of a Kafka record?

Topic name, key, value, headers, timestamp

Key, value, headers, timestamp

O Topic name, headers, timestamp

Key, headers, timestamp

How can we track removal operations with immutable events?

- O We should write removal events to a separate datastore
- O We can remove arbitrary records from a Kafka topic
- We need to create a separate event that marks a removal operation
- We should copy all existing events

 from a Kafka topic, except the removed event

How does Kafka allow to implement transactions across microservices?



It can serve as a WAL for distibuted systems

- O It implements a two-phase commit protocol across multiple databases
- It should be used as the only database in a system that should serve all user queries directly
- O update operations by communicating via Kafka

What is the main benefit of using Kappa architecture over Lambda architecture?



Lambda architecture forces only to use mutable data instead of building on top of immutable events.

- C Kappa architecture allows implementing real-time applications, while Lambda architecture is only about batch processing applications
- O Lambda architecture does not allow to use Kafka



With Kappa architecture, each algorithm should only be implemented once

What is the name of a window type in Kafka Streams that divides records into overlapping groups?



Session window

Tumbling window

Sliding window



Hopping window

What is the benefit of reading WAL to reflect updates from a database in Kafka?

O It allows preserving timestamps of records

It is easier to implement than alternatives

O WAL is supported by more datastores

It allows to keep track of updates and deletions in a database What approach would you use to implement a complex distributed transaction if the goal is to make the result as maintainable as possible?



Choreography

- Neither of these options
- Administration



Orchestration

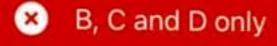
Which of the following are input arguments to the window function in structured streaming?

A: The event time

B: The window interval

C: The sliding interval

D: The batch interval



- A and D only
- A and B only
- A, B and C only
- A, C and D only

Given a high load topic with small message sizes and no requirements from the latency perspective, which property would help to increase the throughput the most?



max.in.flight.request.per.second

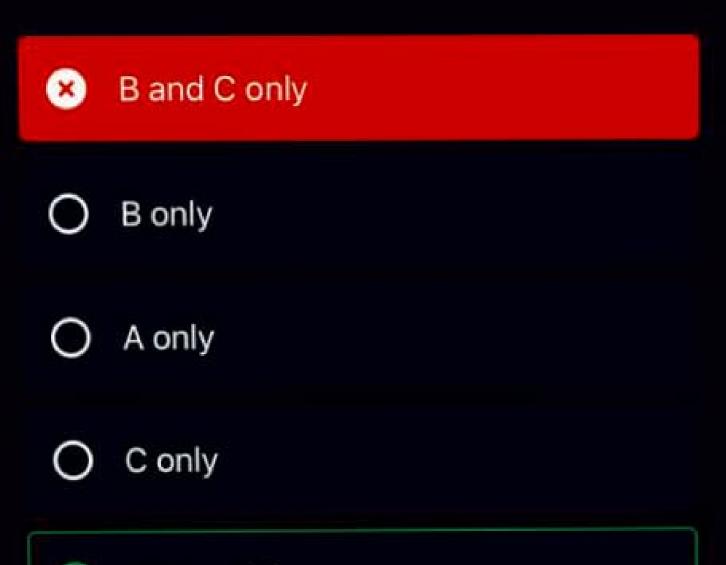
Which of the following are important characteristics of any stream processing system?

A: Ability to work with timesensitive data

A, B and C

B: Manage out of order events

C: Replay data for recalculations



Question 3 of 10

X

Your e-commerce sites sends every transaction with product id as a stream to a Spark application. The data sent can be assumed to always arrive on time.

You want to keep a global count of how many items of each product were sold and save this to an RDBMS.

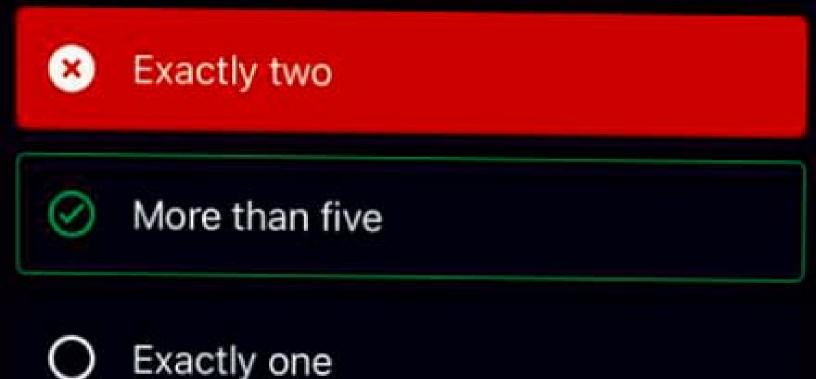
What kind of aggregation would you perform on the incoming stream?



Group by product id and count the number of products

- O Group by day and count the number of products
- O Group by sale price and sum the number of products

Given a topic called "orders," how many producers can you use to produce data on this topic?



O Less than five

Question 8 of 10



Which of the following is an improvement present in Spark 2 but not in Spark 1?

- O Ability to work with dataframes and datasets
- Unified APIs for batch as well as stream processing
- Lazy evaluation of datasets when an action is requested
- Fault-tolerant distributed datasets

What kind of mechanism do Kafka Consumers use for retrieving records?

O Event-based

Ø

Pull

O Push

Command

Which of the following is NOT a valid difference between batch and streaming data?



The order of data received is generally more important in streams

- O Streams are constantly updated, batches may not be
- O Batches are bounded datasets, streams are unbounded datasets



Streams have a single global state that does not change over time

Why are RDDs, the basic Spark programming abstraction, considered resilient?



RDDs keep track of their lineage; they can be reconstructed from the source even when cluster nodes crash.

After every transformation, an RDD

- has its data replicated to the
 memory of multiple nodes and can be retrieved from other machines if nodes crash.
- After every transformation, an RDD has its data replicated to multiple nodes in the file system and it's permanently stored on more than one machine.

Question 9 of 10

You receive log messages at the rate of two messages per second and your batch interval is five seconds.

If the cliding window size is 20

your batch interval is five seconds.

If the sliding window size is 20 seconds, how many RDDs are included in the window at any point in time?



4

O 8

O 20

O 40

O 10

Which option best describes data in streaming?



Unbounded and Continuous

Unbounded and Discrete

Bounded and Continuous

Bounded and Discrete

Question 2 of 10



What is checkpointing in streaming Spark applications?

O Backing up all input data sources on multiple machines in the cluster

Periodically saving data so state can be reconstructed from this intermediate context



Backing up all RDDs on stable storage after every transformation

Question 8 of 10



What do you need to do in order to be able to run SQL queries on your dataframe?



Create a temporary view of your data that registers it as a SQL table

O Set the enable SQL flag to true to use SQL queries



Format your data in rows and columns and write it out to disk

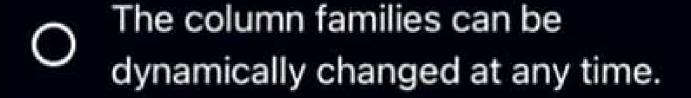
Which of these statements are true about HBase?



Only primitives can be row keys.



The column family is specified when you create the table.



- The column values uniquely identify a row.
- O All rows should have the same column names.

What do you call a logical division of data?



Block

O File

O Division



Partition

Question 5 of 10



If the output data type of the mapper is ImmutableBytesWritable, LongWritable the input arguments to the reduce() method is:



ImmutableBytesWritable, LongWritable

O ImmutableBytesWritable, Collection<LongWritable>



ImmutableBytesWritable,
Iterable<LongWritable>

O LongWritable, ImmutableBytesWritable



Which of these best defines an RDD?

O A table of rows and columns which allow array manipulations

 A collection of records which are held in memory across multiple machines in a cluster

A data type which allows operations that are idempotent

A map of key values pairs which are held in memory and allow fast lookup

Question 4 of 10

X

What are the main features of the Hadoop distributed computing framework?



Storage and processing of data across multiple machines with fault tolerance and recovery

- O Distribution of processes but not storage across multiple machines while handling fault tolerance and recovery
- O Using single powerful machine
 which integrates with other dumb
 terminals all of which send their
 processes to the single machine
- Multiple machines which are connected to each other but each job runs only on one machine at a

Aim -

When you sample data, what occurs if you specify withreplacement = True?



The probability of each item from being selected increases as other items are taken

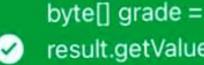


The probability of each item from being selected remains constant as other items are taken

The probability of each item from
 being selected decreases as other items are taken



What is the command to retrieve a single column value from a Result object representing a row in the **HBase table?**



result.getValue(Bytes.toBytes("acad emic"), Bytes.toBytes("grade"))

string grade = (String) result.getValue(Bytes.toBytes("acad emic"), Bytes.toBytes("grade"))

byte[] grade = result.getValue("academic", "grade")

byte[] grade = result.getValue(Bytes.toBytes("acad emic"))

Hadoop has some limitations such that it cannot be used as a database. Which of the following is NOT a limitation of Hadoop?

- Data in Hadoop has a well-defined structure and that structure is wellenforced.
- O Hadoop alone is more suited for long-running batch processes.
- O Hadoop alone does not allow random access to individual records.
- O Hadoop can store files of any kind, text, binary, media files, etc.



How does lazy evaluation function in a Spark RDD?

All transformations write out
 intermediate results to disk and so can be lazily combined later.

Transformations are only applied when a result is requested, they are evaluated on demand.

O All operations on an RDD are performed in memory not on disk.

What is the command to see what partitions exist in a Hive table?



display partitions <tablename>

O describe table <tablename>

O describe formatted <tablename>

Which of the following set operations are allowed in Hive?



Which of these is NOT a feature of a data warehouse?



Loss of data is acceptable and even expected, and lost data is ignored.

- O Data may be available a couple of days after the updates were made.
- They are optimized to serve reads, and certain write operations may not even be allowed.
- O Records may be stored for many years.

Question 6 of 10



What is the name of Cloudera's Platform-as-a-Service offering?

- O Cloudera Director
- O CDH
- Cloudera Manager
- Cloudera Altus

Question 3 of 10

X

What are the default data types for the input to the map phase when reading from an HBase table using a TableMapper class?



Byte, Result

- O Integer, String
- O ImmutableBytesWritable, Row



ImmutableBytesWritable, Result

Which of these is NOT a characteristic of a relation in Pig?



A relation is analogous to a table in traditional databases and has strict schema and data type constraints.

- A relation is immutable, contents of a relation cannot be edited in-place.
- A relation is ephemeral, it survives only for the duration of a Pig session.
- A relation is referenced by a name, O similar to a variable in a programming language.

How does the MapReduce programming model take advantage of multiple machines in a cluster?

The Map operation is performed on multiple machines on different

O subsets of data, and the Reduce operation is performed in parallel with the Map process



The Map operation is performed on a single machine, and the Reduce operation is parallelized across multiple machines



The Map operation is performed on multiple machines on different subsets of data, and the Reduce operation combines the output of various mappers

Which of the following is NOT a way to sample data from Hive tables?



- Random Sampling
- O The LIMIT Keyword
- O Block Sampling
- Row Count Sampling

While implementing the MapReduce classes using Hadoop libraries, the output key and value type of the mapper have to match:



The input key and value types of the mapper

O The output key and value types of the reducer



The input key and value types of the reducer

When would you choose to use broadcast variables in Spark?

O When you want to announce cluster updates to all nodes

When you want to efficiently copy over data to each node to reduce network data transfer

O When you want to perform join operations on 2 dataframes

Which of the following statements about Project Tungsten in Spark 2 is false?

- It optimizes Spark queries like a compiler would rather than like a DBMS would.
- It contains many virtual function
 calls to help model very complex operations.
- O It is faster than the volcano iterator model of Spark 1.
- O It heavily utilizes loop unrolling to make operations faster.

Which of these is NOT true of unions in Pig?



Pig preserves duplicates in the final union

- O Pig does not preserve the order of tuples in the final union
- O Pig requires that the schema of the relations be compatible



Pig will resolve schema incompatibilities by padding the tuples with nulls

Which of the following is true of the flatten function in Pig?

8

It applies to the tuple complex data type and extracts individual fields as separate records.

It applies to the map complex data
O type and extracts key-value pairs as separate records.



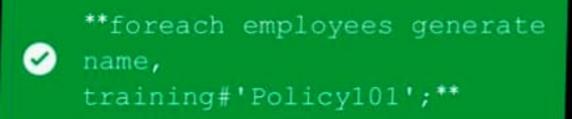
It applies to the bag complex data type and extracts each tuple as separate records.

```
employee_id: chararray
name: chararray
title: chararray
department: chararray
salary: double
address: tuple (apartment: chararray,
    street: chararray, zipcode: int)
training: map[chararray]
```

- **
- **

The training field is a map of all the training sessions the employee has attended. It maps the name of the training session to the date on which it was held.

Access the employee name and the dates on which they attended the 'Policy101' training and store it in a relation



Which of the following attributes is not true of MapReduce processing?



Has a Map step for preprocessing and a Reduce step for aggregating

O Is used by Hadoop, with Java in its native coding environment



Processes relational or dimensional data using SQL or MDX exclusively

Pioneered at Google

What function would you use to apply a transformation to every element of an RDD?

X rdd.mapReduce() rdd.withColumn() rdd.filter() rdd.transform() rdd.map()

What role does ETL, Extract, Transform and Load play when you work with data?

0

Allows you to pull data from original sources and pre-process them in useful ways before it is stored in a data warehouse for analysis

- Allows you to transform data when moving it from one database to another
- Has plugins to multiple databases to allow transforming data and loading into a variety of destination databases
- Allows you to work with multiple sources, even if they are remotely located to bring them all together in

11-- ----- £------

What is the Pig function that is used to descrialize and load data from the file system?

```
***PigDeserializer()**

**PigLoader()**

**PigStorage()**
```

O **PigLoadStore()**

O **PigStore() **

If no data type is specified for a certain field when loading a relation, what is the default data type that Pig assumes for this field?

- × **chars**
 - O **bytes**
 - O **chararray**
 - Ø bytearray
 - O **int**

Which is these is NOT a primitive data type in Hive?

- O Decimal
- String
- Month
 - O Integer



Which of these is NOT a standard scheduling algorithm in YARN?

Capacity Scheduler

First In First Out Scheduler

Fair Scheduler



Shortest Job First Scheduler

Which is the correct syntax for a join operation on Hive tables?

o select students.id, students.name,
courses.title from students, courses
where students.id =
courses.student_id

ourses.title from students.id <
courses where students.id <
courses.student_id

select students.id, students.name, courses.title from students join courses where students.id = courses.student_id

Which of the following vendors does not offer a Hadoop distribution of its own?



O IBM

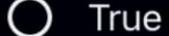
SAP/Sybase

MapR

When writing a Bash script for HDFS, you must specify 'HDFS' in script.



False



What is one major drawback of the Capacity Scheduler?



Underutilization of cluster resources

O Excessive overhead in queue management

O Starvation of jobs with short processing times

Cong wait times

HiveQL is best described as a:

8

Expressive language

Scripting language



Declarative language

O Procedural language

Why does the pseudo-distributed mode require password-less Secure Shell (SSH) login to be set up?

8

SSH is required to encrypt client interactions with the cluster so it is secure, and passwords ensure that this communication is smooth and not interrupted by password requests.

SSH is required for communication between master and slave nodes in the cluster and manual password logins do not work with machine to machine communication.

SSH is required for communication between master and slave nodes in the cluster, and since this communication is very frequent, you



reduce the overhead by configuring

Hadoop can be run on cloud computing infrastructure.



True



What is the command to submit a MapReduce job to a Hadoop cluster using the command line?

O hadoop submit

hadoop jar

O hadoop run

O hadoop fs

Bash scripts can automate what HDFS task?

8

Optimizing YARN applications

Checking HBase job status

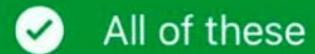
0

Ingesting, exporting, and moving data

Updating Masternode key pairs

Which of the following is a valid definition of big data?

O Into an MPP data warehouse appliance



O Loaded onto Hadoop or NoSQL clusters

O Analysis of petabyte-scale collections of data

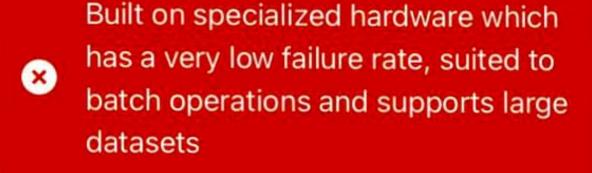
When removing a file using the rm command, files are deleted permanently.

True

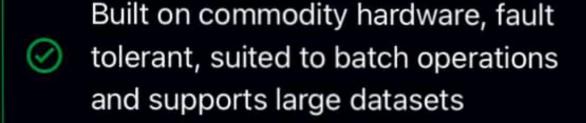


False

What are the main characteristics of the Hadoop Distributed File System?



Built on commodity hardware, fault tolerant, suited to low latency transactional operations and supports large datasets



tolerant, suited to low latency
transactional operations and
supports large structured datasets

..ith ...all dafinad fialda

Built on commodity hardware, fault

Why do reads from the HDFS have to access both the name node and the data nodes?



Reads to the name node check the file permissions and reads to the data nodes read where the file is located and the file contents.

- Reads to the name node update edit

 logs and reads to the data nodes
 retrieve file contents.
- Reads to the name node retrieve the contents as the name node in turn reads from the data nodes.

where a file and the corresponding blocks are located on the cluster, and reads to the data nodes read the actual contents of the file.

Reads to the name node look up

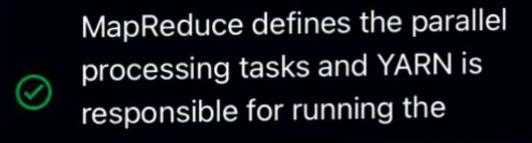
YARN in Hadoop versions after 2.0?

MapReduce helps the various

parallel processes communicate with each other and YARN is responsible for management of resources in the cluster to help run processes.

YARN defines the parallel processing tasks and MapReduce is responsible for running the processes on nodes in the cluster.

MapReduce contains the parallel processing logic and YARN partitions the data so the processes can run efficiently.



Which of the following is a valid justification for use of NoSQL technology for Big Data applications?

- 8
- The Hadoop Distributed File System is explicitly tailored to unstructured data and RESTful APIs
- O Many NoSQL databases are based on open source technology
- O NoSQL's use of scale-out architectures work well for big data
- 0

All NoSQL products use MapReduce in their query interfaces

How do you empty the trash in HDFS dfs?



empty

trash_out

O refresh