

100% Pass Databricks - Professional Reliable Databricks-Certified-Professional-Data-Engineer Exam Question



The FreePdfDump is one of the best platforms that has been helping the Databricks-Certified-Professional-Data-Engineer exam candidates for many years. Over this long time period the countless Databricks Certified Professional Data Engineer Exam Databricks-Certified-Professional-Data-Engineer exam candidates have passed their dream Databricks Databricks-Certified-Professional-Data-Engineer Certification Exam and they have become certified Databricks Databricks-Certified-Professional-Data-Engineer professionals. All the successful Databricks Databricks-Certified-Professional-Data-Engineer certification professionals are doing jobs in small, medium, and large size enterprises.

Databricks Certified Professional Data Engineer (Databricks-Certified-Professional-Data-Engineer) certification exam is designed for professionals who want to demonstrate their expertise in using Databricks to manage big data and create data pipelines. Databricks Certified Professional Data Engineer Exam certification exam is ideal for data engineers, data architects, data scientists, and other professionals who work with big data and want to validate their skills in using Databricks to build data pipelines.

>> **Reliable Databricks-Certified-Professional-Data-Engineer Exam Question** <<

New Databricks-Certified-Professional-Data-Engineer Practice Questions, Valid Databricks-Certified-Professional-Data-Engineer Exam Papers

Before you buy our product, you can download and try out it freely so you can have a good understanding of our Databricks-Certified-Professional-Data-Engineer test prep. In such a way, the client can visit the page of our Databricks-Certified-Professional-Data-Engineer exam questions on the website. So the client can understand our Databricks-Certified-Professional-Data-Engineer Exam Materials well and decide whether to buy our Databricks-Certified-Professional-Data-Engineer training guide or not since that they have checked the quality of our Databricks-Certified-Professional-Data-Engineer exam questions. We provide the best Databricks-Certified-Professional-Data-Engineer learning guide to our client and you will be satisfied.

Databricks Certified Professional Data Engineer Exam Sample Questions (Q19-Q24):

NEW QUESTION # 19

A DLT pipeline includes the following streaming tables:

Raw_1ot ingest raw device measurement data from a heart rate tracking device.

Bgm_stats incrementally computes user statistics based on BPM measurements from raw_1ot.

How can the data engineer configure this pipeline to be able to retain manually deleted or updated records in the raw_1ot table while recomputing the downstream table when a pipeline update is run?

- A. Set the SkipChangeCommits flag to true raw_1ot
- **B. Set the pipelines, reset, allowed property to false on raw_1ot**
- C. Set the pipelines, reset, allowed property to false on bpm_stats
- D. Set the skipChangeCommits flag to true on bpm_stats

Answer: B

Explanation:

In Databricks Lakehouse, to retain manually deleted or updated records in the raw_iot table while recomputing downstream tables when a pipeline update is run, the property `pipelines.reset.allowed` should be set to `false`. This property prevents the system from resetting the state of the table, which includes the removal of the history of changes, during a pipeline update. By keeping this property as `false`, any changes to the raw_iot table, including manual deletes or updates, are retained, and recomputation of downstream tables, such as `bpm_stats`, can occur with the full history of data changes intact.

:

Databricks documentation on DLT pipelines: <https://docs.databricks.com/data-engineering/delta-live-tables/delta-live-tables-overview.html>

NEW QUESTION # 20

A data engineer is designing a data pipeline in Databricks that needs to process records from a Kafka stream where late-arriving data is common. Which approach should the data engineer use?

- A. Use batch processing and overwrite the entire output table each time to ensure late data is incorporated correctly.
- B. Implement a custom solution using Databricks Jobs to periodically reprocess all historical data.
- C. Use an Auto CDC pipeline with batch tables to simplify late data handling.
- **D. Use a watermark to specify the allowed lateness to accommodate records that arrive after their expected window, ensuring correct aggregation and state management.**

Answer: D

Explanation:

Databricks and Apache Spark document watermarks as the standard mechanism for handling late-arriving data in Structured Streaming. A watermark defines how long the engine should continue waiting for out-of-order event-time data and helps manage state for aggregations, joins, and deduplication. (Databricks Documentation) This directly matches Kafka streaming scenarios where lateness is common. The other options are not the standard streaming solution for event-time late data: Auto CDC is for change data capture use cases, and repeated full historical reprocessing is inefficient and unnecessary when watermarking is the built-in feature designed for this problem. (Databricks Documentation)

NEW QUESTION # 21

Which configuration parameter directly affects the size of a spark-partition upon ingestion of data into Spark?

- **A. `spark.sql.files.maxPartitionBytes`**
- B. `spark.sql.adaptive.coalescePartitions.minPartitionNum`
- C. `spark.sql.autoBroadcastJoinThreshold`
- D. `spark.sql.adaptive.advisoryPartitionSizeInBytes`
- E. `spark.sql.files.openCostInBytes`

Answer: A

Explanation:

This is the correct answer because `spark.sql.files.maxPartitionBytes` is a configuration parameter that directly affects the size of a spark-partition upon ingestion of data into Spark. This parameter configures the maximum number of bytes to pack into a single partition when reading files from file-based sources such as Parquet, JSON and ORC. The default value is 128 MB, which means each partition will be roughly 128 MB in size, unless there are too many small files or only one large file. Verified Reference: [Databricks Certified Data Engineer Professional], under "Spark Configuration" section; Databricks Documentation, under "Available Properties - `spark.sql.files.maxPartitionBytes`" section.

NEW QUESTION # 22

A Delta Lake table was created with the below query:

Realizing that the original query had a typographical error, the below code was executed:

```
ALTER TABLE prod.sales_by_stor RENAME TO prod.sales_by_store
```

Which result will occur after running the second command?

- A. The table reference in the metastore is updated and no data is changed.
- B. A new Delta transaction log is created for the renamed table.
- C. The table name change is recorded in the Delta transaction log.
- D. All related files and metadata are dropped and recreated in a single ACID transaction.
- E. The table reference in the metastore is updated and all data files are moved.

Answer: A

Explanation:

Explanation

The query uses the CREATE TABLE USING DELTA syntax to create a Delta Lake table from an existing Parquet file stored in DBFS. The query also uses the LOCATION keyword to specify the path to the Parquet file as /mnt/finance_eda_bucket/tx_sales.parquet. By using the LOCATION keyword, the query creates an external table, which is a table that is stored outside of the default warehouse directory and whose metadata is not managed by Databricks. An external table can be created from an existing directory in a cloud storage system, such as DBFS or S3, that contains data files in a supported format, such as Parquet or CSV.

The result that will occur after running the second command is that the table reference in the metastore is updated and no data is changed. The metastore is a service that stores metadata about tables, such as their schema, location, properties, and partitions. The metastore allows users to access tables using SQL commands or Spark APIs without knowing their physical location or format.

When renaming an external table using the ALTER TABLE RENAME TO command, only the table reference in the metastore is updated with the new name; no data files or directories are moved or changed in the storage system. The table will still point to the same location and use the same format as before. However, if renaming a managed table, which is a table whose metadata and data are both managed by Databricks, both the table reference in the metastore and the data files in the default warehouse directory are moved and renamed accordingly. Verified References:

[Databricks Certified Data Engineer Professional], under "Delta Lake" section; Databricks Documentation, under "ALTER TABLE RENAME TO" section; Databricks Documentation, under "Metastore" section; Databricks Documentation, under "Managed and external tables" section.

NEW QUESTION # 23

When evaluating the Ganglia Metrics for a given cluster with 3 executor nodes, which indicator would signal proper utilization of the VM's resources?

- A. CPU Utilization is around 75%
- B. The five Minute Load Average remains consistent/flat
- C. Total Disk Space remains constant
- D. Network I/O never spikes
- E. Bytes Received never exceeds 80 million bytes per second

Answer: A

Explanation:

In the context of cluster performance and resource utilization, a CPU utilization rate of around 75% is generally considered a good indicator of efficient resource usage. This level of CPU utilization suggests that the cluster is being effectively used without being overburdened or underutilized.

* A consistent 75% CPU utilization indicates that the cluster's processing power is being effectively employed while leaving some headroom to handle spikes in workload or additional tasks without maxing out the CPU, which could lead to performance degradation.

* A five Minute Load Average that remains consistent/flat (Option A) might indicate underutilization or a bottleneck elsewhere.

* Monitoring network I/O (Options B and C) is important, but these metrics alone don't provide a complete picture of resource utilization efficiency.

* Total Disk Space (Option D) remaining constant is not necessarily an indicator of proper resource utilization, as it's more related to storage rather than computational efficiency.

References:

* Ganglia Monitoring System: Ganglia Documentation

* Databricks Documentation on Monitoring: Databricks Cluster Monitoring

NEW QUESTION # 24

.....

