

# Databricks - Databricks-Certified-Professional-Data-Engineer–Newest Real Exam



BTW, DOWNLOAD part of Prep4King Databricks-Certified-Professional-Data-Engineer dumps from Cloud Storage:  
[https://drive.google.com/open?id=11qEFOFLRiFDWv659KbG\\_bt0zL0bJSK](https://drive.google.com/open?id=11qEFOFLRiFDWv659KbG_bt0zL0bJSK)

Once you have any questions about our Databricks-Certified-Professional-Data-Engineer actual exam, you can contact our staff online or send us an email. We have a dedicated all-day online service to help you solve problems. Before purchasing, you may be confused about what kind of Databricks-Certified-Professional-Data-Engineer Guide questions you need. You can consult our staff online. After the consultation, your doubts will be solved and you will choose the Databricks-Certified-Professional-Data-Engineer learning materials that suit you.

Databricks Certified Professional Data Engineer (Databricks-Certified-Professional-Data-Engineer) Exam is a certification program designed for individuals who want to demonstrate their expertise in building, deploying, and maintaining data engineering solutions using Databricks. Databricks-Certified-Professional-Data-Engineer Exam is intended for data engineers, data architects, and other data professionals who work with large-scale data processing systems and want to validate their skills and knowledge in this area.

Databricks Certified Professional Data Engineer certification is a valuable credential for data engineers who work with Databricks. It demonstrates that the candidate has a deep understanding of Databricks and can use it effectively to solve complex data engineering problems. Databricks Certified Professional Data Engineer Exam certification can help data engineers advance their careers, increase their earning potential, and gain recognition as experts in the field of big data and machine learning.

Databricks Certified Professional Data Engineer certification is a valuable credential for data engineers who work with the Databricks platform. It validates their skills and expertise and demonstrates to employers that they have the knowledge and experience needed to work with Databricks effectively. By passing the exam and earning the certification, data engineers can enhance their career prospects and gain a competitive advantage in the job market.

>> **Databricks-Certified-Professional-Data-Engineer Real Exam** <<

## **Databricks-Certified-Professional-Data-Engineer Valid Exam Vce Free & Databricks-Certified-Professional-Data-Engineer Valid Vce Dumps**

We would like to provide our customers with different kinds of Databricks-Certified-Professional-Data-Engineer practice torrent to learn, and help them accumulate knowledge and enhance their ability. Besides, we guarantee that the questions of all our users can be answered by professional personal in the shortest time with our Databricks-Certified-Professional-Data-Engineer study guide. One more to mention, we can help you make full use of your sporadic time to absorb knowledge and information. In a word, compared to other similar companies aiming at Databricks-Certified-Professional-Data-Engineer Test Prep, the services and quality of our Databricks-Certified-Professional-Data-Engineer exam questions are highly regarded by our customers and potential clients.

## **Databricks Certified Professional Data Engineer Exam Sample Questions (Q70-Q75):**

**NEW QUESTION # 70**

A query is taking too long to run. After investigating the Spark UI, the data engineer discovered a significant amount of disk spill. The compute instance being used has a core-to-memory ratio of 1:2.

What are the two steps the data engineer should take to minimize spillage? (Choose 2 answers)

- A. Reduce `spark.sql.files.maxPartitionBytes`.
- B. Choose a compute instance with a higher core-to-memory ratio.
- C. Choose a compute instance with more disk space.
- D. Choose a compute instance with more network bandwidth.
- E. Increase `spark.sql.files.maxPartitionBytes`.

**Answer: A,B**

Explanation:

Comprehensive and Detailed Explanation From Exact Extract of Databricks Data Engineer Documents:

Databricks recommends addressing disk spilling—which occurs when Spark tasks run out of memory—by increasing memory per core and controlling partition size. Selecting an instance type with a higher memory-to-core ratio (A) provides each task with more available RAM, directly reducing the chance of spilling to disk. Additionally, reducing `spark.sql.files.maxPartitionBytes` (D) creates smaller partitions, preventing any single task from holding too much data in memory. Increasing partition size (C) or disk capacity (B) does not solve memory bottlenecks, and bandwidth (E) affects network I/O, not spill behavior. Therefore, the correct actions are A and D.

### NEW QUESTION # 71

A Structured Streaming job deployed to production has been experiencing delays during peak hours of the day. At present, during normal execution, each microbatch of data is processed in less than 3 seconds. During peak hours of the day, execution time for each microbatch becomes very inconsistent, sometimes exceeding

30 seconds. The streaming write is currently configured with a trigger interval of 10 seconds.

Holding all other variables constant and assuming records need to be processed in less than 10 seconds, which adjustment will meet the requirement?

- A. Use the trigger once option and configure a Databricks job to execute the query every 10 seconds; this ensures all backlogged records are processed with each batch.
- B. Increase the trigger interval to 30 seconds; setting the trigger interval near the maximum execution time observed for each batch is always best practice to ensure no records are dropped.
- C. Decrease the trigger interval to 5 seconds; triggering batches more frequently may prevent records from backing up and large batches from causing spill.
- D. The trigger interval cannot be modified without modifying the checkpoint directory; to maintain the current stream state, increase the number of shuffle partitions to maximize parallelism.
- E. Decrease the trigger interval to 5 seconds; triggering batches more frequently allows idle executors to begin processing the next batch while longer running tasks from previous batches finish.

**Answer: A**

Explanation:

The scenario presented involves inconsistent microbatch processing times in a Structured Streaming job during peak hours, with the need to ensure that records are processed within 10 seconds. The trigger once option is the most suitable adjustment to address these challenges:

\* Understanding Triggering Options:

\* Fixed Interval Triggering (Current Setup): The current trigger interval of 10 seconds may contribute to the inconsistency during peak times as it doesn't adapt based on the processing time of the microbatches. If a batch takes longer to process, subsequent batches will start piling up, exacerbating the delays.

\* Trigger Once: This option allows the job to run a single microbatch for processing all available data and then stop. It is useful in scenarios where batch sizes are unpredictable and can vary significantly, which seems to be the case during peak hours in this scenario.

\* Implementation of Trigger Once:

\* Setup: Instead of continuously running, the job can be scheduled to run every 10 seconds using a Databricks job. This scheduling effectively acts as a custom trigger interval, ensuring that each execution cycle handles all available data up to that point without overlapping or queuing up additional executions.

\* Advantages: This approach allows for each batch to complete processing all available data before the next batch starts, ensuring consistency in handling data surges and preventing the system from being overwhelmed.

\* Rationale Against Other Options:

\* Option A and E (Decrease Interval): Decreasing the trigger interval to 5 seconds might exacerbate the problem by increasing the frequency of batch starts without ensuring the completion of previous batches, potentially leading to higher overhead and less efficient processing.

\* Option B (Increase Interval): Increasing the trigger interval to 30 seconds could lead to latency issues, as the data would be processed less frequently, which contradicts the requirement of processing records in less than 10 seconds.

\* Option C (Modify Partitions): While increasing parallelism through more shuffle partitions can improve performance, it does not address the fundamental issue of batch scheduling and could still lead to inconsistency during peak loads.

\* Conclusion:

\* By using the trigger once option and scheduling the job every 10 seconds, you ensure that each microbatch has sufficient time to process all available data thoroughly before the next cycle begins, aligning with the need to handle peak loads more predictably and efficiently.

References

\* Structured Streaming Programming Guide - Triggering

\* Databricks Jobs Scheduling

### NEW QUESTION # 72

A Delta Lake table with Change Data Feed (CDF) enabled in the Lakehouse named `customer_churn_params` is used in churn prediction by the machine learning team. The table contains information about customers derived from a number of upstream sources. Currently, the data engineering team populates this table nightly by overwriting the table with the current valid values derived from upstream data sources. The churn prediction model used by the ML team is fairly stable in production. The team is only interested in making predictions on records that have changed in the past 24 hours. Which approach would simplify the identification of these changed records?

- A. Apply the churn model to all rows in the `customer_churn_params` table, but implement logic to perform an upsert into the predictions table that ignores rows where predictions have not changed.
- B. Modify the overwrite logic to include a field populated by calling `current_timestamp()` as data are being written; use this field to identify records written on a particular date.
- C. Replace the current overwrite logic with a MERGE statement to modify only those records that have changed; write logic to make predictions on the changed records identified by the Change Data Feed.
- D. Convert the batch job to a Structured Streaming job using the complete output mode; configure a Structured Streaming job to read from the `customer_churn_params` table and incrementally predict against the churn model.

**Answer: C**

Explanation:

Comprehensive and Detailed Explanation From Exact Extract:

\* Exact extract: "Change data feed (CDF) provides row-level change information for Delta tables."

\* Exact extract: "Use `table_changes` to query the set of rows that were inserted, updated, or deleted between two versions (or timestamps)." References: Delta Lake Change Data Feed; Delta Lake MERGE INTO.

### NEW QUESTION # 73

A user new to Databricks is trying to troubleshoot long execution times for some pipeline logic they are working on. Presently, the user is executing code cell-by-cell, using `display()` calls to confirm code is producing the logically correct results as new transformations are added to an operation. To get a measure of average time to execute, the user is running each cell multiple times interactively.

Which of the following adjustments will get a more accurate measure of how code is likely to perform in production?

- A. The only way to meaningfully troubleshoot code execution times in development notebooks is to use production-sized data and production-sized clusters with Run All execution.
- B. Scala is the only language that can be accurately tested using interactive notebooks; because the best performance is achieved by using Scala code compiled to JARs, all PySpark and Spark SQL logic should be refactored.
- C. The Jobs UI should be leveraged to occasionally run the notebook as a job and track execution time during incremental code development because Photon can only be enabled on clusters launched for scheduled jobs.
- D. Production code development should only be done using an IDE; executing code against a local build of open source Spark and Delta Lake will provide the most accurate benchmarks for how code will perform in production.
- E. Calling `display()` forces a job to trigger, while many transformations will only add to the logical query plan; because of caching, repeated execution of the same logic does not provide meaningful results.

**Answer: E**

Explanation:

This is the correct answer because it explains which of the following adjustments will get a more accurate measure of how code is likely to perform in production. The adjustment is that calling `display()` forces a job to trigger, while many transformations will only add to the logical query plan; because of caching, repeated execution of the same logic does not provide meaningful results. When developing code in Databricks notebooks, one should be aware of how Spark handles transformations and actions. Transformations are operations that create a new `DataFrame` or `Dataset` from an existing one, such as `filter`, `select`, or `join`. Actions are operations that trigger a computation on a `DataFrame` or `Dataset` and return a result to the driver program or write it to storage, such as `count`, `show`, or `save`. Calling `display()` on a `DataFrame` or `Dataset` is also an action that triggers a computation and displays the result in a notebook cell. Spark uses lazy evaluation for transformations, which means that they are not executed until an action is called. Spark also uses caching to store intermediate results in memory or disk for faster access in subsequent actions. Therefore, calling `display()` forces a job to trigger, while many transformations will only add to the logical query plan; because of caching, repeated execution of the same logic does not provide meaningful results. To get a more accurate measure of how code is likely to perform in production, one should avoid calling `display()` too often or clear the cache before running each cell. Verified References: [Databricks Certified Data Engineer Professional], under "Spark Core" section; Databricks Documentation, under "Lazy evaluation" section; Databricks Documentation, under "Caching" section.

#### NEW QUESTION # 74

The Delta Live Tables Pipeline is configured to run in Development mode using the Triggered Pipeline Mode. what is the expected outcome after clicking Start to update the pipeline?

- A. All datasets will be updated continuously and the pipeline will not shut down. The compute resources will persist with the pipeline
- B. All datasets will be updated at set intervals until the pipeline is shut down. The compute resources will persist after the pipeline is stopped to allow for additional development and testing
- C. All datasets will be updated at set intervals until the pipeline is shut down. The compute resources will be deployed for the update and terminated when the pipeline is stopped
- D. All datasets will be updated once and the pipeline will shut down. The compute resources will be terminated
- E. All datasets will be updated once and the pipeline will shut down. The compute resources will persist to allow for additional development and testing

**Answer: A**

Explanation:

Explanation

The answer is All datasets will be updated once and the pipeline will shut down. The compute re-sources will persist to allow for additional testing.

DLT pipeline supports two modes Development and Production, you can switch between the two based on the stage of your development and deployment lifecycle.

Development and production modes

When you run your pipeline in development mode, the Delta Live Tables system:

\*Reuses a cluster to avoid the overhead of restarts.

\*Disables pipeline retries so you can immediately detect and fix errors.

In production mode, the Delta Live Tables system:

\*Restarts the cluster for specific recoverable errors, including memory leaks and stale credentials.

\*Retries execution in the event of specific errors, for example, a failure to start a cluster.

Use the buttons in the Pipelines UI to switch between development and production modes. By default, pipelines run in development mode.

Switching between development and production modes only controls cluster and pipeline execution behavior.

Storage locations must be configured as part of pipeline settings and are not affected when switching between modes.

Please review additional DLT concepts using below link

<https://docs.databricks.com/data-engineering/delta-live-tables/delta-live-tables-concepts.htm#delta-live-tables-c>

#### NEW QUESTION # 75

.....

In this cut-throat competitive world of Databricks, the Databricks Databricks-Certified-Professional-Data-Engineer certification is the most desired one. But what creates an obstacle in the way of the aspirants of the Databricks Databricks-Certified-Professional-

