

# Dumps Databricks-Certified-Professional-Data-Engineer Free Download | Test Databricks-Certified-Professional-Data-Engineer Questions Answers



Customizable Databricks Databricks-Certified-Professional-Data-Engineer practice exams (desktop and web-based) of ActualtestPDF are designed to give you the best learning experience. You can attempt these Databricks-Certified-Professional-Data-Engineer practice tests multiple times till the best preparation for the Databricks-Certified-Professional-Data-Engineer test. On every take, our Databricks-Certified-Professional-Data-Engineer Practice Tests save your progress so you can view it to see and strengthen your weak concepts easily. Customizable Databricks-Certified-Professional-Data-Engineer practice exams allow you to adjust the time and Databricks-Certified-Professional-Data-Engineer questions numbers according to your practice needs.

To prepare for the DCPDE exam, candidates should have a solid understanding of data engineering concepts, such as data modeling, data integration, data transformation, and data quality. They should also have experience working with big data technologies, such as Apache Spark, Apache Kafka, and Apache Hadoop.

Databricks Certified Professional Data Engineer exam is a hands-on exam that requires the candidate to complete a set of tasks using Databricks. Databricks-Certified-Professional-Data-Engineer Exam evaluates the candidate's ability to design and implement data pipelines, work with data sources and sinks, and perform transformations using Databricks. Databricks-Certified-Professional-Data-Engineer exam also tests the candidate's ability to optimize and tune data pipelines for performance and reliability.

**>> Dumps Databricks-Certified-Professional-Data-Engineer Free Download <<**

## 2026 Newest Dumps Databricks-Certified-Professional-Data-Engineer Free Download | 100% Free Test Databricks-Certified-Professional-Data-Engineer Questions Answers

The customers can prepare from the actual Databricks-Certified-Professional-Data-Engineer and can clear Databricks Certified Professional Data Engineer Exam exam with ease and if they failed to do it despite all of their efforts they can get a full refund of their money according to terms and conditions. The Databricks-Certified-Professional-Data-Engineer exam solutions is packed with a lot of premium features and it is getting updated on the daily basis according to the syllabus. Databricks Databricks-Certified-Professional-Data-Engineer updates real questions so the students can easily prepare for it and clear Databricks Databricks-

Certified-Professional-Data-Engineer exam.

Databricks Certified Professional Data Engineer (Databricks-Certified-Professional-Data-Engineer) exam is a certification program designed to validate the skills and expertise of data engineers in developing and managing big data pipelines using Databricks. Databricks-Certified-Professional-Data-Engineer Exam is ideal for data engineers, ETL developers, and data architects who work with Databricks and want to showcase their skills and proficiency.

## Databricks Certified Professional Data Engineer Exam Sample Questions (Q17-Q22):

### NEW QUESTION # 17

Newly joined data analyst requested read-only access to tables, assuming you are owner/admin which section of Databricks platform is going to facilitate granting select access to the user

- A. Azure Databricks control pane IAM
- B. Azure RBAC
- C. Data explorer
- D. Admin console
- E. User settings

**Answer: C**

Explanation:

Explanation

Answer is Data Explorer

<https://docs.databricks.com/sql/user/data/index.html>

Data explorer lets you easily explore and manage permissions on databases and tables. Users can view schema details, preview sample data, and see table details and properties. Administrators can view and change owners, and admins and data object owners can grant and revoke permissions.

To open data explorer, click Data in the sidebar.

### NEW QUESTION # 18

A security analytics pipeline must enrich billions of raw connection logs with geolocation data. The join hinges on finding which IPv4 range each event's address falls into.

Table 1: network\_events ( $\approx$  5 billion rows)

```
event_id ip_int
42 3232235777
```

Table 2: ip\_ranges ( $\approx$  2 million rows)

```
start_ip_int end_ip_int country
3232235520 3232236031 US
```

The query is currently very slow:

```
SELECT n.event_id, n.ip_int, r.country
FROM network_events n
```

```
JOIN ip_ranges r
```

```
ON n.ip_int BETWEEN r.start_ip_int AND r.end_ip_int;
```

Which change will most dramatically accelerate the query while preserving its logic?

- A. Increase `spark.sql.shuffle.partitions` from 200 to 10000.
- B. Add a range-join hint `/*+ RANGE JOIN(r, 65536) */`.
- C. Force a sort-merge join with `/*+ MERGE(r) */`.
- D. Add a broadcast hint: `/*+ BROADCAST(r) */` for `ip_ranges`.

**Answer: B**

Explanation:

Comprehensive and Detailed Explanation from Databricks Documentation:

The query joins billions of rows (`network_events`) with millions of rows (`ip_ranges`) using a range predicate (`BETWEEN`). Unlike equality joins (`=`), range joins are not efficiently handled by broadcast or sort-merge joins because:

Broadcast Join (D): Effective for small tables but only for equality joins. Since this query uses a range condition, broadcast will not reduce the complexity of scanning billions of records across non-equality conditions.

Sort-Merge Join (C): Works for ordered joins but is inefficient on range conditions. Sorting billions of records adds excessive overhead and will not resolve the bottleneck.

Increasing Shuffle Partitions (A): Only spreads out shuffle work but does not address the fundamental inefficiency of range-based lookups at scale.

Range Joins in Spark (RANGE\_JOIN hint):

Databricks provides range join optimizations specifically for conditions such as BETWEEN. By applying a RANGE\_JOIN hint, Spark can build optimized data structures (such as interval indexes or partition pruning strategies) that map billions of input rows to ranges much faster. This avoids brute-force scans and unnecessary shuffle costs.

Thus, Option B is the correct solution because:

It leverages range-join optimization, which is purpose-built for queries joining massive event logs to smaller lookup tables with IP ranges.

This ensures Spark can evaluate billions of rows against millions of ranges with optimized matching logic, drastically improving query performance while preserving correctness.

## NEW QUESTION # 19

Which of the following operations are not supported on a streaming dataset view?

`spark.readStream.format("delta").table("sales").createOrReplaceTempView("streaming_view")`

- A. `SELECT sum(unitssold) FROM streaming_view`
- B. `SELECT id, sum(unitssold) FROM streaming_view GROUP BY id ORDER BY id`
- C. `SELECT max(unitssold) FROM streaming_view`
- D. `SELECT id, count(*) FROM streaming_view GROUP BY id`
- E. `SELECT * FROM streaming_view ORDER BY id`

**Answer: E**

Explanation:

Explanation

The answer is `SELECT * FROM streaming_view order by id` Please Note: Sorting with Group by will work without any issues see below explanation for each option of the options, Graphical user interface, text, application Description automatically generated

```
1 %sql
2 SELECT *
3 FROM streaming_view
4 ORDER BY id

Error in SQL statement: AnalysisException: Sorting is not supported on streaming DataFrames/Datasets, unless it is on aggregated DataFrame/Dataset in Complete output mode; line 3 pos 0;
Sort [id#4123L ASC NULLS FIRST], true
+- Project [timestamp#38, value#39L, unitssold#4118L, id#4123L]
   +- SubqueryAlias streaming_view
```

Certain operations are not allowed on streaming data, please see highlighted in bold.

<https://spark.apache.org/docs/latest/structured-streaming-programming-guide.html#unsupported-operations>

- \* Multiple streaming aggregations (i.e. a chain of aggregations on a streaming DF) are not yet supported on streaming Datasets.
- \* Limit and take the first N rows are not supported on streaming Datasets.
- \* Distinct operations on streaming Datasets are not supported.
- \* Deduplication operation is not supported after aggregation on a streaming Datasets.
- \* Sorting operations are supported on streaming Datasets only after an aggregation and in Complete Output Mode.

Note: Sorting without aggregation function is not supported.

Here is the sample code to prove this,

Setup test stream

Graphical user interface, text, application, email Description automatically generated

Cmd 1

```
1 from pyspark.sql import functions as F
```

Command took 0.03 seconds -- by akhil.vangala@avanade.com at

Cmd 2

```
1 spark.conf.set("spark.sql.shuffle.partitions",4
```

Command took 0.04 seconds -- by akhil.vangala@avanade.com at

Cmd 3

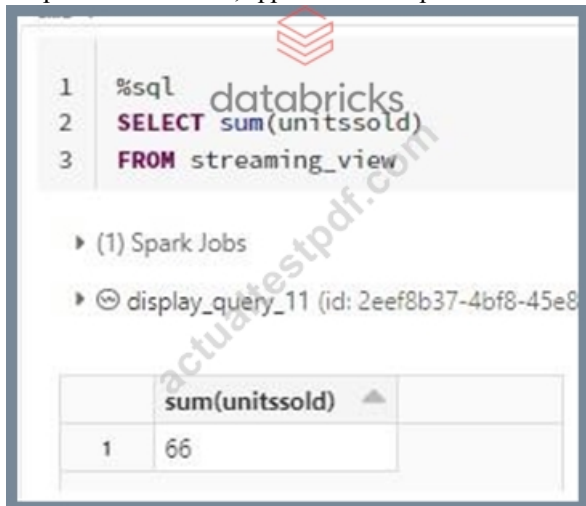
```
1 df = (spark.readStream.format("rate")
2     .option("rowPerSecond", 100)
3     .load())
4 df = df.withColumn("unitssold",F.col("value"))
5         .withColumn("id",F.col("value"))
6 df.createOrReplaceTempView("streaming_view")
```

df: pyspark.sql.dataframe.DataFrame = [timestamp: timestamp]

Command took 0.07 seconds -- by akhil.vangala@avanade.com at

Sum aggregation function has no issues on stream

Graphical user interface, application Description automatically generated



The screenshot shows the Databricks SQL interface. At the top, there is a code editor with the following SQL query:

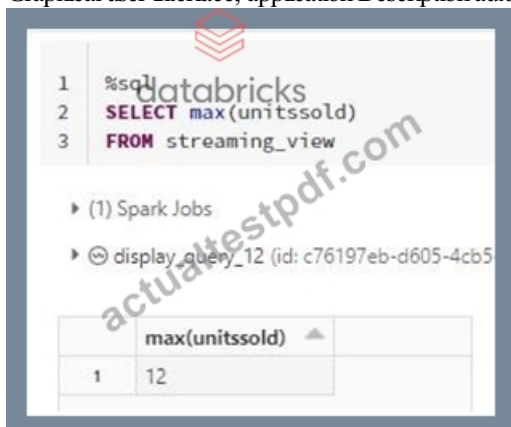
```
1 %sql
2 SELECT sum(unitssold)
3 FROM streaming_view
```

Below the code editor, it indicates "(1) Spark Jobs" and "display\_query\_11 (id: 2eef8b37-4bf8-45e8)". The results are displayed in a table:

	sum(unitssold)
1	66

Max aggregation function has no issues on stream

Graphical user interface, application Description automatically generated



The screenshot shows the Databricks SQL interface. At the top, there is a code editor with the following SQL query:

```
1 %sql
2 SELECT max(unitssold)
3 FROM streaming_view
```

Below the code editor, it indicates "(1) Spark Jobs" and "display\_query\_12 (id: c76197eb-d605-4cb5)". The results are displayed in a table:

	max(unitssold)
1	12

Group by with Order by has no issues on stream



```

1 %sql
2 SELECT id, count(*)
3 FROM streaming_view
4 GROUP BY id
5 ORDER BY id

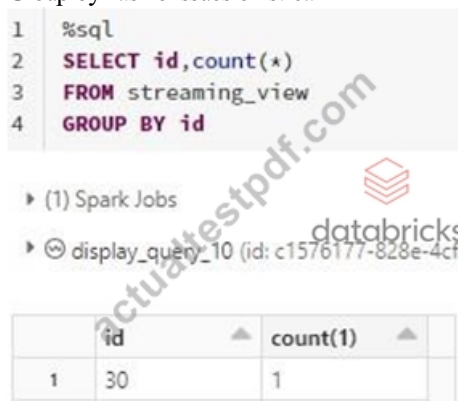
```

▶ (1) Spark Jobs

▶ display\_query\_13 (id: 0e8112ed-bd8e-4...

	id	count(1)
1	0	1

Group by has no issues on stream



```

1 %sql
2 SELECT id, count(*)
3 FROM streaming_view
4 GROUP BY id

```

▶ (1) Spark Jobs

▶ display\_query\_10 (id: c1576177-828e-4cf...

	id	count(1)
1	30	1

Table Description automatically generated

Order by without group by fails.

Graphical user interface, text, application Description automatically generated



```

1 %sql
2 SELECT *
3 FROM streaming_view
4 ORDER BY id

```

Error in SQL statement: AnalysisException: Sorting is not supported on streaming DataFrames/Datasets, unless it is on aggregated DataFrame/Dataset in Complete output mode; line 3 pos 0;

Sort [id#4123L ASC NULLS FIRST], true

-- Project [timestamp#38, value#39L, unitssold#4118L, id#4123L]

-- SubqueryAlias streaming\_view

## NEW QUESTION # 20

A junior data engineer has been asked to develop a streaming data pipeline with a grouped aggregation using DataFrame df. The pipeline needs to calculate the average humidity and average temperature for each non-overlapping five-minute interval. Incremental state information should be maintained for 10 minutes for late-arriving data.

Streaming DataFrame df has the following schema:

"device\_id INT, event\_time TIMESTAMP, temp FLOAT, humidity FLOAT"

Code block:

Choose the response that correctly fills in the blank within the code block to complete this task.

- A. await("event\_time + '10 minutes'")
- B. slidingWindow("event\_time", "10 minutes")
- C. withWatermark("event\_time", "10 minutes")
- D. awaitArrival("event\_time", "10 minutes")
- E. delayWrite("event\_time", "10 minutes")

**Answer: C**

Explanation:

The correct answer is A. withWatermark("event\_time", "10 minutes"). This is because the question asks for incremental state information to be maintained for 10 minutes for late-arriving data. The withWatermark method is used to define the watermark for late data. The watermark is a timestamp column and a threshold that tells the system how long to wait for late data. In this case, the watermark is set to 10 minutes. The other options are incorrect because they are not valid methods or syntax for watermarking in

Structured Streaming. References:

\* Watermarking: <https://docs.databricks.com/spark/latest/structured-streaming/watermarks.html>

\* Windowed aggregations:

<https://docs.databricks.com/spark/latest/structured-streaming/window-operations.html>

### NEW QUESTION # 21

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. Decrease the trigger interval to 5 seconds; triggering batches more frequently may prevent records from backing up and large batches from causing spill.
- B. 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.
- C. 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.
- 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. 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.

**Answer: C**

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



• • • • •

<https://www.actualtestpdf.com/Databricks/Databricks-Certified-Professional-Data-Engineer-practice-exam-dumps.html>

- [illegible]