# Exam Questions DP-203

Data Engineering on Microsoft Azure

## https://www.2passeasy.com/dumps/DP-203/

**NEW QUESTION 1**
- (Exam Topic 1)
You need to implement the surrogate key for the retail store table. The solution must meet the sales transaction dataset requirements.
What should you create?

A. a table that has an IDENTITY property
B. a system-versioned temporal table
C. a user-defined SEQUENCE object
D. a table that has a FOREIGN KEY constraint

**Answer:** A

**Explanation:**
Scenario: Implement a surrogate key to account for changes to the retail store addresses.
A surrogate key on a table is a column with a unique identifier for each row. The key is not generated from the table data. Data modelers like to create surrogate keys on their tables when they design data warehouse models. You can use the IDENTITY property to achieve this goal simply and effectively without affecting load performance.
Reference:
https://docs.microsoft.com/en-us/azure/synapse-analytics/sql-data-warehouse/sql-data-warehouse-tables-identity

**NEW QUESTION 2**
- (Exam Topic 2)
Which Azure Data Factory components should you recommend using together to import the daily inventory data from the SQL server to Azure Data Lake Storage?
To answer, select the appropriate options in the answer area.
NOTE: Each correct selection is worth one point.

Integration runtime type:
- Azure integration runtime
- Azure-SSIS integration runtime
- Self-hosted integration runtime

Trigger type:
- Event-based trigger
- Schedule trigger
- Tumbling window trigger

Activity type:
- Copy activity
- Lookup activity
- Stored procedure activity

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**
Box 1: Self-hosted integration runtime
A self-hosted IR is capable of running copy activity between a cloud data stores and a data store in private network.
Box 2: Schedule trigger
Schedule every 8 hours Box 3: Copy activity Scenario:

≫ Customer data, including name, contact information, and loyalty number, comes from Salesforce and can be imported into Azure once every eight hours. Row modified dates are not trusted in the source table.

≫ Product data, including product ID, name, and category, comes from Salesforce and can be imported into Azure once every eight hours. Row modified dates are not trusted in the source table.

**NEW QUESTION 3**
- (Exam Topic 3)
You are monitoring an Azure Stream Analytics job.
The Backlogged Input Events count has been 20 for the last hour. You need to reduce the Backlogged Input Events count.
What should you do?

A. Drop late arriving events from the job.
B. Add an Azure Storage account to the job.
C. Increase the streaming units for the job.
D. Stop the job.

**Answer:** C

**Explanation:**

General symptoms of the job hitting system resource limits include:

≫ If the backlog event metric keeps increasing, it's an indicator that the system resource is constrained (either because of output sink throttling, or high CPU).
Note: Backlogged Input Events: Number of input events that are backlogged. A non-zero value for this metric implies that your job isn't able to keep up with the number of incoming events. If this value is slowly increasing or consistently non-zero, you should scale out your job: adjust Streaming Units.
Reference:
https://docs.microsoft.com/en-us/azure/stream-analytics/stream-analytics-scale-jobs https://docs.microsoft.com/en-us/azure/stream-analytics/stream-analytics-monitoring

**NEW QUESTION 4**
- (Exam Topic 3)
You use Azure Data Factory to prepare data to be queried by Azure Synapse Analytics serverless SQL pools. Files are initially ingested into an Azure Data Lake Storage Gen2 account as 10 small JSON files. Each file contains the same data attributes and data from a subsidiary of your company.
You need to move the files to a different folder and transform the data to meet the following requirements: ≫ Provide the fastest possible query times.

≫ Automatically infer the schema from the underlying files.
How should you configure the Data Factory copy activity? To answer, select the appropriate options in the answer area.
NOTE: Each correct selection is worth one point.

Copy behavior:
- Flatten hierarchy
- Merge files
- Preserve hierarchy

Sink file type:
- CSV
- JSON
- Parquet
- TXT

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**
Box 1: Preserver herarchy
Compared to the flat namespace on Blob storage, the hierarchical namespace greatly improves the performance of directory management operations, which improves overall job performance.
Box 2: Parquet
Azure Data Factory parquet format is supported for Azure Data Lake Storage Gen2. Parquet supports the schema property.
Reference:
https://docs.microsoft.com/en-us/azure/storage/blobs/data-lake-storage-introduction https://docs.microsoft.com/en-us/azure/data-factory/format-parquet

**NEW QUESTION 5**
- (Exam Topic 3)
You are building an Azure Stream Analytics job to identify how much time a user spends interacting with a feature on a webpage.
The job receives events based on user actions on the webpage. Each row of data represents an event. Each event has a type of either 'start' or 'end'.
You need to calculate the duration between start and end events.
How should you complete the query? To answer, select the appropriate options in the answer area. NOTE: Each correct selection is worth one point.

```
SELECT
    [user],
    feature,
    ┌─────────────────┐▼
    │ DATEADD(        │
    │ DATEDIFF(       │
    │ DATEPART(       │
    └─────────────────┘
        second,
    ┌──────────────┐▼   (Time) OVER (PARTITION BY [user], feature LIMIT DURATION(hour, 1) WHEN Event = 'start'),
    │ ISFIRST      │
    │ LAST         │
    │ TOPONE       │
    └──────────────┘
        Time) as duration
FROM input TIMESTAMP BY Time
WHERE
        Event = 'end'
```

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**
Box 1: DATEDIFF
DATEDIFF function returns the count (as a signed integer value) of the specified datepart boundaries crossed between the specified startdate and enddate.
Syntax: DATEDIFF ( datepart , startdate, enddate ) Box 2: LAST
The LAST function can be used to retrieve the last event within a specific condition. In this example, the condition is an event of type Start, partitioning the search by PARTITION BY user and feature. This way, every user and feature is treated independently when searching for the Start event. LIMIT DURATION limits the search back in time to 1 hour between the End and Start events.
Example: SELECT
[user], feature, DATEDIFF(
second,
LAST(Time) OVER (PARTITION BY [user], feature LIMIT DURATION(hour,
1) WHEN Event = 'start'), Time) as duration
FROM input TIMESTAMP BY Time
WHERE
Event = 'end' Reference:
https://docs.microsoft.com/en-us/azure/stream-analytics/stream-analytics-stream-analytics-query-patterns

**NEW QUESTION 6**
- (Exam Topic 3)
You have an Azure Data Lake Storage account that has a virtual network service endpoint configured.
You plan to use Azure Data Factory to extract data from the Data Lake Storage account. The data will then be loaded to a data warehouse in Azure Synapse Analytics by using PolyBase.
Which authentication method should you use to access Data Lake Storage?

A. shared access key authentication
B. managed identity authentication
C. account key authentication
D. service principal authentication

**Answer:** B

**Explanation:**
Reference:
https://docs.microsoft.com/en-us/azure/data-factory/connector-azure-sql-data-warehouse#use-polybase-to-load-d

**NEW QUESTION 7**
- (Exam Topic 3)
You have an enterprise-wide Azure Data Lake Storage Gen2 account. The data lake is accessible only through an Azure virtual network named VNET1.
You are building a SQL pool in Azure Synapse that will use data from the data lake.
Your company has a sales team. All the members of the sales team are in an Azure Active Directory group named Sales. POSIX controls are used to assign the Sales group access to the files in the data lake.
You plan to load data to the SQL pool every hour.
You need to ensure that the SQL pool can load the sales data from the data lake.
Which three actions should you perform? Each correct answer presents part of the solution. NOTE: Each area selection is worth one point.

A. Add the managed identity to the Sales group.
B. Use the managed identity as the credentials for the data load process.
C. Create a shared access signature (SAS).
D. Add your Azure Active Directory (Azure AD) account to the Sales group.
E. Use the snared access signature (SAS) as the credentials for the data load process.
F. Create a managed identity.

**Answer:** ADF

**Explanation:**
The managed identity grants permissions to the dedicated SQL pools in the workspace.
Note: Managed identity for Azure resources is a feature of Azure Active Directory. The feature provides Azure services with an automatically managed identity in Azure AD Reference:
https://docs.microsoft.com/en-us/azure/synapse-analytics/security/synapse-workspace-managed-identity

**NEW QUESTION 8**
- (Exam Topic 3)
What should you recommend using to secure sensitive customer contact information?

A. data labels
B. column-level security
C. row-level security
D. Transparent Data Encryption (TDE)

**Answer:** B

**Explanation:**
Scenario: All cloud data must be encrypted at rest and in transit.
Always Encrypted is a feature designed to protect sensitive data stored in specific database columns from access (for example, credit card numbers, national identification numbers, or data on a need to know basis). This includes database administrators or other privileged users who are authorized to access the database to perform management tasks, but have no business need to access the particular data in the encrypted columns. The data is always encrypted, which means the encrypted data is decrypted only for processing by client applications with access to the encryption key.
References:
https://docs.microsoft.com/en-us/azure/sql-database/sql-database-security-overview

**NEW QUESTION 9**
- (Exam Topic 3)
You have a Microsoft SQL Server database that uses a third normal form schema.
You plan to migrate the data in the database to a star schema in an A?\ire Synapse Analytics dedicated SQI pool.
You need to design the dimension tables. The solution must optimize read operations.
What should you include in the solution? to answer, select the appropriate options in the answer area. NOTE: Each correct selection is worth one point.

**Answer Area**

Transform data for the dimension tables by: ▼

For the primary key columns in the dimension tables, use:
New IDENTITY columns
A new computed column
The business key column from the source sys

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**

**Answer Area**

Transform data for the dimension tables by: ▼

For the primary key columns in the dimension tables, use:
New IDENTITY columns
A new computed column
The business key column from the source sys

**NEW QUESTION 10**
- (Exam Topic 3)
You have an Azure Synapse Analytics dedicated SQL pool that contains a table named Table1.
You have files that are ingested and loaded into an Azure Data Lake Storage Gen2 container named container1.
You plan to insert data from the files into Table1 and azure Data Lake Storage Gen2 container named container1.
You plan to insert data from the files into Table1 and transform the data. Each row of data in the files will produce one row in the serving layer of Table1.
You need to ensure that when the source data files are loaded to container1, the DateTime is stored as an additional column in Table1.
Solution: You use a dedicated SQL pool to create an external table that has a additional DateTime column. Does this meet the goal?

A. Yes
B. No

**Answer:** A

**NEW QUESTION 10**
- (Exam Topic 3)
You develop a dataset named DBTBL1 by using Azure Databricks. DBTBL1 contains the following columns:
• SensorTypeID
• GeographyRegionID
• Year
• Month
• Day
• Hour
• Minute
• Temperature
• WindSpeed
• Other

You need to store the data to support daily incremental load pipelines that vary for each GeographyRegionID. The solution must minimize storage costs.
How should you complete the code? To answer, select the appropriate options in the answer area. NOTE: Each correct selection is worth one point.

**Answer Area**

```
df.write
```

| .bucketBy | ("*") |
| .format | ("GeographyRegionID") |
| .partitionBy | ("GeographyRegionID","Year","Month","Day") |
| .sortBy | ("Year","Month","Day","GeographyRegionID") |

| .csv("/DBTBL1") |
| .json("/DBTBL1") |
| .parquet("/DBTBL1") |
| .saveAsTable("/DBTBL1") |

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**

**Answer Area**

```
df.write
```

| .bucketBy | ("*") |
| .format | ("GeographyRegionID") |
| .partitionBy | ("GeographyRegionID","Year","Month","Day") |
| .sortBy | ("Year","Month","Day","GeographyRegionID") |

| .csv("/DBTBL1") |
| .json("/DBTBL1") |
| .parquet("/DBTBL1") |
| .saveAsTable("/DBTBL1") |

**NEW QUESTION 12**
- (Exam Topic 3)
You have a self-hosted integration runtime in Azure Data Factory.
The current status of the integration runtime has the following configurations:
➤ Status: Running
➤ Type: Self-Hosted
➤ Version: 4.4.7292.1
➤ Running / Registered Node(s): 1/1
➤ High Availability Enabled: False
➤ Linked Count: 0
➤ Queue Length: 0
➤ Average Queue Duration. 0.00s
The integration runtime has the following node details:
➤ Name: X-M
➤ Status: Running
➤ Version: 4.4.7292.1
➤ Available Memory: 7697MB
➤ CPU Utilization: 6%

≫ Network (In/Out): 1.21KBps/0.83KBps

≫ Concurrent Jobs (Running/Limit): 2/14

≫ Role: Dispatcher/Worker

≫ Credential Status: In Sync

Use the drop-down menus to select the answer choice that completes each statement based on the information presented.

NOTE: Each correct selection is worth one point.

If the X-M node becomes unavailable, all executed pipelines will:

| ▼ |
|---|
| fail until the node comes back online |
| switch to another integration runtime |
| exceed the CPU limit |

The number of concurrent jobs and the CPU usage indicate that the Concurrent Jobs (Running/Limit) value should be:

| ▼ |
|---|
| raised |
| lowered |
| left as is |

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**
Box 1: fail until the node comes back online We see: High Availability Enabled: False
Note: Higher availability of the self-hosted integration runtime so that it's no longer the single point of failure in your big data solution or cloud data integration with Data Factory.
Box 2: lowered We see:
Concurrent Jobs (Running/Limit): 2/14 CPU Utilization: 6%
Note: When the processor and available RAM aren't well utilized, but the execution of concurrent jobs reaches a node's limits, scale up by increasing the number of concurrent jobs that a node can run
Reference:
https://docs.microsoft.com/en-us/azure/data-factory/create-self-hosted-integration-runtime

**NEW QUESTION 13**
- (Exam Topic 3)
You are designing a statistical analysis solution that will use custom proprietary1 Python functions on near real-time data from Azure Event Hubs.
You need to recommend which Azure service to use to perform the statistical analysis. The solution must minimize latency.
What should you recommend?

A. Azure Stream Analytics
B. Azure SQL Database
C. Azure Databricks
D. Azure Synapse Analytics

**Answer:** A

**NEW QUESTION 14**
- (Exam Topic 3)
You are designing a real-time dashboard solution that will visualize streaming data from remote sensors that connect to the internet. The streaming data must be aggregated to show the average value of each 10-second interval. The data will be discarded after being displayed in the dashboard.
The solution will use Azure Stream Analytics and must meet the following requirements:

≫ Minimize latency from an Azure Event hub to the dashboard.

≫ Minimize the required storage.

≫ Minimize development effort.

What should you include in the solution? To answer, select the appropriate options in the answer area. NOTE: Each correct selection is worth one point

Azure Stream Analytics input type:

| Azure Event Hub |
| Azure SQL Database |
| Azure Stream Analytics |
| Microsoft Power BI |

Azure Stream Analytics output type:

| Azure Event Hub |
| Azure SQL Database |
| Azure Stream Analytics |
| Microsoft Power BI |

Aggregation query location:

| Azure Event Hub |
| Azure SQL Database |
| Azure Stream Analytics |
| Microsoft Power BI |

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**
Reference:
https://docs.microsoft.com/en-us/azure/stream-analytics/stream-analytics-power-bi-dashboard

**NEW QUESTION 18**
- (Exam Topic 3)
You are developing a solution using a Lambda architecture on Microsoft Azure. The data at test layer must meet the following requirements:
Data storage:
•Serve as a repository (or high volumes of large files in various formats.
•Implement optimized storage for big data analytics workloads.
•Ensure that data can be organized using a hierarchical structure. Batch processing:
•Use a managed solution for in-memory computation processing.
•Natively support Scala, Python, and R programming languages.
•Provide the ability to resize and terminate the cluster automatically. Analytical data store:
•Support parallel processing.
•Use columnar storage.
•Support SQL-based languages.
You need to identify the correct technologies to build the Lambda architecture.
Which technologies should you use? To answer, select the appropriate options in the answer area NOTE: Each correct selection is worth one point.

| Architecture requirement | Technology |
| --- | --- |
| Data storage | ▼ |
| | Azure SQL Database |
| | Azure Blob Storage |
| | Azure Cosmos DB |
| | Azure Data Lake Store |
| Batch processing | ▼ |
| | HDInsight Spark |
| | HDInsight Hadoop |
| | Azure Databricks |
| | HDInsight Interactive Query |
| Analytical data store | ▼ |
| | HDInsight HBase |
| | Azure SQL Data Warehouse |
| | Azure Analysis Services |
| | Azure Cosmos DB |

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**
Data storage: Azure Data Lake Store
A key mechanism that allows Azure Data Lake Storage Gen2 to provide file system performance at object storage scale and prices is the addition of a hierarchical namespace. This allows the collection of objects/files within an account to be organized into a hierarchy of directories and nested subdirectories in the same way that the file system on your computer is organized. With the hierarchical namespace enabled, a storage account becomes capable of providing the scalability and cost-effectiveness of object storage, with file system semantics that are familiar to analytics engines and frameworks.
Batch processing: HD Insight Spark
Aparch Spark is an open-source, parallel-processing framework that supports in-memory processing to boost the performance of big-data analysis applications. HDInsight is a managed Hadoop service. Use it deploy and manage Hadoop clusters in Azure. For batch processing, you can use Spark, Hive, Hive LLAP, MapReduce.
Languages: R, Python, Java, Scala, SQL Analytic data store: SQL Data Warehouse
SQL Data Warehouse is a cloud-based Enterprise Data Warehouse (EDW) that uses Massively Parallel Processing (MPP).
SQL Data Warehouse stores data into relational tables with columnar storage. References:
https://docs.microsoft.com/en-us/azure/storage/blobs/data-lake-storage-namespace https://docs.microsoft.com/en-us/azure/architecture/data-guide/technology-choices/batch-processing https://docs.microsoft.com/en-us/azure/sql-data-warehouse/sql-data-warehouse-overview-what-is

**NEW QUESTION 23**
- (Exam Topic 3)
You have an Azure Databricks workspace named workspace1 in the Standard pricing tier.
You need to configure workspace1 to support autoscaling all-purpose clusters. The solution must meet the following requirements:

≫ Automatically scale down workers when the cluster is underutilized for three minutes.

≫ Minimize the time it takes to scale to the maximum number of workers.

≫ Minimize costs.
What should you do first?

A. Enable container services for workspace1.
B. Upgrade workspace1 to the Premium pricing tier.
C. Set Cluster Mode to High Concurrency.
D. Create a cluster policy in workspace1.

**Answer:** B

**Explanation:**
For clusters running Databricks Runtime 6.4 and above, optimized autoscaling is used by all-purpose clusters in the Premium plan
Optimized autoscaling:
Scales up from min to max in 2 steps.
Can scale down even if the cluster is not idle by looking at shuffle file state. Scales down based on a percentage of current nodes.
On job clusters, scales down if the cluster is underutilized over the last 40 seconds.
On all-purpose clusters, scales down if the cluster is underutilized over the last 150 seconds.
The spark.databricks.aggressiveWindowDownS Spark configuration property specifies in seconds how often a cluster makes down-scaling decisions. Increasing the value causes a cluster to scale down more slowly. The maximum value is 600.
Note: Standard autoscaling
Starts with adding 8 nodes. Thereafter, scales up exponentially, but can take many steps to reach the max. You can customize the first step by setting the spark.databricks.autoscaling.standardFirstStepUp Spark configuration property.

Scales down only when the cluster is completely idle and it has been underutilized for the last 10 minutes. Scales down exponentially, starting with 1 node.
Reference:

**NEW QUESTION 25**
- (Exam Topic 3)
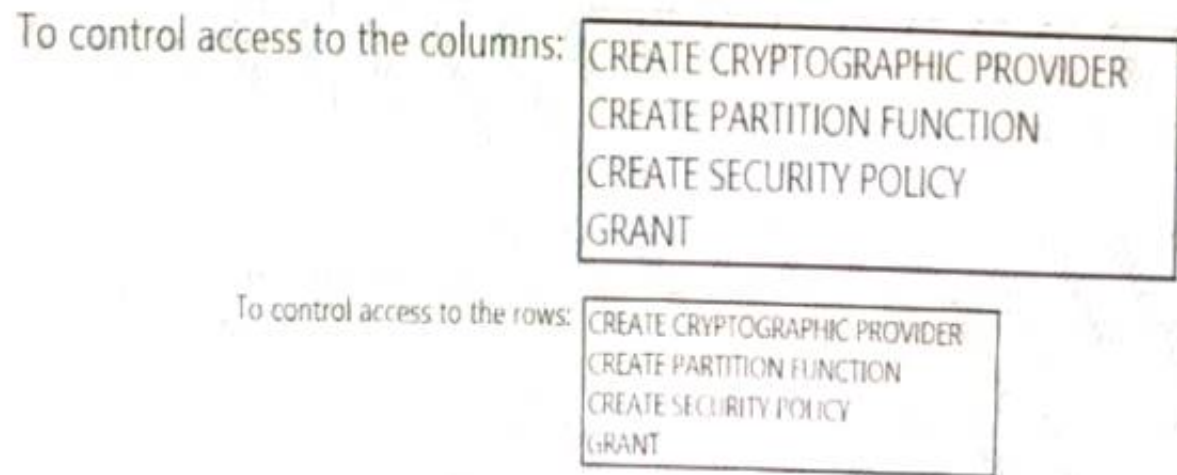You have an Azure subscription that contains the following resources:
* An Azure Active Directory (Azure AD) tenant that contains a security group named Group1.
* An Azure Synapse Analytics SQL pool named Pool1.
You need to control the access of Group1 to specific columns and rows in a table in Pool1
Which Transact-SQL commands should you use? To answer, select the appropriate options in the answer area. NOTE: Each appropriate options in the answer area.

**Answer Area**

To control access to the columns:
- CREATE CRYPTOGRAPHIC PROVIDER
- CREATE PARTITION FUNCTION
- CREATE SECURITY POLICY
- GRANT

To control access to the rows:
- CREATE CRYPTOGRAPHIC PROVIDER
- CREATE PARTITION FUNCTION
- CREATE SECURITY POLICY
- GRANT

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**

**Answer Area**

To control access to the columns:
- CREATE CRYPTOGRAPHIC PROVIDER
- CREATE PARTITION FUNCTION
- CREATE SECURITY POLICY
- GRANT

To control access to the rows:
- CREATE CRYPTOGRAPHIC PROVIDER
- CREATE PARTITION FUNCTION
- CREATE SECURITY POLICY
- GRANT

**NEW QUESTION 30**
- (Exam Topic 3)
You plan to implement an Azure Data Lake Gen2 storage account.
You need to ensure that the data lake will remain available if a data center fails in the primary Azure region. The solution must minimize costs.
Which type of replication should you use for the storage account?

A. geo-redundant storage (GRS)
B. zone-redundant storage (ZRS)
C. locally-redundant storage (LRS)
D. geo-zone-redundant storage (GZRS)

**Answer:** A

**Explanation:**
Geo-redundant storage (GRS) copies your data synchronously three times within a single physical location in the primary region using LRS. It then copies your data asynchronously to a single physical location in the secondary region.
Reference:
https://docs.microsoft.com/en-us/azure/storage/common/storage-redundancy

**NEW QUESTION 32**
- (Exam Topic 3)
You have an Azure Data Lake Storage Gen2 account that contains a JSON file for customers. The file contains two attributes named FirstName and LastName.

You need to copy the data from the JSON file to an Azure Synapse Analytics table by using Azure Databricks. A new column must be created that concatenates the FirstName and LastName values.
You create the following components:

≫ A destination table in Azure Synapse

≫ An Azure Blob storage container

≫ A service principal

Which five actions should you perform in sequence next in is Databricks notebook? To answer, move the appropriate actions from the list of actions to the answer area and arrange them in the correct order.

**Actions**

Mount the Data Lake Storage onto DBFS.

Write the results to a table in Azure Synapse.

Perform transformations on the file.

Specify a temporary folder to stage the data.

Write the results to Data Lake Storage.

Read the file into a data frame.

Drop the data frame.

Perform transformations on the data frame.

**Answer Area**

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**
Step 1: Read the file into a data frame.
You can load the json files as a data frame in Azure Databricks. Step 2: Perform transformations on the data frame.
Step 3: Specify a temporary folder to stage the data
Specify a temporary folder to use while moving data between Azure Databricks and Azure Synapse. Step 4: Write the results to a table in Azure Synapse.
You upload the transformed data frame into Azure Synapse. You use the Azure Synapse connector for Azure Databricks to directly upload a dataframe as a table in a Azure Synapse.
Step 5: Drop the data frame
Clean up resources. You can terminate the cluster. From the Azure Databricks workspace, select Clusters on the left. For the cluster to terminate, under Actions, point to the ellipsis (...) and select the Terminate icon.
Reference:
https://docs.microsoft.com/en-us/azure/azure-databricks/databricks-extract-load-sql-data-warehouse

**NEW QUESTION 35**
- (Exam Topic 3)
You store files in an Azure Data Lake Storage Gen2 container. The container has the storage policy shown in the following exhibit.

```
{
    "rules": [
        {
            "enabled": true,
            "name": "contosorule",
            "type": "Lifecycle",
            "definition": {
                "actions": {
                    "version": {
                        "delete": {
                            "daysAfterCreationGreaterThan": 60
                        }
                    },
                    "baseBlob": {
                        "tierToCool": {
                            "daysAfterModificationGreaterThan": 30
                        },
                    },
                }
            },
            "filters": {
                "blobTypes": [
                    "blockBlob"
                ],
                "prefixMatch": [
                    "container1/contoso"
                ]
            }
        }
    ]
}
```

Use the drop-down menus to select the answer choice that completes each statement based on the information presented in the graphic.
NOTE: Each correct selection Is worth one point.

**Answer Area**

The files are **[answer choice]** after 30 days.

| deleted from the container |
| moved to archive storage |
| moved to cool storage |
| moved to hot storage |

The storage policy applies to **[answer choice]**.

| container1/contoso1.csv |
| container1/docs/contoso.json |
| container1/mycontoso/contoso.csv |

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**

**Answer Area**

The files are **[answer choice]** after 30 days.

| deleted from the container |
| moved to archive storage |
| moved to cool storage |
| moved to hot storage |

The storage policy applies to **[answer choice]**.

| container1/contoso1.csv |
| container1/docs/contoso.json |
| container1/mycontoso/contoso.csv |

**NEW QUESTION 38**
- (Exam Topic 3)
You have an Azure Synapse Analytics serverless SQL pool named Pool1 and an Azure Data Lake Storage Gen2 account named storage1. The AllowedBlobpublicAccess porperty is disabled for storage1.
You need to create an external data source that can be used by Azure Active Directory (Azure AD) users to access storage1 from Pool1.
What should you create first?

A. an external resource pool
B. a remote service binding
C. database scoped credentials
D. an external library

**Answer:** C

**NEW QUESTION 42**
- (Exam Topic 3)
You are designing an application that will store petabytes of medical imaging data
When the data is first created, the data will be accessed frequently during the first week. After one month, the data must be accessible within 30 seconds, but files will be accessed infrequently. After one year, the data will be accessed infrequently but must be accessible within five minutes.
You need to select a storage strategy for the data. The solution must minimize costs.
Which storage tier should you use for each time frame? To answer, select the appropriate options in the answer area.
NOTE: Each correct selection is worth one point.

First week:
| Archive |
| Cool |
| Hot |

After one month:
| Archive |
| Cool |
| Hot |

After one year:
| Archive |
| Cool |
| Hot |

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**
First week: Hot
Hot - Optimized for storing data that is accessed frequently. After one month: Cool
Cool - Optimized for storing data that is infrequently accessed and stored for at least 30 days.
After one year: Cool

**NEW QUESTION 47**
- (Exam Topic 3)
You plan to create an Azure Synapse Analytics dedicated SQL pool.
You need to minimize the time it takes to identify queries that return confidential information as defined by the company's data privacy regulations and the users who executed the queries.
Which two components should you include in the solution? Each correct answer presents part of the solution. NOTE: Each correct selection is worth one point.

A. sensitivity-classification labels applied to columns that contain confidential information
B. resource tags for databases that contain confidential information
C. audit logs sent to a Log Analytics workspace
D. dynamic data masking for columns that contain confidential information

**Answer:** AC

**Explanation:**
A: You can classify columns manually, as an alternative or in addition to the recommendation-based classification:

≫ Select Add classification in the top menu of the pane.

≫ In the context window that opens, select the schema, table, and column that you want to classify, and the information type and sensitivity label.

≫ Select Add classification at the bottom of the context window.

C: An important aspect of the information-protection paradigm is the ability to monitor access to sensitive data. Azure SQL Auditing has been enhanced to include a new field in the audit log called data_sensitivity_information. This field logs the sensitivity classifications (labels) of the data that was returned by a query. Here's an example:



Reference:
https://docs.microsoft.com/en-us/azure/azure-sql/database/data-discovery-and-classification-overview

**NEW QUESTION 49**
- (Exam Topic 3)
You are designing an Azure Stream Analytics job to process incoming events from sensors in retail environments.
You need to process the events to produce a running average of shopper counts during the previous 15 minutes, calculated at five-minute intervals.
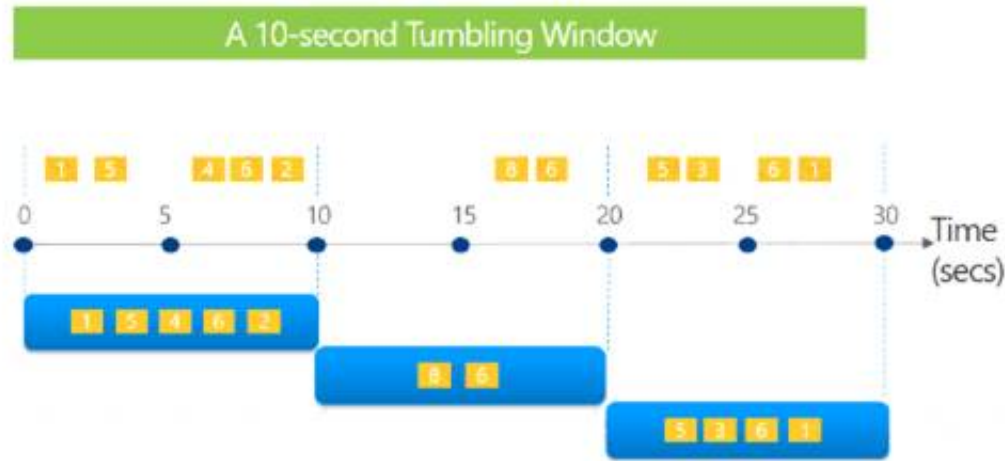Which type of window should you use?

A. snapshot
B. tumbling
C. hopping
D. sliding

**Answer:** B

**Explanation:**
Tumbling windows are a series of fixed-sized, non-overlapping and contiguous time intervals. The following diagram illustrates a stream with a series of events and how they are mapped into 10-second tumbling windows.

## Tell me the count of tweets per time zone every 10 seconds



```
SELECT TimeZone, COUNT(*) AS Count
FROM TwitterStream TIMESTAMP BY CreatedAt
GROUP BY TimeZone, TumblingWindow(second,10)
```

Reference:
https://docs.microsoft.com/en-us/stream-analytics-query/tumbling-window-azure-stream-analytics

**NEW QUESTION 52**
- (Exam Topic 3)
You implement an enterprise data warehouse in Azure Synapse Analytics. You have a large fact table that is 10 terabytes (TB) in size.
Incoming queries use the primary key SaleKey column to retrieve data as displayed in the following table:

| SaleKey | CityKey | CustomerKey | StockItemKey | InvoiceDateKey | Quantity | UnitPrice | TotalExcludingTax |
|---|---|---|---|---|---|---|---|
| 49309 | 90858 | 70 | 69 | 10/22/13 | 8 | 16 | 128 |
| 49313 | 55710 | 126 | 69 | 10/22/13 | 2 | 16 | 32 |
| 49343 | 44710 | 234 | 68 | 10/22/13 | 10 | 16 | 160 |
| 49352 | 66109 | 163 | 70 | 10/22/13 | 4 | 16 | 64 |
| 49488 | 65312 | 230 | 70 | 10/22/13 | 8 | 16 | 128 |
| 49646 | 85877 | 271 | 70 | 10/24/13 | 1 | 16 | 16 |
| 49798 | 41238 | 288 | 69 | 10/24/13 | 1 | 16 | 16 |

You need to distribute the large fact table across multiple nodes to optimize performance of the table. Which technology should you use?

A. hash distributed table with clustered index
B. hash distributed table with clustered Columnstore index
C. round robin distributed table with clustered index
D. round robin distributed table with clustered Columnstore index
E. heap table with distribution replicate

**Answer:** B

**Explanation:**
Hash-distributed tables improve query performance on large fact tables.
Columnstore indexes can achieve up to 100x better performance on analytics and data warehousing workloads and up to 10x better data compression than traditional rowstore indexes.
Reference:
https://docs.microsoft.com/en-us/azure/sql-data-warehouse/sql-data-warehouse-tables-distribute https://docs.microsoft.com/en-us/sql/relational-databases/indexes/columnstore-indexes-query-performance

**NEW QUESTION 55**
- (Exam Topic 3)
You have an Azure Factory instance named DF1 that contains a pipeline named PL1.PL1 includes a tumbling window trigger.
You create five clones of PL1. You configure each clone pipeline to use a different data source.
You need to ensure that the execution schedules of the clone pipeline match the execution schedule of PL1. What should you do?

A. Add a new trigger to each cloned pipeline
B. Associate each cloned pipeline to an existing trigger.
C. Create a tumbling window trigger dependency for the trigger of PL1.
D. Modify the Concurrency setting of each pipeline.

**Answer:** B

**NEW QUESTION 57**
- (Exam Topic 3)

You have a SQL pool in Azure Synapse.

A user reports that queries against the pool take longer than expected to complete. You need to add monitoring to the underlying storage to help diagnose the issue.

Which two metrics should you monitor? Each correct answer presents part of the solution. NOTE: Each correct selection is worth one point.

A. Cache used percentage
B. DWU Limit
C. Snapshot Storage Size
D. Active queries
E. Cache hit percentage

**Answer:** AE

**Explanation:**
A: Cache used is the sum of all bytes in the local SSD cache across all nodes and cache capacity is the sum of the storage capacity of the local SSD cache across all nodes.
E: Cache hits is the sum of all columnstore segments hits in the local SSD cache and cache miss is the columnstore segments misses in the local SSD cache summed across all nodes
Reference:
https://docs.microsoft.com/en-us/azure/synapse-analytics/sql-data-warehouse/sql-data-warehouse-concept-resou

**NEW QUESTION 58**
- (Exam Topic 3)
You are planning a streaming data solution that will use Azure Databricks. The solution will stream sales transaction data from an online store. The solution has the following specifications:
* The output data will contain items purchased, quantity, line total sales amount, and line total tax amount.
* Line total sales amount and line total tax amount will be aggregated in Databricks.
* Sales transactions will never be updated. Instead, new rows will be added to adjust a sale.
You need to recommend an output mode for the dataset that will be processed by using Structured Streaming. The solution must minimize duplicate data.
What should you recommend?

A. Append
B. Update
C. Complete

**Answer:** C

**NEW QUESTION 59**
- (Exam Topic 3)
You have the following table named Employees.

| first_name | last_name | hire_date | employee type |
|---|---|---|---|
| Jane | Doe | 2019-08-23 | new |
| Ben | Smith | 2017-12-15 | Standard |

You need to calculate the employee _type value based on the hire date value.

How should you complete the Transact-SQL statement? To answer, drag the appropriate values to the correct targets. Each value may be used once, more than once, or not at all. You may need to drag the split bar between panes or scroll to view content
NOTE: Each correct selection is worth one point.



A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**

Values

Answer Area

CASE

ELSE

OVER

PARTITION

ROW_NUMBER

```
SELECT
    *,
    CASE
        WHEN hire_date >= '2019-01-01' THEN
    'New'  PARTITION  'Standard'
    END AS employee_type
FROM
    employees;
```

**NEW QUESTION 61**
- (Exam Topic 3)
You have two Azure Data Factory instances named ADFdev and ADFprod. ADFdev connects to an Azure DevOps Git repository.
You publish changes from the main branch of the Git repository to ADFdev. You need to deploy the artifacts from ADFdev to ADFprod.
What should you do first?

A. From ADFdev, modify the Git configuration.
B. From ADFdev, create a linked service.
C. From Azure DevOps, create a release pipeline.
D. From Azure DevOps, update the main branch.

**Answer:** C

**Explanation:**
In Azure Data Factory, continuous integration and delivery (CI/CD) means moving Data Factory pipelines from one environment (development, test, production) to another.
Note:
The following is a guide for setting up an Azure Pipelines release that automates the deployment of a data factory to multiple environments.

> In Azure DevOps, open the project that's configured with your data factory.

> On the left side of the page, select Pipelines, and then select Releases.

> Select New pipeline, or, if you have existing pipelines, select New and then New release pipeline.

> In the Stage name box, enter the name of your environment.

> Select Add artifact, and then select the git repository configured with your development data factory.
Select the publish branch of the repository for the Default branch. By default, this publish branch is adf_publish.

> Select the Empty job template. Reference:
https://docs.microsoft.com/en-us/azure/data-factory/continuous-integration-deployment

**NEW QUESTION 64**
- (Exam Topic 3)
You use Azure Data Lake Storage Gen2.
You need to ensure that workloads can use filter predicates and column projections to filter data at the time the data is read from disk.
Which two actions should you perform? Each correct answer presents part of the solution. NOTE: Each correct selection is worth one point.

A. Reregister the Microsoft Data Lake Store resource provider.
B. Reregister the Azure Storage resource provider.
C. Create a storage policy that is scoped to a container.
D. Register the query acceleration feature.
E. Create a storage policy that is scoped to a container prefix filter.

**Answer:** BD

**NEW QUESTION 65**
- (Exam Topic 3)
You are implementing Azure Stream Analytics windowing functions.
Which windowing function should you use for each requirement? To answer, select the appropriate options in the answer area.
NOTE: Each correct selection is worth one point.

**Answer Area**

Segment the data stream into distinct time segments that repeat but do not overlap:

Hopping
Sliding
Tumbling

Segment the data stream into distinct time segments that repeat and can overlap:

Hopping
Sliding
Tumbling

Segment the data stream to produce an output only when an event occurs:

Hopping
Sliding
Tumbling

A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**

**Answer Area**

Segment the data stream into distinct time segments that repeat but do not overlap:

Hopping
Sliding
Tumbling

Segment the data stream into distinct time segments that repeat and can overlap:

Hopping
Sliding
Tumbling

Segment the data stream to produce an output only when an event occurs:

Hopping
Sliding
Tumbling

**NEW QUESTION 66**
- (Exam Topic 3)
You have an Azure Storage account and a data warehouse in Azure Synapse Analytics in the UK South region. You need to copy blob data from the storage account to the data warehouse by using Azure Data Factory. The solution must meet the following requirements:
≫ Ensure that the data remains in the UK South region at all times.
≫ Minimize administrative effort.
Which type of integration runtime should you use?

A. Azure integration runtime
B. Azure-SSIS integration runtime
C. Self-hosted integration runtime

**Answer:** A

**Explanation:**

| IR type | Public network | Private network |
|---|---|---|
| Azure | Data Flow<br>Data movement<br>Activity dispatch | |
| Self-hosted | Data movement<br>Activity dispatch | Data movement<br>Activity dispatch |
| Azure-SSIS | SSIS package execution | SSIS package execution |

Reference:
https://docs.microsoft.com/en-us/azure/data-factory/concepts-integration-runtime

**NEW QUESTION 68**
- (Exam Topic 3)
You have an Azure Synapse Analytics dedicated SQL pool that contains a table named Table1. You have files that are ingested and loaded into an Azure Data Lake Storage Gen2 container named
container1.
You plan to insert data from the files into Table1 and azure Data Lake Storage Gen2 container named container1.
You plan to insert data from the files into Table1 and transform the data. Each row of data in the files will produce one row in the serving layer of Table1.
You need to ensure that when the source data files are loaded to container1, the DateTime is stored as an additional column in Table1.
Solution: In an Azure Synapse Analytics pipeline, you use a data flow that contains a Derived Column transformation.

A. Yes
B. No

**Answer:** B

**NEW QUESTION 71**
- (Exam Topic 3)
You plan to implement an Azure Data Lake Storage Gen2 container that will contain CSV files. The size of the files will vary based on the number of events that occur per hour.
File sizes range from 4.KB to 5 GB.
You need to ensure that the files stored in the container are optimized for batch processing. What should you do?

A. Compress the files.
B. Merge the files.
C. Convert the files to JSON
D. Convert the files to Avro.

**Answer:** D

**NEW QUESTION 74**
- (Exam Topic 3)
You are building an Azure Analytics query that will receive input data from Azure IoT Hub and write the results to Azure Blob storage.
You need to calculate the difference in readings per sensor per hour.
How should you complete the query? To answer, select the appropriate options in the answer area.
NOTE: Each correct selection is worth one point.



A. Mastered
B. Not Mastered

**Answer:** A

**Explanation:**
Box 1: LAG
The LAG analytic operator allows one to look up a "previous" event in an event stream, within certain constraints. It is very useful for computing the rate of growth of a variable, detecting when a variable crosses a threshold, or when a condition starts or stops being true.
Box 2: LIMIT DURATION
Example: Compute the rate of growth, per sensor: SELECT sensorId,
growth = reading
LAG(reading) OVER (PARTITION BY sensorId LIMIT DURATION(hour, 1)) FROM input
Reference:
https://docs.microsoft.com/en-us/stream-analytics-query/lag-azure-stream-analytics

**NEW QUESTION 76**
- (Exam Topic 3)
You need to design an Azure Synapse Analytics dedicated SQL pool that meets the following requirements:

➢ Can return an employee record from a given point in time.

➢ Maintains the latest employee information.

➢ Minimizes query complexity.
How should you model the employee data?

A. as a temporal table
B. as a SQL graph table
C. as a degenerate dimension table
D. as a Type 2 slowly changing dimension (SCD) table

**Answer:** D

**Explanation:**
A Type 2 SCD supports versioning of dimension members. Often the source system doesn't store versions, so the data warehouse load process detects and manages changes in a dimension table. In this case, the dimension table must use a surrogate key to provide a unique reference to a version of the dimension member. It also includes columns that define the date range validity of the version (for example, StartDate and EndDate) and possibly a flag column (for example, IsCurrent) to easily filter by current dimension members.
Reference:
https://docs.microsoft.com/en-us/learn/modules/populate-slowly-changing-dimensions-azure-synapse-analytics

**NEW QUESTION 81**
- (Exam Topic 3)
You have an Azure Synapse workspace named MyWorkspace that contains an Apache Spark database named mytestdb.
You run the following command in an Azure Synapse Analytics Spark pool in MyWorkspace. CREATE TABLE mytestdb.myParquetTable(
EmployeeID int, EmployeeName string, EmployeeStartDate date) USING Parquet
You then use Spark to insert a row into mytestdb.myParquetTable. The row contains the following data.

| EmployeeName | EmployeeID | EmployeeStartDate |
| --- | --- | --- |
| Alice | 24 | 2020-01-25 |

One minute later, you execute the following query from a serverless SQL pool in MyWorkspace. SELECT EmployeeID
FROM mytestdb.dbo.myParquetTable WHERE name = 'Alice';
What will be returned by the query?

A. 24
B. an error
C. a null value

**Answer:** A

**Explanation:**
Once a database has been created by a Spark job, you can create tables in it with Spark that use Parquet as the storage format. Table names will be converted to lower case and need to be queried using the lower case name. These tables will immediately become available for querying by any of the Azure Synapse workspace Spark pools. They can also be used from any of the Spark jobs subject to permissions.
Note: For external tables, since they are synchronized to serverless SQL pool asynchronously, there will be a delay until they appear.
Reference:
https://docs.microsoft.com/en-us/azure/synapse-analytics/metadata/table

**NEW QUESTION 83**
- (Exam Topic 3)
You have an Azure Synapse Analytics dedicated SQL pool that contains a table named Table1. You have files that are ingested and loaded into an Azure Data Lake Storage Gen2 container named
container1.
You plan to insert data from the files into Table1 and azure Data Lake Storage Gen2 container named container1.
You plan to insert data from the files into Table1 and transform the data. Each row of data in the files will produce one row in the serving layer of Table1.
You need to ensure that when the source data files are loaded to container1, the DateTime is stored as an additional column in Table1.
Solution: In an Azure Synapse Analytics pipeline, you use a Get Metadata activity that retrieves the DateTime of the files.
Does this meet the goal?

A. Yes
B. No

**Answer:** B

**NEW QUESTION 85**
......

# THANKS FOR TRYING THE DEMO OF OUR PRODUCT

Visit Our Site to Purchase the Full Set of Actual DP-203 Exam Questions With Answers.

We Also Provide Practice Exam Software That Simulates Real Exam Environment And Has Many Self-Assessment Features. Order the DP-203 Product From:

## https://www.2passeasy.com/dumps/DP-203/

# Money Back Guarantee

## DP-203 Practice Exam Features:

* DP-203 Questions and Answers Updated Frequently

* DP-203 Practice Questions Verified by Expert Senior Certified Staff

* DP-203 Most Realistic Questions that Guarantee you a Pass on Your FirstTry

* DP-203 Practice Test Questions in Multiple Choice Formats and Updatesfor 1 Year