

1) (C) modify the settings of the Research division workspace to use an Azure Repos repository

2) Domain → The Fabric Admin portal

3) (D) an Azure Data factory Pipeline That executes dataflow to retrieve the maximum value of orderID column in destination

dataflow = Because procedure does not return value

maximum value = to identify the latest record already loaded in the Table

4) (B) Spark-SQL ("Select \* From Lakehouse1.ResearchProduct")

Spark.read.format("delta").load("Tables/Research Product")

external-table(ResearchProduct)

I take 100

64) (D) F All fabric capacities start with f

65) delta → tables / Productline1

- CSV is not compatible with managed tables (not supported)  
which is crucial in maintaining metadata & data life cycle in lakehouse

- Parquet: only delta is available for querying via SQL endpoint

- If you use saveAsTable() you don't need to specify the path "table", if you use save() you specify the full Path

66) (A), (B) and  (D) are correct: Microsoft Power BI Desktop  
Power BI Service  
Tabular Editor

67) Viewer access → Read lakehouse & Data warehouse and shortcuts with T-SQL Through TDS endpoint

Contributor → Read lakehouse through lakehouse explorer

(5) Read All Apache Spark  
Build Reports on default dataset  
Read All Apache Spark

1) ReadAll Permission (Read all data using Apache spark):

- The shared recipient has read access to the underlying Parquet files in onelake, which can be consumed using Spark
- Should be provided only if the shared recipient wants complete access to your warehouse's files using a Spark engine

2) Build permission (Build reports on default dataset):

The shared recipient can build reports on top of default semantic model that is connected to your warehouse.

(6) Average('Survey'[ResponseValue]),  
↳ only logical numeric value  
to aggregate

Period → The period to aggregate on.

(7) View → As it shows data as soon as the data is updated (immediately) in data store  
Case when (ListPrice > 50 And ListPrice ≤ 1000 Then 'Medium')

Note Within the warehouse, the user can add warehouse objects:  
tables or views to the default Power BI semantic model

(8) dataflow (not optimal as we can use the copy data tool to do the ingestion but it's not mentioned).

stored procedure → not making sense

xql → used for real time data

spark notebook → mentioned use low code as possible

(4) (C) Lakehouse

- Suitable for semi-structured & unstructured
- mentioned That Files loaded by The data engineers to one lake will be stored in The Parquet format and will meet Delta Lake Spec.

~~Hive~~ metastore:

centralized repository which store metadata for tables, partitions, columns, schemas and other data structure used by hive & spark

(60) (B) and (E) focussing on security & Efficiency

(61) SSO , Direct query.

Direct query Provide fresh data as it query directly from source

Direct Lake Provide almost fresh data, slightly lagging as it depends on frequency of updates to The Lakehouse files .

Single Sign on (SSO): preferred Authentication method, it allow users to authenticate using organization credentials .

- Simplify access across multiple systems
- " security by reducing credential need to manage

BASIC : less robust Authentication using <sup>username</sup> <sub>password</sub>

service Principal : <sup>used</sup> <sub>work</sub> with <sup>non-</sup> interactive and automated tools, Applications

(62) (A) & (D) Data flows are ideal for creating date dimensions from scratch, They are suitable for calculations and transformations before storage

(63) Create a Pipeline That has dependencies between activities and schedule The Pipeline

This solution is ideal as it ensure all data load activities are executed one after another in sequence

Note (dataflow lacks The comprehensive orchestration capabilities required That handle dependencies for multiple workflows or different operations)

(10) (A) \* Where can Come with group by  
\* Having Clause Can not be used by non-aggregate columns

(11) Row\_Number(), where  $x=1$

(12) display(df, summary=true) → Check The statistics of a given Apache spark dataframe

(13) - Sql analytics endpoint operates in read-only mode over lakehouse Delta tables. You can read only data from Delta tables using The sql analytics endpoint, They can save functions, views and set sql object-level security

- External tables created by Spark can not be accessed by Sql Endpoint

(13) Managed tables created by are fully-controlled by Lakehouse including storage, metadata management meaning These tables are fully supported for both read & write operations via SQL endpoints

(14) `Table.Profile`. Provide a statistical summary of table's columns (min, max, Avg, std, count)

`Table.Range` Extracts range of rows from a table starting from specified index

`Table.Max` Returns the maximum value in a column based on specified comparison criteria

`Table.MaxN` Returns the top N rows with the maximum values in a column, based on the specified comparison criteria

(15) `SparkSQL`  
`Pyspark` You can invoke the predict function using one of the following 3 ways:

- 1) Transformer API from SynapseML
- 2) Spark SQL API
- 3) Pyspark user-defined function (UDF)

(16) `display` by interactive visualization in notebook, displays data in barchart (and other various rich format data visualizations)

(2) `displayHTML` → show HTML content

(3) `show` → show/print off in the console

(17) to make duplicate rows displayed in Power BI report , we can add unique field to each row

(18) Define - introduce one or more entity definitions that can be applied to one or more Evaluate statement

Evaluate - introduce a statement containing a table expression required to execute a DAX query

Summarize - returns a summary table for the requested totals over set of groups

measure - to define a measure

Var - store the result of an expression in a variable for reuse in calculations.

(19) (#)

Distinct doesn't mean it occurs only once  
to enable the XMLA read-write, you need to modify the capacity settings as XMLA must be enabled at the capacity level

essential for creating & publishing semantic models  
(new way to author, collaborate and save projects)  
- store semantic model using TMDL format  
in visual studio code (Microsoft) for bulk changes  
- essential for regular report development as it contains whole model & report file on Desktop

PBIT →

- format used for Power BI Templates it includes the report layout & inquiries . Not include actual data

Power BI → - useful for creating reusable templates

(Q2) PBIDS → File format used for Power BI data source files, it define the data source & their connection details, but does not contain data model or report

(Q3) Grant Alter : give the permissions to perform DDL operations including creating, Truncating, Altering

Execute Permission : Grant executing Functions & stored procedures

(Q4) Automatic aggregations  
use state of the art machine learning to continuously optimize Direct query Semantic models for maximum report query performance (self training & optimizing) created by Power BI

User-defined aggregations:

Aggregations explicitly created & managed by model designer to optimize query performance over large Direct query Semantic models.  
By using These aggregations, you cache data at the aggregated level in-memory

Query Caching : is not applicable for Direct query  
or LiveConnect Semantic models

## (26) - Query folding :

Process:

Process of PowerQuery query is translated into native query  
that can be executed by the underlying data source , such  
as relational database ; allowing most of data transforma-  
logic to be pushed back to the data source , leveraging its  
Computational Power

Note - Queries that rely solely on unstructured data  
sources or that don't have a compute engine  
such as CSV or Excel files , don't have query folding

only refresh complete days setting :

ensure all rows of entire day are included in the  
refresh operation

- loading data takes an excessive amount of time and resources  
cause : Data type mismatch

The data type & format must be the same . If not :  
query folding cannot occur

(27) To scale out semantic model: Set large dataset storage format to on

(28)- Row-level security only applies to queries on a warehouse or sql analytics endpoint in fabric.

- Power BI queries on a warehouse in Direct Lake mode will fall back to Direct Query mode to abide by row-level security

### (30) Vertipaq Analyzer

is a tool specifically designed to analyze the memory usage of tables and columns in Power BI Models

# system.Discover\_storage\_Table\_Column\_segments dynamic Management View (DMV)

- This DMV provides detailed information about the storage and memory usage of table columns within a Power BI Models
- Return information about the column segments
  - used for storing data for in-memory tables

(32) - splitting Date/time Column into separate date & time Columns is beneficial in compression

- Each Component date & time has fewer unique Values than a Combined date/time column leading to reduce memory usage and potentially fast query Performance ; If you frequently query by only one component

- Remember Calculated Columns are Computed & stored in the model leading to additional memory consumption  
Unlike measure that is Computed in-memory (calculated at query time)

(33) In Power BI Desktop and tabular Editor, from settings, we can unEnable Implicit measures

(36) - to manage and utilize sql endpoint . its essential to make shortcut in table's section to work with structured data

only tables in delta format are available in SQL analytics endpoint. Parquet, CSV and other formats can not be queried using SQL Analytics endpoint

### (37) Load to tables option capability:

tables are always loaded using the delta lake table format with N-order optimization enabled

### Optimize feature

used to improve the performance of existing delta table by compacting files and optimizing their layout

(38) To enable specify the partition column in the destination settings:

- 1) from destination tab : Set Mode to Overwrite
- 2) from destination tab : Enable Partition
- 3) " " : Specify the 'Partition Column'

(39) Create Table As Clone of :

only the metadata of the table is copied. The underlying data of the table, stored as parquet files, is not copied.

(Zero-Copy Clone)

(40) Data stores that support delta-table and v-order optimized lakehouse & warehouse

Azure db ~~X~~ , MySQL db ~~X~~

(41) Spark employs compression by default and its partition can be read in parallel across multiple nodes

(45) In spark : on the existing of read & select, The spark read all the columns first, then select the specified columns, but on existing option of infer Schema, It will read only the selecting columns and infer datatype of them.

(47) The function Value.NativeQuery in Power Query is used to execute a native SQL query against a database.

Enable folding a feature in Power Query ensure query folding

a process where Power Query pushes transformations back to the source system so, they are executed in the database or data source rather than locally.

Improve Performance by offloading heavy computations to data Source

(48) In Delta Tables:

- 1) Bin-Compaction is achieved by the optimize Command  
it merges all changes into bigger, consolidating Parquet files

2) Dereferenced Storage Clean-up is achieved by the Vacuum Command.

(49) Type 1 SCD:

always reflect the latest values, and when changes in source data are detected, the dimension table data is overwritten

note managing effective end date of rows is characteristic of Type 2 SCD

For Azure Data Lake Storage (ADLS) Gen2 account,

If we need to make a shortcut to it:

Connection: https

Endpoint: dfs

(51) In Broadcast Join: It minimize data shuffling

- is an optimization technique used in Spark SQL engine.
- It utilized when one of data frames is small enough to be stored in the memory of all executor nodes
- Performance improving by of join operation by minimizing data shuffling across the network
- It's advantageous when working with large dataset

(52) Automatic: allow Direct Lake with fall back to Direct Query (default behavior)

(53) df.explain() - Print (logical & Physical) Plans to the Console for debugging purposes.

df.show() - Prints the first n rows of df to console

df.summary() - Computes specified statistics for numeric and string columns (mean, std, min - max - -)

(57) Retrieve Delta table History: → Help in maintenance, user, and reviewing

You can retrieve information including the operations, user, timestamp for each write to a Delta table by running the history command for (30 days)

Describe History table\_name - get the full history of the table

Describe History table\_name limit 1 - get the last operation only

Describe Detail -- get information about schema, partitioning, table size and so on.

Refresh Table\* table\_name

Command to refresh the metadata and date of a specific table. It ensures that the table's metadata and data in spark's catalog are synchronized with the underlying data source, especially after external changes

Explain table table\_name

Used to provide logical/Physical Plans for an input statement

(73) Dataflow Support Power Query M formula language expressions that has transformations capabilities during ingestion (copy data activity don't)

Error if exists :

on Saving dataframe , If data already exists : an exception to be thrown

Ignore :

" " "

" : The save operation is expected to not to save the contents of the df and not to change existing data

APPend

overwRite

(75) In sql : Function LEAST (expr [,...exp])

Return the minimum value from a list of one or more expressions

(79) Fuzzy merge: Type of data join that allows for inexact matching between columns in 2 dataset

useful when keys don't match because formatting or types

## Fuzzy Append

Similar to merge, but focuses on appending new rows to a dataset base on exact matching

## (7a) In SQL analytic Endpoint:

You can modify the default behaviour to not automatically add tables and views to the default Power BI semantic model

## (83)

lookup activity in data factory and Synapse analytics:

- read and returns the contents of a configuration file or table (data from a source)
- also returns the result of executing a query or stored procedure

It can be consumed in a subsequent copy, transformation, or control flow activities (ForEach activity)

## switch activity:

evaluates a set of activities corresponding to a case that matches the condition evaluation

## Append Variable activity :

to Add a Value to an existing array variable defined in a data factory or synapse Analytic Pipeline

### (qo) Calculate

Function evaluates an expression in a modified filter context

Calculate( sum(sales[Amount]), sales[Region] = "North")  
                  ↑  
                  expression  
                  ↓  
                  filter 1

### Generate

used to Create a Cartesian Product between 2 tables  
Generate(Products,Customers)

### CombineValues

Concatenate Values from multiple Columns or expressions into a Single text string , Separated by a specified delimiter

CombineValues ( "", Employees[FirstName] , Employees[lastName])  
                  ↓  
                  expression 1  
                  ↓  
                  delimiter  
                  expression 2

## SelectedMeasure

Function returns The measure Currently being evaluated in a Calculation group • used within Calculation items in a Calculation group  
SelectedMeasure()

## SelectedValue

returns The Value when There is a Single Value in The Context; otherwise, it return an alternative result (Blank)

SelectedValue( Products [ProductName] , "No Products selected")  
alternative Result

(q4)

## Boolean Expression

- True or False Value for each row

## Table Expression

- A Table with rows & Columns
- Operates on entire table or Subsets of it

- Direct filtering in calculate, keep, filter, sum
- Functions requiring a Table

Example  
Sales[Region] = "East"

Filter(Sales, Sales[Amount] > 100)

The Calculate Function accepts table expression by filter Dax , which evaluates its filter expression for each row of the table

However, it could be achieved much more efficiently by using Boolean expression

Here is code that use Boolean expression using keepFilters Dax function which ensures any existing filters applied and preserved , not overwritten

```
RedSales =  
Calculate(  
[Sales],  
keepFilters('Product'[color] = "Red")  
)
```

It's recommended to pass filter arguments as Boolean expression

It's because import model tables are in-memory column stores . They are explicitly optimized to efficiently filter columns in this way

Can not reference columns from multiple tables

QUESTION

### (99) DateTrunc

Used to truncate a date or datetime value to a specified precision, such as year, quarter, month, ... etc  
(remove smaller components of date)

Select DateTrunc (month,  $\downarrow$  orderDate);  
to first of that month in that orderDate

### DateName

Returns a character string representing the specified date part of specified date

DateName (weekday, orderdate);

↳ name of weekday (sunday) in that orderdate  
day → number of day (not name)

- Dataflow and Pipelines can not be shared at item level
- 1- Read all data using sql (**ReadData**)
  - ⇒ Read all objects within warehouse using T-sql
- 2- Read all oneLake data (**ReadAll**)
  - ⇒ Read all the warehouse's underlying oneLake files using Apache Spark, Pipelines, or other apps that access the oneLake data directly
- 3- Build reports on the default semantic model connected to warehouse (**Build**)
- to implement Row-level security, user should have access
  - either → **ReadData**
  - **object-level security**

\* to Create a RLS:

- 1) create schema
  - 2) Create a function that filter a table based on logged-user
  - 3) Create a security policy to apply The filtering function
- \* to implement OLS & CLS
- ⇒ The user should have access Read Permission only to warehouse
- ⇒ Can be done at warehouse or lakehouse T-SQL endpoint

\* using Manage onelake Data Access (PREVIEW)

- You can:
- 1) Define roles
  - 2) For each role give specific tables / files access to
- \* to Configure access, you have to be: Admin Member Contributor of workspace

## Data Governance

### \* Sensitivity Labels :

Labels needs to be created and managed in Microsoft Purview (before they become accessible in fabric)

In some industries, labelling data and information with a sensitivity label is necessary for compliance with information protection regulations

### \* Endorsement :

Endorses quality Content to help Colleagues discover it  
Promoted  $\Rightarrow$  recommended for others to use it  
Certified  $\Rightarrow$  Certified by your org as trusted source  
Master Data  $\Rightarrow$  Core source (single truth) of organizational data records

## Reusable Power BI assets

- 1) Power BI template can act as a starting point for a new report's layout, data model, and queries
    - It's a reusable asset which can improve efficiency & consistency
    - Export an existing report (.pbix file) as a Power BI Template
    - If there are any parameters in the template, you are asked to input these at the time of creating a new report
  - 2) The (.pbids file) is a Power BI data source file, which include connection details to a data source
    - Can help to quickly transfer all the data connection from one project to another
    - Pbids file can be exported from an existing report file
    - File can be edited in a text editor accordingly
- \* You can share semantic model from either :
- 1) Semantic model's options menu on OneLake data hub (stats)
  - 2) Data details page options "publish" button.

- we make version control at workspace level (settings)
  - Git in Azure Repos → same tenant as Fabric tenant
  - GitHub
  - GitHub Enterprise
- The semantic model can be represented
  - 1) TMSL
  - 2) THDL

## Semantic Model deployment (XNLA endpoint)

XNLA Endpoint: Support open-platform Connectivity from Microsoft and third-party client applications and tools to the semantic models workload

By default, read only connectivity using the endpoint is enabled

Read-write Provides more semantic model management, governance, advanced semantic modeling, debugging and monitoring

They are available to Premium or Embedded only

- \* Ingest batch with Dataflows Gen2 , Pipelines or notebooks
- \* Can ingest data for time series analysis in GUI-based way with eventstreams
- \* lakehouse → storing raw files and transforming to delta tables, that can be manipulated by spark (PySpark - sparkSQL)
  - warehouse → use T-SQL to work with large volumes of data
  - EventHouse → storing timeseries data that can be queried by kql
- \* You can convert column data types:
  - 1) Dataflow Gen2
  - 2) PySpark in notebooks

Direct Lake for Power BI :

Prerequisite: Fabric Capacity / Power BI Premium  
Lake house + SQL Endpoint / warehouse

Delta tables

### Limitations

- Querying one single warehouse or warehouse
- T-SQL views are not supported (will fall back to DAX)
- DAX queries exceeding limits or using unsupported features fall back to DQ mode
- No DAX Calculated cols
- No Composite model
- No Date time relationships

\* Relationship filter direction

Single

Both

## Dax Iterator Function

- Iterate over a table, evaluate the expression  
For each row, and then aggregate result
  - Ends with X
  - must whenever more than 1 column is used in the expression

## Aggregation

When you create Aggregation in composite model

original table → must be direct query

with Dim tables → You can set Dual storage mode  
to overcome limited relationships

Aggregate tables → in import mode gives the maximum performance

\* Unlike CSV files, Parquet files have the schema and data types embedded within them

Parquet files are the format supported by Delta Lake which allows for a transactional / data warehousing style capability to be used alongside the files

Parquet are better for compression