

— NO.

Title

Enjoy the most efficient handwriting experience! Taking notes on the go, whether for inspiration, ideas, knowledge learning, business insights, or even sketches...



Creative notes

By reading we enrich the mind;
by writing we polish it.

NOTE IN Created by Notein

In General

Data Warehouse

Characteristics:

- 1) Data integration: Ingest data from multiple sources into a unified store/repository
- 2) Schema-on-write: Data is transformed & structured before loading, adhering to a predefined schema for optimized querying
- 3) Security & governance: provide Role-based access control (RBAC), data encryption, and Compliance features for enterprise security.
- 4) ACID Compliance
- 5) High performance
- 6) Historical data management

Architecture

- 1) bottom layer: where data flow into DWH through ETL.
 - In many DWs, data is stored in RDBMs or similar system
- 2) The middle layer: build around an analytics engine (OLAP)
- 3) The Top layer: Includes user interfaces & reporting tools.

Data lake

- low-cost data storage solutions.
- handle massive volumes of data
- uses a schema-on-read approach → They don't apply standard format to incoming data but on accessing it
- Data lakes stores data on its native format; structured, unstructured, semi-structured

Architecture

- Early ones are built on HDFS, modern ones often use a cloud object store (Microsoft blob storage)
- Data lakes separate data storage from compute resources, which make them more cost-effective and scalable than DWH
- To process data in datalake: users can connect external data processing tools such as spark
- used: to maintain backups & archive old unused data | store datasets for ML, AI and big data analytics workloads

Data lake house

- Schema-on-Read & schema-on-write
- Diverse data types :
- unified storage: combines the scalability of data lakes with the performance of data warehouse
- cost in governance & security | Cost efficiency

Fabric lakehouse

- * is a data architecture platform for storing, managing, and analyzing **structured** and **unstructured** data in a single location
- * flexible & scalable solution that allows organizations to handle large volumes of data using various tools and frameworks to process and analyze the data.
- * lakehouse combine → **Scalability of a data lake** with → **performance and structure of a data warehouse** providing a unified platform for data storage, management, and analytics
- * you can perform sql-analytics directly or top of the delta tables in the lake to provides a frictionless and performant experience (ingestion → reporting)
through **sql analytics endpoint** (read-only)
- Note only the tables in delta format are available in sql analytics endpoint
- + lakehouse provide a fully managed file to table experience
through **Automatic table discovery and registration**
 - 1- You drop a file into the managed area of lakehouse
 - 2- system auto. validates it for the **structured formats** and registers it into the metastore with necessary metadata such as **columns by names**.

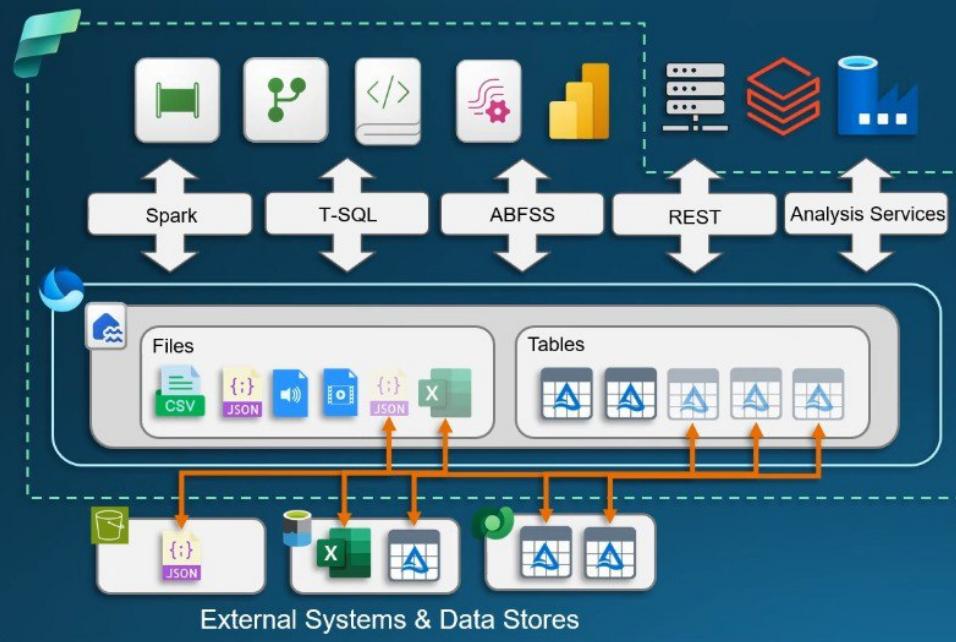
only support
delta table

formats, Compression and more.

3- You can after That reference The file as a table and use like spark to interact

* You can interact with lakehouse through pipelines, Notebooks
lakehouse explorer , Data flow Gen2 , Apace spark jobs

Lakehouse



Consume

- Most Fabric tools support interacting with Lakehouse data
- Can be consumed with external tools
 - Custom applications
 - Azure Databricks and Data Factory

Compute

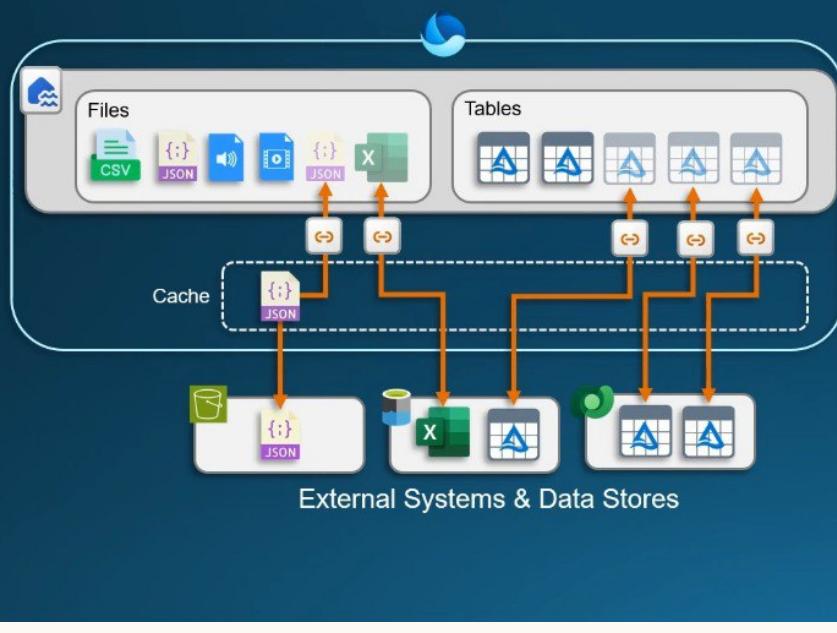
- Different interfaces with different capabilities and limitations
- Powered by Microsoft Fabric Capacity
 - Capacity needs to be turned on for data access

Store

- Can store files and tables
- Data is stored to Microsoft Fabric's OneLake
 - Unified data lake for entire organization
 - Only one OneLake per Fabric
 - SaaS version of Azure Data Lake Storage
- Supports storing files in any format
- Stores tables in Delta Parquet format
- Shortcuts
 - Allows virtualizing data via Lakehouse
 - Physical data stays in the external location
 - Data doesn't need to be duplicated over
 - Data is always up to date
 - Supports files and tables

Shortcuts in OneLake allow you quickly & easily source data from external cloud providers and use it across all fabric workloads

Lakehouse Shortcuts



Core Features of Shortcuts

- No data duplication

- Instant access

- Supports multiple sources

- Other Fabric Data Stores
- ADLS Gen2
- Amazon S3 (or compatible)
- Dataverse
- Google Cloud Storage
- Apache Iceberg

- Mount to Files or Tables section of Lakehouse

- Shortcut caching

- Workspace setting
- Stores frequently accessed data locally in OneLake
- Improves performance by avoiding repeated reads from external sources
- Invisible to users
- Cached files are kept for 1–28 days (configurable)
- Files larger than 1 GB are not cached
- Currently supported for GCS, S3, S3 compatible, (and on-premises data gateway shortcuts)

Is behind The scenes Copy of data

However:

Managed entirely by fabric

each time these workloads read data from cross-cloud sources charges additional egress fees on the data.

Thankfully, shortcut caching allows the data to be sourced once and then used across all fabric workloads without additional egress fees.

stores local copy of the data

Inside OneLake after the first need.

retention period

defines: How long data is stored and available before it's automatically removed

Basics of Delta Tables



dbo.person

OneLake file structure

```
.../dbo/person/
  - _delta_log
  - 0001.parquet
  - 0002.parquet
  - 0003.parquet
  - 0004.parquet
```



Table Representation

first_name	last_name
Aleksi	Partanen
John	Smith
Jane	Doe
Alex	Wang
Lisa	Bush
Jack	Doe

OneLake file structure

- Delta log keeps track of the parquet files
- Parquet files store the actual data

Operations

- INSERT
 - Adds new files
 - Updates delta log
- UPDATE / DELETE / MERGE
 - Doesn't modify files directly
 - Marks old files as invalid
 - Writes new files
 - Updates delta log

Some of the special features

- Time travel
- ACID transactions
- Schema evolution

Schema evolution

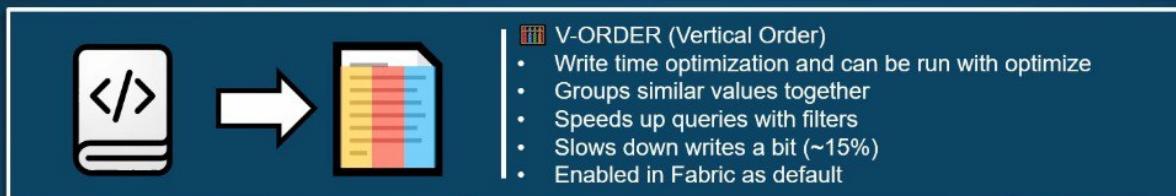
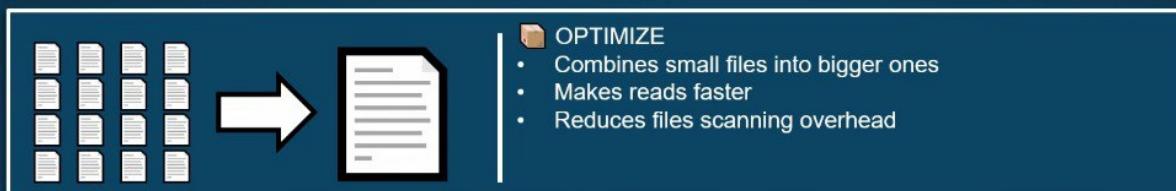
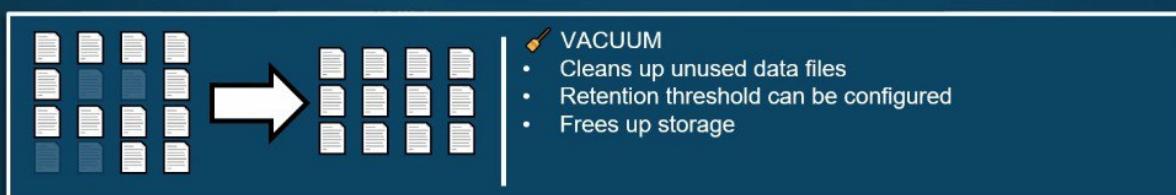


- ability of data system to handle changes in data structure over time without breaking existing pipelines, queries, or reports

- allow systems to adapt safely to changes in schema:

new fields added, data type changes, columns are removed

Maintenance & Optimization



NOTE

Created by Notein

Optimize Command

- consolidate multiple smaller parquet files into fewer larger files
- Reading lots of small files leads to more overhead & reduced performance, while larger files have better data distribution & compression
- we can add applying V-order beside it Optimize in UI
- Optimize Command takes time to run as it's potentially quite heavy on resources
- after applying optimize Command:
 - 1) in delta-log, we find new added json file state that previous Parquet files are removed and packed in N file(s) and a metric column with json data is returned, where it shows what happened during optimization (How many files added, removed)
 - 2) previous Parquet files are not actually removed to use time travel feature on need
- Optimize Command is Idempotent → nothing happens on executing it a second time.
- You can add a predicate with the command to optimize only part of table. However, the predicate needs to be applied to a Partitioning Column.

Vacuum Command

- As you see using optimize command, The old files are not removed and essentially total storage increased
- We can then use Vacuum to remove old unneeded files
- Vacuum can break the "time travel" feature.
- We can use retention threshold to control the deletion based on it
specify age of files that should be deleted based on it (maximum 7 days)

V-Order

- Write time optimization technique
- It applies special sorting, encoding, row group distribution, and compression to achieve about 50% more compression at the cost of some write overhead (15%)
- V-order is compatible with open parquet format, so any engine can read parquet files, can read it with V-order.
- disabled by default

- ↓
- table property reduces number of parquet files written when delta table is created & Increases file size
 - Comes with extra processing cost on writing, but improves reading performance as fewer files must be read
 - boosts Vacuum & Optimize Commands
 - advised to be used on partitioned table
 - enabled by default in Microsoft Fabric Runtime for Apache spark

Merge optimization

- use Merge statement to update/modify delta tables with advanced options
- It update data from a source table, view or dataframe into a target table by using Merge Command
- The open-source distribution of delta is not fully optimized for handling unmodified data (even matched rows are shuffled across spark nodes)
with this feature, an optimization called "Low Shuffle Merge" is done to exclude unmodified rows from an expensive shuffling operation
- enabled by default.

Fabric Datawarehouse

- **Lake centric warehouse** → Fabric datawarehouse stores data directly in the data lake using open formats like parquet which enables:
 - Easily handles massive data volumes
 - Different engineers can access the same data with different tools
 - Data doesn't need to be copied
(No data duplication)
- The datawarehouse uses a high-performance distributed processing engine.
- Data stored in Delta-parquet format
- Datawarehousing items:
Fabric Datawarehouse is a lake warehouse that supports 2 distinct warehousing items:
 - The Fabric data warehouse &
 - The SQL analytics endpoint

Warehouse

Key Points

- Feels like a traditional relational database
- Stores data in Delta Parquet format and stores data to Fabric OneLake under the hood
- Built for structured, tabular data
- Not designed for raw or unstructured files
- Ideal for the Silver and Gold layers of the Medallion architecture
- Best suited for cleaned and transformed data, not raw source files
- Familiar to anyone who's used SQL Server or other relational databases
- Powered by the cloud and an open storage format

Objects

- **Schemas**: Group related objects (like folders)
- **Tables**: Store actual data
- **Views**: Saved SQL queries that behave like virtual tables
- **Functions**: Return values from input (e.g., fiscal quarter logic)

The screenshot shows the Fabric Datawarehouse interface. The left sidebar is the Explorer pane, which includes sections for Workspaces, Catalog, Monitor, Real-Time, Workloads, and a list of databases (wh_dp700, wh_dp700_e004, wh_dp700_e005). The 'wh_dp700' database is selected. Under 'Schemas', there is a 'dbo' schema containing a 'person' table. The 'Tables' section also lists 'Views', 'Functions', 'Stored Procedures', and security-related items. To the right of the Explorer is a 'Data preview - person' window showing a table with three rows of data:

ID	Name
1	Alice
2	Bob
3	Charlie

- **Stored Procedures:** Encapsulate SQL logic for tasks and automation



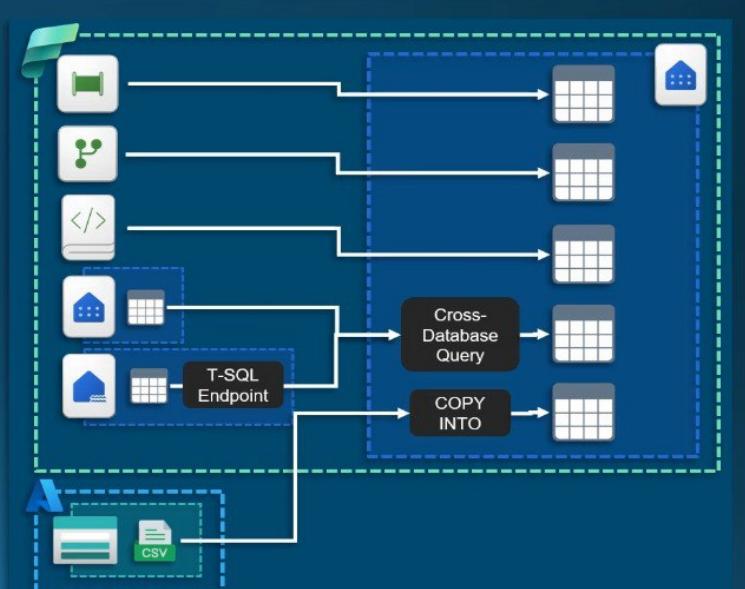
- The same feature in lakehouse:

A Warehouse, SQL database, and Fabric One lake all automatically provision a SQL analytics endpoint when created.

with it, T-SQL Commands can define and query data objects only (~~Read-only~~)

Ingesting Data to Warehouse

- Data pipelines
- Dataflow Gen2
- Notebooks
 - T-SQL Notebooks
 - Spark Connector
 - ODBC and other more generic interfaces
- T-SQL
 - Cross-Database Queries
 - COPY INTO statement
- Other ways also that are not in the scope of DP-700
 - Fabric Copy Job
 - Azure Data Factory
 - ODBC, JDBC etc.



Reading Data from Warehouse

- Familiar T-SQL experience
- Power BI
- SQL Server Management Studio (SSMS)
- Data Pipelines
- Dataflows Gen2
- Notebooks
- Cross-Database Queries from another Warehouse
- Many other tools, services and apps
 - Copy Job
 - Azure Data Factory
 - ODBC, JDBC etc.



Time Travel & Table Cloning



⌚ Time Travel

- View table as it existed previously
- Recover data after accidental DELETE or TRUNCATE
- Less need for manual backups or versioning
- 30-day history window
- Great for audits, comparisons, or rollbacks



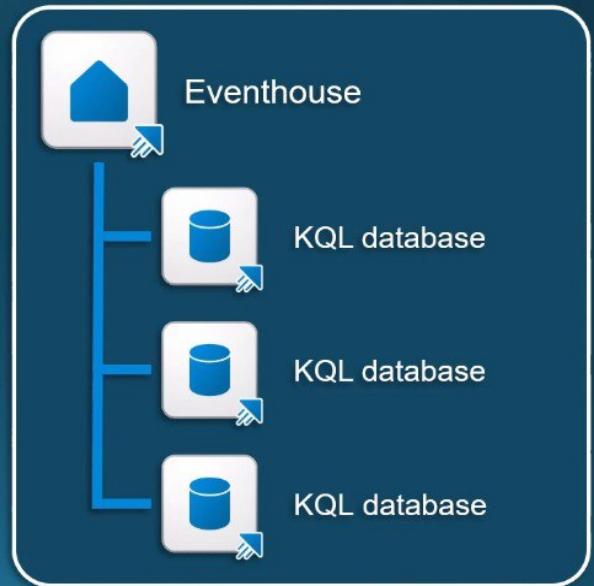
✍ Table Clones

- Instantly create a copy of a table
- Can clone current or past version (uses time travel)
- Zero-copy clone -> Only metadata is copied
- Inherits properties like RLS, dynamic data masking etc.
- Currently only within the same Warehouse

Event House

Eventhouse Overview

- Eventhouse is a data store optimized for a real-time data and analytics
- Eventhouse is a logical grouping of KQL databases
- Each KQL database is optimized for fast event ingestion and high-speed analytics
- Ideal for working with streaming data, telemetry, clickstreams, and IoT events
- Kusto Query Language (KQL) is the primary way to query KQL databases
- T-SQL can be also used
- Data can be made available in OneLake for open access across Fabric services
- Supports flexible data ingestion from many sources, including OneLake storage, Eventstream, cloud storage like Azure Blob and Amazon S3, and Fabric tools like Pipelines and Dataflows.



Eventhouse Strengths and Weaknesses

Strengths

- Optimized for real-time streaming and telemetry data ingestion
- Designed for low-latency analytics with fast query responses
- Ideal for semi-structured or event-based data (e.g., JSON logs, IoT data)
- Highly scalable to handle massive data volumes (millions of events/second)
- Supports real-time dashboards, anomaly detection, and operational monitoring



Weaknesses

- Not designed for traditional relational data modeling
- Transaction support is limited compared to relational engines
- For more shorter-term data than multi-year archives
- Limited support for complex multi-table joins
- Mastering Kusto Query Language (KQL) is necessary to unlock full capabilities

Key Objects and Concepts

• KQL Queryset

- A Fabric item that organizes multiple saved KQL queries

• Tables

- Core storage unit where data is ingested and queried

• Materialized Views

- Precomputed query results to boost performance for dashboards and reports

• Functions

- Saved reusable KQL query logic for simplifying complex queries

• Shortcuts

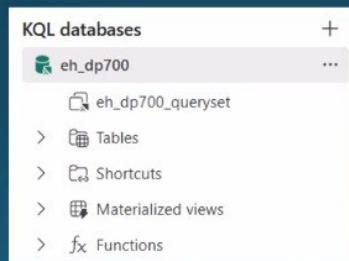
- Query external data without duplicating it inside Eventhouse

The screenshot shows the Eventhouse interface. On the left, a code editor window titled "eh_dp700" contains the following KQL query:

```
1 aleksi_test
2 | take 100
3
4 employees
5 | getschema
6
7 departments
8 | summarize LastIngestionTime = max(ingestion_time())
```

On the right, there are two configuration panels:

- Retention**: A toggle switch is set to "Unlimited". Below it is a field labeled "Retention period (days)" with a placeholder "Unlimited". A note at the bottom states "Min 1, max 36500 days".
- Caching**: A toggle switch is set to "Unlimited". Below it is a field labeled "Caching period (days)" with a placeholder "Unlimited". A note at the bottom states "Caching must be lower or equal to retention".



• Retention Policies

- Automatically manage how long data is kept before deletion

• Caching Policies

- Improve query performance by caching frequently accessed data

• OneLake Availability

- Make Eventhouse data accessible in OneLake

Materialized Views

- Precomputed query results stored as managed tables
- Automatically updated as new data is ingested
- Improve performance for dashboards and recurring queries
- Ideal for aggregation, filtering, or light transformations on streaming data
- Important properties to remember
 - backfill (boolean)
 - Include historical data
 - Defaults to False
 - effectiveDateTime (datetime)
 - Starting point for historical backfill

```
.create async materialized-view with (
    backfill = true,
    effectiveDateTime = datetime(2024-12-01T00:00:00Z),
    docString = "Daily website visits with backfill"
) e015_mv_daily_visits_with_backfill
on table e015_website_visits
{
    e015_website_visits
    | extend visit_day = bin(visit_timestamp, 1d)
    | summarize visit_count = count() by visit_day
}
```

OneLake Availability

- Default Eventhouse storage is optimized for fast querying inside Eventhouse, not automatically accessible by other Fabric items without using the KQL database engine
- OneLake availability enables Eventhouse data to be stored in OneLake as Delta Lake tables
- Can be enabled at the database level or individual table level
- Allows querying the data in your KQL database in Delta Lake format via other Fabric engines such as Direct Lake mode in Power BI, Warehouse, Lakehouse, Notebooks, and more
- The Data retention policy of your KQL database is also applied to the data in OneLake
- When OneLake availability is enabled, you cannot rename tables, alter schemas, apply Row-Level Security, or delete, truncate, or purge data

