

Using Lakehouse data at scale with Power BI. Featuring Power BI Direct Lake mode!

Benni De Jagere



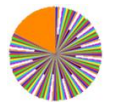
Slides





Benni De Jagere

Senior Program Manager | Fabric Customer Advisory Team (FabricCAT)



Fabric CAT

.be Member

@BenniDeJagere

/bennidejagere

/bennidejagere

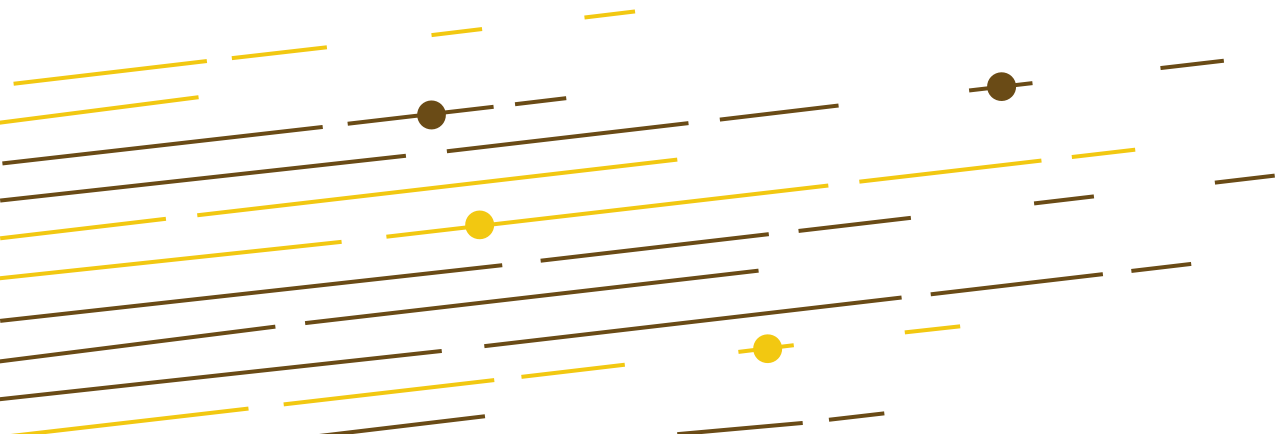
/bennidejagere

#SayNoToPieCharts



Disclaimer: We're not benchmarking

Session Objectives



Session Objectives

- Introduce Fabric and OneLake
- Set the scene for Direct Lake
- Take it for spin.. 😊

Introducing Fabric



Microsoft Fabric

The unified data platform for the era of AI



Data
Factory



Synapse Data
Engineering



Synapse Data
Science



Synapse Data
Warehousing



Synapse Real
Time Analytics



Power BI



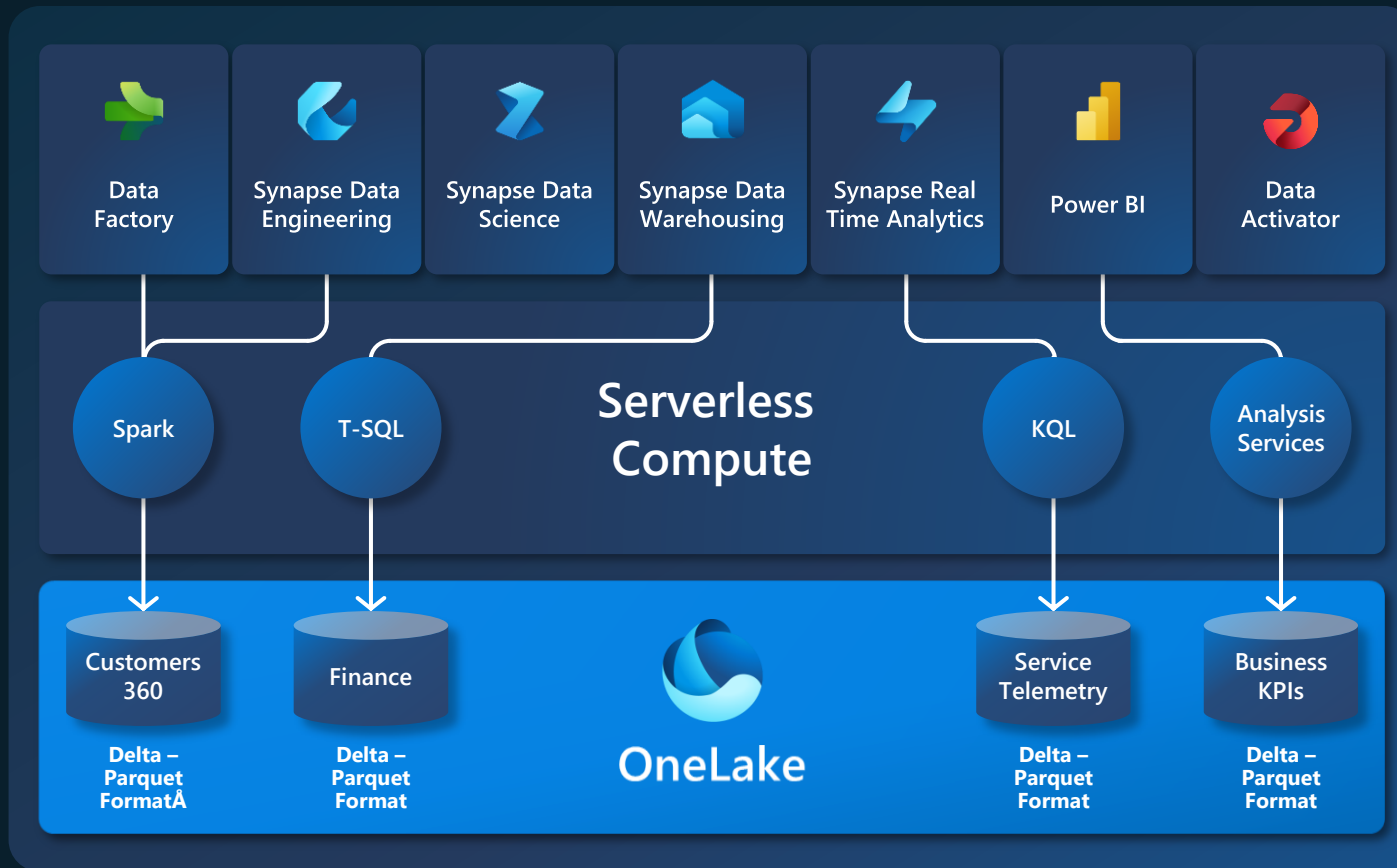
Data
Activator



OneLake

One Copy for all computes

Real separation of compute and storage



All the compute engines store their data automatically in OneLake

The data is stored in a single common format

Delta – Parquet, an open standards format, is the storage format for all tabular data in Analytics vNext

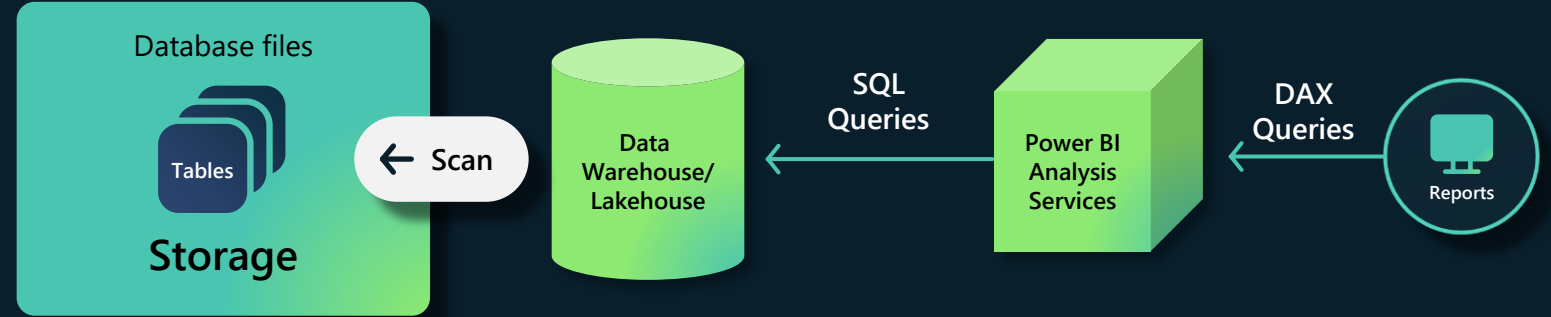
Once data is stored in the lake, it is directly accessible by all the engines without needing any import/export

All the compute engines have been fully optimized to work with Delta Parquet as their native format

Shared universal security model is enforced across all the engines

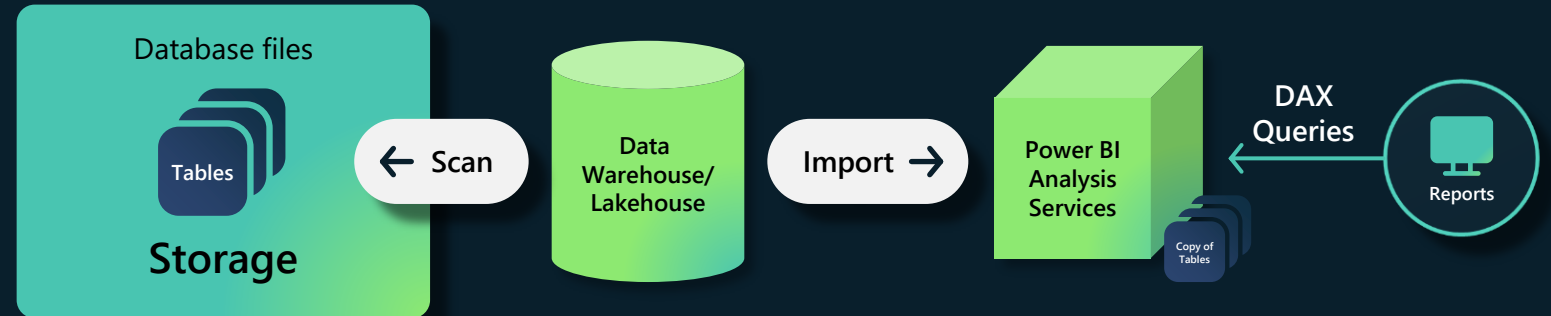
"Direct Query Mode"

Slow, but real time



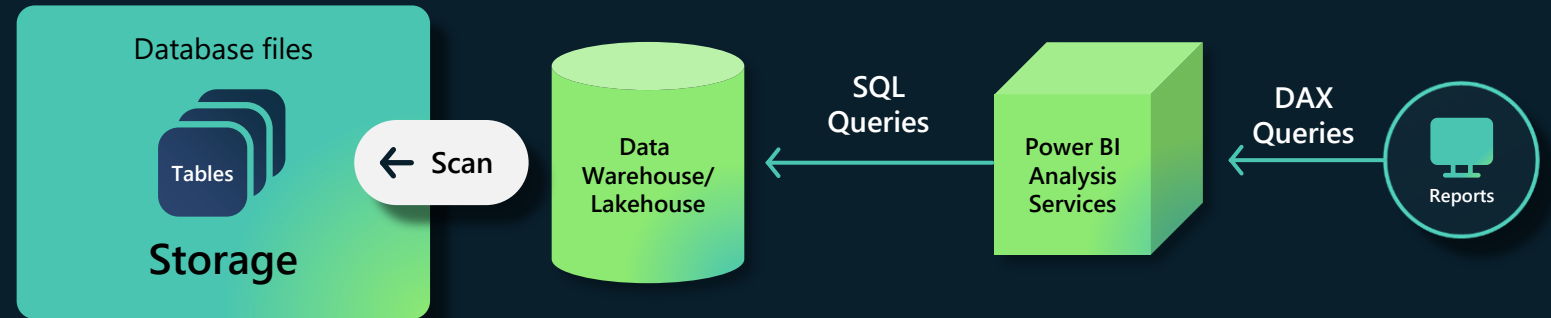
"Import Mode"

Latent & duplicative but fast



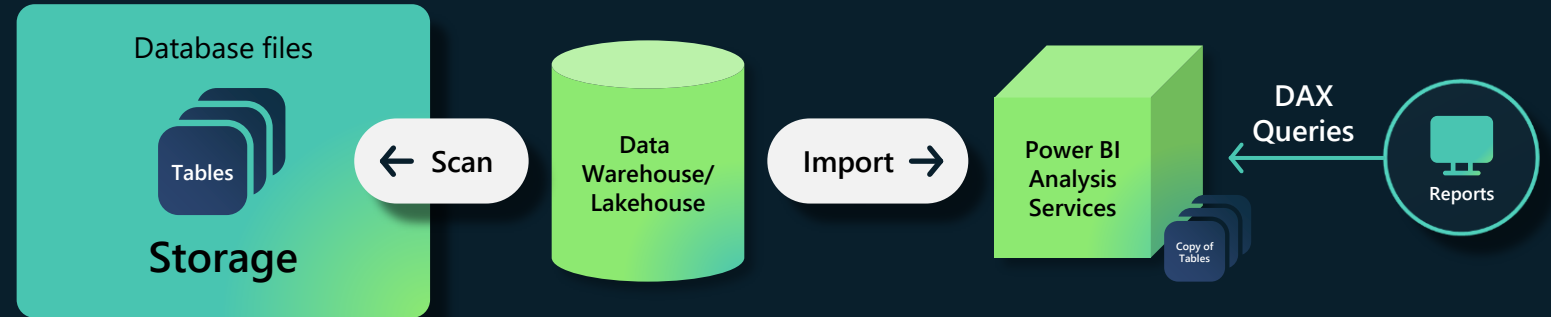
"Direct Query Mode"

Slow, but real time



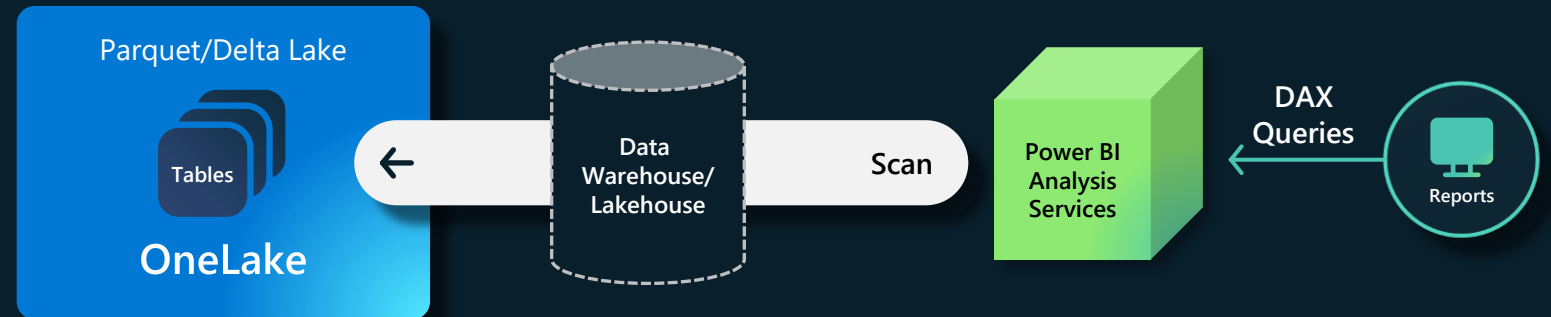
"Import Mode"

Latent & duplicative but fast



"Direct Lake Mode"

Perfect!



Why Delta?



Why Delta (Parquet)?

Open Standard for data format

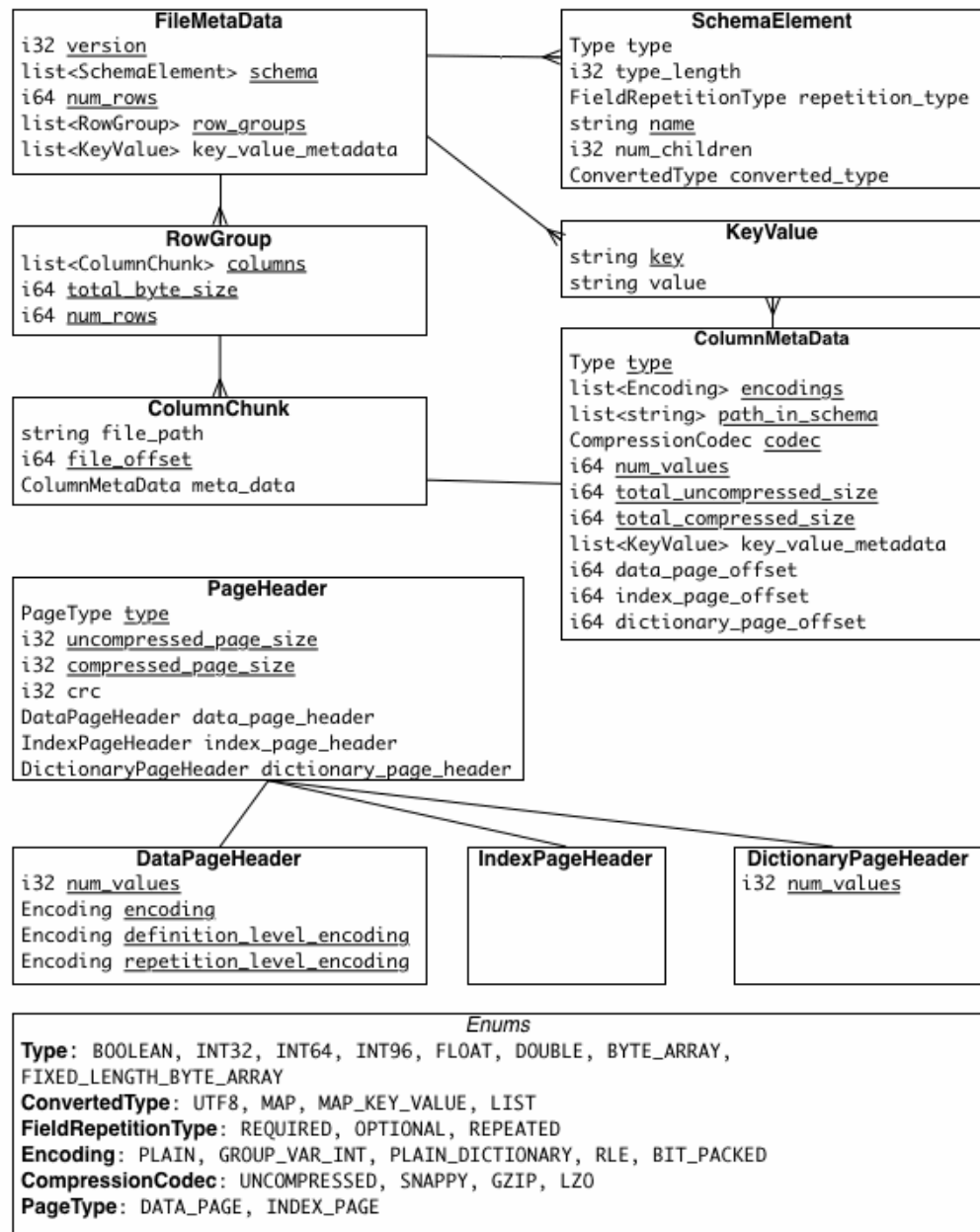
Column oriented, efficient data storage and retrieval

Efficient Data Compression and Encoding

Becoming the Industry Standard

Well suited for pruning (Column, rowgroup)

Thrives on bulk operations



Inside Delta (Parquet)

Header:

RowGroup1:

StoreID: StoreA, StoreA, StoreA

DateTime : 2023-01-01, 2023-01-02, 2023-01-03

ProductID : SKU001, SKU001, SKU001

Value: 10, 15, 12

RowGroup2:

....

Footer:

Inside Delta (Parquet) – Dictionary IDs

Header:

RowGroup1:

StoreID: 1, 1, 1

DateTime : 1, 2, 3

ProductID : 1, 1, 1

Value: 1, 2, 3

RowGroup2:

....

Footer:

Introducing V-Ordering

Write time optimization to parquet files

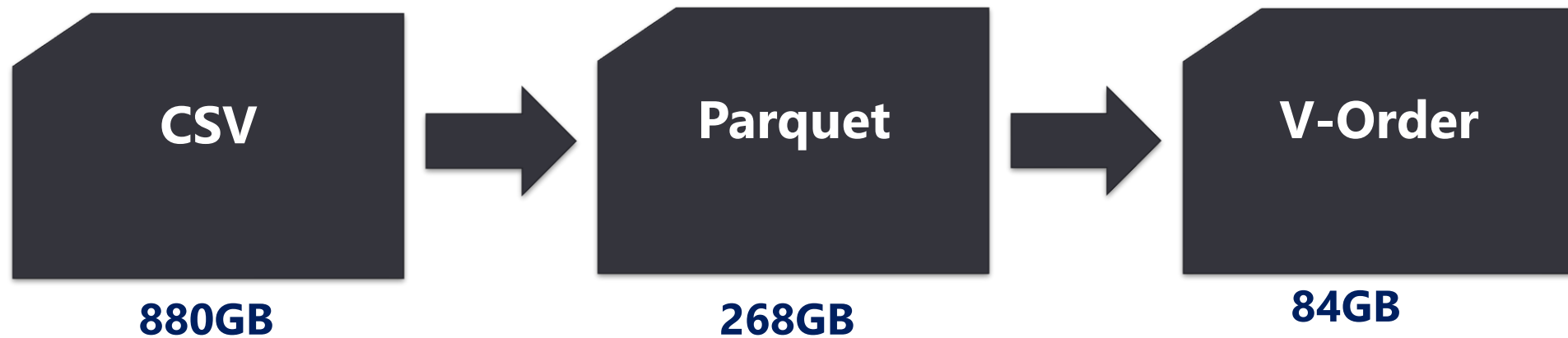
Sorting, row group distribution, dictionary encoding, and compression (Shuffling)

Complies to the open standard

Z-Order, compaction, vacuum, time travel, etc. are orthogonal to V-Order

V-ordering in action

Microsoft Internal DB (162 tables)



x3.2

Reduced IO for workloads

V-ordering in our demo case

CSV Properties

General

Sharing

Security

Previous Versions

Customise

CSV

Type: File folder

Location: C:\OneDrive\OneDrive - dataMinds vzw\Documents\IF

Size: 32.1 GB (34,544,780,233 bytes)

Size on disk: 32.1 GB (34,545,246,208 bytes)

Contains: 242 Files, 7 Folders

	TABLE_NAME	SCHEMA_NAME	Rows	TotalReservedSpaceMB	UsedDataSpaceMB	FreeUnusedSpaceMB
1	Trips_FA	Analytical	181940575	6413	6412	0
2	Time_DI	Analytical	86400	15	15	15
3	Bike_DI	Analytical	35553	1	1	1
4	Date_DI	Analytical	7304	19	19	19
5	Date_DI	Analytical	7304	0	0	0
6	Station_DI	Analytical	3430	0	0	0
7	Gender_DI	Analytical	59	0	0	0
8	Region_DI	Analytical	8	0	0	0
9	RideType_DI	Analytical	4	0	0	0
10	UserType_DI	Analytical	3	0	0	0
11	TripType_DI	Analytical	3	0	0	0
12	MemberType_DI	Analytical	3	0	0	0
13	FileType_DI	Analytical	3	0	0	0

	Name	Date modified	Type	Size
	4b7c39a4-613d-445a-9533-e4c2a08ab671.parquet	6/5/2023 4:00:58 PM	PARQUET	5.50 GB
	_delta_log	6/5/2023 4:00:58 PM	Folder	1 items

Copy data details

Copy_mns

Source

Azure SQL Database

→

Destination

Lakehouse

Data read: ⓘ

24.016 GB

Rows read:

181,940,575

Data written: ⓘ

5.909 GB

Files written: ⓘ

1

Rows written: ⓘ

181,940,575

Status

✔ Succeeded

Start time

6/5/2023, 3:22:03 PM

Pipeline run activity ID

09aa4ccc-d8ae-4e70-aec7-c76a020ddd3c

Throughput

10.321 MB/s

Total duration

00:38:53

▼ Duration breakdown

Start time

6/5/2023, 3:22:04 PM

Optimized throughput ⓘ

Standard

Used parallel copies ⓘ

1



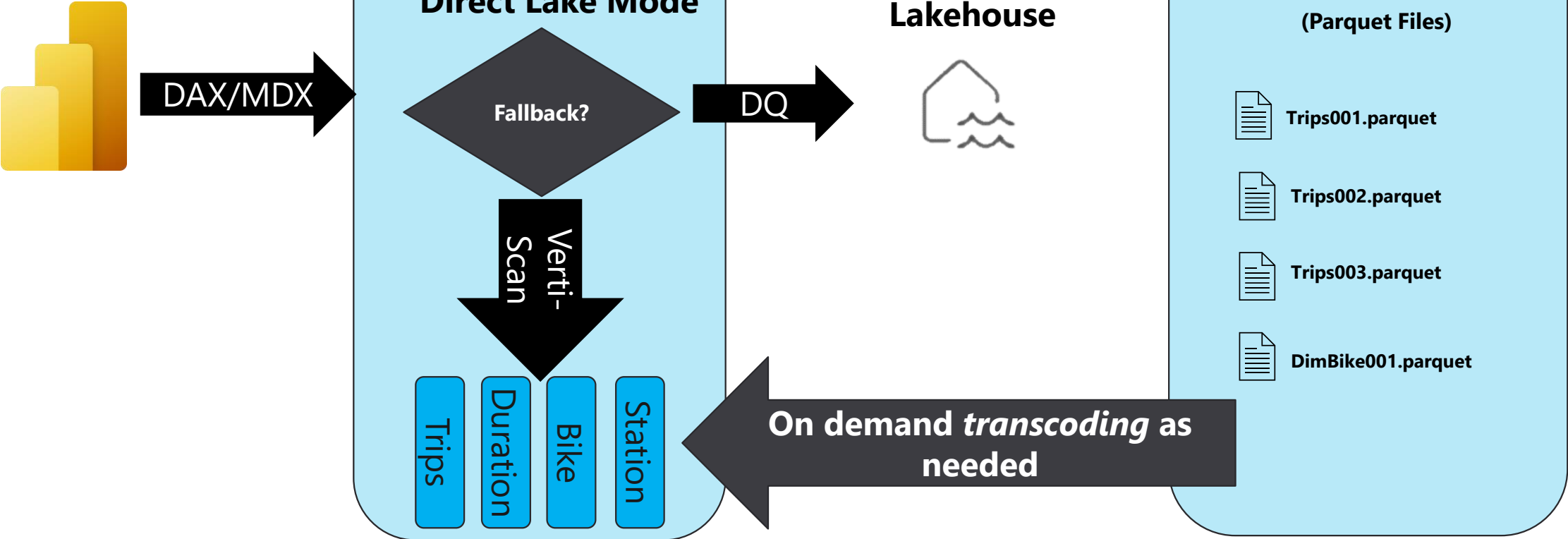
DirectLake Mode

- On start, no data is loaded in-memory
- Column data is transcoded from Parquet files when queried
- Multi-column tables can have mix of transcoded (resident) and non-resident
- Column data can get evicted over time
- DirectLake fallback to SQL Server for suitable sub-queries
- “Framing” of dataset determines what gets loaded from DeltaLake

STOP! Demo time!

Using Direct Lake mode over a Lakehouse

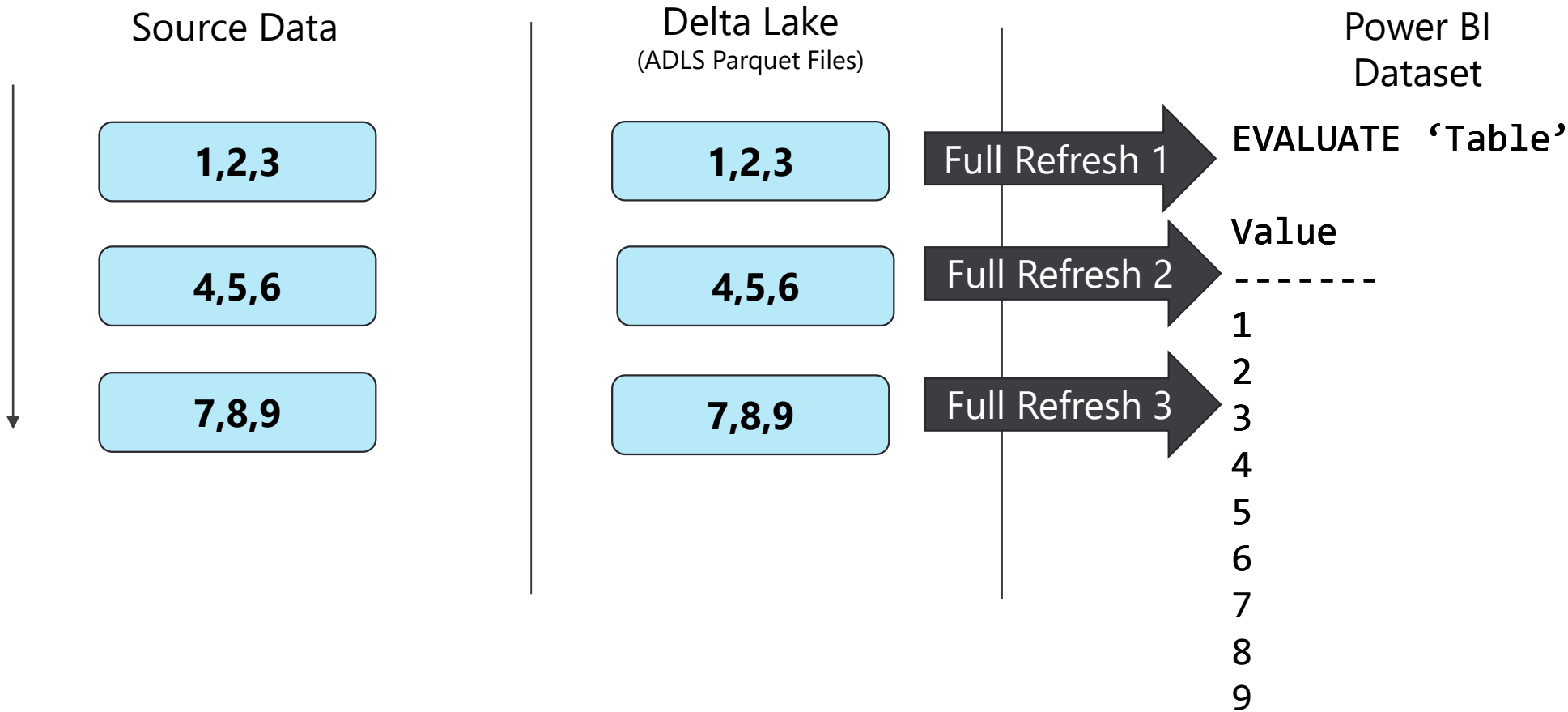
DQ Fallback



Framing

- What is framing
 - "point in time" way of tracking what data can be queried by DirectLake
- Why is this important
 - Delta-lake data is transient for many reasons
- ETL Process
 - Ingest data to delta lake tables
 - Transform as needed using preferred tool
 - When ready, perform *Framing* operation on dataset
- Framing is near instant and acts like a cursor
 - Determines the set of .parquet files to use/ignore for *transcoding* operations

Framing



STOP! Demo time!

Let's look at Framing

Optimizing Delta for Direct Lake mode

Optimizing Delta for Direct Lake mode

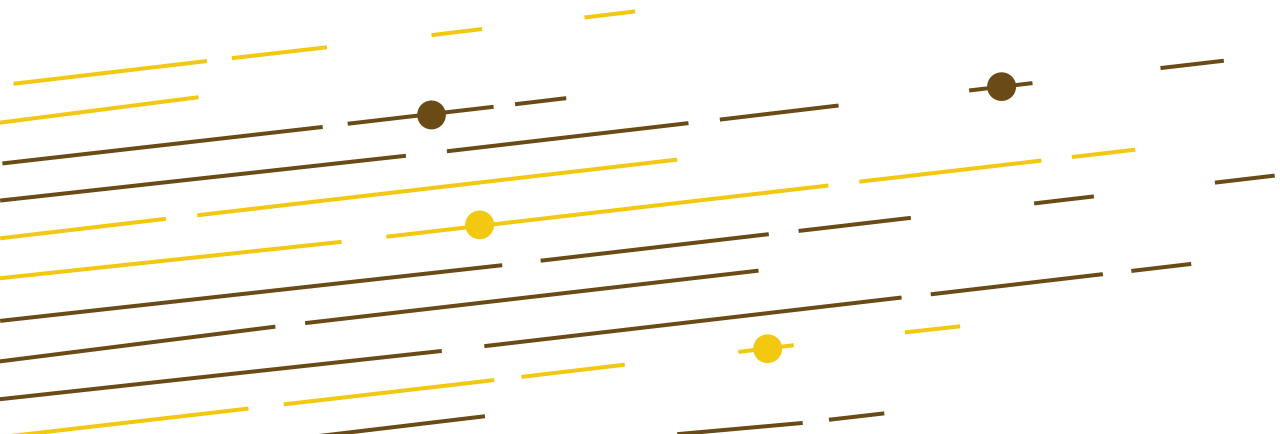
- V-Order makes a big difference, as it's tailored for Verti-Scan
- Direct Lake will work over Shortcuts to external data
 - Expect a performance impact, because reasons ..
- Direct Lake thrives on fewer, larger .parquet files
 - Physical structure will always be crucial
 - OPTIMIZE (bin-compaction) and VACUUM in the Data Engineering process will be key
 - Especially with streaming/small batch architectures, keep this in mind
- Principle of lean models will still apply
 - Only include what's needed for the reports and datasets

Common Answers to Common Questions

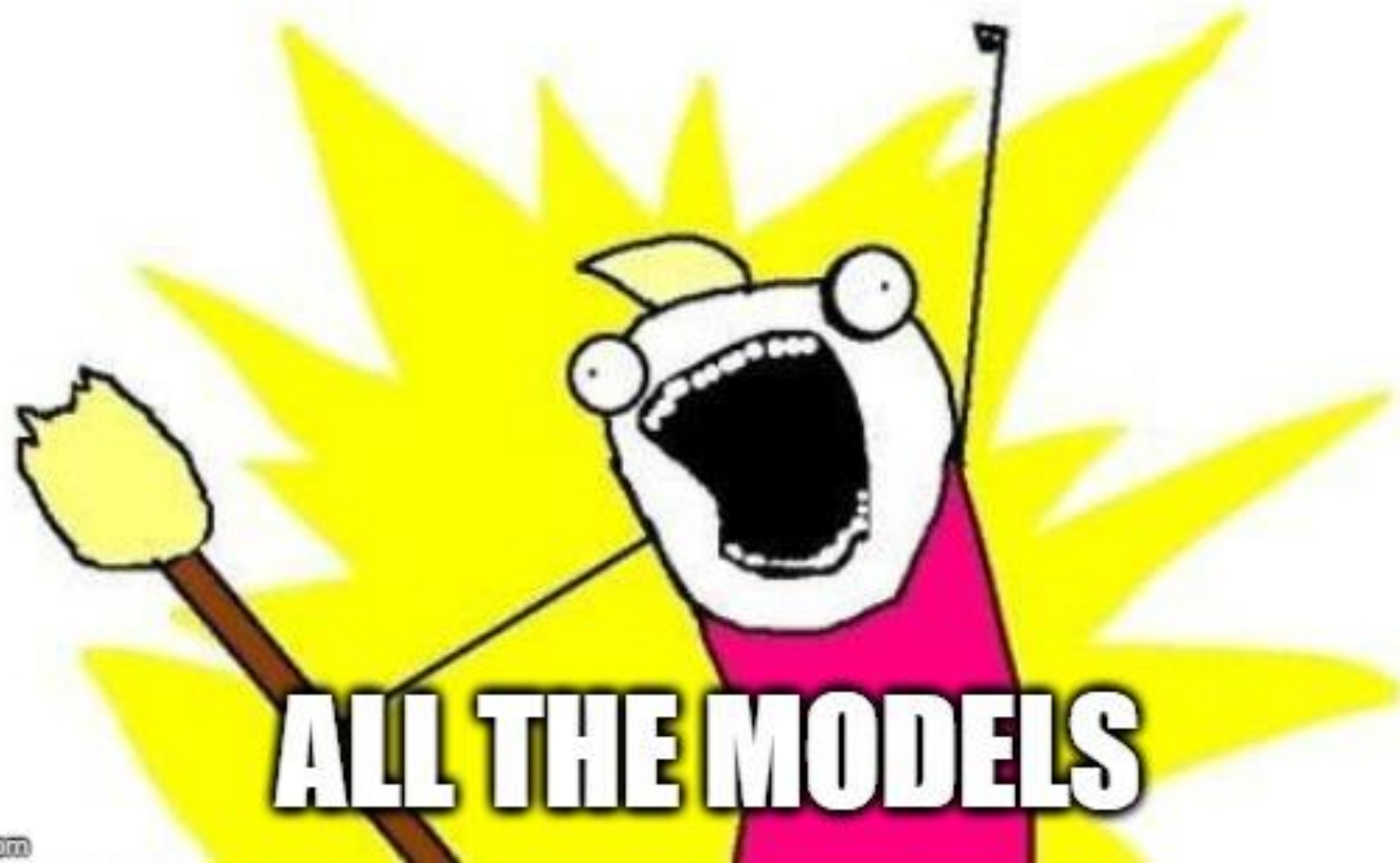
"Greatest Hits"

- Delta doesn't like spaces in object names 😊
- Delta Tables are a hard requirement for Direct Lake mode
 - Dataflows Gen2, Pipelines, Notebooks can create them for you in the lakehouse
- Web modelling is the only way to use DirectLake for now
- XMLA Read/Write is not yet supported
 - No External Tools, Calc Groups, ..
- DirectLake doesn't have unique DAX limitations
 - DQ does ..
- No confirmed plans right now to support Apache Iceberg, HUDI, ..
- No, you can't have Copilot yet

What does this mean for my data modelling?



STAR SCHEMA



Thanks, @KoVer!

Data should be transformed as far upstream as possible, and as far downstream as necessary.

Matthew Roche, 2021

(The purple haired sword aficionado)

<https://ssbipolar.com/2021/05/31/roches-maxim>

Resources

- <https://learn.microsoft.com/en-us/power-bi/enterprise/directlake-overview>
- <https://learn.microsoft.com/en-us/power-bi/enterprise/directlake-analyze-qp>
- <https://learn.microsoft.com/en-us/fabric/data-engineering/lakehouse-pbi-reporting>
- <https://learn.microsoft.com/en-us/fabric/data-engineering/delta-optimization-and-v-order?tabs=sparksql>
- <https://fabric.guru/power-bi-direct-lake-mode-frequently-asked-questions>
- <https://www.fourmoo.com/2023/05/24/using-power-bi-directlake-in-microsoft-fabric/>



Slides



https://github.com/BenniDeJagere/Presentations/{Year}/{YYYYMMDD}_{Event}





Thank you

