# Dave Ruijter

Solution Architect Data & Analytics
Blue Rocket Data Consulting & Solutions

✉ dave@blue-rocket.it

🐦 @DaveRuijter

in linkedin.com/in/DaveRuijter

🌐 ModernData.ai

**BLUE ROCKET**
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Agenda

- Challenges
  - Data Platform
  - Power BI
- Deep Dive
  - Data Platform
  - Power BI
- Hybrid Tables & Incremental Refresh
- Refresh Challenges & Orchestration
- Scaling

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# After this session

## Design and implement

Better design and implement complex data models, including hybrid tables, aggregations, and combined storage modes (import, DirectQuery , dual).

## Orchestrate

Orchestrate the end-to-end data processing, with a pipeline chain from data ingest in the data lake house to the incremental Power BI dataset refresh.

## Performance

Use different techniques to identify performance bottlenecks in your solutions and how to solve those ("does it fold"?).
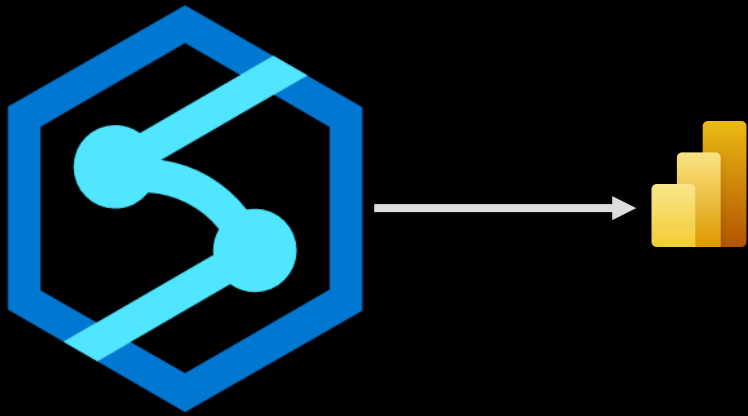
## Cost

Implement a cost-efficient solution, that still meets the scalability demands.
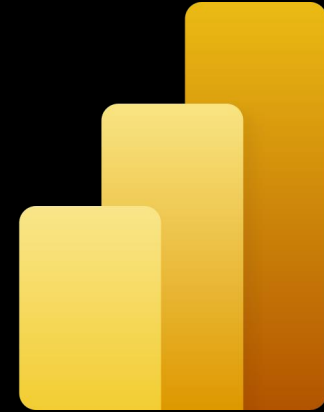
BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Solution challenges

BLUE ROCKET
DATA CONSULTING & SOLUTIONS
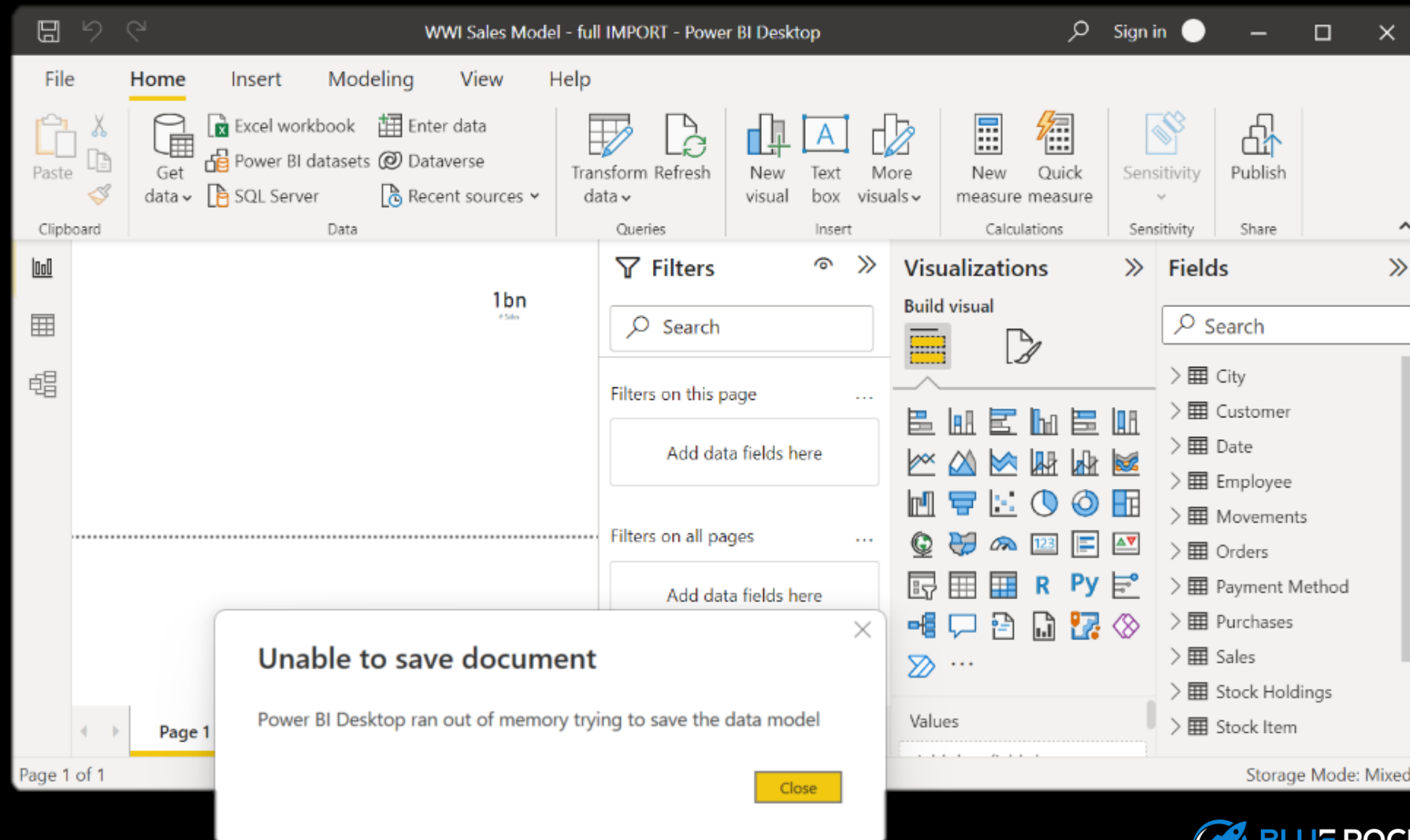
Data-Marc.com

two separate worlds

# Data Platform Solution Challenges

- Power BI report can't handle the volume of data
- Showing near real-time data in Power BI report

# Power BI report can't handle the volume

Should I put everything on DirectQuery  instead?
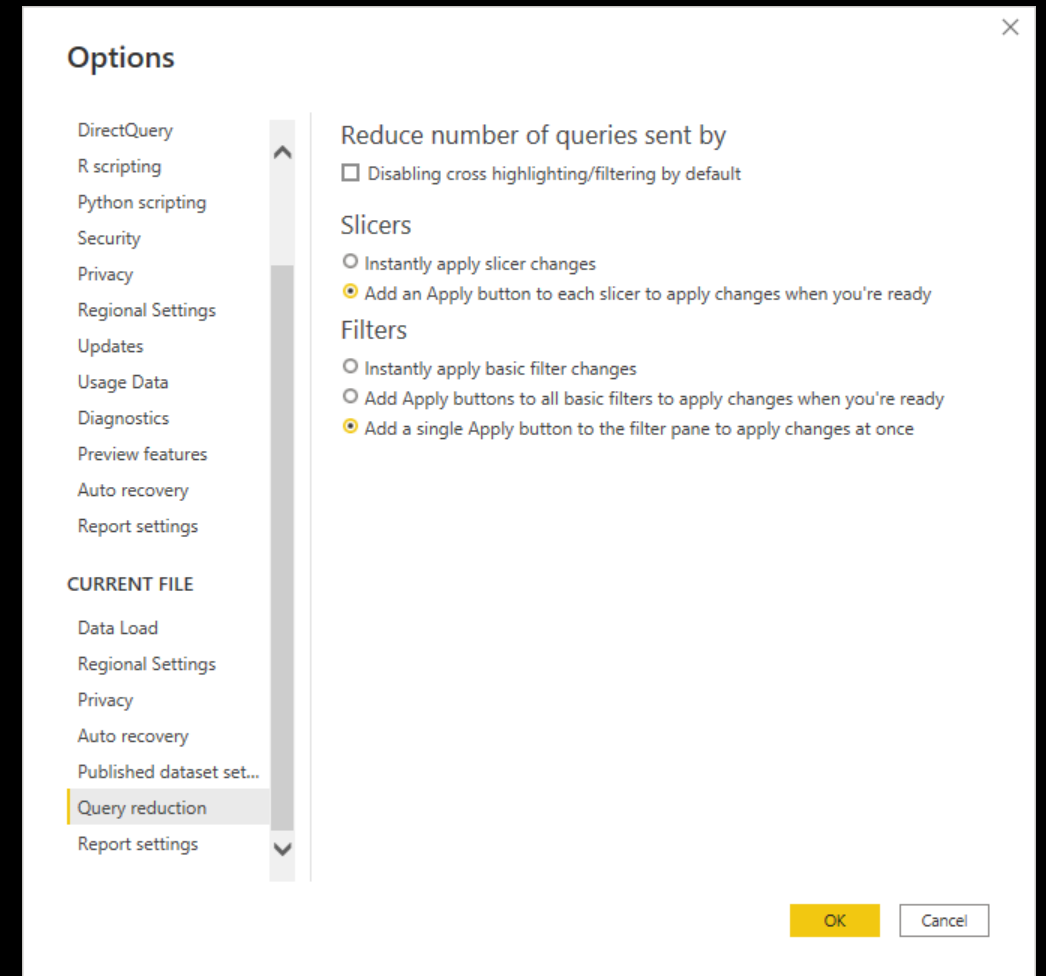
# DirectQuery limitations

- Limited Power Query capabilities

- DirectQuery != streaming / live! Front-end still caches data

- No built-in date hierarchy (year/quarter/month/day)

- Lowest granularity data is seconds (no milliseconds)

- No parent-child support in DAX with *PATH()*

- Slower end user performance

- 1M row per query

- DAX limitations, only simple calculations possible

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# DirectQuery query reduction

Consider requesting to click **Apply** before queries are executed to the source
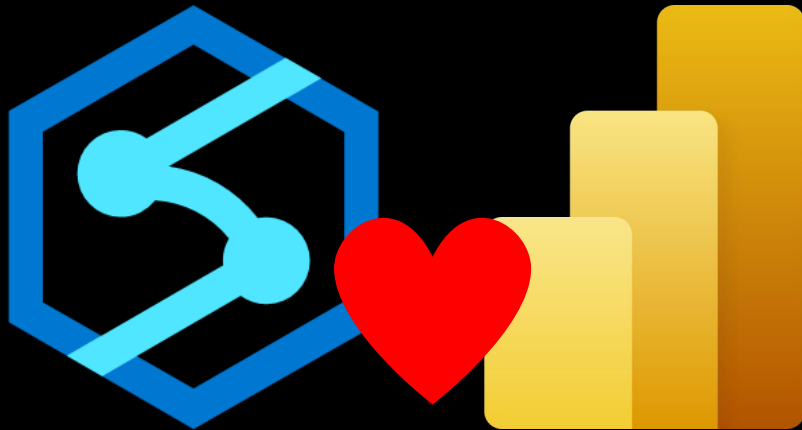
Applies to

- Slicers

- Filters (filter pane)



Options

DirectQuery
R scripting
Python scripting
Security
Privacy
Regional Settings
Updates
Usage Data
Diagnostics
Preview features
Auto recovery
Report settings

**CURRENT FILE**

Data Load
Regional Settings
Privacy
Auto recovery
Published dataset set...
Query reduction
Report settings

Reduce number of queries sent by

☐ Disabling cross highlighting/filtering by default

Slicers

○ Instantly apply slicer changes
⦿ Add an Apply button to each slicer to apply changes when you're ready

Filters

○ Instantly apply basic filter changes
○ Add Apply buttons to all basic filters to apply changes when you're ready
⦿ Add a single Apply button to the filter pane to apply changes at once

OK    Cancel

# Showing near real-time data in Power BI

- Refresh takes to long

- Poor end-user performance on DirectQuery

- Streaming datasets only allow one table

- Potentially queries are not foldable, therefore incremental refresh does not work (depending on source)

# Data Platform Solution ~~Challenges~~

- Optimize Power BI model
- Use Aggregations
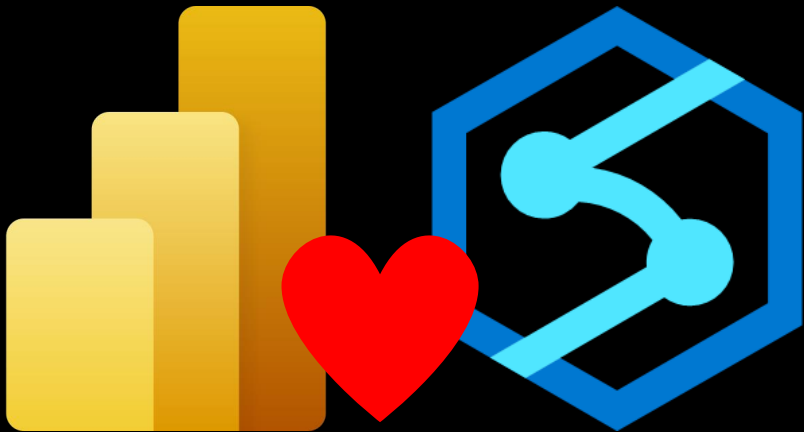- Use Hybrid Table
- End-2-end orchestration

# Power BI Solution Challenges

- Loading data from challenging sources (flat files, APIs)

- Some sources are manually maintained (like mapping tables)

- Data stored on decentralized storages, like SharePoint pages.

- Data from source systems are exported, rather than connected to analytical systems

- Store historical data in Power BI for trend analysis

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

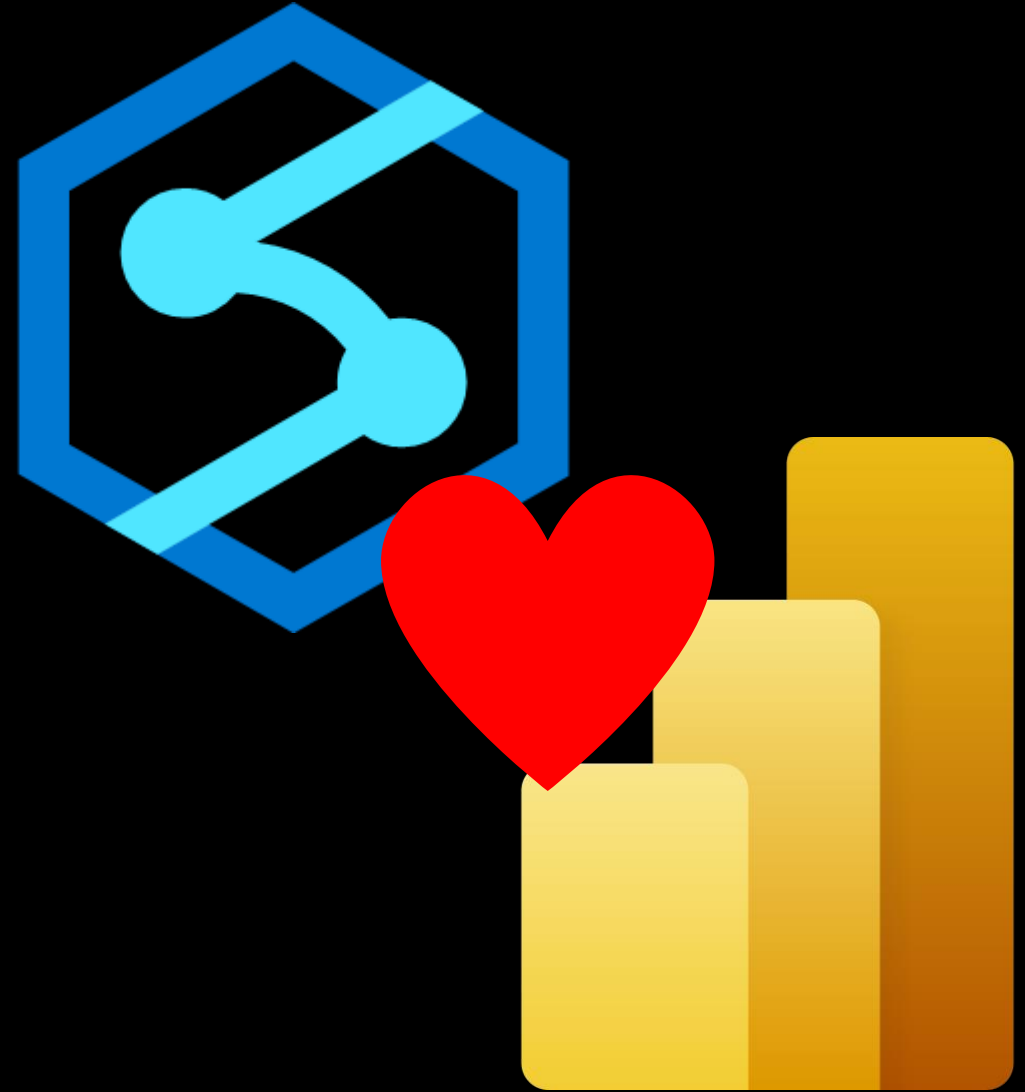Data-Marc.com

# Power BI Solution ~~Challenges~~

- ## Use Synapse Analytics to ingest data
  - ### Easily connect to various types of data sources

- ## Use Data Lakehouse: Bronze, Silver & Gold layers
  - ### Easily store historical data

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Better together

Data platform deep dive

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com
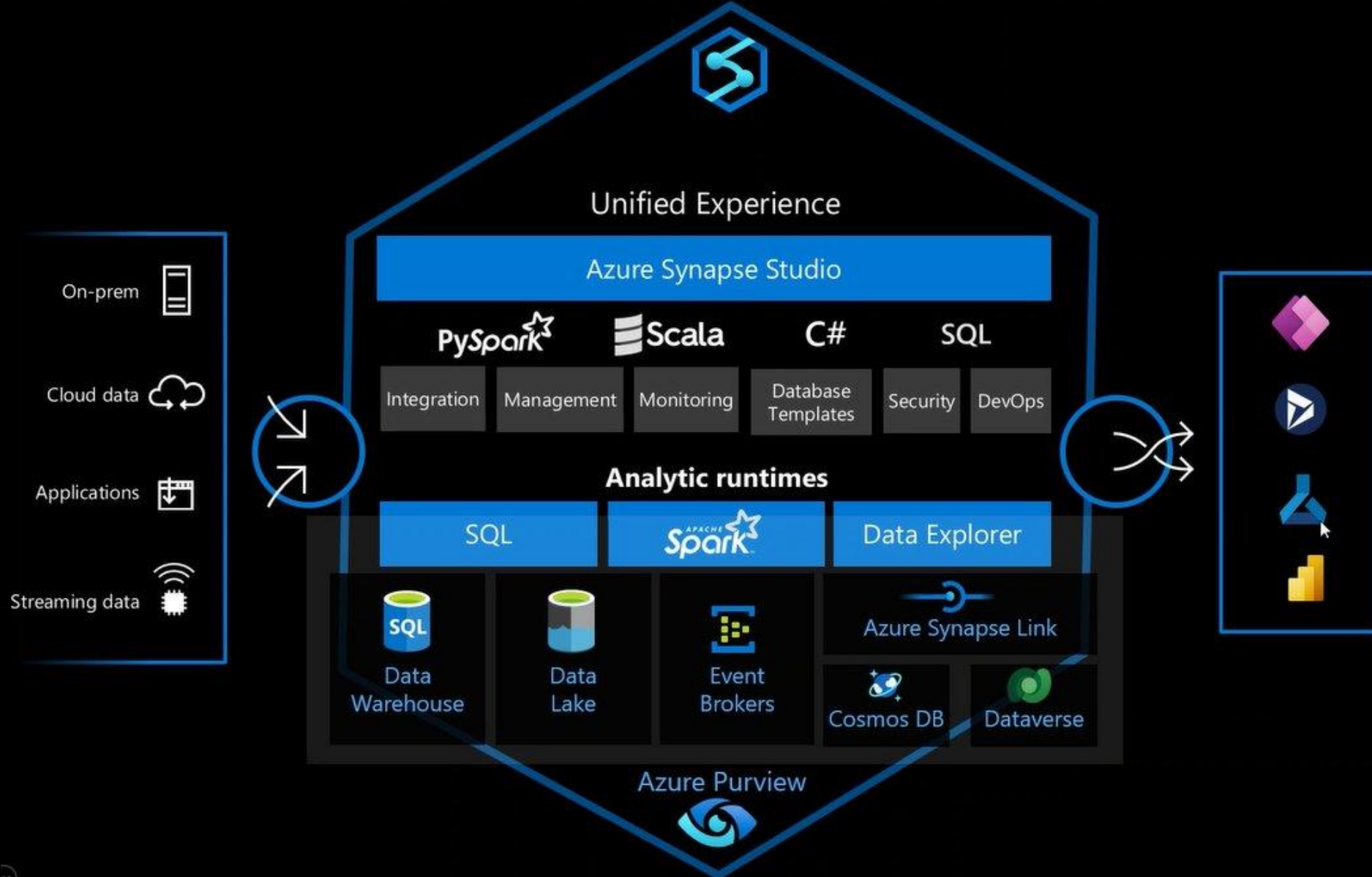
# Improvement areas

- Ingesting data from APIs using Synapse Analytics
- Store (historical!) data in the delta lakehouse architecture

# Ingesting data from APIs using Synapse Analytics

- Pipeline (same as ADF)

- Data Flow (same as ADF) *not to be confused with dataflows in PBI*

- Spark Notebook (4 languages available)

- Wrangling Dataflows (Same as ADF)

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Store (historical!) data in the lake

- The Layered approach: Bronze, Silver & Gold

- Keep original raw data, build up history -> **bronze**
- Cleanse and refine data, standard file format -> **silver**
- Aggregate, prepare, transform, merge, make start schema -> **gold**

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Store (historical!) data in the lake

- The Layered approach:
  Bronze, Silver & Gold

- Keep original raw data, build up history -> **bronze**

- Cleanse and refine data, standard file format -> **silver**

- Aggregate, prepare, transform, merge, make start schema -> **gold**

| Bronze | Silver | Gold |
|---|---|---|
| Raw data | Apply metadata | Implement business rules |
| All history, system replayable | Protect data (GDPR) | Fit for purpose |
| | Current & historical view | |

# Using Delta Lakehouse with layered approach

Batch data

Streaming data

Ingested / raw data
(Bronze)

Refined data
(Silver)
**DELTA LAKE**

Star Schemas
(Gold)
**DELTA LAKE**

SQL Serverless Views

Power BI

Azure
Data Lake Storage Gen2

**BLUE ROCKET**
DATA CONSULTING & SOLUTIONS

Data-Marc.com

Demo Data Platform
Lakehouse architecture

# Warm-up time of Serverless SQL pools

# Best practices for serverless SQL pools

- Azure AD Pass-through Authentication performance <= shared access signature credentials

- Colocate

- Same region

- Convert large CSV and JSON files to Parquet

- Try to optimize storage layout by using partitioning and keeping your files in the range between 100 MB and 10 GB

- Use appropriate data types (smallest, integer-based, varchar)

- Use filename and filepath functions to target specific partitions

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Better together

Power BI deep dive

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Improvement areas

- Data model
- Query Folding
- Aggregations
- Storage modes
- Hybrid tables

BLUE ROCKET
DATA CONSULTING & SOLUTIONS
Data-Marc.com

# But before we start changing our solution, let's measure…

- Refresh durations
- Model Size
- Vertipaq Analyzer
- Performance Analyzer
- Query folding applied?

## Perfect E2E.pbix

| Total Size | Last Data Refresh | Analysis Date |
|---|---|---|
| **74,89 MB** ⓘ | 1-3-2022 20:05:54 +01:00 | 1-3-2022 20:05:55 +01:00 |

| Compatibility | Tables | Columns | Server |
|---|---|---|---|
| 1550 | 7 | 124 | localhost:60032 |

# Performance analyzer in Power BI Desktop

| Performance analyzer | | Sync slicers |
|---|---|---|

**Performance analyzer** ✕

▶ Start recording   ↻ Refresh visuals   ◉ Stop

◇ Clear   ⬓ Export

| Name | Duration (ms) ↓ |
|---|---|
| ⊞ 0.0% | 2279 |
| ⊞ | 1440 |
| ⊞ Simple Image | 4331 |
| ⊞ Net Sales vs "What If" Analysis | 2391 |
| ⊞ OneNote | 2543 |
| ⊞ "What If" Analysis Forecast | 50 |
| ▽ Changed a slicer | - |
| ⊞ What If... | 128 |
| ⊞ Return Rate | 1028 |
| ⊞ Net Sales (Forecast) | 1618 |
| ⊞ Extra Profit | 2046 |
| ⊞ Card | 1425 |
| ⊞ "What If" Analysis Forecast | 1890 |
| ⊞ Returns | 2249 |
| ⊞ OneNote | 1155 |
| ⊞ | 1723 |
| ⊞ | 1722 |
| ⊞ | 1723 |
| ⊞ Button | 398 |
| ⊞ Last Refresh: Jun 30th, 2019 / ... | 397 |
| ⊞ $30,772 | 2107 |
| ⊞ +17.1% | 1528 |
| ⊞ | 1715 |
| ⊞ "What If" Analysis Forecast | 388 |
| ⊞ Simple Image | 2721 |

Learn more about optimizing your report's performance on our support site. Find help tuning your report from specialist Power BI partners on AppSource.

- **DAX Query**
- **Visual Display**
- **Other**
  - Preparing queries
  - Waiting for other visuals to complete
  - Other background processes

# Vertipaq analyzer

## See where your data volume is

VertiPaq Analyzer Metrics

Tables | Columns | Relationships | Partitions | Summary

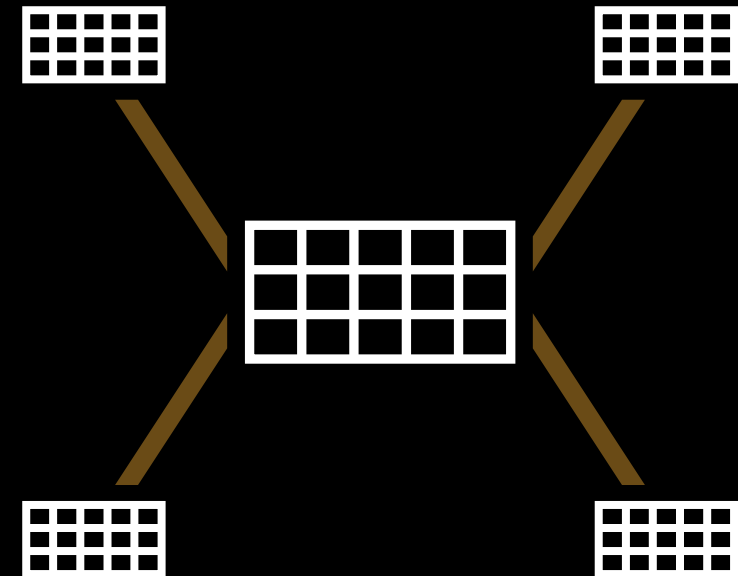| Name | Cardinality | Table Size | Col Size | Data | Dictionary | Hier Size | Encoding | Data Type |
|---|---|---|---|---|---|---|---|---|
| ▲ **Sales Agg** | **1.933.444** | **55.428.208** | **55.422.984** | **29.854.056** | **21.688.736** | **3.880.192 Many** | | **-** |
| Total Including Tax | 119.940 | 55.428.208 | 11.252.400 | 4.973.680 | 5.319.152 | 959.568 HASH | | Double |
| Tax Amount | 119.137 | 55.428.208 | 11.241.312 | 4.972.888 | 5.315.288 | 953.136 HASH | | Double |
| Total Excluding Tax | 118.952 | 55.428.208 | 11.241.120 | 4.973.680 | 5.315.776 | 951.664 HASH | | Double |
| Profit | 113.830 | 55.428.208 | 11.139.328 | 4.975.824 | 5.252.816 | 910.688 HASH | | Double |
| Invoice Date Key | 1.444 | 55.428.208 | 3.109.056 | 3.013.344 | 84.112 | 11.600 HASH | | DateTime |
| Delivery Date Key | 1.443 | 55.428.208 | 3.109.016 | 3.013.344 | 84.072 | 11.600 HASH | | DateTime |
| Count Rows | 9.722 | 55.428.208 | 2.356.128 | 1.974.104 | 304.200 | 77.824 HASH | | Int64 |
| Customer Key | 403 | 55.428.208 | 1.929.700 | 1.916.016 | 10.420 | 3.264 HASH | | Int64 |
| Salesperson Key | 101 | 55.428.208 | 44.804 | 41.176 | 2.780 | 848 HASH | | Int64 |
| RowNumber-2662979B-1795-4F74-8F37-6A1BA8059B61 | 0 | 55.428.208 | 120 | 0 | 120 | 0 VALUE | | Int64 |
| ▷ **Date** | **3.287** | **23.045.462** | **22.981.430** | **73.096** | **22.805.342** | **102.992 Many** | | **-** |
| ▷ **Customer** | **403** | **6.436.260** | **6.436.260** | **2.448** | **6.419.380** | **14.432 Many** | | **-** |
| ▷ **Employee** | **213** | **2.157.356** | **2.157.356** | **1.120** | **2.151.460** | **4.776 Many** | | **-** |
| ▷ **StockItem** | **0** | **8.576** | **8.576** | **160** | **8.416** | **0 HASH** | | **-** |
| ▷ **Sales** | **0** | **8.392** | **8.352** | **176** | **8.176** | **0 HASH** | | **-** |
| ▷ **sales_model Employee** | **0** | **8.352** | **8.352** | **176** | **8.176** | **0 HASH** | | **-** |
| ▷ **City** | **0** | **6.496** | **6.496** | **112** | **6.384** | **0 HASH** | | **-** |

# Star schema all the things!

**Facts**

▪ Contains <span style="color:orange">numerical information</span> about a business process or items to be aggregated

▪ Aggregations provide totals, averages, etc.
Power BI implements these using <span style="color:orange">Measures</span>

▪ Usefulness limited without context
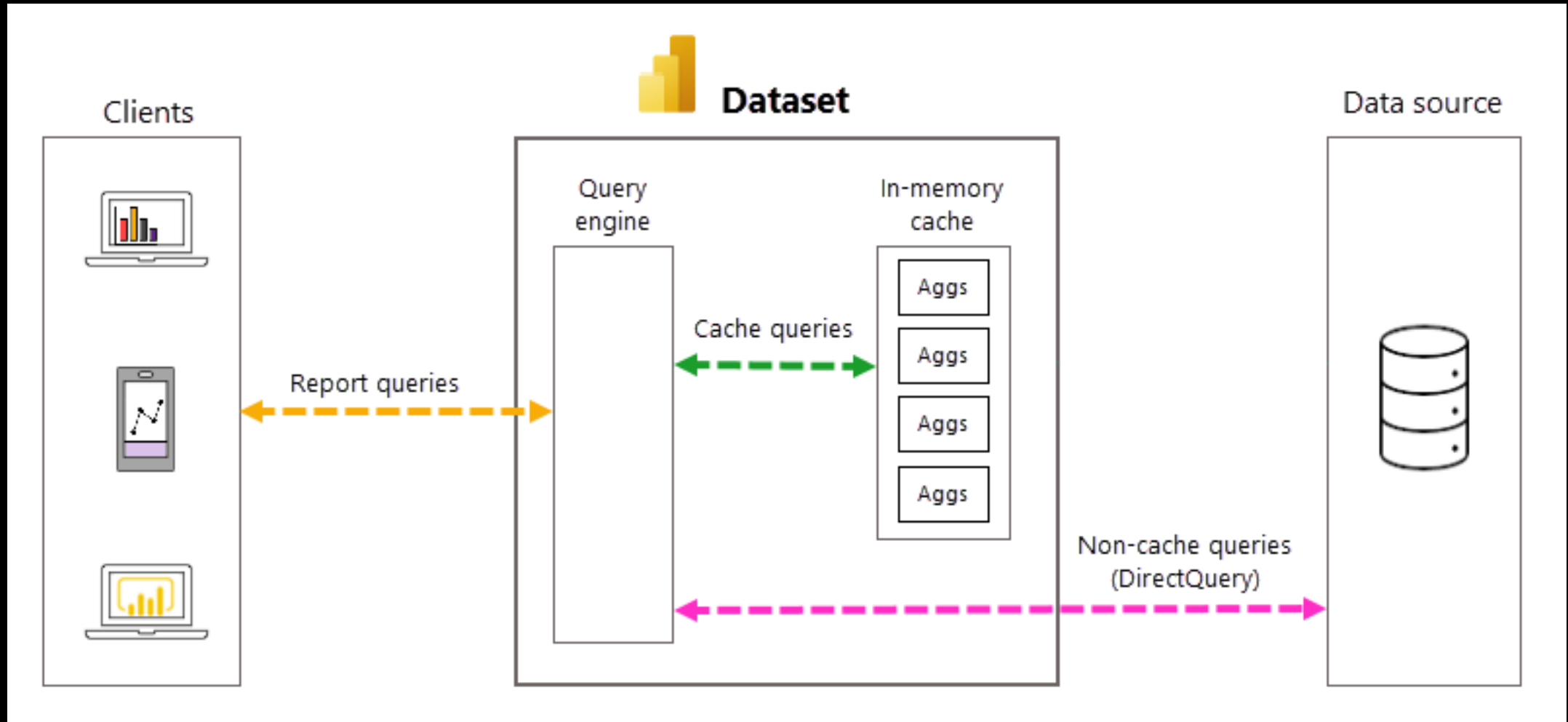Context is provided by <span style="color:orange">dimensions</span> that slice the data

**Dimensions**

▪ Contains <span style="color:orange">descriptive information</span> that define how a fact should roll up.

▪ Examples: Date, Month, Customer, Geography, Product, Payment type.

▪ Without dimensions there is no context.

▪ Also called: Lookup table on steroids.

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Aggregations

# Relationships and storage modes

**A model relationship is _limited_ when there's no guaranteed "one" side. It can be the case for three reasons:**

- The relationship uses a Many-to-many cardinality type (even if one or both columns contain unique values)

- The storage mode combination is Import and DirectQuery

- The relationship is cross source group

# Automatic Aggregations

**Power BI Premium per User, Premium Capacity and Embedded datasets**

**Automatic aggregations based on Query logs (7 days)**

**Supported sources:**

- Azure SQL Database

- Azure Synapse Dedicated SQL pool

- Google BigQuery

- Snowflake

Demo Aggregations

# Hybrid tables & incremental refresh

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Incremental refresh

- Incremental refresh is supported for Power BI Premium, Premium per user, **Power BI Pro**, and Power BI Embedded datasets.

- Getting the latest data in **real time** with DirectQuery is **only supported for Power BI Premium**, Premium per user, and Power BI Embedded datasets.

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Incremental refresh config

# Hybrid tables

- Live / Realtime data in Power BI

- Combines different storage modes on partition level in a single table

- Goes hand-in-hand with Incremental Refresh

| Granularity | Name | Row Count | |
|---|---|---|---|
| Year | 2011 | 295,489,717 | |
| Year | 2012 | 297,678,498 | |
| Year | 2013 | 295,575,442 | |
| Year | 2014 | 292,477,875 | |
| Year | 2015 | 297,780,469 | |
| Year | 2016 | 294,060,081 | Archived: **Import** |
| Year | 2017 | 300,419,682 | |
| Year | 2018 | 296,541,108 | |
| Year | 2019 | 292,787,420 | |
| Year | 2020 | 299,273,979 | |
| Quarter | 2021Q1 | 74,135,277 | |
| Month | 2021Q104 | 24,939,498 | |
| Day | 2021Q10501 | 820,805 | |
| Day | 2021Q10502 | 826,885 | Incremental refresh: **Import** |
| Day | 2021Q10503 | 821,043 | |
| Day-DirectQuery | 2021Q10504-DQ | 271,110 | Real time: **DirectQuery** |
| **Total** | | **3,063,898,887** | |

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Hybrid tables

- Implementation with Incremental Refresh

- Customizable via 3rd party tooling like Tabular Editor

>> Limitation: Only 1 DQ partition per table allowed at the moment.



Incremental refresh and real-time data

ⓘ These settings will apply when you publish the dataset to the Power BI service. Once you do that, you won't be able to download it back to Power BI Desktop. Learn more

**1. Select table**

Sales Agg

**2. Set import and refresh ranges**

🔘 Incrementally refresh this table

Archive data starting | 10 | Years | before refresh date

Data imported from 3/2/2012 to 2/20/2022.

Incrementally refresh data starting | 10 | Days | before refresh date

Data will be incrementally refreshed from 2/20/2022 to 3/2/2022.

**3. Choose optional settings**

☑ Get the latest data in real time with DirectQuery (Premium only) Learn more

☑ Only refresh complete days Learn more

☐ Detect data changes Learn more

**4. Review and apply**

| Archived | Incremental Refresh | Real-time |

10 years before refresh date · 10 days before refresh date · Refresh date
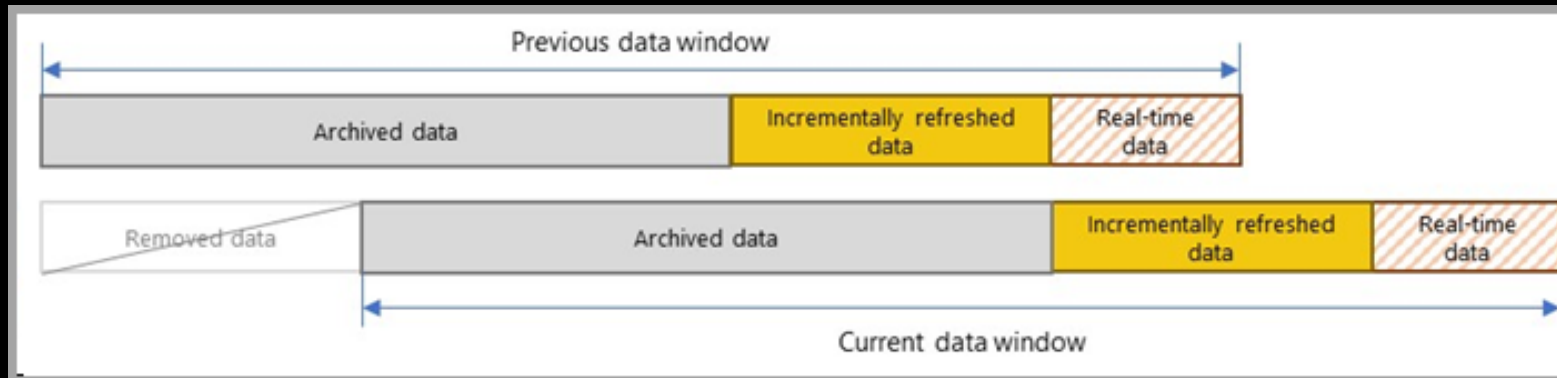
Apply | Cancel

# Hybrid tables – what challenge does it solve?

- Realtime scenarios without full tables on DQ mode

- No complex refresh mechanisms needed with partition refresh and queries over XMLA

- No more multiple tables and complex DAX to combine to achieve the same goal

# Hybrid tables – Keep in mind that…

- Premium feature

- DAX restrictions for DirectQuery apply

- Limited Power Query capabilities (due to DQ)

- Requires Large Dataset Format (storage) in workspace

- Performance hit on upstream data sources

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

Demo Hybrid Tables
Latest data in real time

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Refresh challenges

# Refresh options

- Scheduled in the service

- Manual trigger

- Power Automate

- PowerShell

- API

# Effective refreshing

**Considerations**

- Refreshing the entire model takes too long with high load on sources

- Can we only refresh certain tables?

- Can we only refresh certain partitions?

- Can we use DQ tables/partitions (Hybrid Tables)

**What do we need?**

- Async refresh API

- XMLA Endpoints

# Async refresh API

- Specify the objects to refresh

POST
*https://api.powerbi.com/v1.0/myorg/groups/f089354e-8366-4e18-aea3-4cb4a3a50b48/datasets/cfafbeb1-8037-4d0c-896e-a46fb27ff229/refreshes*

```
{
    "type": "Full",
    "commitMode": "transactional",
    "maxParallelism": 2,
    "retryCount": 2,
    "objects": [
        {
            "table": "DimCustomer",
            "partition": "DimCustomer"
        },
        {
            "table": "DimDate"
        }
    ]
}
```

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Wait a sec...

## Enhanced refresh with the Power BI REST API is now generally available

**Ogbemi Ekwejunor-Etchie**
Program Manager

June 2, 2022

in Share   Tweet   f Like

We're excited to move Enhanced refresh with the Power BI REST API (formerly asynchronous refresh) ...
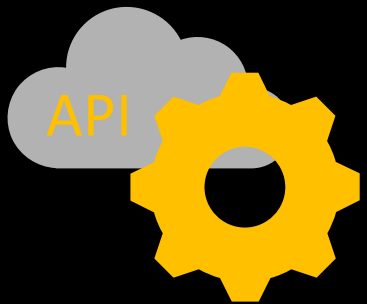
examples.

- Cancel an in-progress refresh operation.
- Check the status of historical, current, and pending refreshes.
- Refresh individual tables and individual partitions.

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Central E2E orchestration

- Combine pipelines from Data Platform with Power BI
- Lowest latency between source and report
- Consider including backup operations for Power BI
- Incremental loading where possible



BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

Demo end-2-end orchestration

# Scaling

# Scaling data platform

- Spark Cluster:
  - Use multiple cluster configs
  - Autoscale -> *It can take 1 to 5 minutes for a scaling operation to complete*
  - Dynamic allocation of executors
  - Automatic pause

| Size | vCore | Memory |
|------|-------|--------|
| Small | 4 | 32 GB |
| Medium | 8 | 64 GB |
| Large | 16 | 128 GB |
| XLarge | 32 | 256 GB |
| XXLarge | 64 | 512 GB |
| XXX Large (Isolated Compute) | 80 | 504 GB |

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Auto-scale (gen2 only)

Auto-scale adds:

- Additional vCores

- Applies for at least 24h

Configured through

- Max. number of scalable vCores

- Azure subscription – Pay as you go

# Power BI Premium – Auto scale

# Power BI Premium – Auto scale

# Wrap-up

| Design and implement | Orchestrate | Performance | Cost |
|---|---|---|---|
| Better design and implement complex data models, including hybrid tables, aggregations, and combined storage modes (import, DirectQuery , dual). | Orchestrate the end-to-end data processing, with a pipeline chain from data ingest in the data lake house to the incremental Power BI dataset refresh. | Use different techniques to identify performance bottlenecks in your solutions and how to solve those ("does it fold"?). | Implement a cost-efficient solution, that still meets the scalability demands. |

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com

# Thanks for attending!

Please fill in the evaluation

Evals.datagrillen.com

BLUE ROCKET
DATA CONSULTING & SOLUTIONS

Data-Marc.com