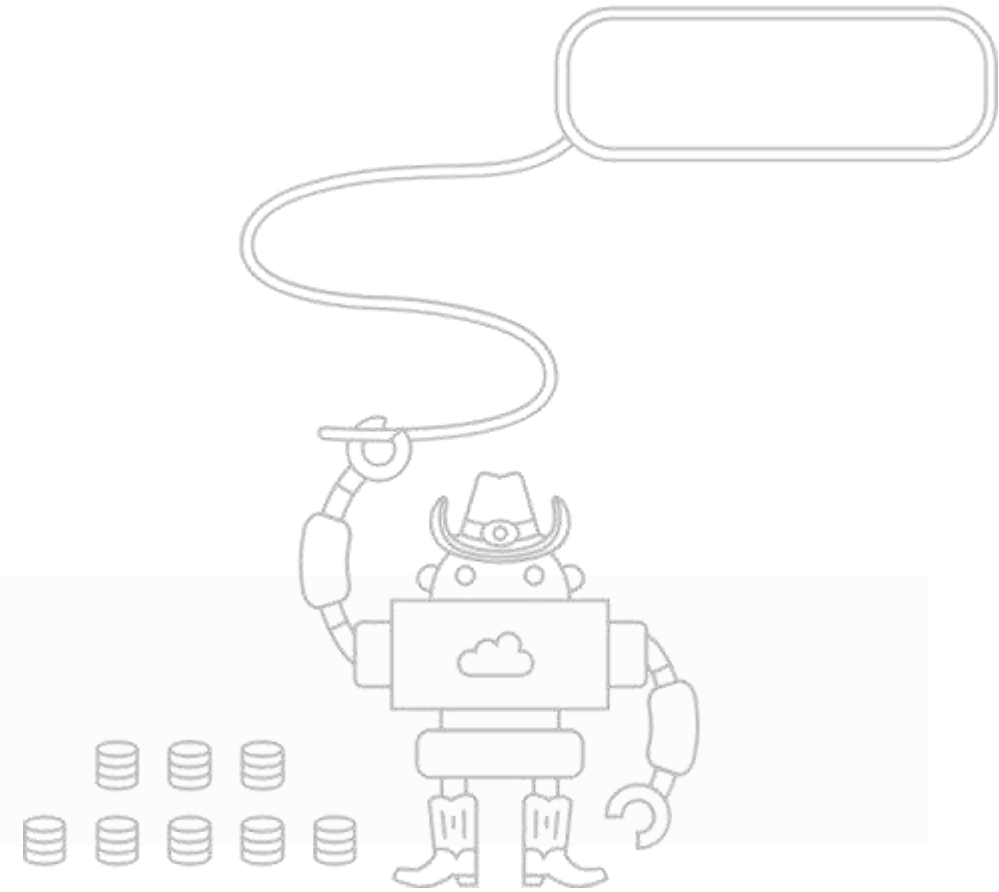


ADF WRANGLING DATA FLOWS

POWER QUERY GOES SPARK

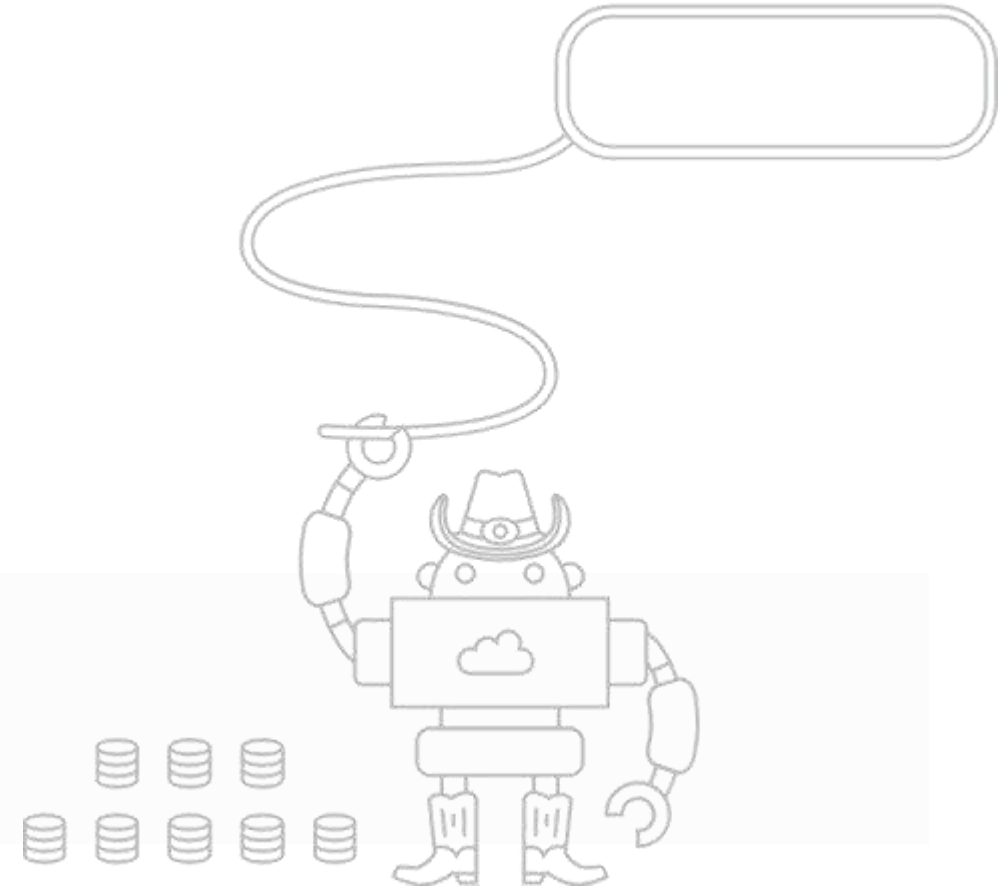


Christoph Seck | KI Group

c.seck@kigroup.de

~~ADF~~ WRANGLING DATA FLOWS

POWER QUERY GOES ~~SPARK~~ DATA LAKE



Christoph Seck | KI Group

c.seck@kigroup.de

>280
EMPLOYEES
> 150 ENGINEERS

FOUNDED
1999

>300
PROJECTS

5 LOCATIONS

COLOGNE

STUTT GART

BERLIN



MUNICH

LISBON



KI group

HOME FOR ENTREPRENEURS,
CREATORS AND SOLVERS

PASSION FOR
TECH AND
INNOVATIONS



BUSINESS MODELS & PROCESSES

KI mobility **KI** finance
KI challenge.rs **KI** growth
KI chemicals



EXPERIENCES , DESIGN AND LIFESTYLE



SOFTWARE & DATA



KI performance
KI decentralized
KI analytics
KI labs

MOBILAB

KI INTRODUCING group

HOME FOR ENTREPRENEURS,
CREATORS AND SOLVERS

INVESTMENTS, COMPANYBUILDING AND ECOSYSTEM

KI capital



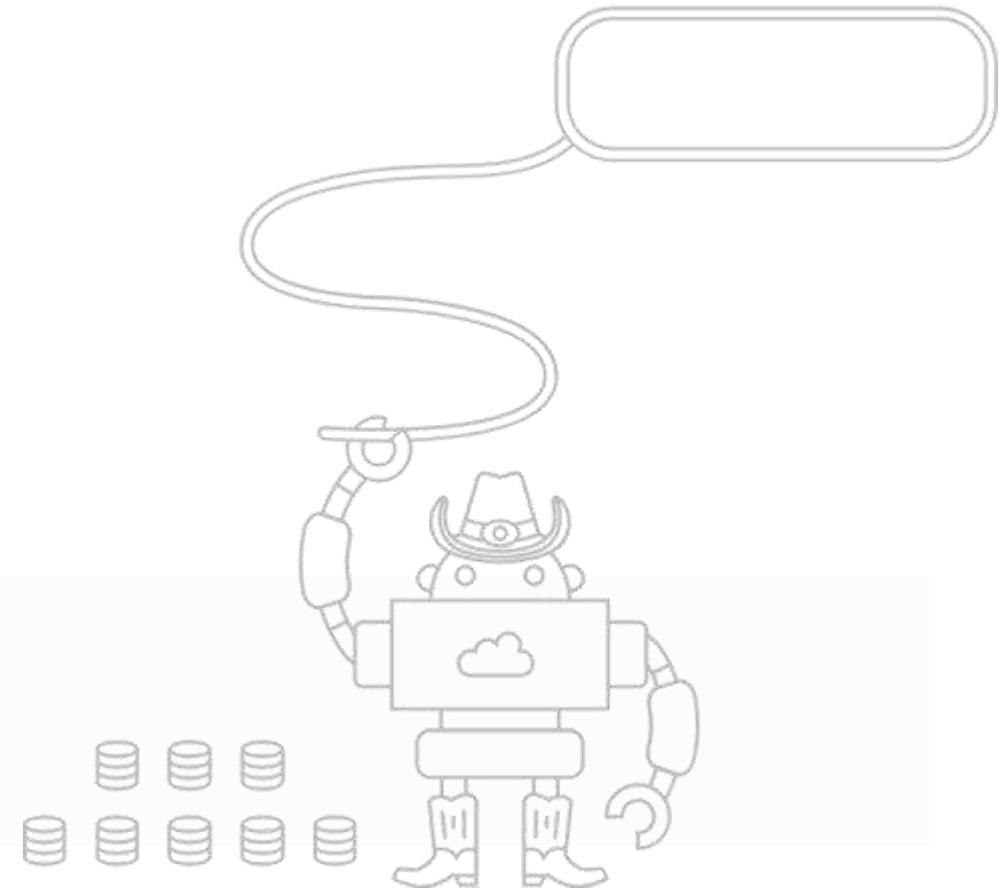
HUMAN RESOURCES

KI professionals **KI** connect



WRANGLING DATA FLOWS

POWER QUERY GOES DATA LAKE

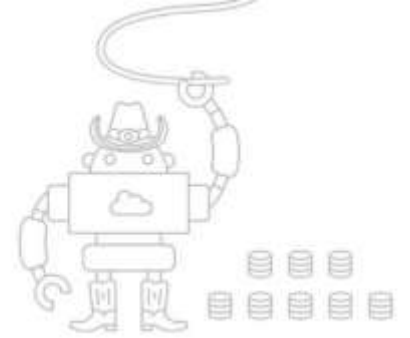


Christoph Seck | KI Group

c.seck@kigroup.de



Agenda



Wrangling in Context



Refresher Power Query & M



ADF Data Wrangling



The swotty Power Cousin



Resumé



Agenda



Wrangling in Context



Refresher Power Query & M



ADF Data Wrangling



The swotty Power Cousin



Resumé

Modern Data Warehouse (MDW)



On-premises data

Oracle, SQL,, Teradata, fileshares, SAP



Cloud data

Azure, AWS, GCP



SaaS data

Salesforce, Dynamics

INGEST



Azure Data Factory

PREPARE



Azure Data Factory



Azure Databricks

TRANSFORM & ENRICH



Azure Data Factory



Azure Databricks

SERVE



Azure SQL Data Warehouse

VISUALIZE



Power BI

STORE

Azure Data Lake Storage Gen2



Data Pipeline Orchestration & Monitoring

Azure Data Factory



Modern Data Warehouse (MDW)



On-premises data

Oracle, SQL,, Teradata, fileshares, SAP



Cloud data

Azure, AWS, GCP



SaaS data

Salesforce, Dynamics

INGEST



Azure Data Factory

PREPARE



Azure Data Factory



Azure Databricks

TRANSFORM & ENRICH



Azure Data Factory



Azure Databricks

SERVE



Azure SQL Data Warehouse

VISUALIZE



Power BI

STORE

Azure Data Lake Storage Gen2



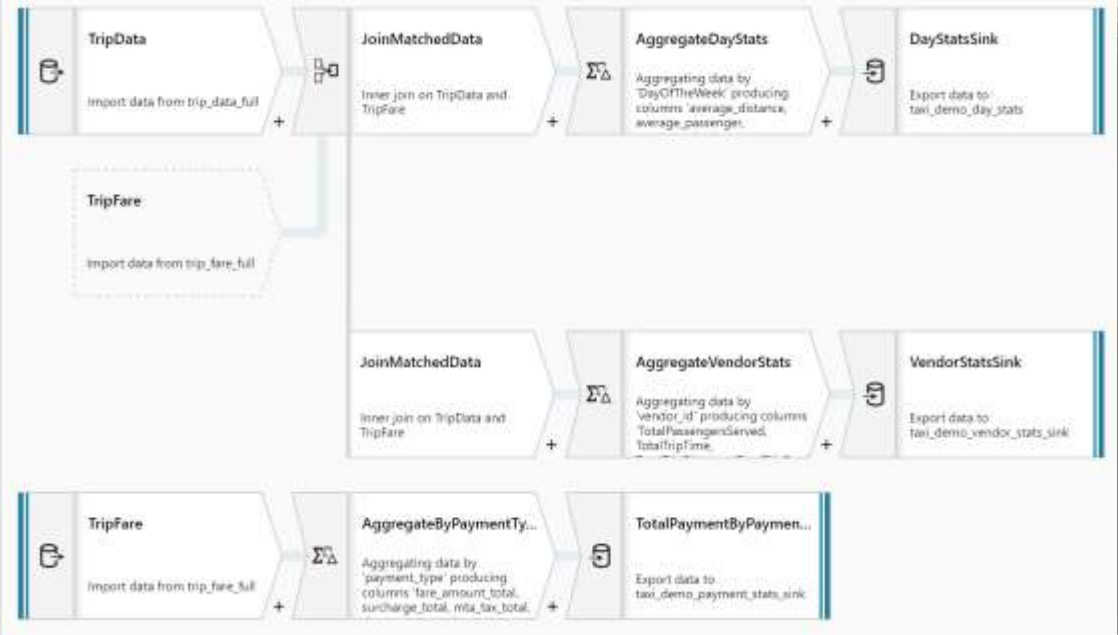
Data Pipeline Orchestration & Monitoring

Azure Data Factory



MAPPING DATAFLOW

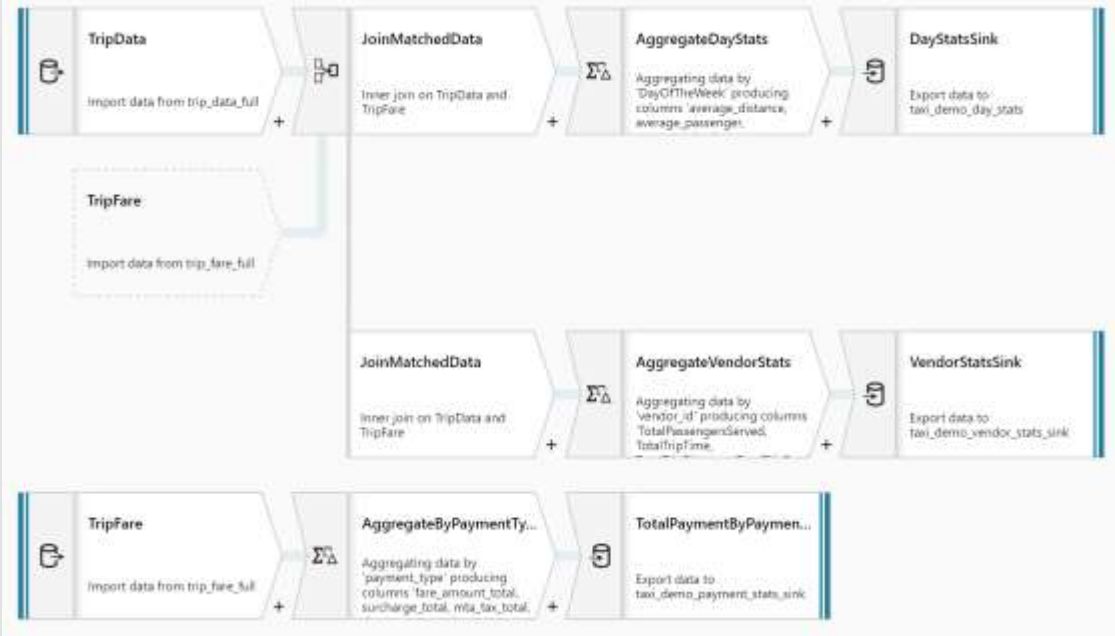
Code-free data transformation @scale



PUBLIC
PREVIEW

MAPPING DATAFLOW

Code-free data transformation @scale



WRANGLING DATAFLOW

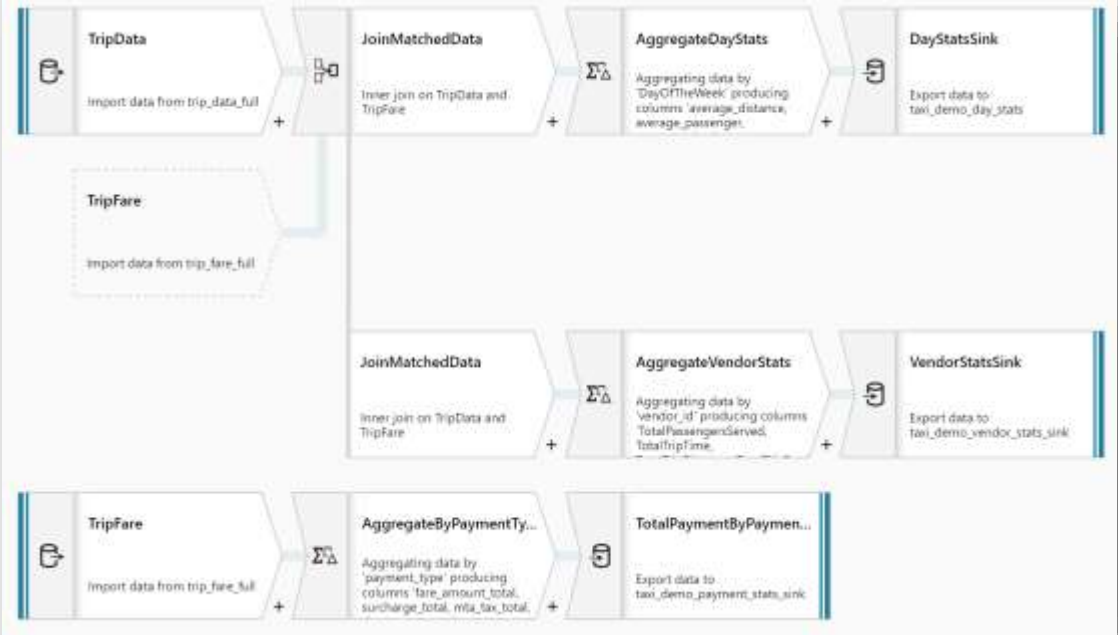
Code-free data preparation @scale

Customer ID	Customer Name	City	Lat	Lon	Email	State	Birthdate
1	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
2	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
3	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
4	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
5	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
6	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
7	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
8	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
9	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
10	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
11	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
12	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
13	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
14	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
15	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
16	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
17	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
18	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
19	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01
20	Henry	Seattle	47.61	-122.33	henry@dataflow.com	WA	1980-01-01

LIMITED PREVIEW

MAPPING DATAFLOW

Code-free data transformation @scale



WRANGLING DATAFLOW

Code-free data preparation @scale

CustomerId	FirstName	LastName	City	Lat	Long	State	BirthDate
1	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
2	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
3	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
4	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
5	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
6	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
7	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
8	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
9	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
10	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
11	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
12	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
13	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
14	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
15	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
16	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
17	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
18	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
19	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
20	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
21	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
22	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
23	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
24	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
25	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
26	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
27	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
28	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
29	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01
30	Henry	Robert	Baltimore	39.2904	-76.6122	MD	1980-01-01

PUBLIC PREVIEW

Data Preparation in today's BI world

- **Finding** & **Connecting** to data is easy for users
- Experiences for data connectivity are **common across different tools**
- Data is **always in the desired shape** for analysis/consumption
- **Reshaping data is easy and quick**, so can be done multiple times without effort
- **Combining** data from multiple sources is straightforward
- **Data Volumes are small and manageable**, data usually comes from a **single type of data source**, **data refresh** is never a problem

In Reality...

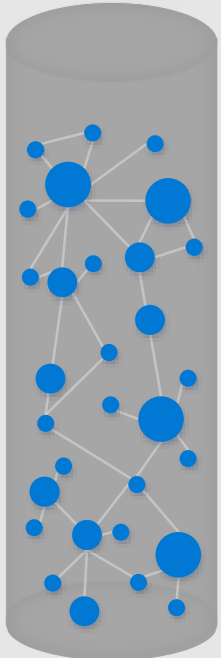
- **Finding** & **Connecting** to data is too difficult
- Experiences for data connectivity are too fragmented
- Data often needs to be **reshaped** before consumption
- Any shaping is one-off and not **repeatable**
- **Combining** data from multiple sources is difficult
- Volume, Velocity & Variety

"Analysts spend up to 80% of their time on data preparation delaying the time to analysis and decision making."

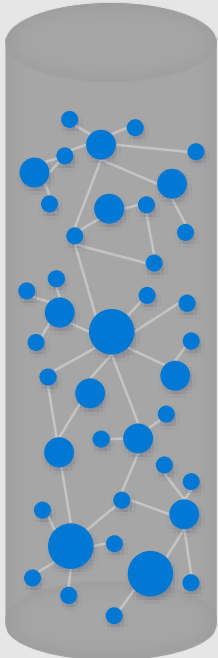
- Gartner

Data is stored within internal and external silos

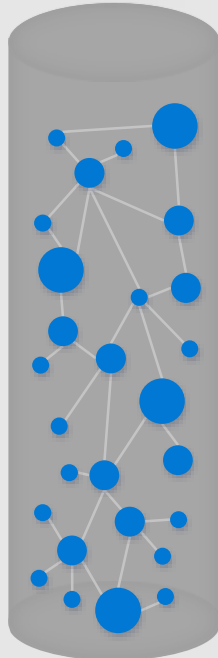
Mobile/Web



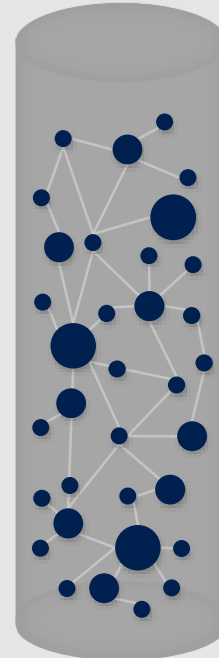
Transactions



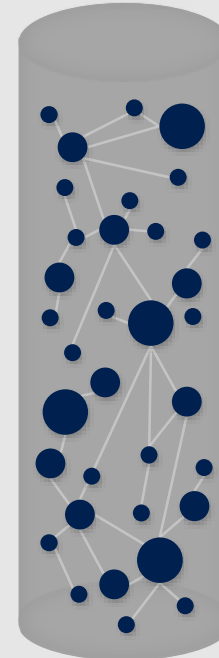
IoT



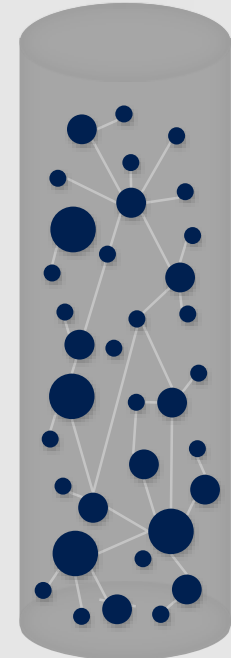
Social



Advertising



Marketplace



Need a shared data model with rich semantics

Mobile/Web



Transactions



IoT



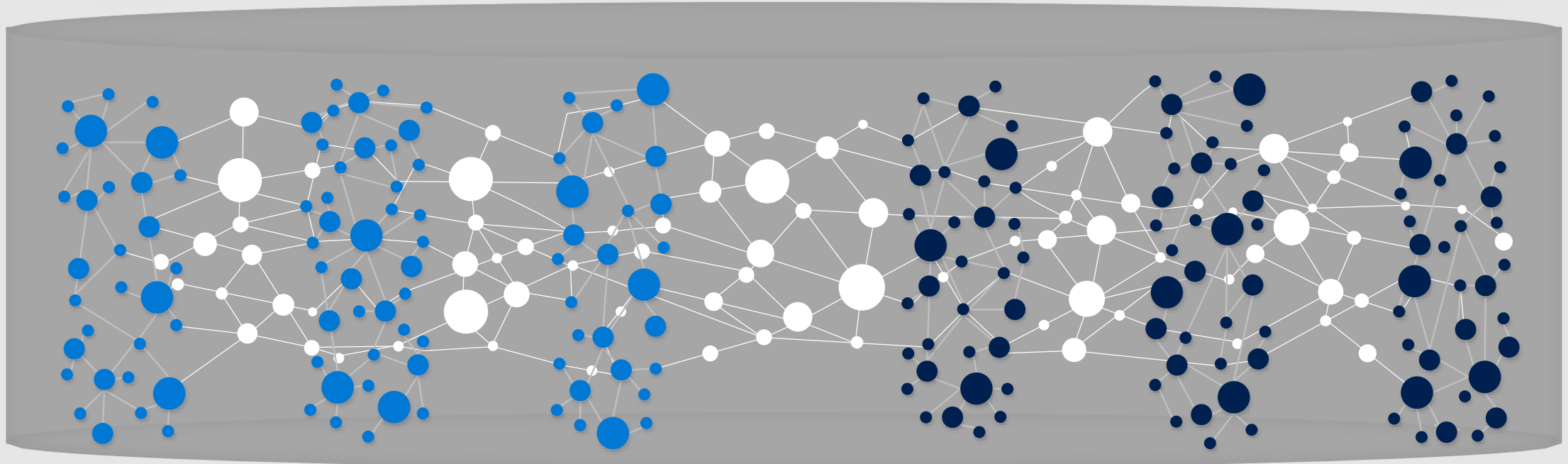
Social



Advertising

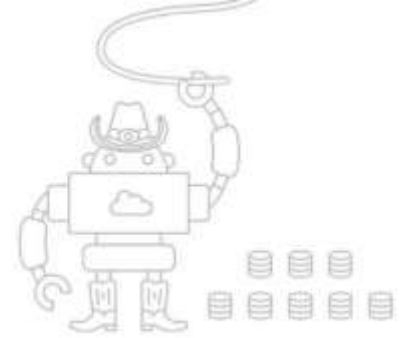


Marketplace





Agenda



Wrangling in Context



Refresher Power Query & M



ADF Data Wrangling



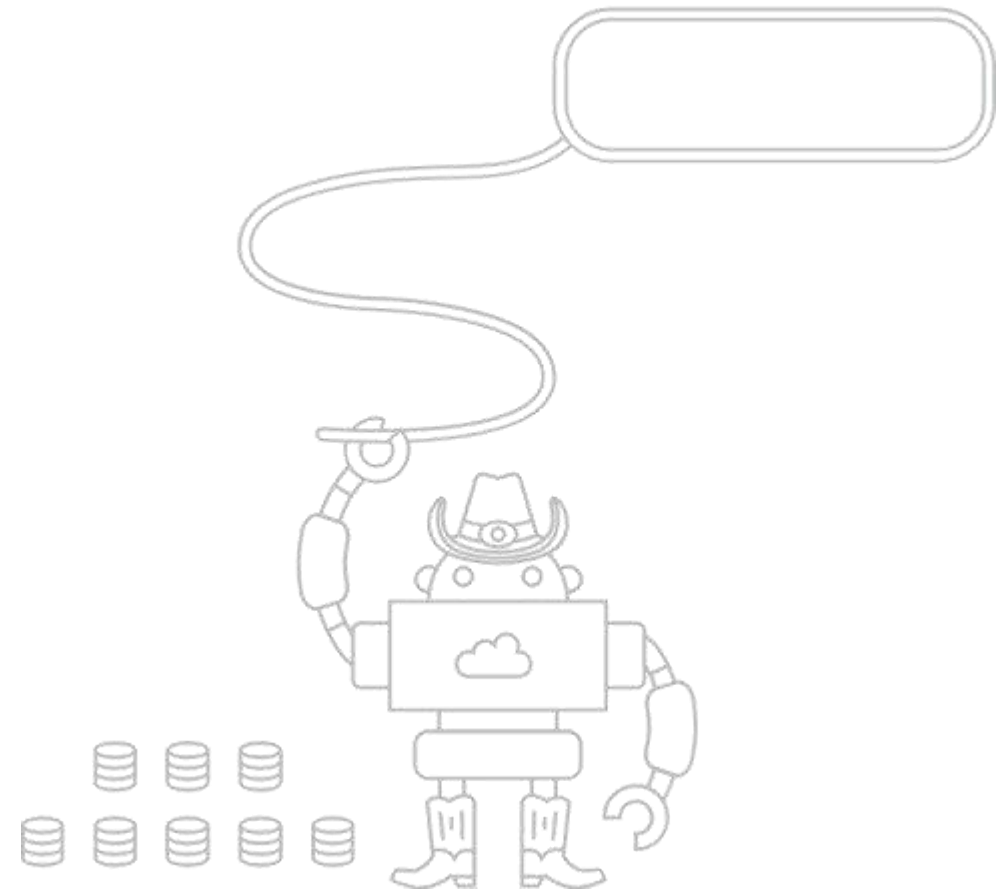
The swotty Power Cousin



Resumé

DEMO

Power Query & M





Agenda



Wrangling in Context



Refresher Power Query & M



ADF Data Wrangling

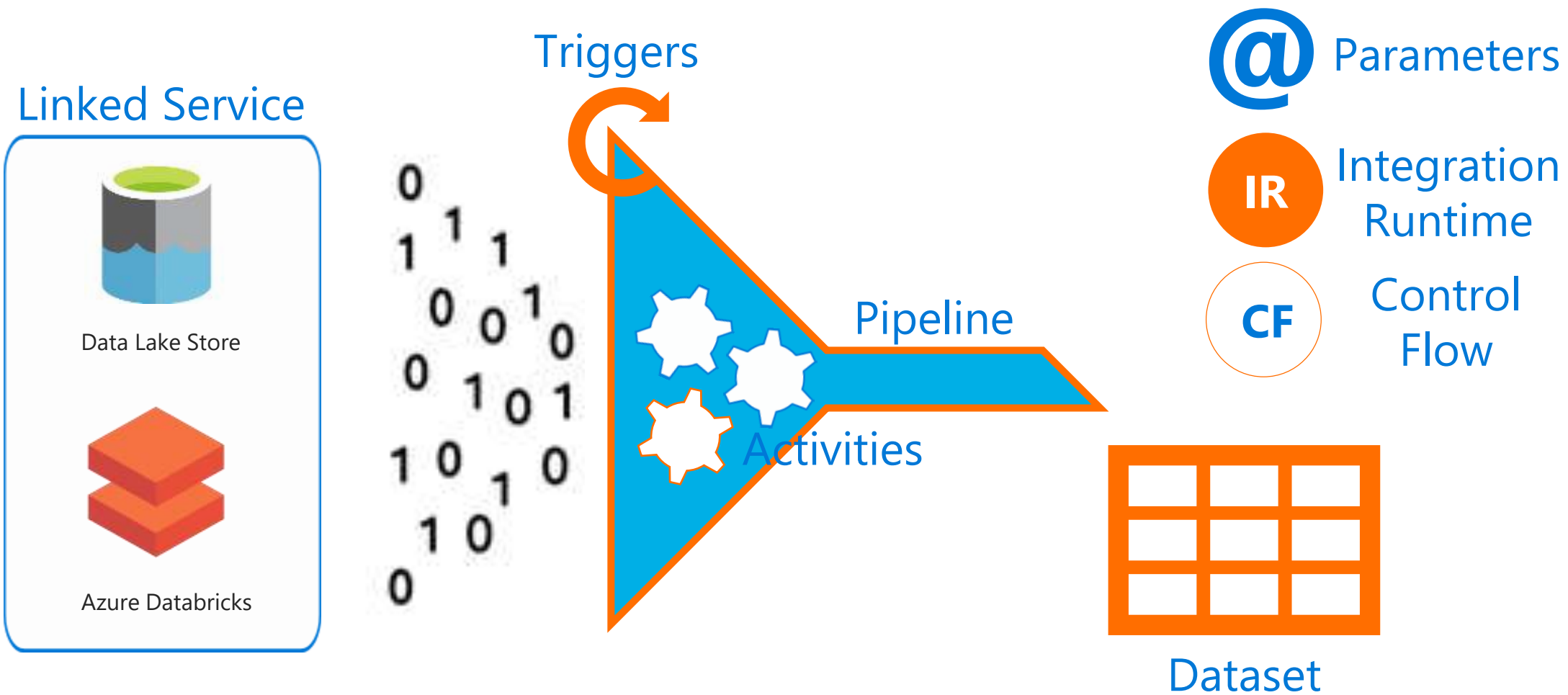


The swotty Power Cousin

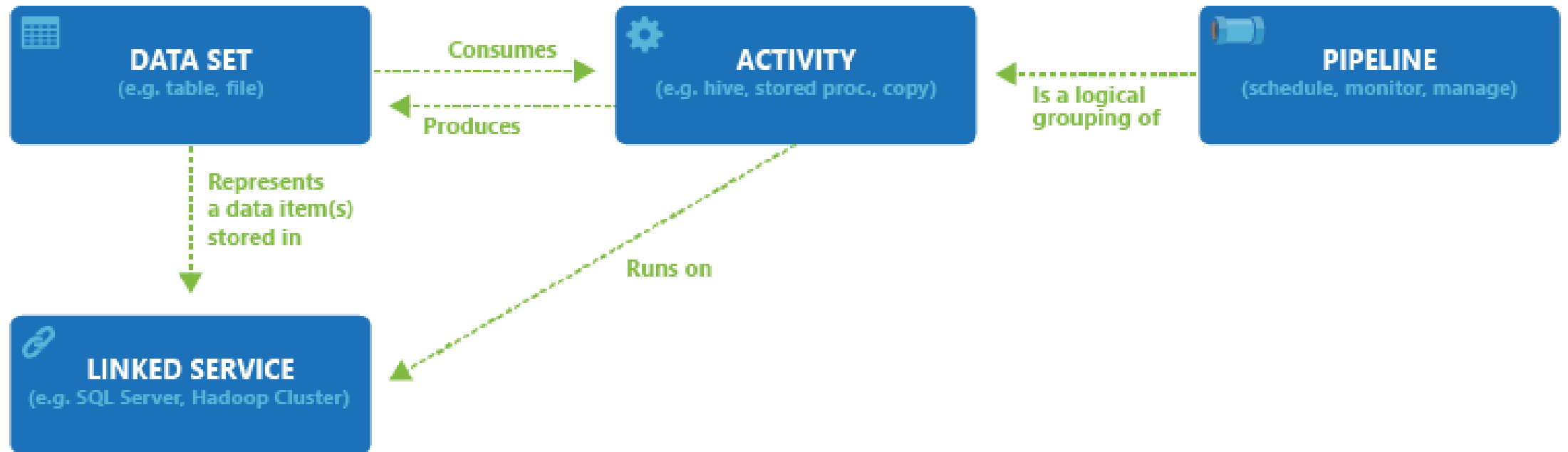


Resumé

Azure Data Factory Components



Component dependencies



Methods for transforming in Azure Data Factory

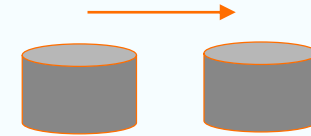
Compute
resources



SSIS Packages



Data Flow



Methods for transforming data in Azure Data Factory

Code free data transformation at scale

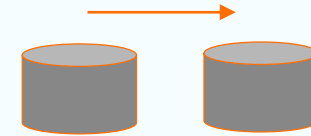
Compute
resources



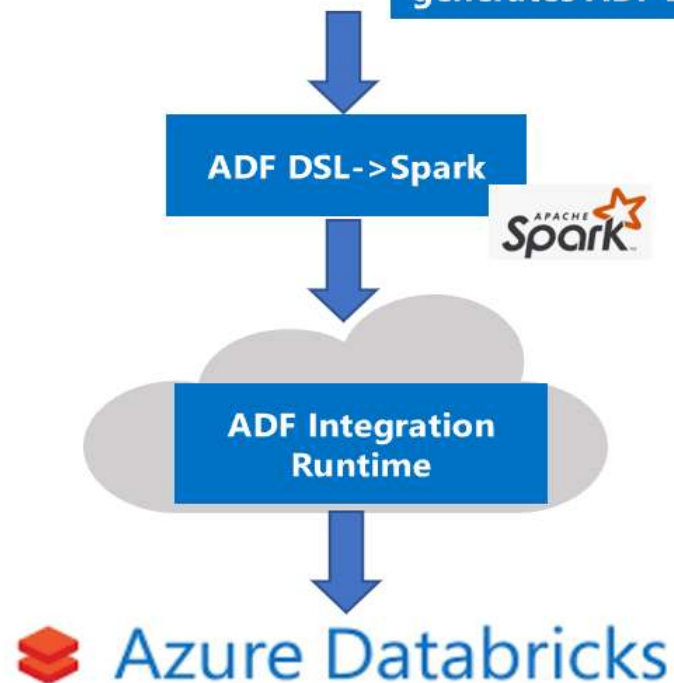
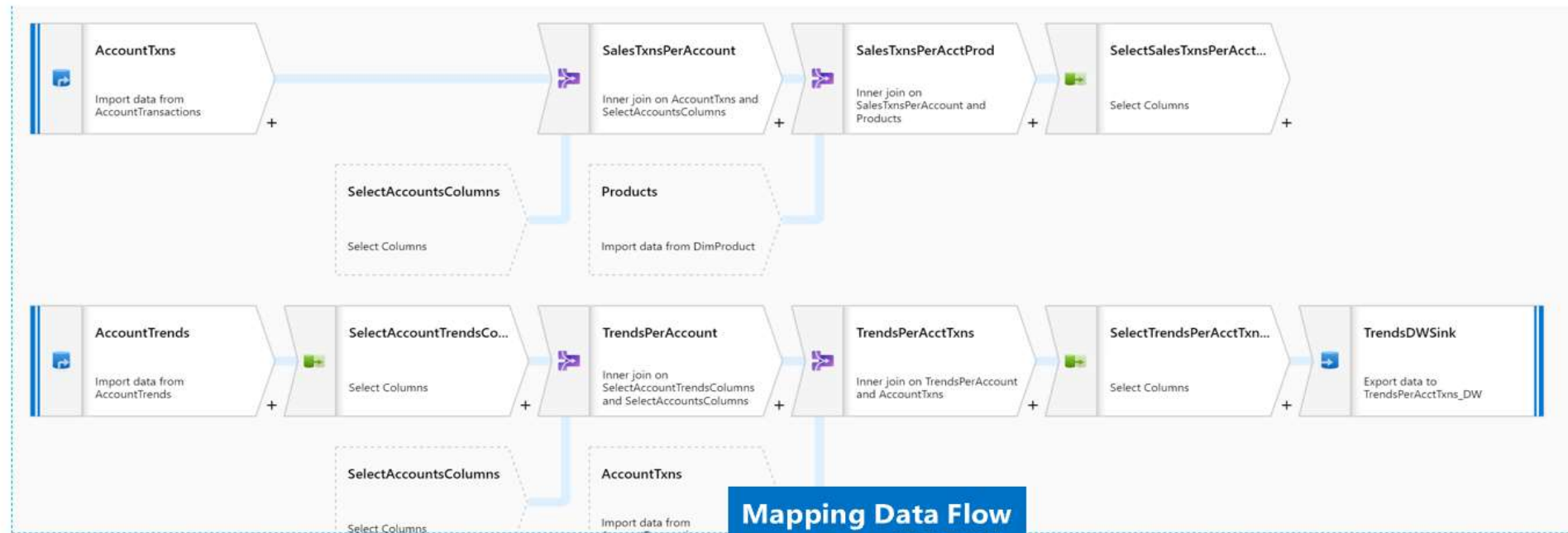
SSIS Packages



Data Flow



Mapping Data Flow: High Level Architecture



Wrangling Data Flow: High Level Architecture

Microsoft Azure

New entities from data (Technical Preview)

Get Data Refresh Options Manage Columns Transform Table Reduce Rows Add Column Add Conditional Column Combine Tables

SalesLT.SalesOrder...

Table.ReverseRows(#"Navigation 1")

	1 SalesOrderID	2 SalesOrderDetailID	3 OrderQty	4 ProductID	5 UnitPrice	6 UnitPriceD
1	71,946	113,406	1	916	31.584	
2	71,938	113,315	3	800	672.294	
3	71,938	113,314	2	675	5.394	
4	71,938	113,313	1	881	32.394	
5	71,938	113,312	3	715	29.984	
6	71,938	113,311	6	939	37.254	
7	71,938	113,310	7	794	1,466.01	
8	71,938	113,309	5	801	672.294	
9	71,938	113,308	4	977	323.994	
10	71,938	113,307	7	880	32.994	
11	71,938	113,306	5	797	672.294	
12	71,938	113,305	3	793	1,466.01	
13	71,938	113,304	5	975	1,020.594	
14	71,938	113,303	2	859	14.694	
15	71,938	113,302	4	874	5.394	
16	71,938	113,301	6	795	1,466.01	
17	71,938	113,300	8	870	2.994	
18	71,938	113,299	3	997	323.994	
19	71,938	113,298	3	883	32.394	
20	71,938	113,297	3	738	202.332	
21	71,938	113,296	5	865	38.1	
22	71,938	113,295	5	708	20.984	
23	71,938	113,294	5	708	20.984	

Name: SalesLT.SalesOrderDetail

Applied Steps: Source, Navigation 1, Reversed Rows

Create

PQO Exp in
ADF UX

PQO gen M &
convert to ADF DSL

ADF DSL->Spark

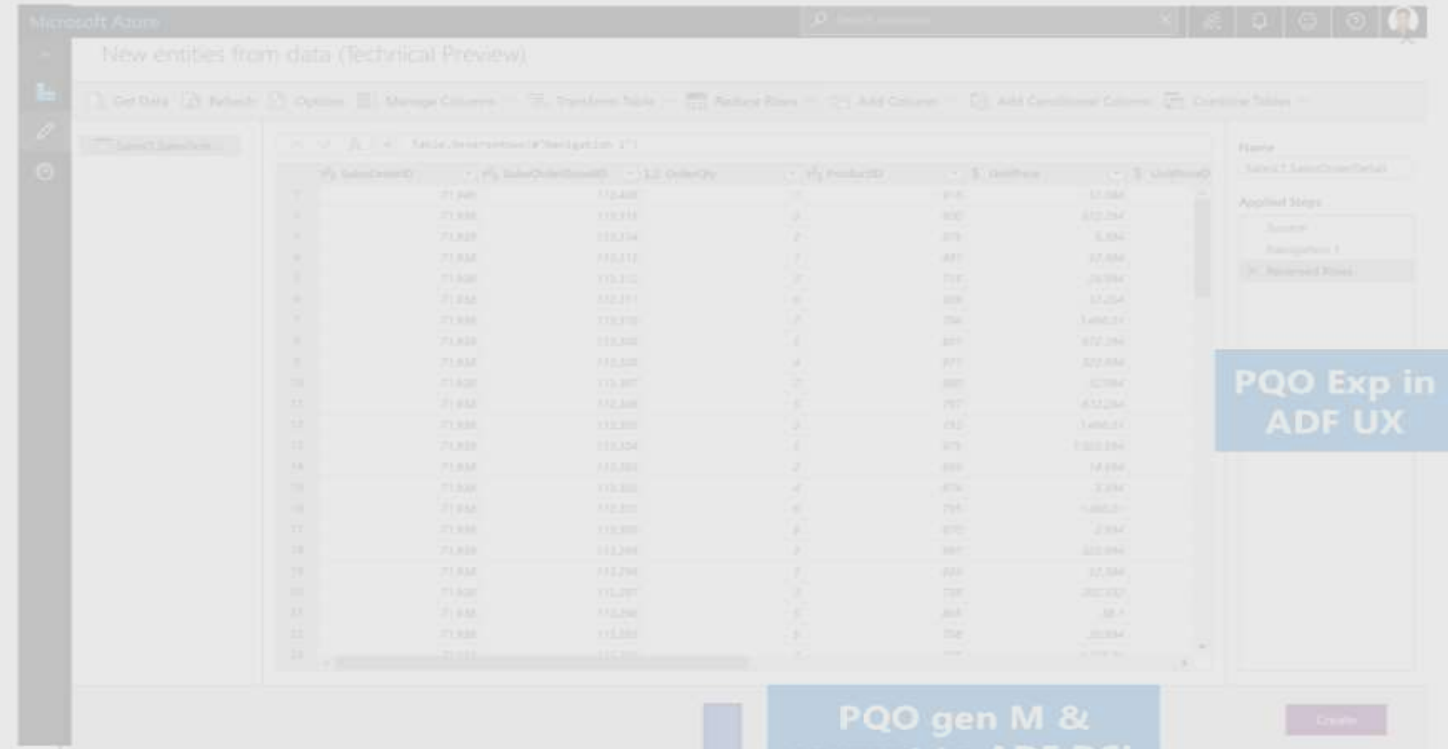


ADF Integration
Runtime



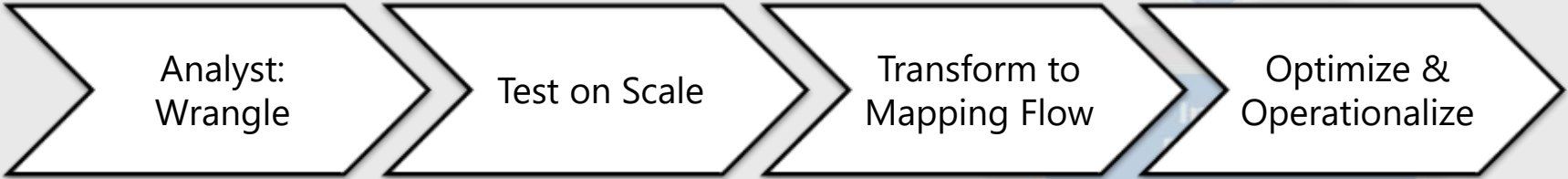
Azure Databricks

Wrangling Data Flow: High Level Architecture



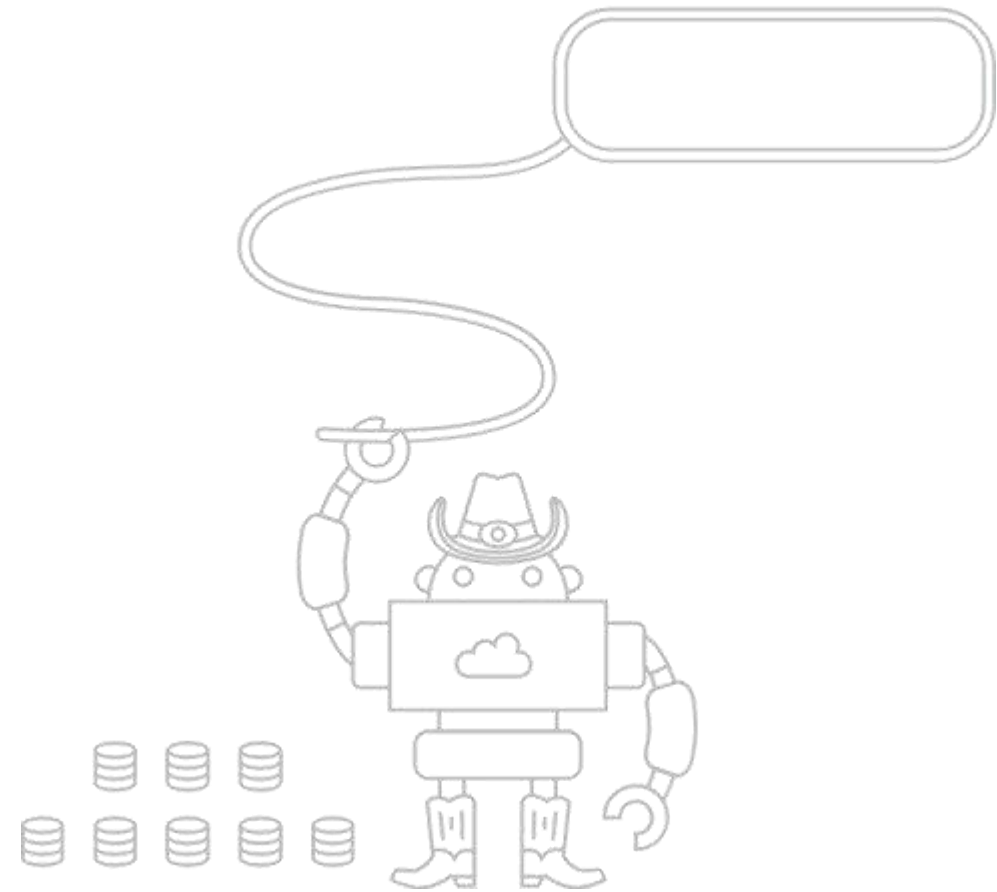
PQO gen M &
convert to ADF DSL

ADF DSL->Spark



Demo

ADF Wrangling Data Flows



Not supportet (yet)

- Merge columns (can be achieved with AddColumn)
- Split column
- Append queries
- 'Use first row as headers' and 'Use headers as first row'

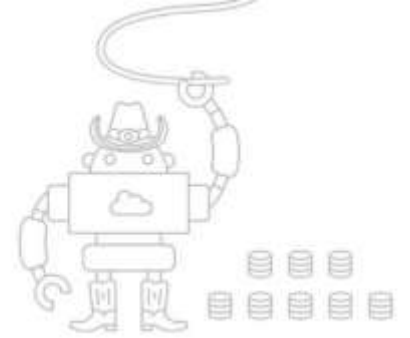
Details at <https://docs.microsoft.com/en-us/azure/data-factory/wrangling-data-flow-functions>

Data Sources

- ADL Gen2
- CSV
- Azure SQL Server



Agenda



Wrangling in Context



Refresher Power Query & M



ADF Data Wrangling

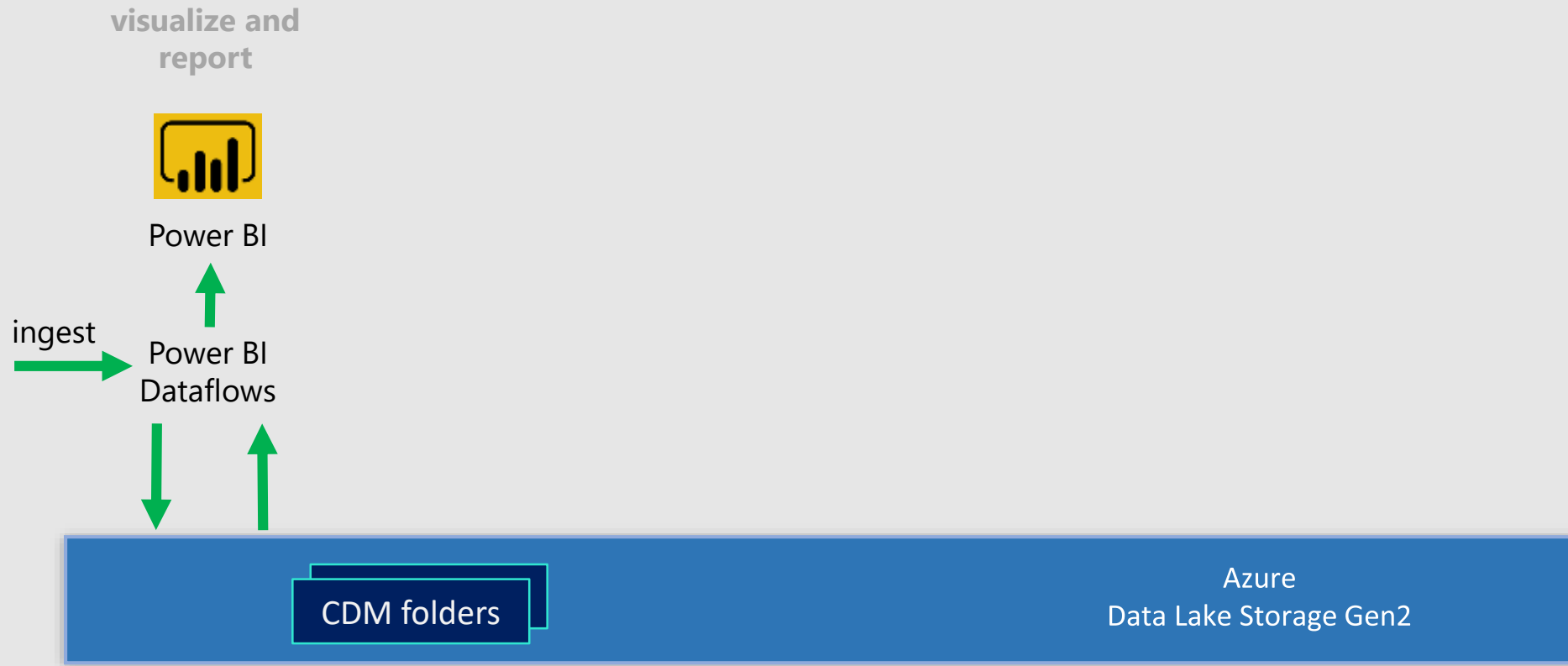


The swotty Power Cousin



Resumé

The idea ...

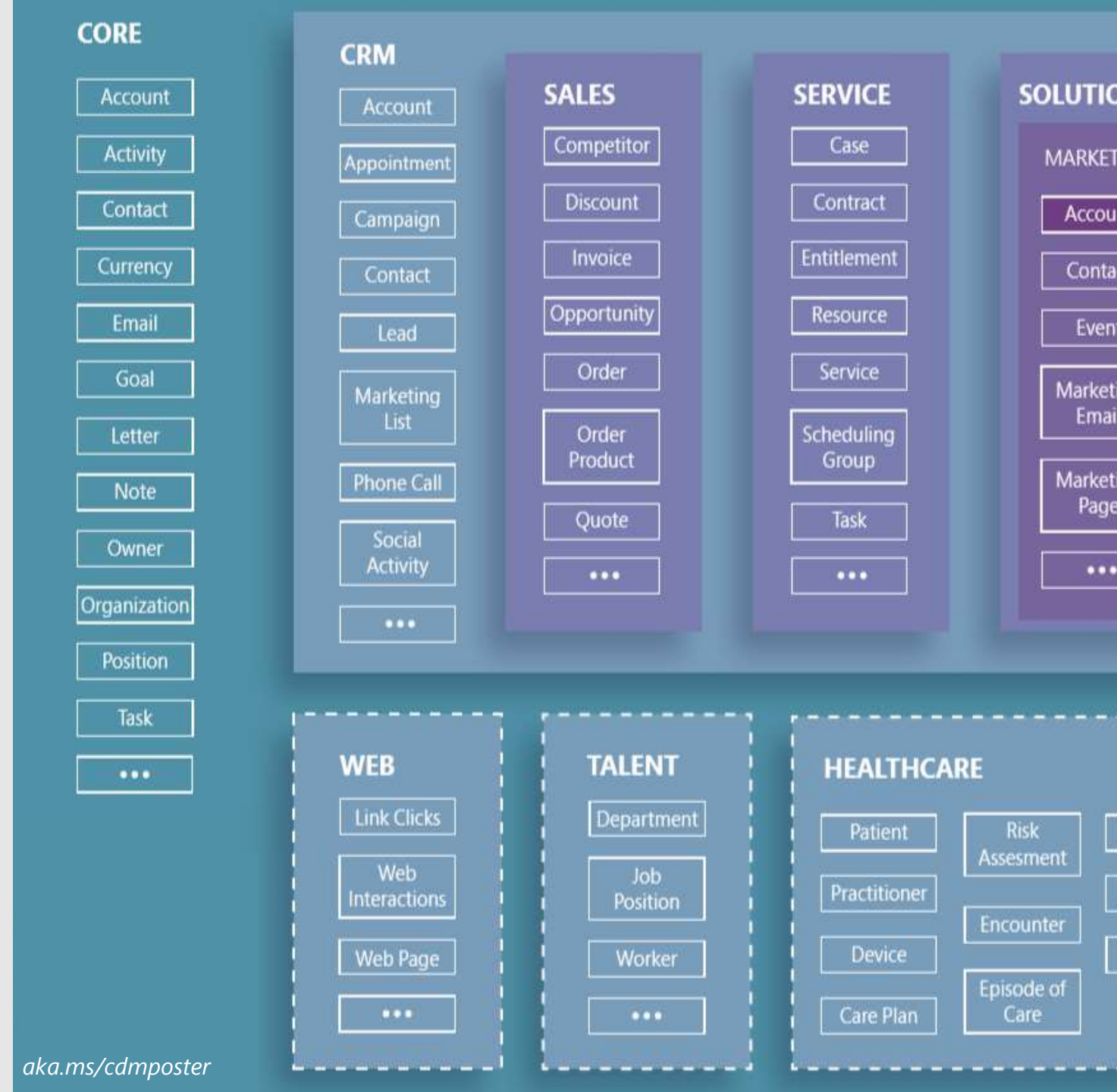


Common Data Model (CDM)

The **Common Data Model (CDM)** provides modular, and extensible business entities (Account, Lead, Opportunity, etc.)

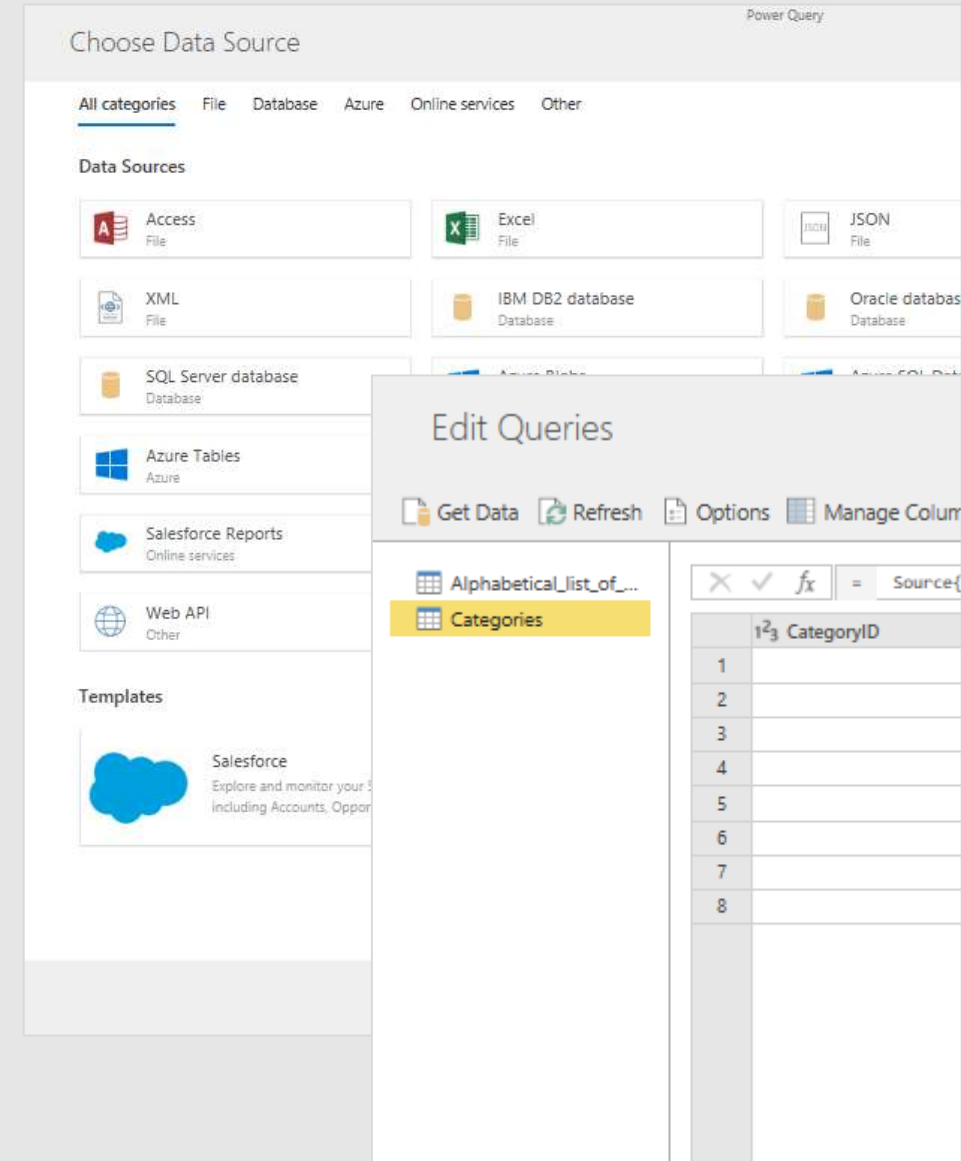
It **unifies data in a well-known schema** with semantic consistency across applications and deployments.

Customers, system integrators, and ISVs can develop turnkey business and intelligence applications, that **share data** based on these CDM entities. They can also extend entities to capture additional business-specific ideas.



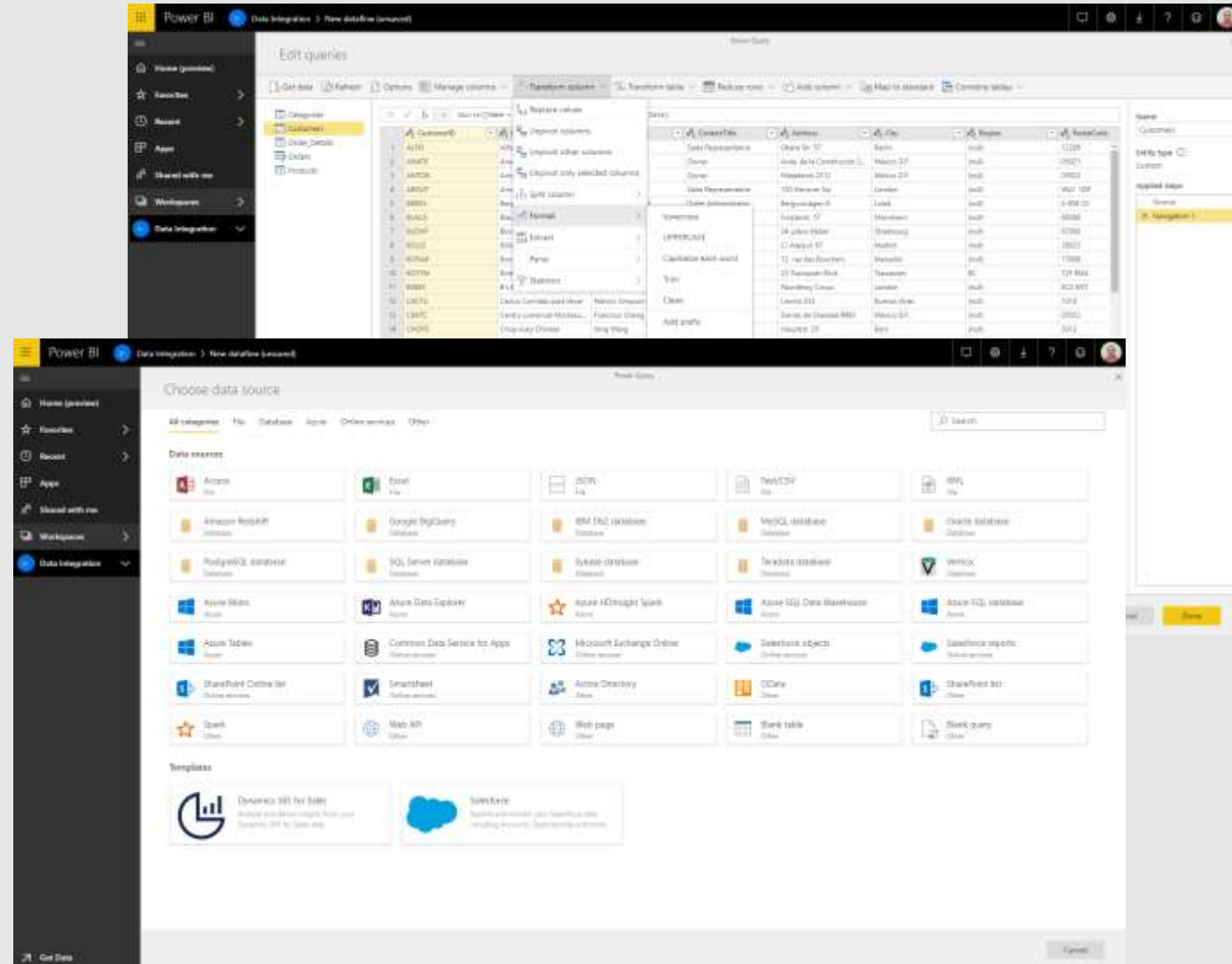
CDM in dataflows

- Dataflows allow customers to bring in data from a variety of disparate sources
- Leveraging **Power Query Online**, same technology in Excel and Power BI, visually transform and combine data sources in a no-code/low-code experience
- **Map** the resulting data to a standard **Common Data Model entity**, or **create a new entity** that can be related back to the standard



Dataflows use Power Query Online

- Power Query available as **web-based self-service data prep experience**
- Supports the same set of **300+ transformations** as desktop Power Query (M Engine parity level)
- Currently **~45 connectors**, including cloud & on-prem data sources via the **On-premises data gateway**



Recently Shipped Features – November 2019

New Data Connectors

- PDF Files
- Local Folder
- SharePoint Folder
- Google BigQuery
- Teradata
- MySQL
- PostgreSQL
- ODBC
- HDInsight Spark
- Apache Spark
- Impala

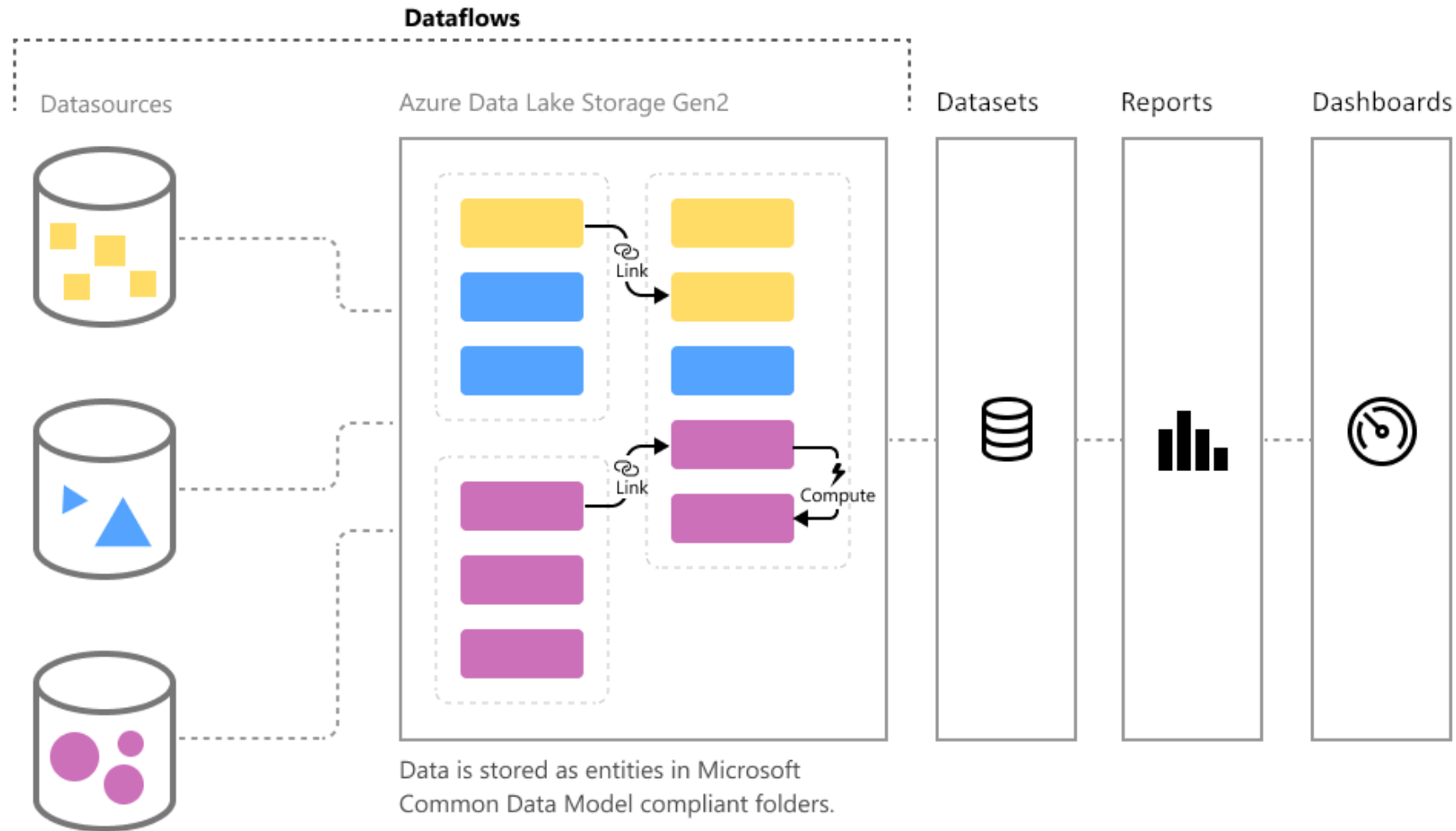
New Data Transformations

- Combine Files UX
- Merge Queries – Visual join kind selection
- Additional Number/Date/DateTime/Duration transformations UX
- List Transforms: Statistics, Sort, Keep/Remove/Reverse items
- Fill Up/Down
- Move Columns left/right/beginning/end
- Replace Errors

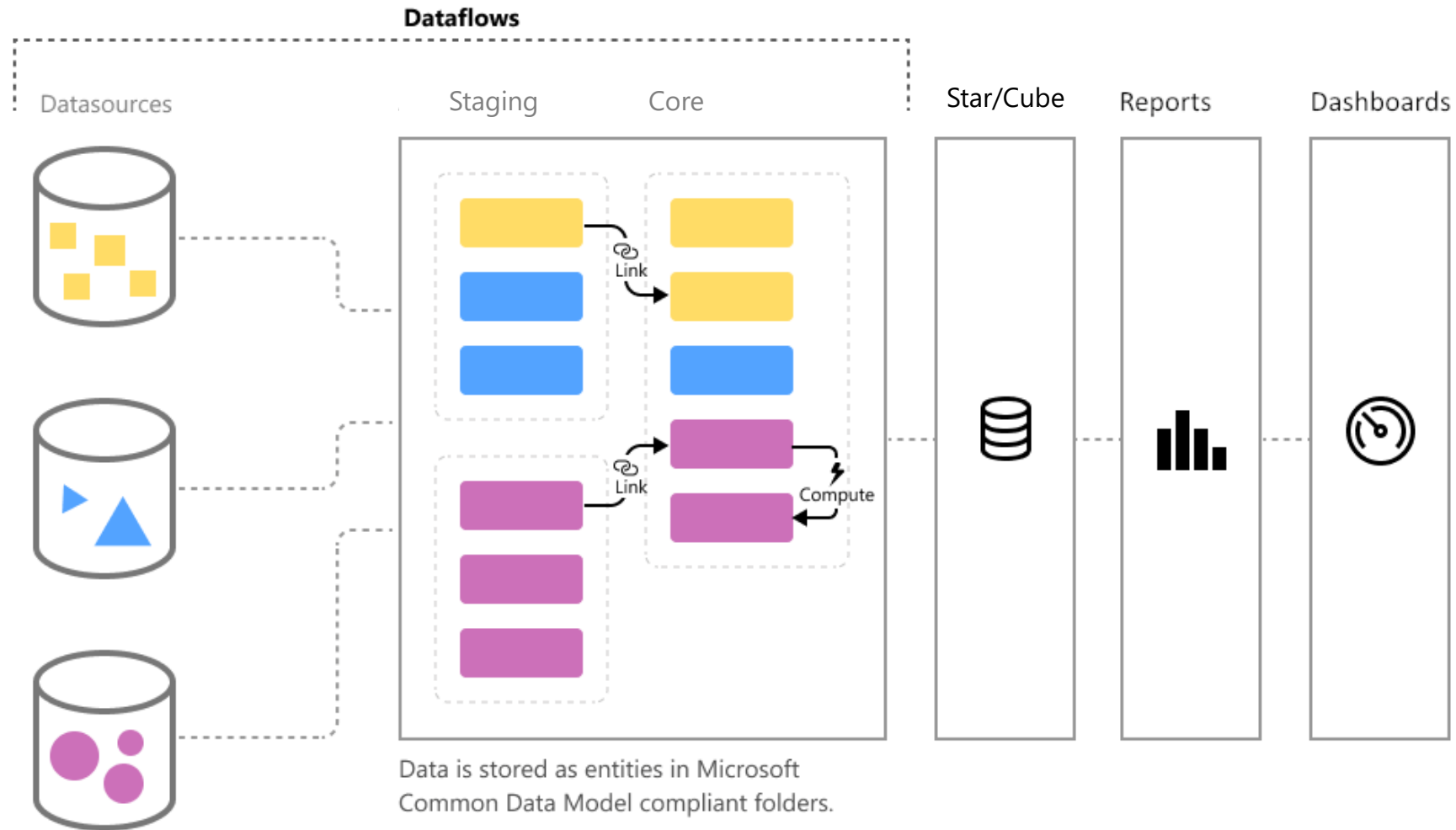
Other Enhancements

- Data Profiling
- Query Parameters
- Function Authoring
- M Intellisense support in Advanced Query Editor & Formula Bar
- Select Related Tables as part of Get Data UX
- More descriptive error messages within Power Query Online

Dataflows are composable, just like Excel



Speaking old architecture ...



Power BI dataflows “FAQ”

- ✓ • A new capability for self-service data preparation in Power BI
- ✓ • Delivered in a familiar Power Query experience
- ✓ • Built on the foundation of Azure Data Lake Storage gen2
- ✓ • Utilize the CDM folder format for data storage & enable mapping to CDM
- ✓ • A tool for business users to drive data reuse without requiring IT involvement
- ✓ • Enable Excel-like data lineage and orchestration

- ✗ • NOT a replacement for datasets
- ✗ • NOT a replacement for a data warehouse
- ✗ • NOT a replacement for Azure Data Factory or SSIS
- ✗ • NOT a Premium-only feature
- ✗ • NOT an additional cost or fee
- ✗ • NOT spelled with a space or any capital letters

Power BI dataflows “FAQ”

- ✓ • A new capability for self-service data preparation in Power BI
- ✓ • Delivered in a familiar Power Query experience
- ✓ • Built on the foundation of Azure Data Lake Storage gen2
- ✓ • Utilize the CDM folder format for data storage & enable mapping to CDM
- ✓ • A tool for business users to drive data reuse without requiring IT involvement
- ✓ • Enable Excel-like data lineage and orchestration

- ✗ • NOT a replacement for datasets
- ✗ • NOT a replacement for a data warehouse
- ✗ • NOT a replacement for Azure Data Factory or SSIS
- ✗ • NOT a Premium-only feature
- ✗ • NOT an additional cost or fee
- ✗ • NOT spelled with a space or any capital letters

Key takeaways from this session

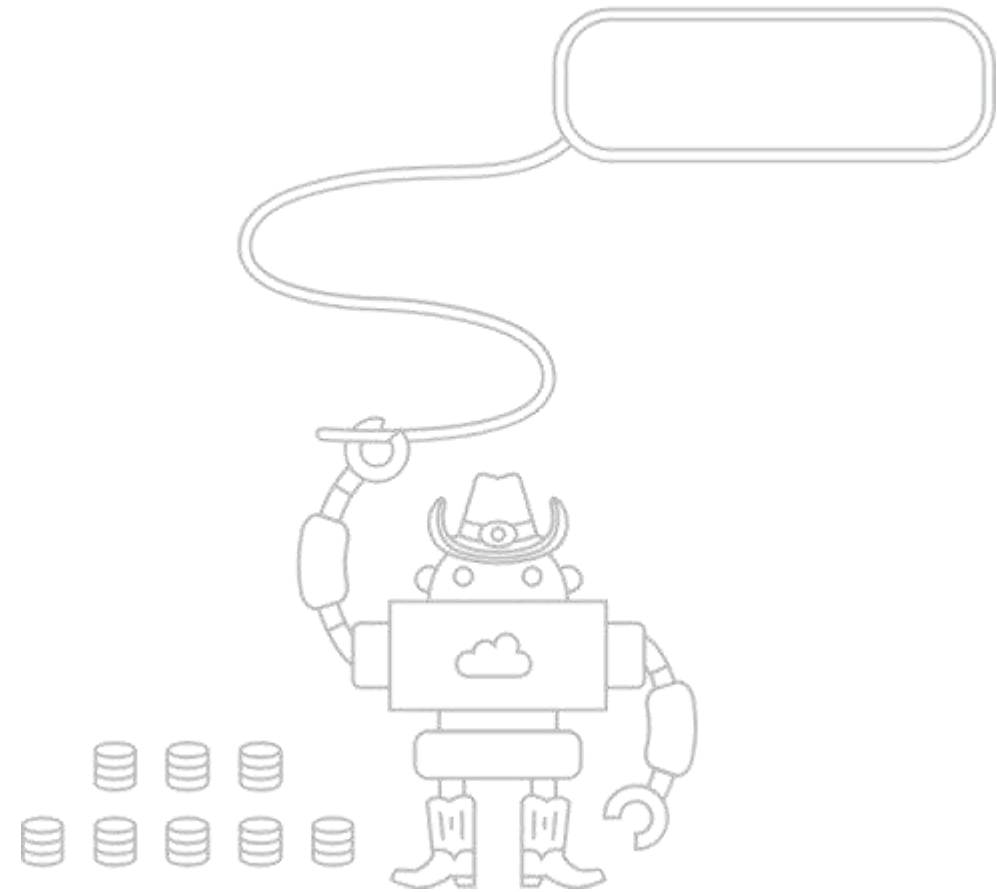
Move your ETL to Dataflows and benefit from reusability and sharing

Try out the new Compute engine inside dataflows to greatly improve refresh times of your ETL Jobs

Have a deeper understanding of how Power BI and dataflows work

Demo

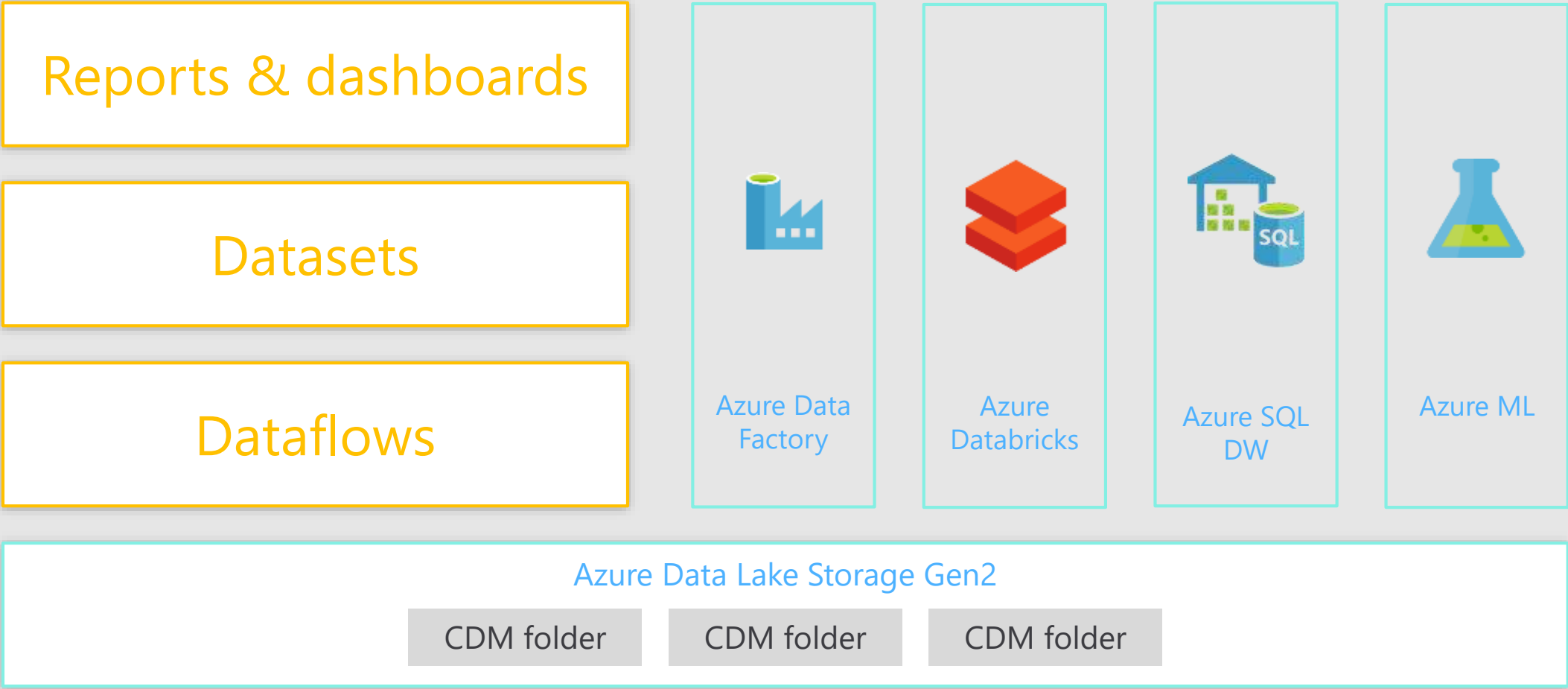
Power BI Dataflows



Use cases

- There are two distinct use cases for Dataflows and ADLS Gen2:
 1. Prep data with Dataflows in Power BI, and all dataflow data will be stored in your organization's ADLS Gen2
 - Data will be stored in CDM Folder format.
 - Power BI, Azure Data services or your custom LOB applications can read the data from the lake

Using an organizational ADLSg2 resource with dataflows



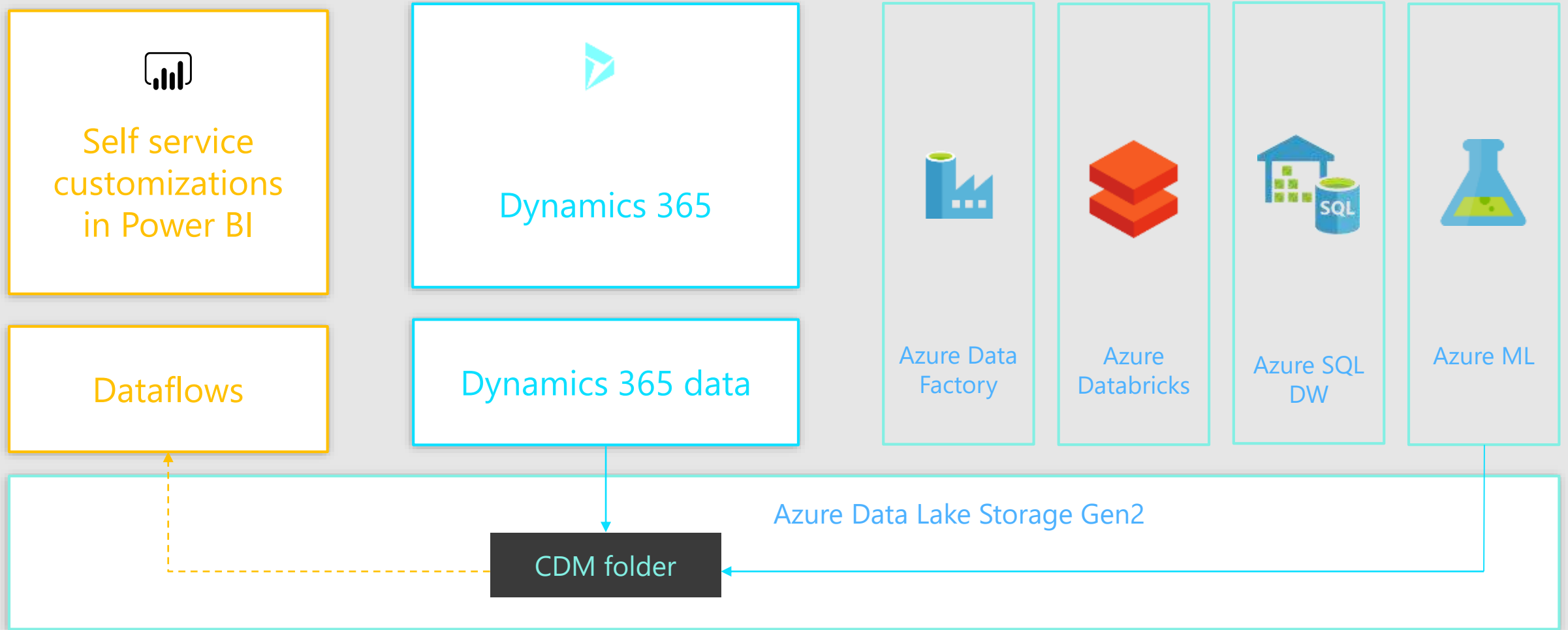
Business analysts
Low/no code

Data scientists
Data engineers

Use cases

- There are two distinct use cases for Dataflows and ADLS Gen2:
 1. Prep data with Dataflows, and all dataflow data will be stored in your organization's ADLS Gen2
 - Data will be stored in CDM Folder format.
 - Power BI, Azure Data services or your custom LOB applications can read the data from the lake
 2. Prep data in the lake, and add it as an external dataflow
 - Data stored in CDM Folder format, can easily be added as a dataflow
 - Business Analysts in Power BI Desktop can easily build reports and dashboards using the Dataflows connector

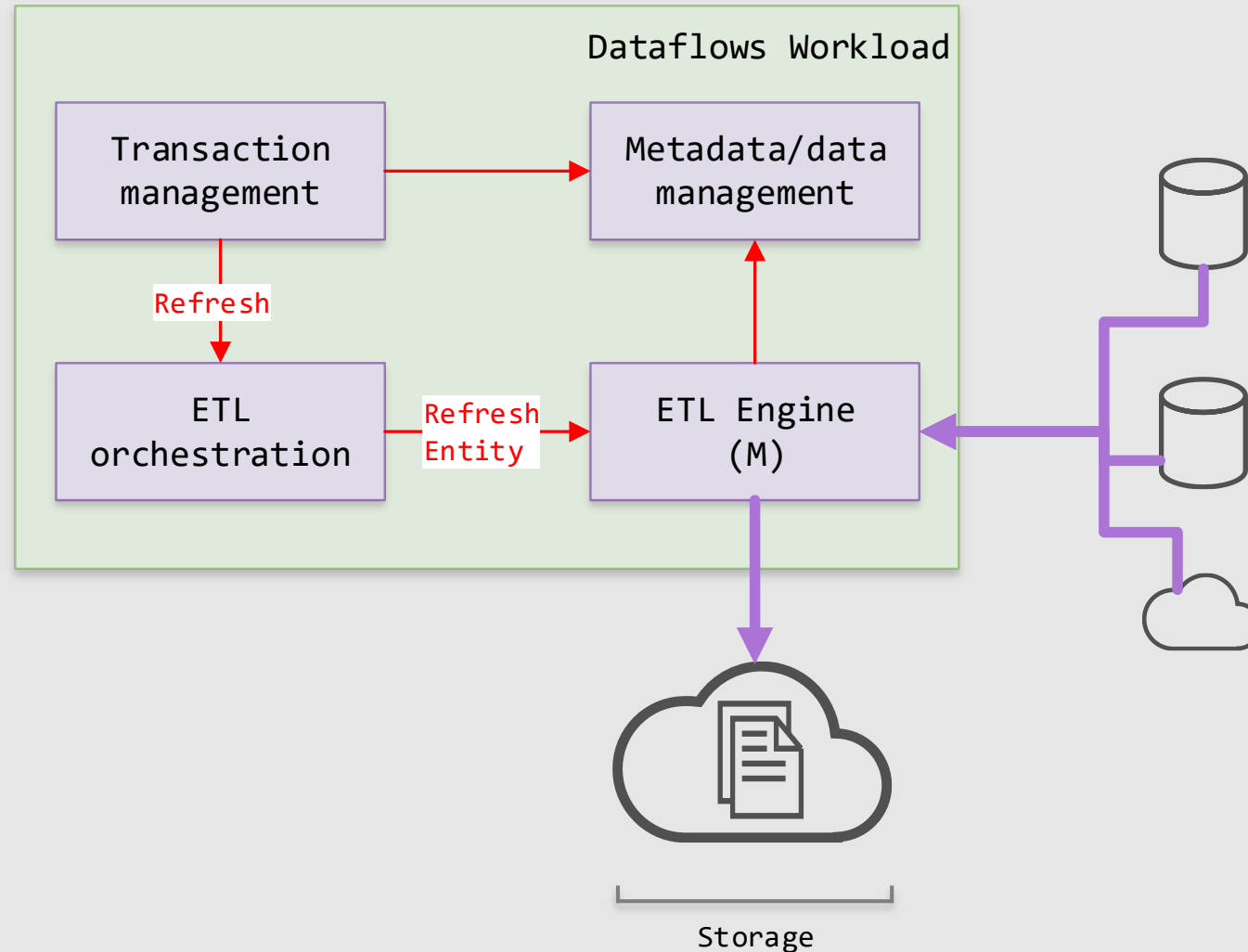
Deliver ready-made insights to Power BI users from Azure



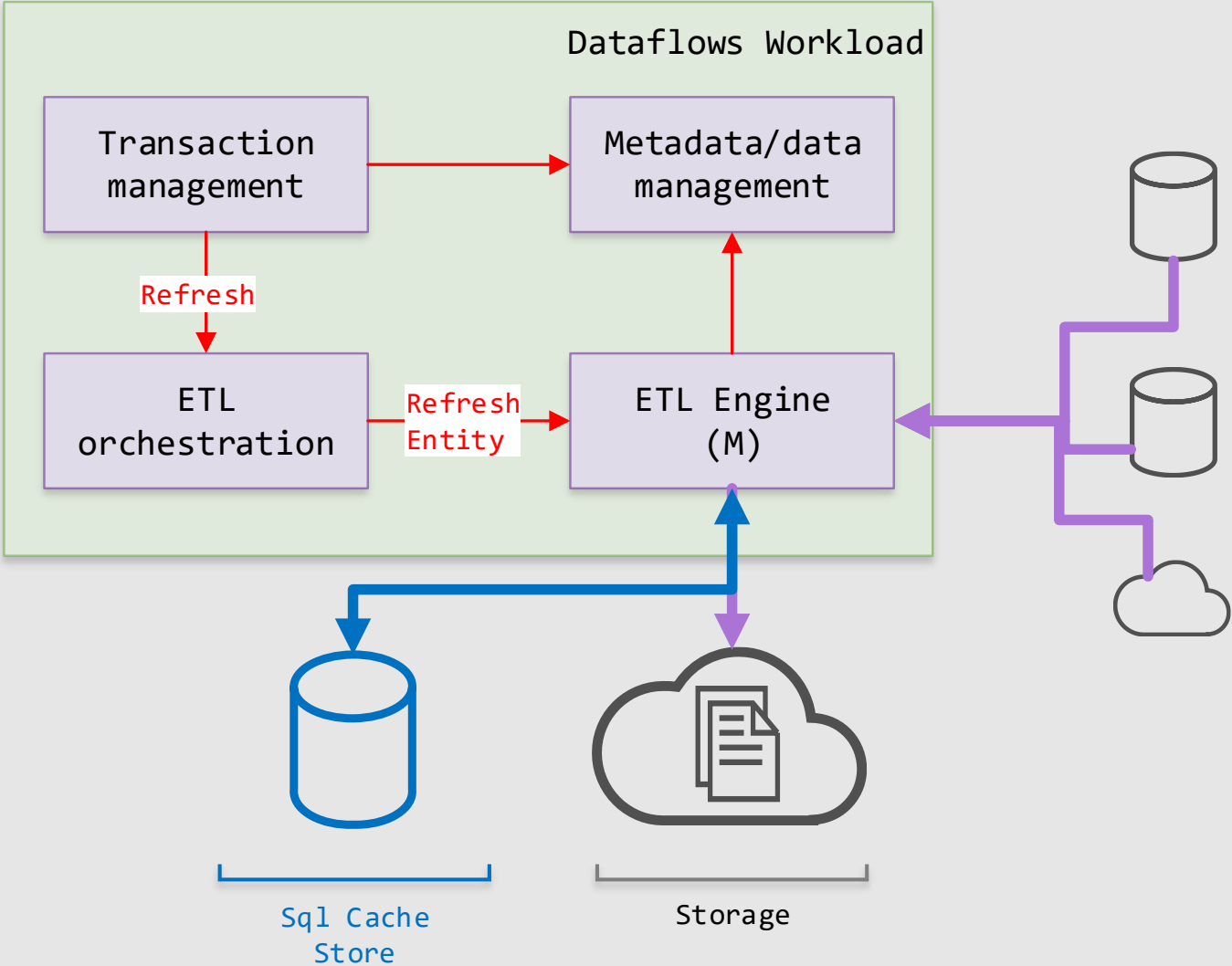
Feature	Pro	Premium
Create dataflows	✓	✓
Connect to dataflows from power bi desktop	✓	✓
Work with on prem and cloud sources	✓	✓
Use ADLS Gen 2 instead of built in storage	✓	✓
Use External CDM folders	✓	✓
Use linked entities (link to other dataflows)		✓
Use computed entities (link to other dataflows and compute over it)		✓
Enhanced compute engine for faster ETL		✓
Direct Query over dataflows		✓
Automated ML and Cognitive Services		✓

Feature	Pro	Premium
Create dataflows	✓	✓
Connect to dataflows from power bi desktop	✓	✓
Work with on prem and cloud sources	✓	✓
Use ADLS Gen 2 instead of built in storage	✓	✓
Use External CDM folders	✓	✓
Use linked entities (link to other dataflows)		✓
Use computed entities (link to other dataflows and compute over it)		✓
Enhanced compute engine for faster ETL		✓
Direct Query over dataflows		✓
Automated ML and Cognitive Services		✓

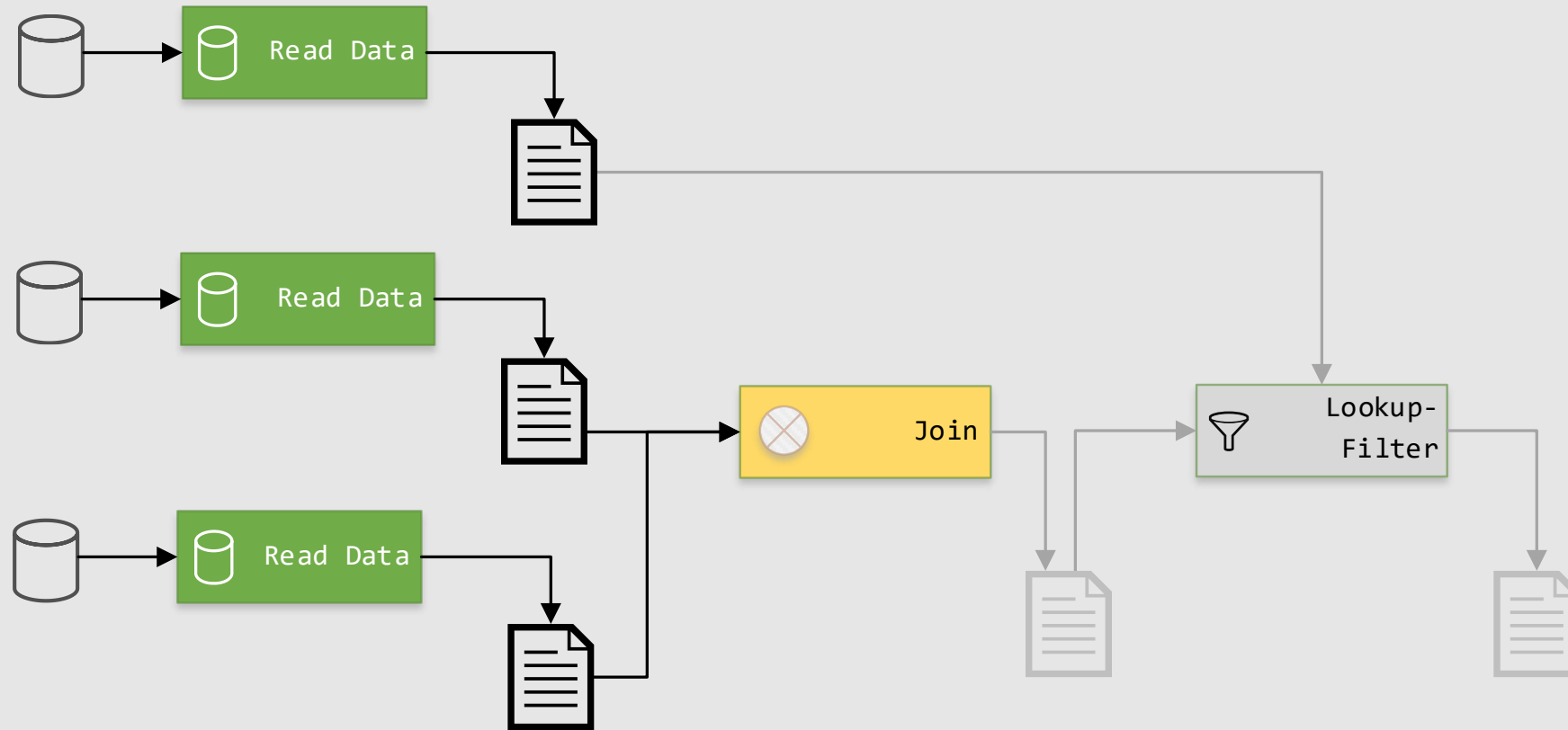
Current Dataflows Compute Engine (GA)



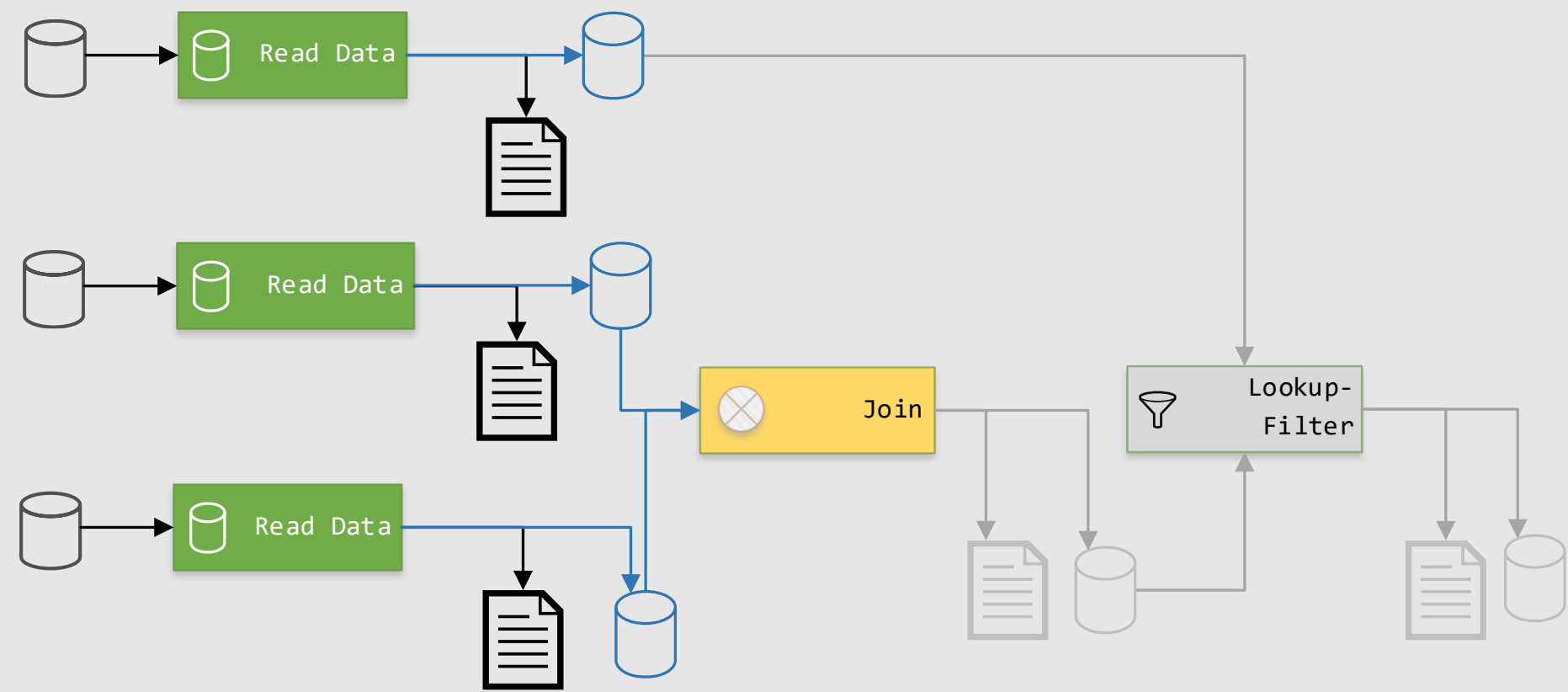
Enhanced Dataflows Compute Engine (Preview)



Current Dataflows Compute Engine (GA)



Enhanced Dataflows Compute Engine (Preview)



10GB Calc Time

● Old Engine ● New Engine



10GB Total Execution Time (minutes)

● Old Engine ● New Engine

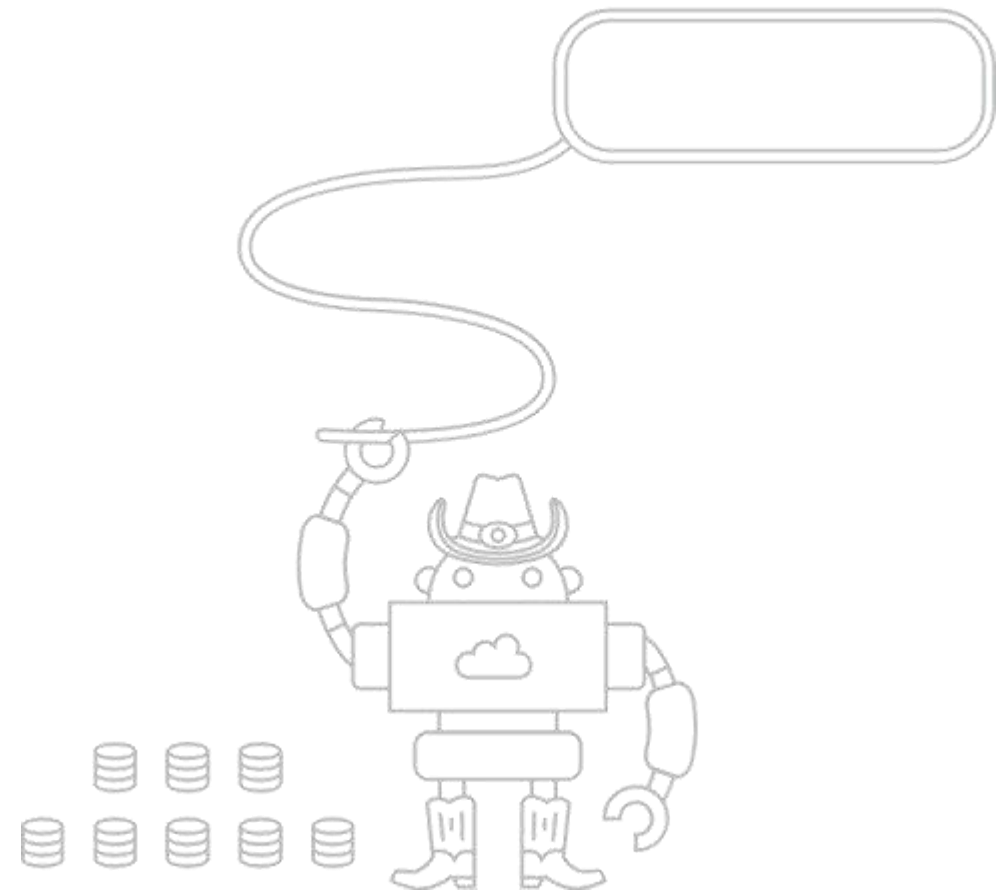


Up to 25X times faster

**Note : Loading data is slower as the data is being written twice*

Demo

Configure the Enhanced Compute Engine



Admin portal

- Usage metrics
- Users
- Audit logs
- Tenant settings
- Capacity settings
- Embed Codes
- Organizational visuals
- Dataflow settings
- Workspaces
- Custom branding

Power BI Premium Power BI Embedded

CAPACITY NAME	CAPACITY ADMINS	ACTIONS	SKU	REGION	STATUS
[REDACTED]			A3	West Europe	Active

[Set up new capacity in Azure](#)

Admin portal

Usage metrics

Users

Audit logs

Tenant settings

Capacity settings

Embed Codes

Organizational visuals

Dataflow settings

Workspaces

Custom branding

West Europe

USER PERMISSIONS

- Users with assignment permissions
Disabled for the entire organization

MORE OPTIONS

- Workloads
- Advanced options

WORKSPACES (3)

🔍 Search content...

✕ Remove all + Assign workspaces

Showing 3 items Sort Workspaces by... ▼

WORKSPACE NAME	WORKSPACE ADMINS	ACTIONS	STATUS
----------------	------------------	---------	--------

Admin portal

Usage metrics

Users

Audit logs

Tenant settings

Capacity settings

Embed Codes

Organizational visuals

Dataflow settings

Workspaces

Custom branding

Max Online Dataset Size (Gb)

Automatic page refresh (Preview)

☐ Off

Minimum refresh interval

Minutes ▼

DATAFLOWS - *Active*

Your workload is ready to use.

☒ On

Max Memory (%)

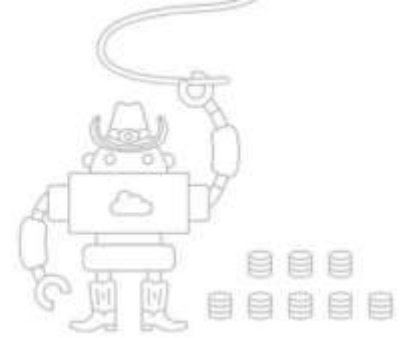
Enhanced Dataflows Compute Engine (Preview)

☒ On

Container Size (Mb)



Agenda



Wrangling in Context



Refresher Power Query & M



ADF Data Wrangling



The swotty Power Cousin



Resumé

Feature	ADF	PBI
Maturity		✓
Source Control	✓	
Modern DWH	✓	
Logical DWH		✓
Agility		✓
Time To Value		✓
Scalability (Data)	✓	
Performance		?
Connectivity		✓
Scalability (Development)	✓	

GET IN TOUCH

!

Christoph Seck
Chief Architect BI

Mail: c.seck@kigroup.de
Mobil: +49 151 12 11 10 95

KI performance GmbH
Mittelstraße 12-14
50672 Köln