

# Azure Synapse Analytics Service

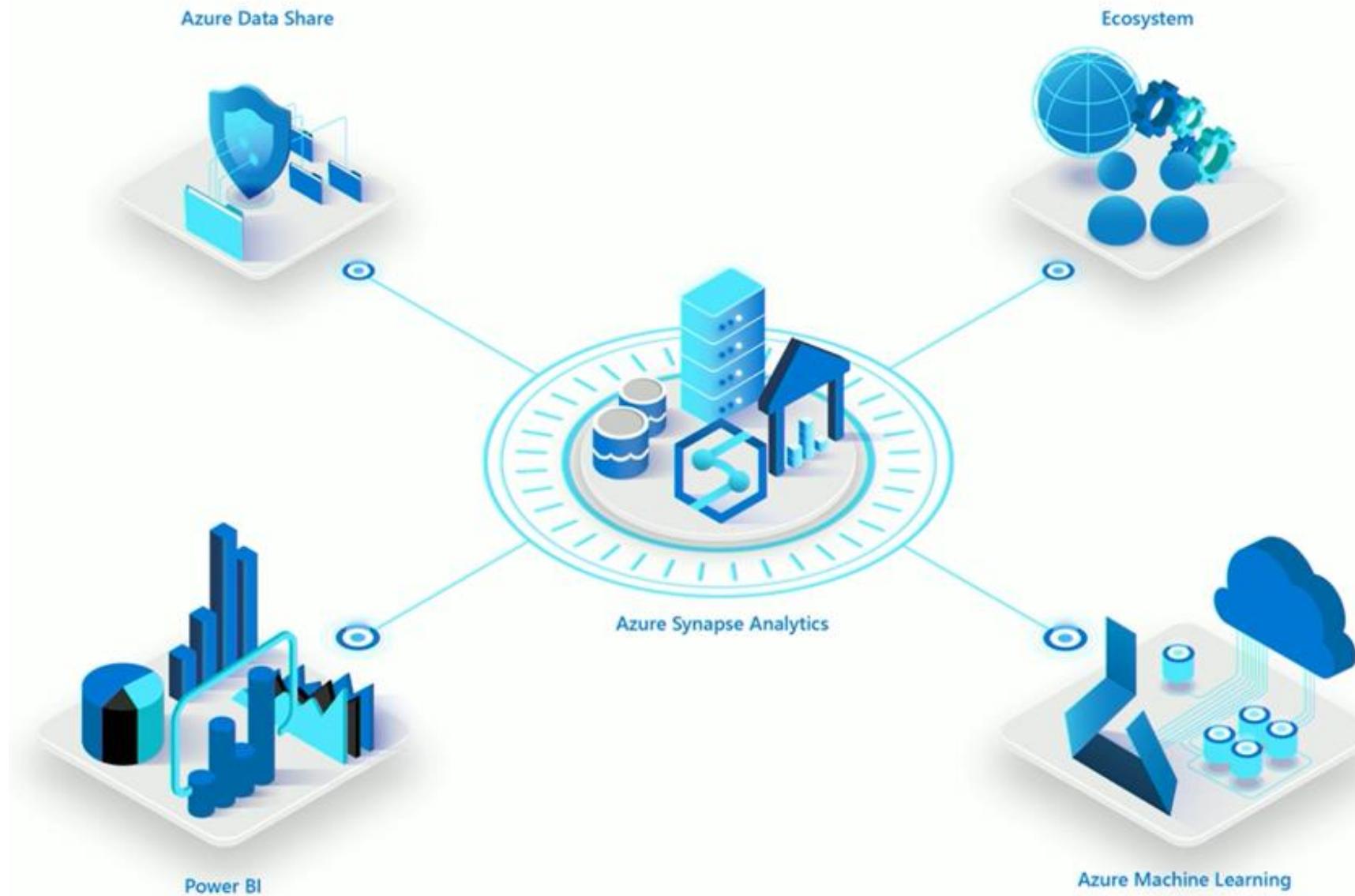
Eshant Garg

Data Engineer, Architect, Advisor

[eshant.garg@gmail.com](mailto:eshant.garg@gmail.com)



Azure  
Synapse  
Analytics



# Introduction

Why Warehousing in Cloud

Azure Synapse Service – a Game Changer

Traditional vs Modern DW Architecture

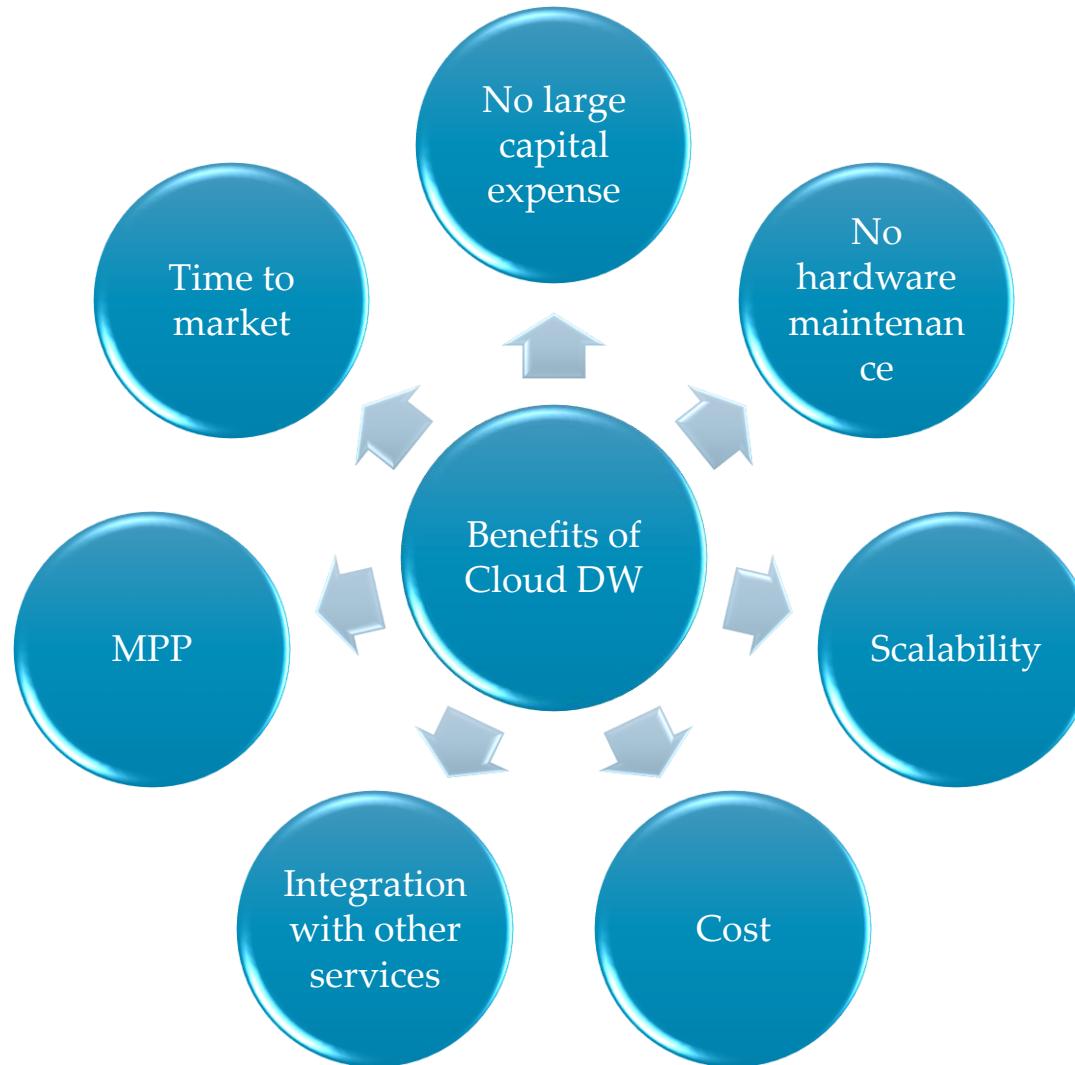
Modern vs Synapse DW Architecture

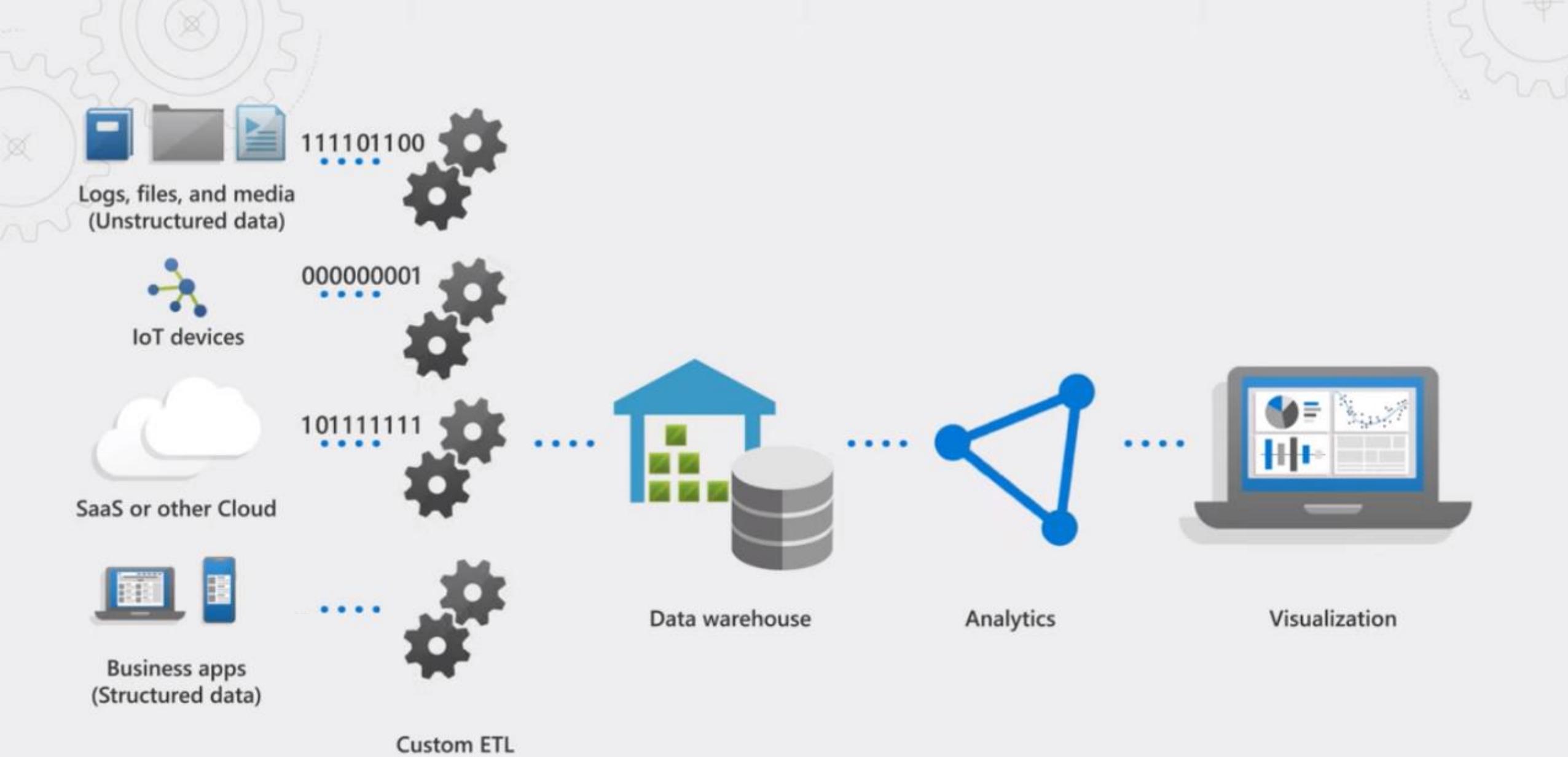
Azure Synapse Studio – Unified Experience

Demo – Provision, Scale, Pause, Firewall settings



# Why Warehousing in Cloud?

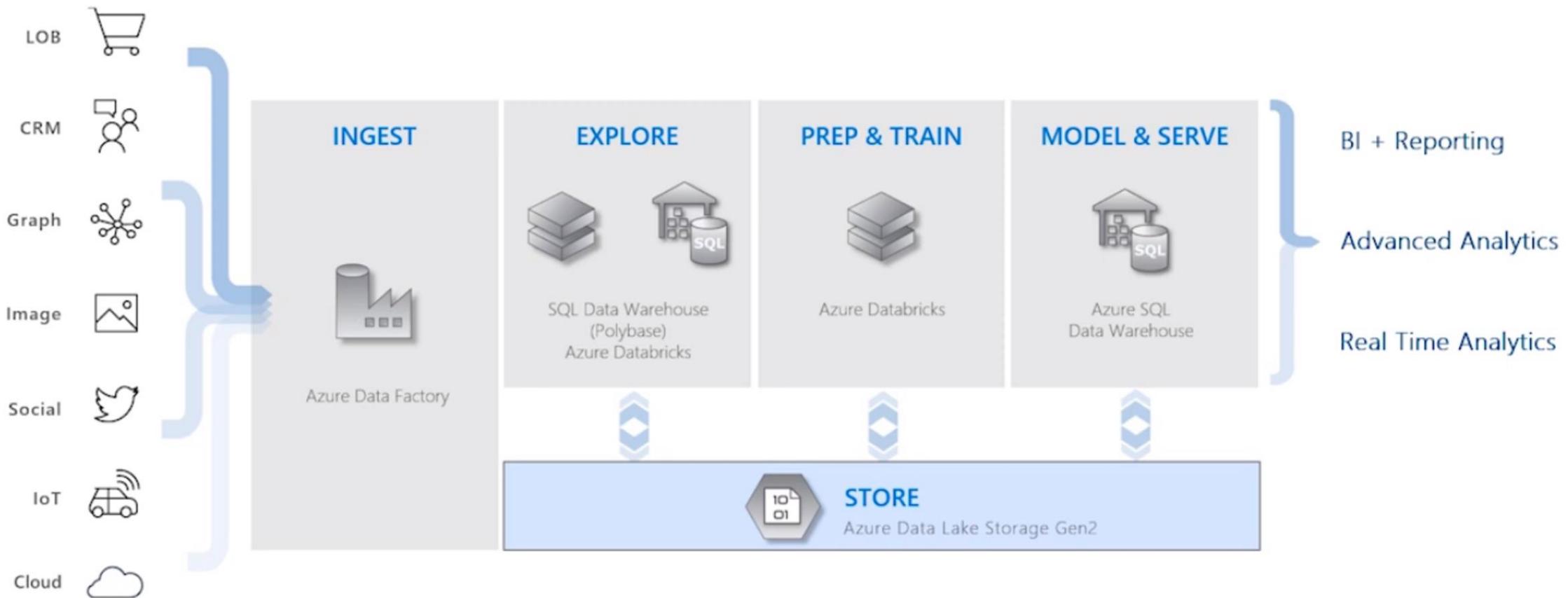




Old data warehouse

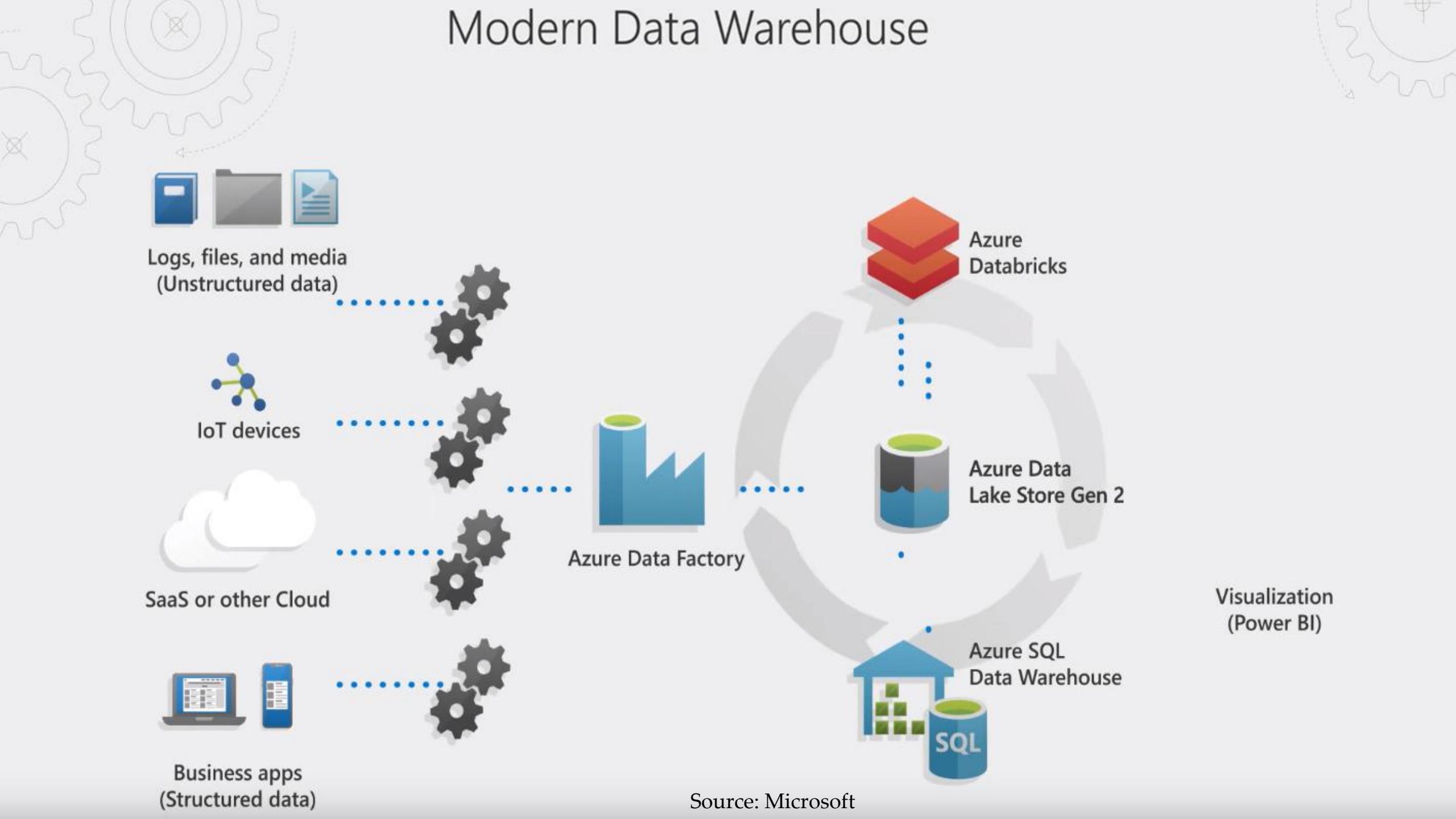
Source: Microsoft

# Modern Data Warehousing

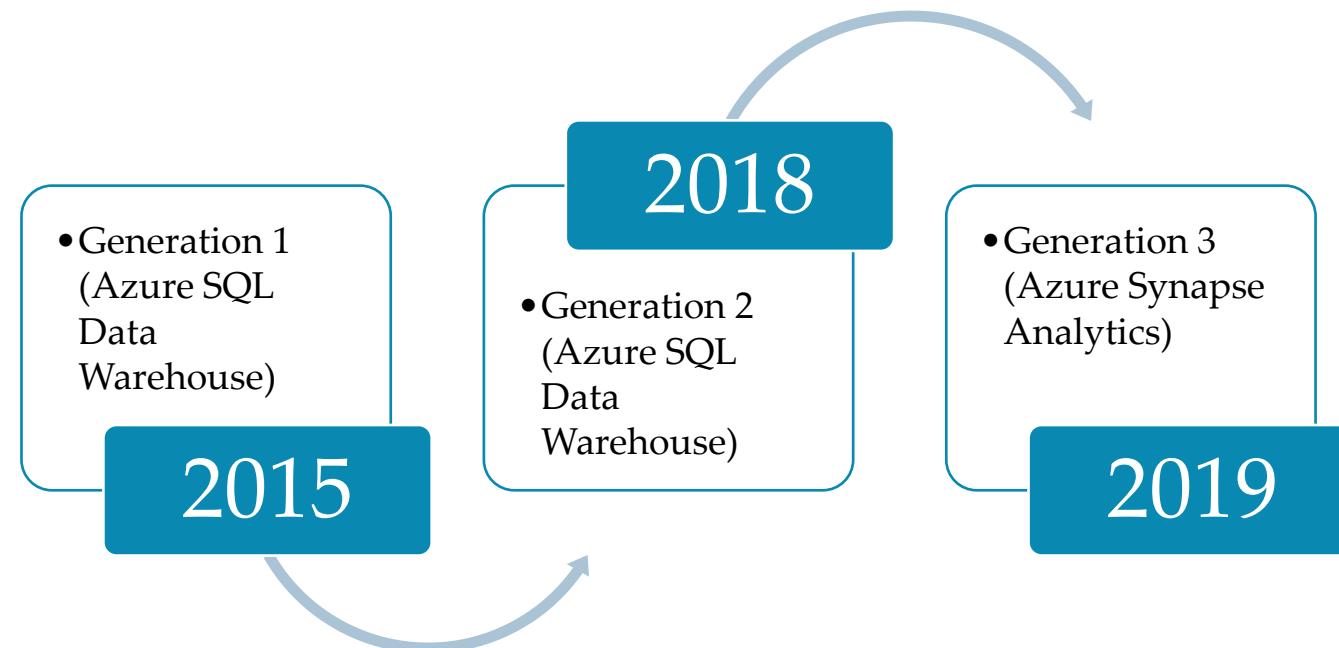


Source: Microsoft

# Modern Data Warehouse



# Azure SQL Data warehouse



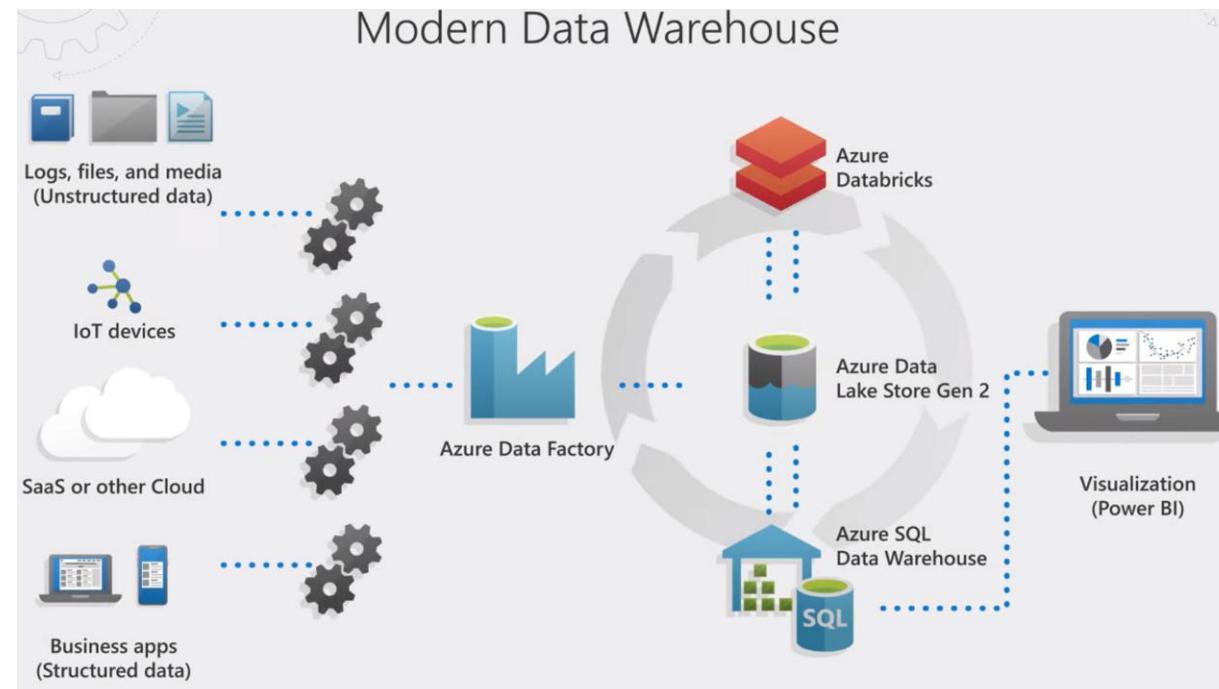
# Azure Synapse Analytics

*"Synapse is the next generation of Azure SQL Data Warehouse, blending big data analytics, data warehousing, and data integration into a single unified service that provides end-to-end analytics with limitless scale."*

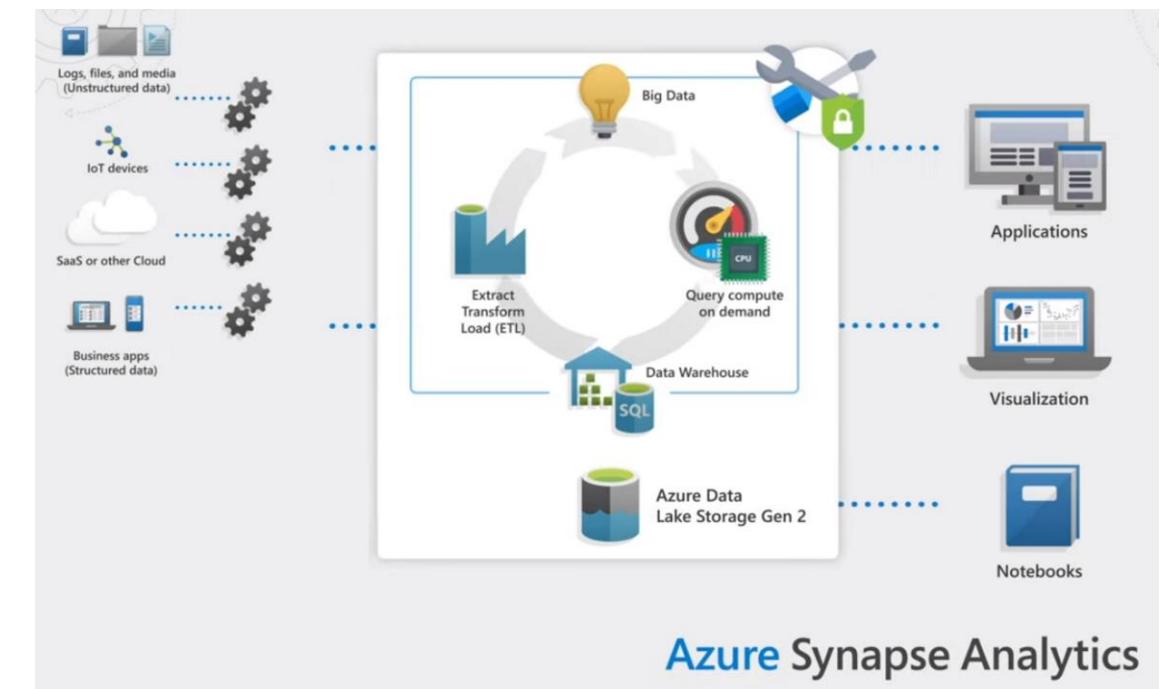
# Azure Synapse Analytics



# Modern vs Synapse Architecture

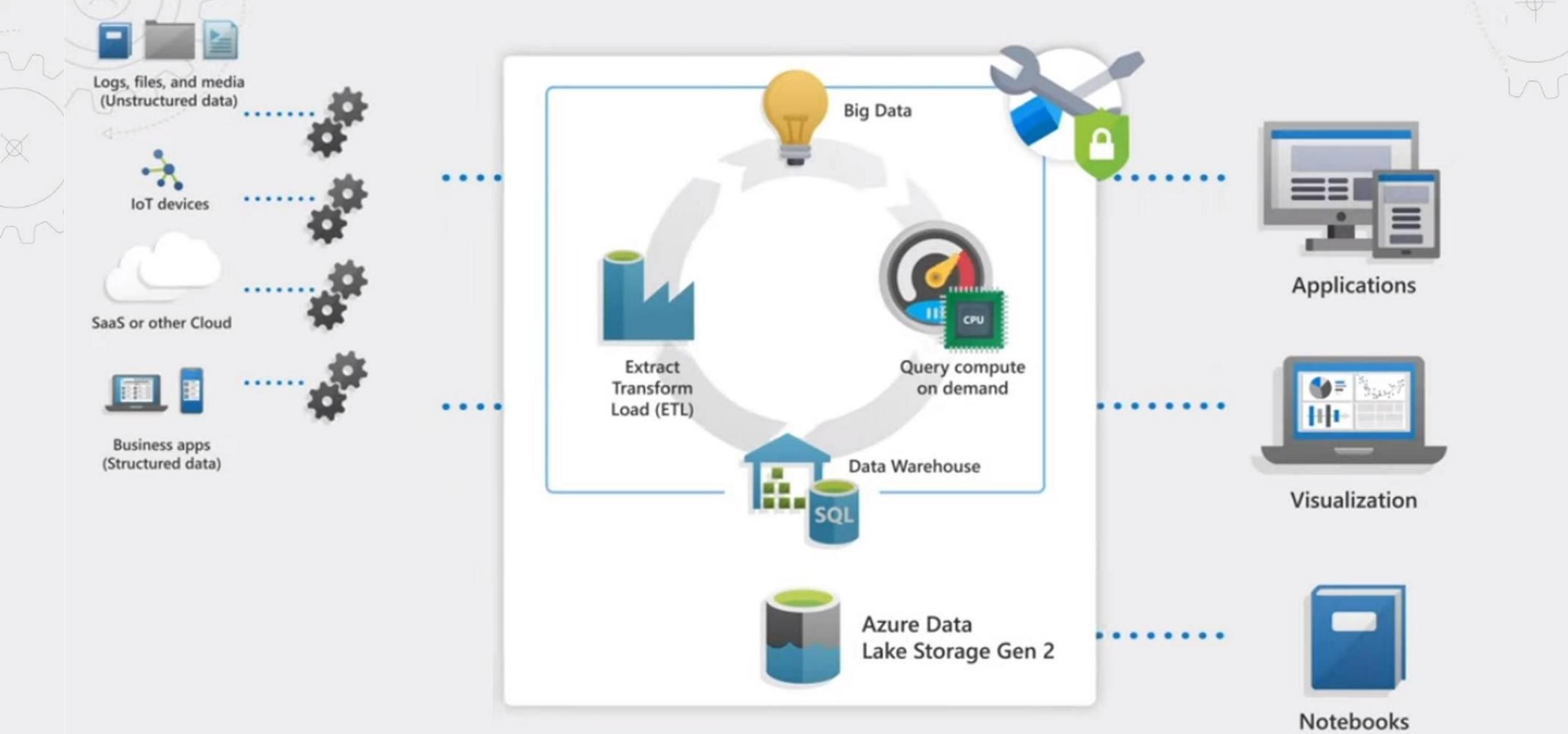


Modern Data Warehouse Architecture

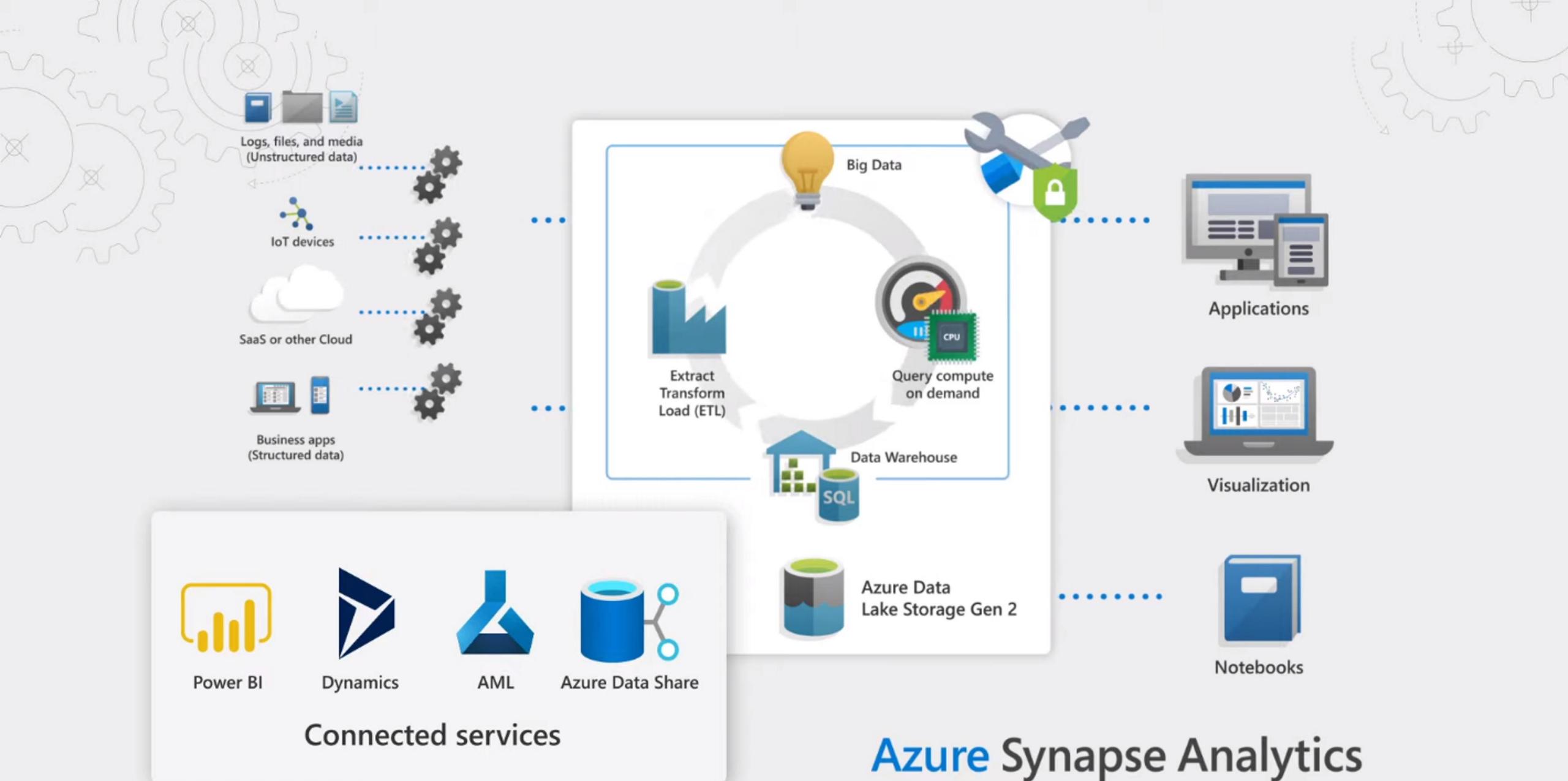


Azure Synapse Analytics

Synapse Analytics Architecture



# Azure Synapse Analytics



# Azure Synapse Analytics



# Demo – Provision Azure Synapse Service

1. Create SQL Server
2. Create Synapse SQL Pool (Azure SQL Data Warehouse)
3. Pause/Resume Compute Node
4. Create Firewall Rule
5. Connect with Microsoft SQL Server Management Studio

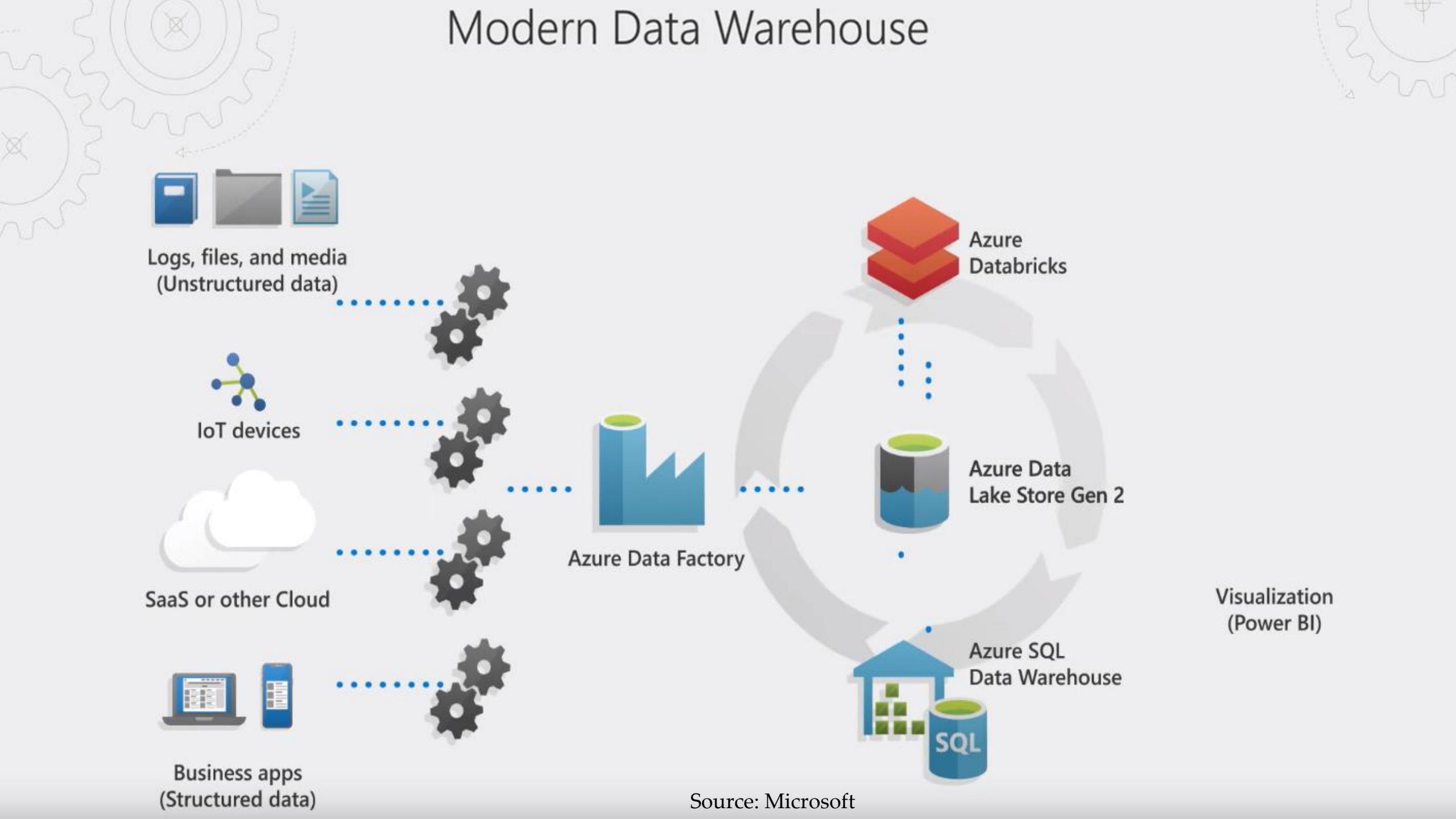


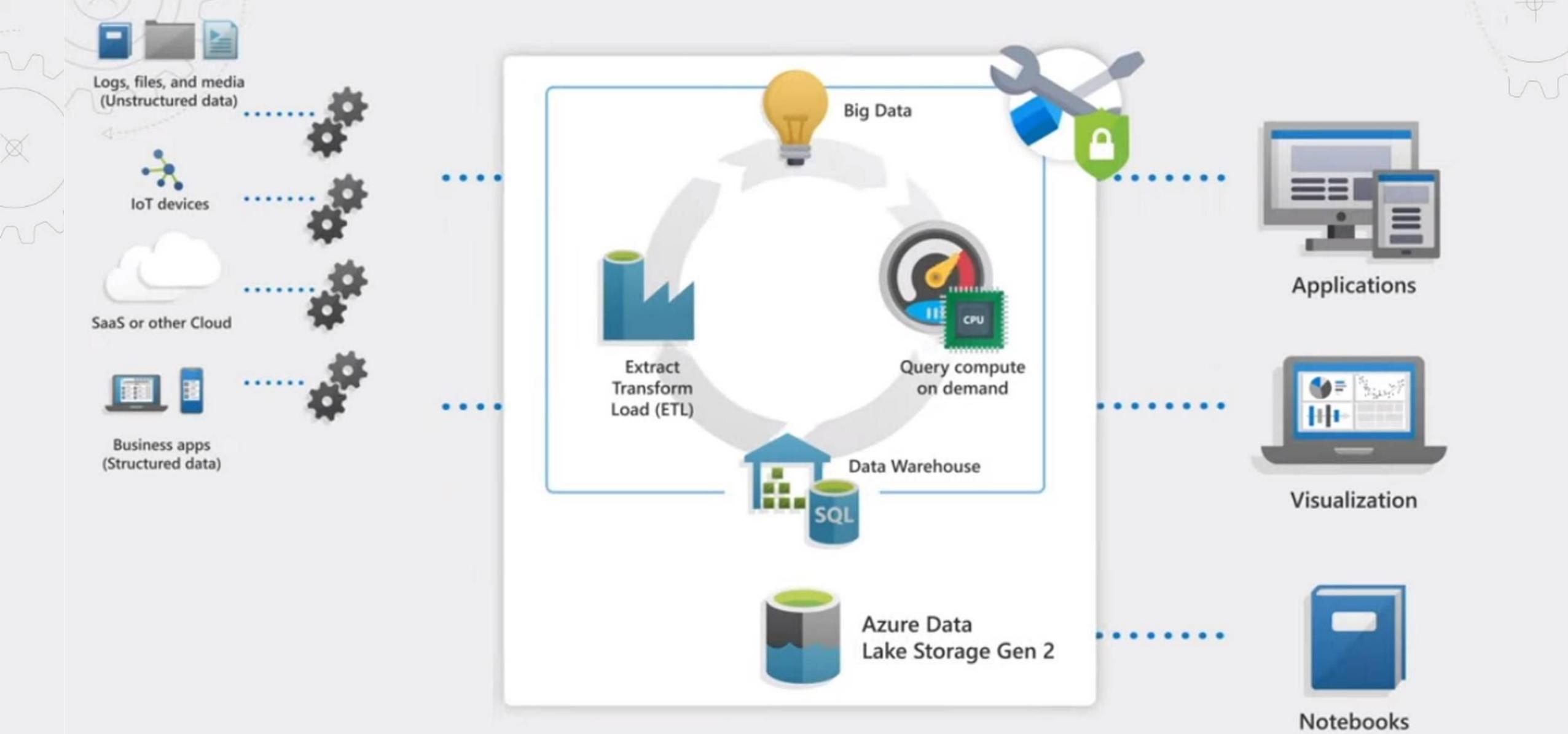


Azure  
Synapse  
Analytics

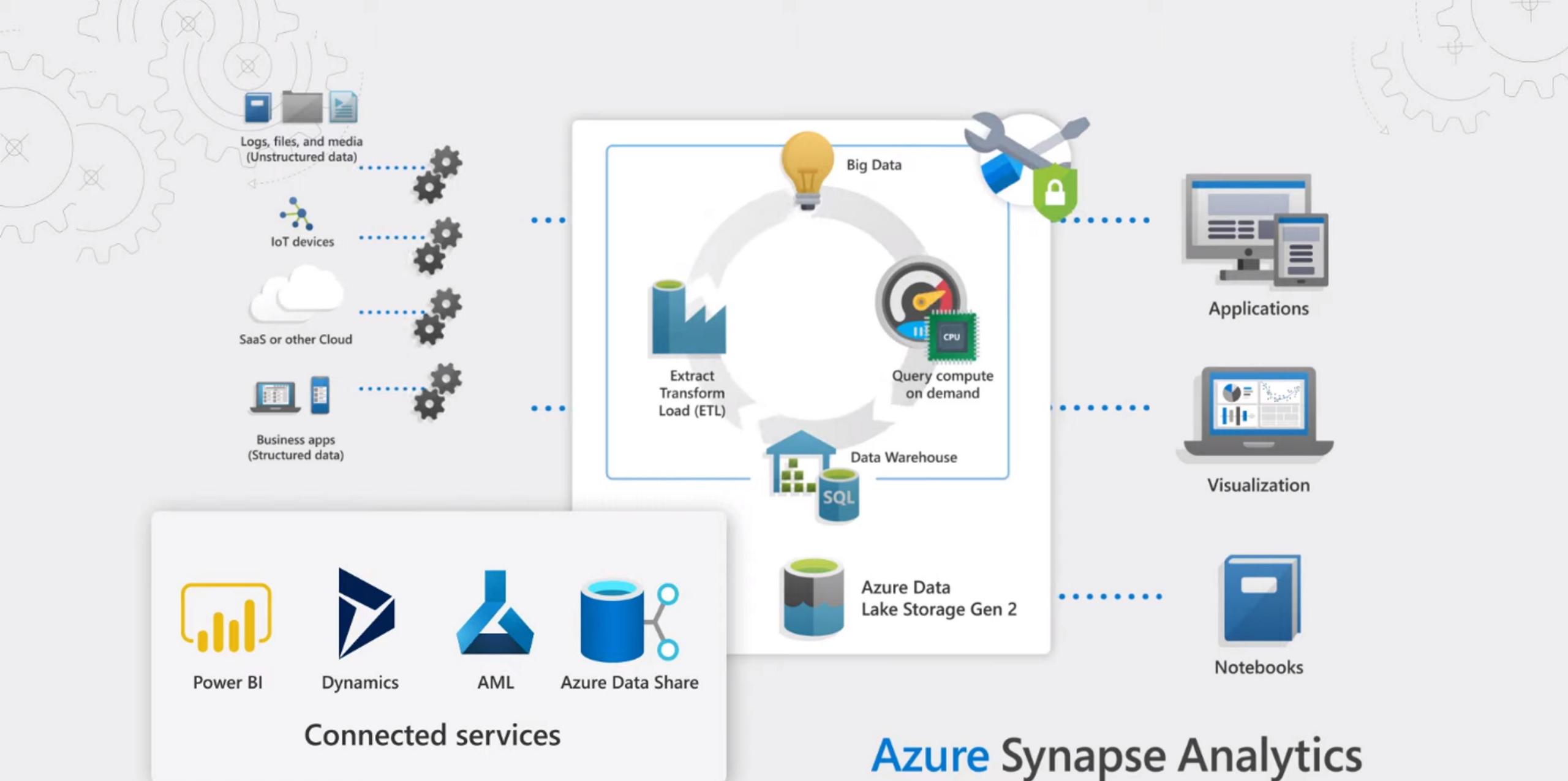


# Modern Data Warehouse



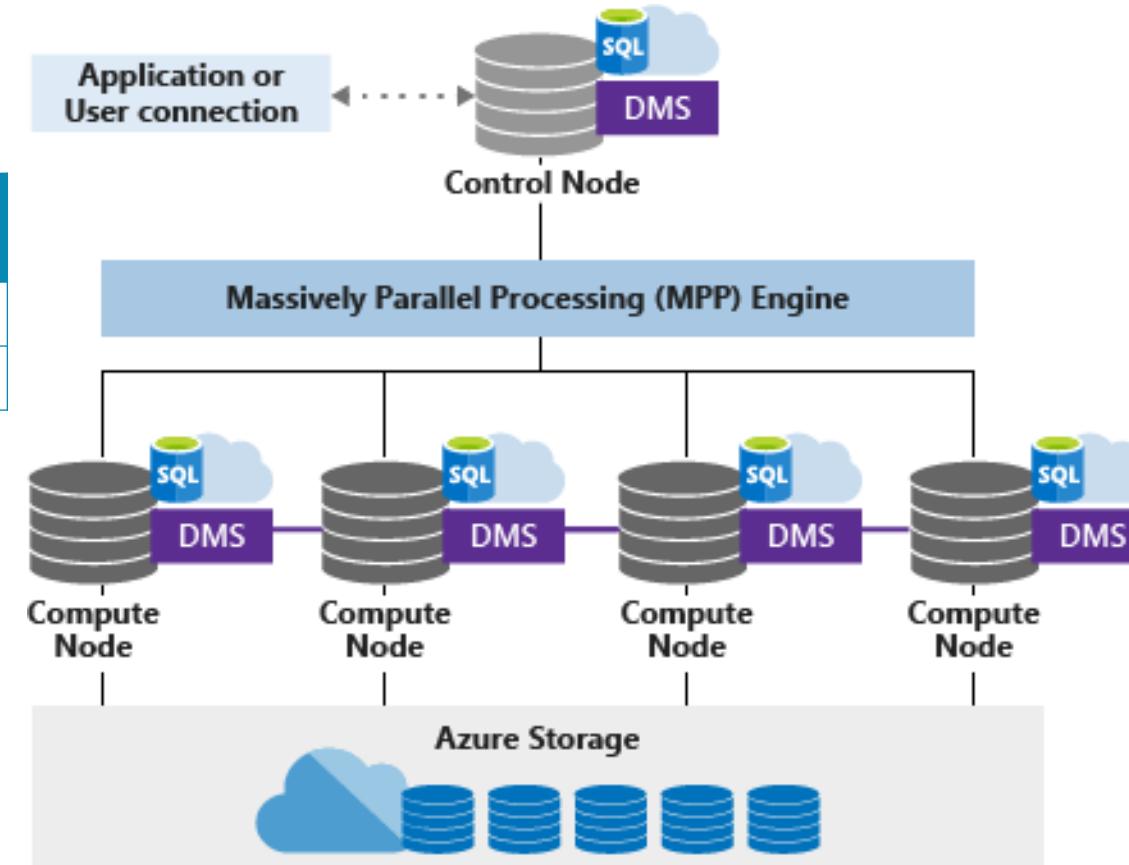


# Azure Synapse Analytics

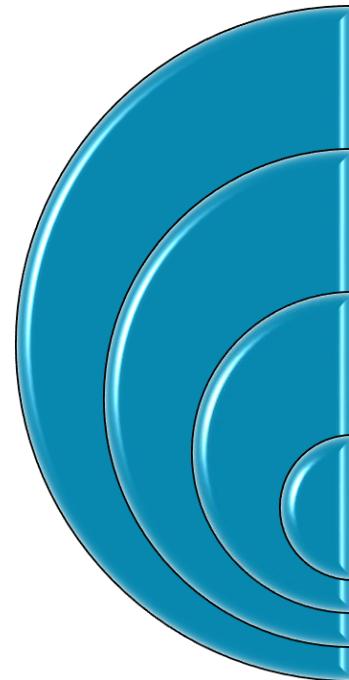


# Azure Synapse MPP Architecture

| DWU | Loading<br>3 Tables | Ran<br>Report |
|-----|---------------------|---------------|
| 100 | 15                  | 20            |
| 500 | 3                   | 4             |



# Azure Storage and Distribution



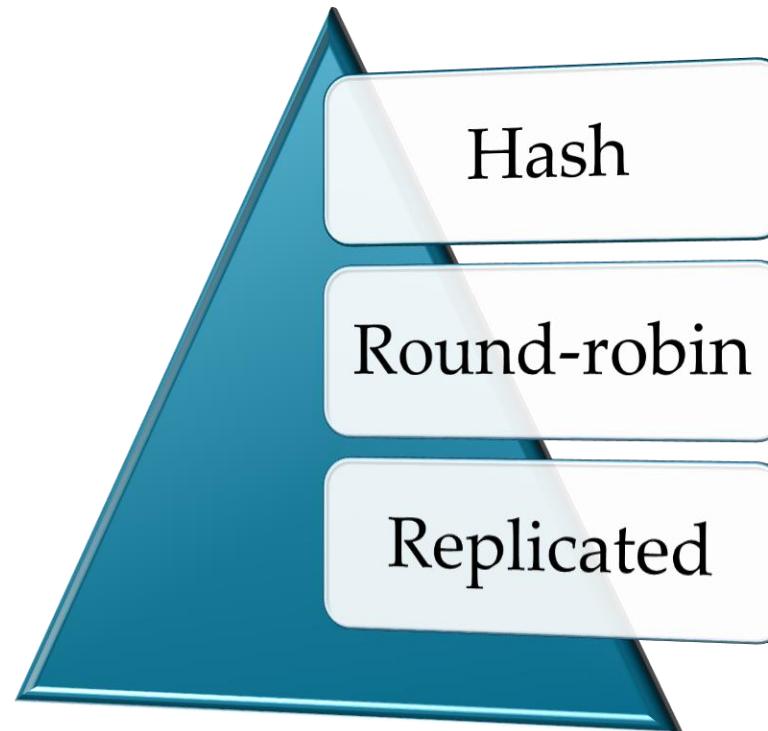
SQL DW charges separately for storage consumption

A distribution is the basic unit of storage and processing for parallel queries

Rows are stored across 60 distributions which are run in parallel

Each compute node manages one or more of the 60 distribution

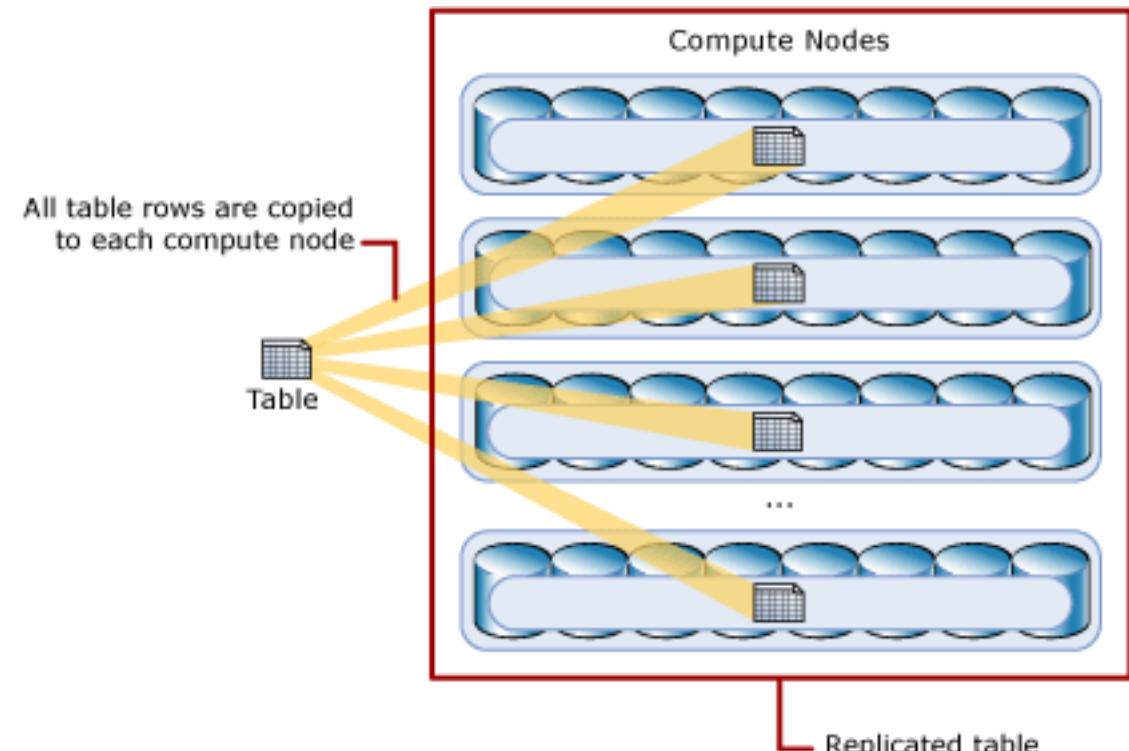
# Sharding Patterns



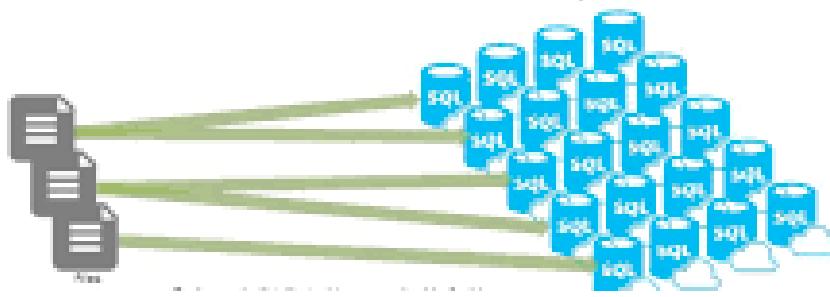
# Replicated Tables

- Caches a full copy on each compute node.
- Used for small tables

```
CREATE TABLE [dbo].[BusinessHierarchies](  
    [BookId] [nvarchar](250) ,  
    [Division] [nvarchar](100) ,  
    [Cluster] [nvarchar](100) ,  
    [Desk] [nvarchar](100) ,  
    [Book] [nvarchar](100) ,  
    [Volcker] [nvarchar](100) ,  
    [Region] [nvarchar](100)  
)  
WITH  
(  
    CLUSTERED COLUMNSTORE INDEX  
    DISTRIBUTION = REPLICATE  
)  
;
```



# Round Robin tables

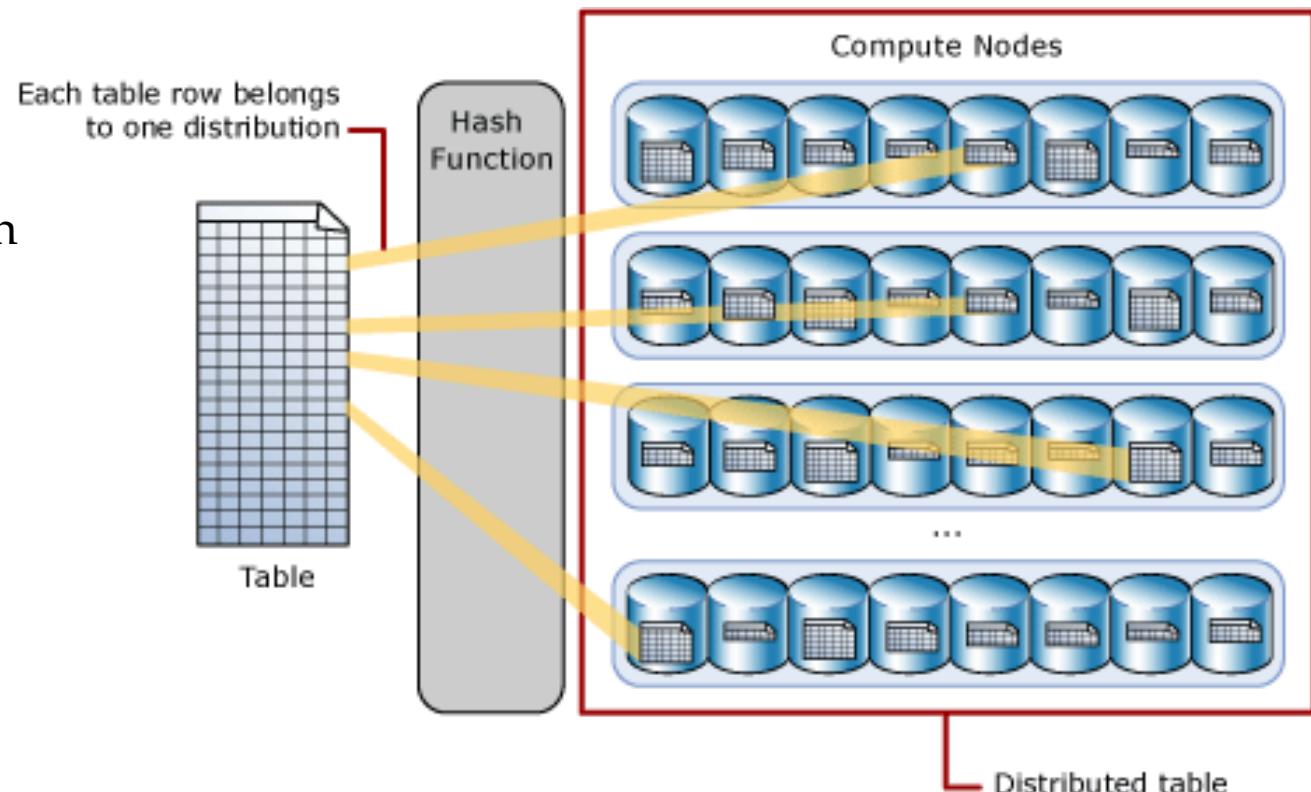


- Generally use to load staging tables
- Distribute data evenly across the table without additional optimization
- Joins are slow, because it requires to reshuffle data
- Default distribution type

```
CREATE TABLE [dbo].[Dates](
    [Date] [datetime2](3) ,
    [DateKey] [decimal](38, 0) ,
    ...
    ...
    [WeekDay] [nvarchar](100) ,
    [Day Of Month] [decimal](38, 0)
)
WITH
(
    CLUSTERED COLUMNSTORE INDEX
    DISTRIBUTION = ROUND_ROBIN
)
;
```

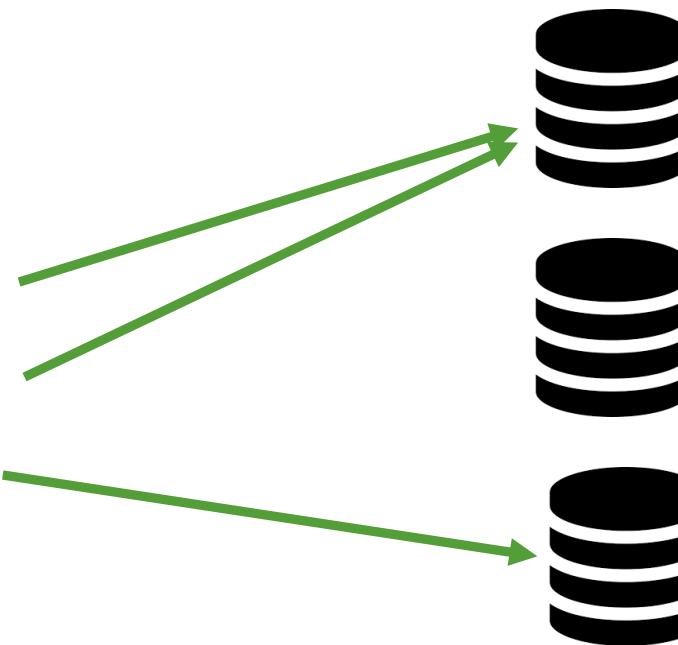
# Hash Distribution Tables

- Highest performance for large tables
- Each row belongs to one particular distribution
- It is used mostly for larger tables



# Hash Distribution Tables

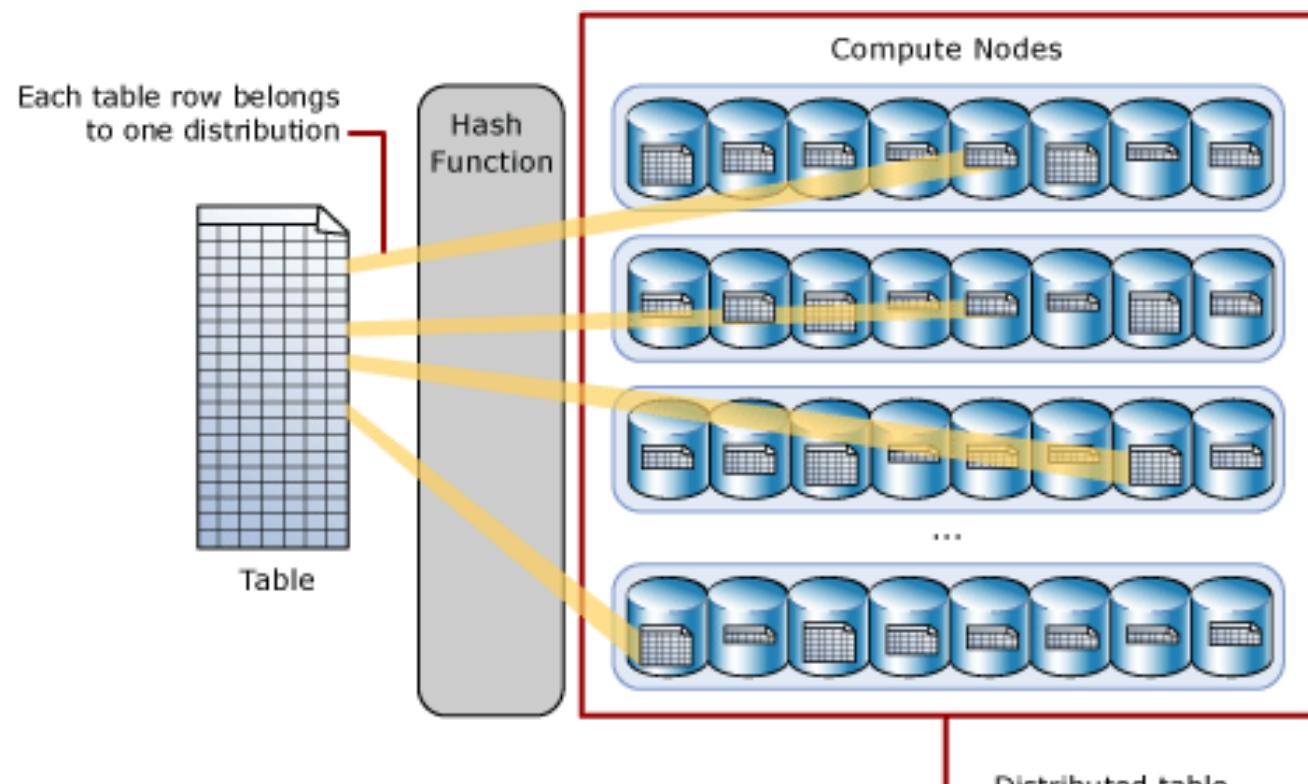
| Record | Product  | Store       |
|--------|----------|-------------|
| 1      | Soccer   | New York    |
| 2      | Soccer   | Los Angeles |
| 3      | Football | Phoenix     |



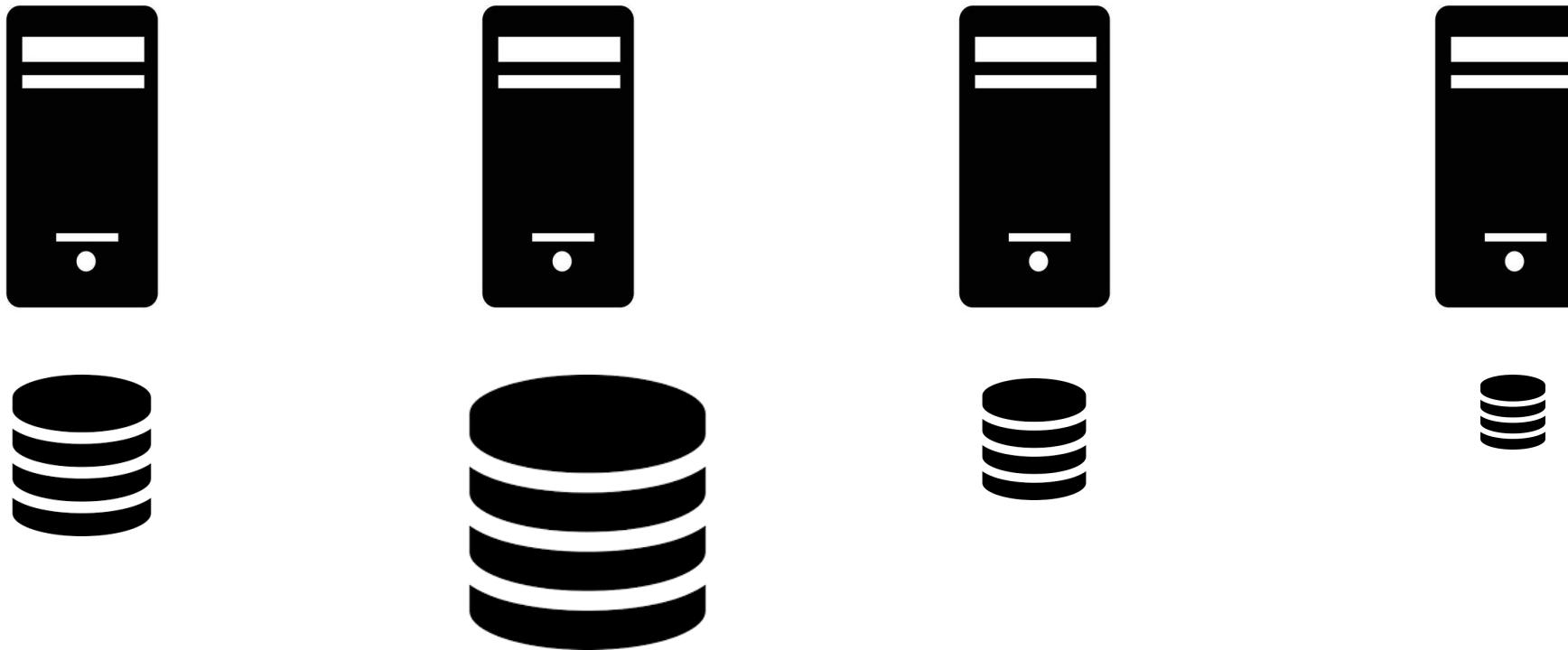
# Hash Distribution Tables

- Highest performance for large tables
- Each row belongs to one particular distribution
- It is used mostly for larger tables

```
CREATE TABLE [dbo].[EquityTimeSeriesData](
    [Date] [varchar](30) ,
    [BookId] [decimal](38, 0) ,
    [P&L] [decimal](31, 7) ,
    [VaRLower] [decimal](31, 7)
)
WITH
(
    CLUSTERED COLUMNSTORE INDEX
    , DISTRIBUTION = HASH([P&L])
)
;
```



# Avoid Data Skew



# Even Distribution

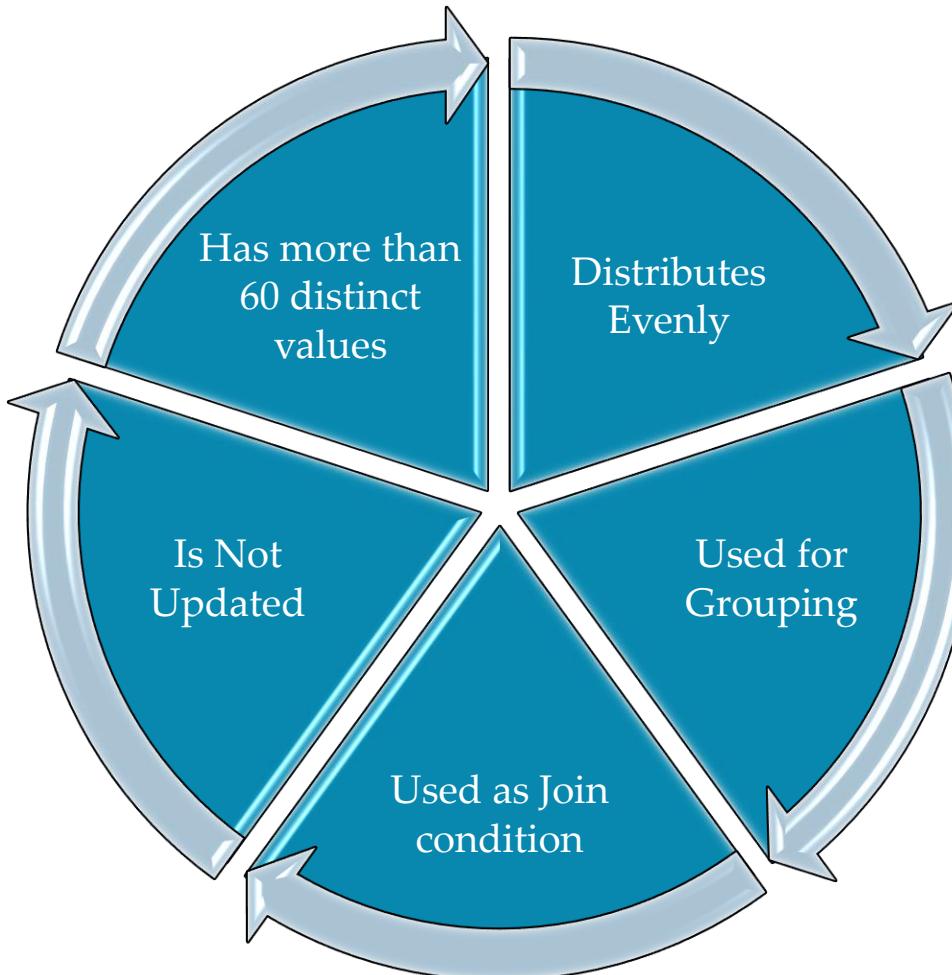


# Distribution key

Determines the method in which Azure SQL Data Warehouse spreads the data across multiple nodes.

Azure SQL Data Warehouse uses up to 60 distributions when loading data into the system.

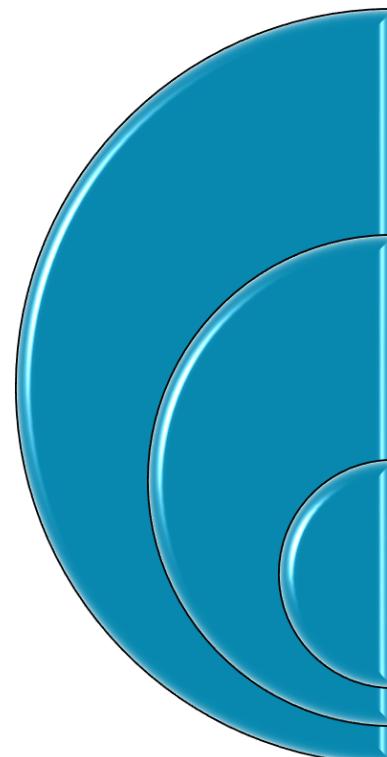
# Good Hash Key



# What Data Distribution to Use?

| Type                  | Great fit for  | Watch out if...  |
|-----------------------|--|--|
| Replicated            | Small-dimension tables in a star schema with less than 2GB of storage after compression  | <ul style="list-style-type: none"><li>Many write transaction are on the table (insert/update/delete)</li><li>You change DWU provisioning frequently</li><li>You use only 2-3 columns, but your table has many columns</li><li>You index a replicated table</li></ul> |
| Round-robin (default) | <ul style="list-style-type: none"><li>Temporary/Staging table</li><li>No obvious joining key or good candidate column.</li></ul> | Performance is slow due to data movement   |
| hash                  | <ul style="list-style-type: none"><li>Fact tables</li><li>Large dimension tables</li></ul>                                       | The distribution key can't be updated  |

# Data types



Use the smallest data type which will support your data

Avoid defining all character columns to a large default length

Define columns as VARCHAR rather than NVARCHAR if you don't need Unicode

# Data types



The goal is to not only save space but also move data as efficiently as possible.

# Data types



Some complex data types (XML, geography, etc)  
are not supported on Azure SQL Data  
Warehouse yet.

# Table types

Clustered  
columnstore

- Updateable primary storage method
- Great for read-only

Heap

- Data is not in any particular order.
- Use when data has no natural order.

Clustered Index

- An index that is physically stored in the same order as the data being indexed

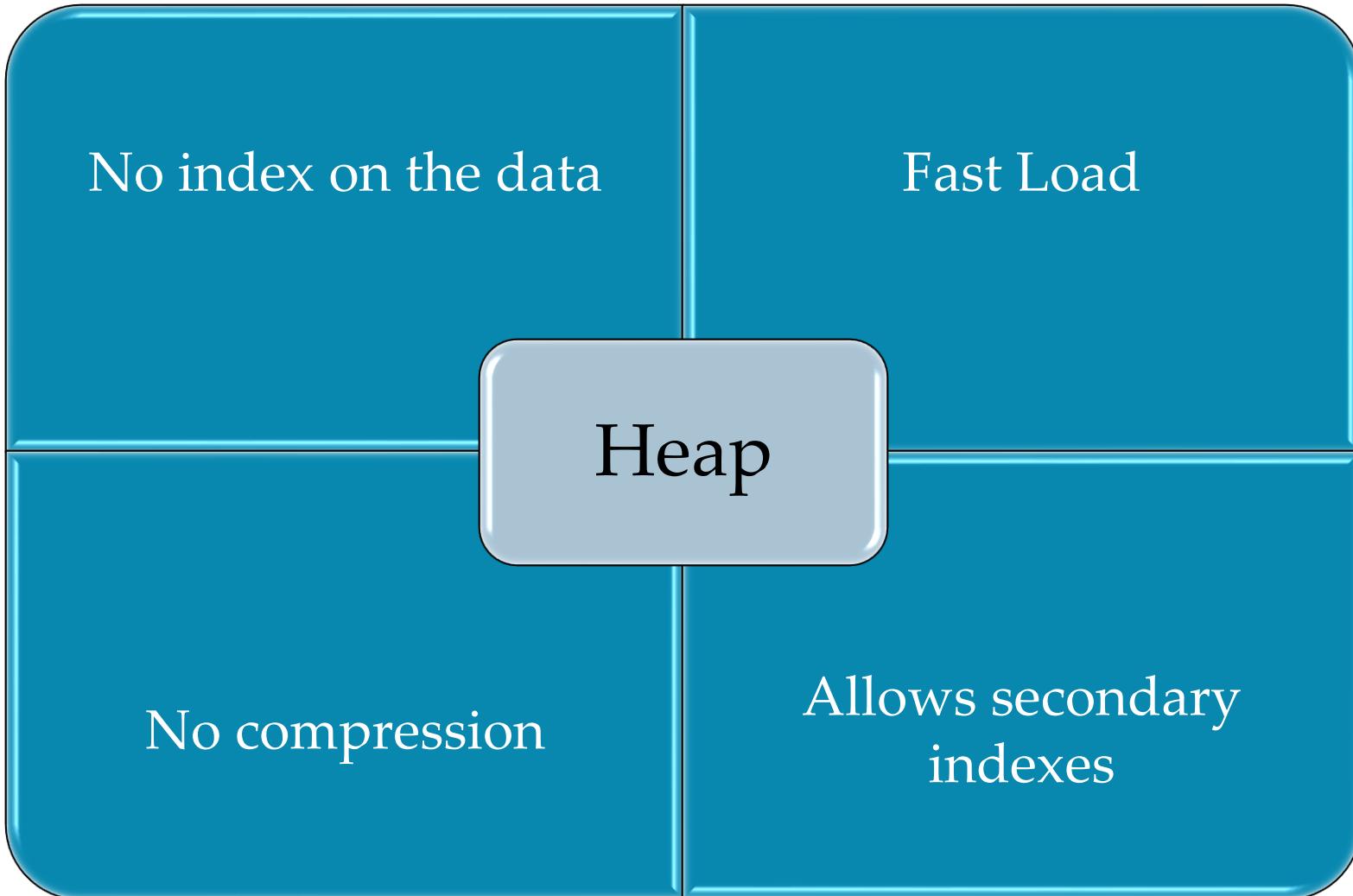
Default table type

High compression  
ratio

Clustered  
columnstore

Ideally segments of  
1M rows

No Secondary  
Indexes



## Clustered B-Tree

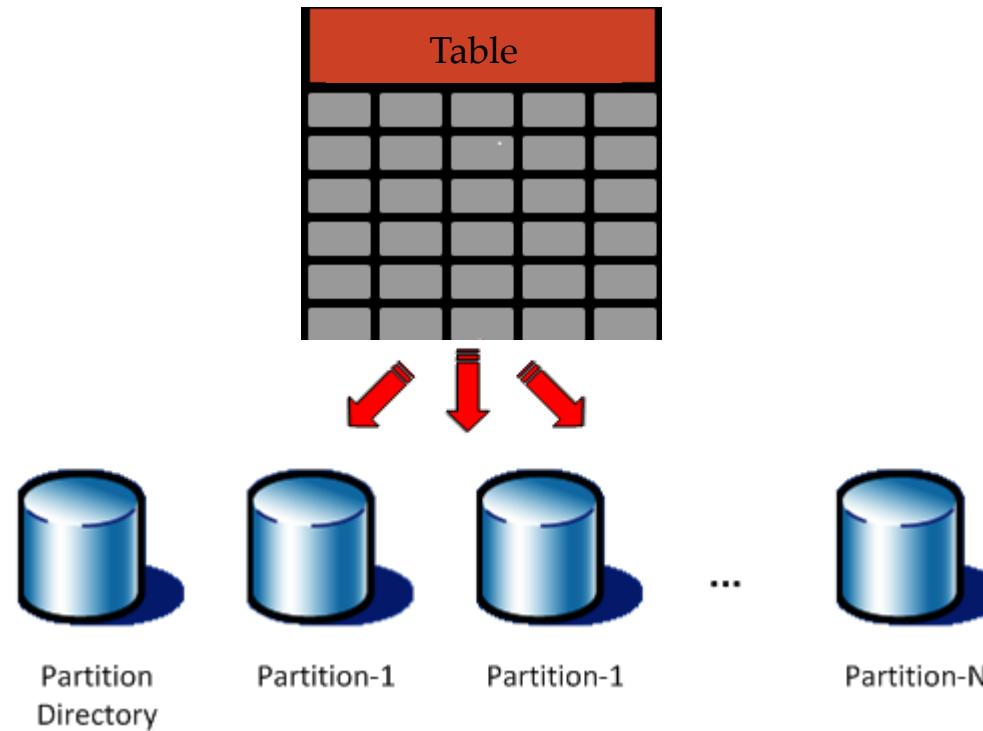
Sorted index on the data

Fast singleton lookup

No compression

Allows secondary  
indexes

# Table Partitioning



# Partitioning



Table partitions enable you to divide your data into smaller groups of data

Improve the efficiency and performance of loading data by use of partition deletion, switching and merging

Usually data is partitioned on a date column tied to when the data is loaded into the database

Can also be used to improve query performance

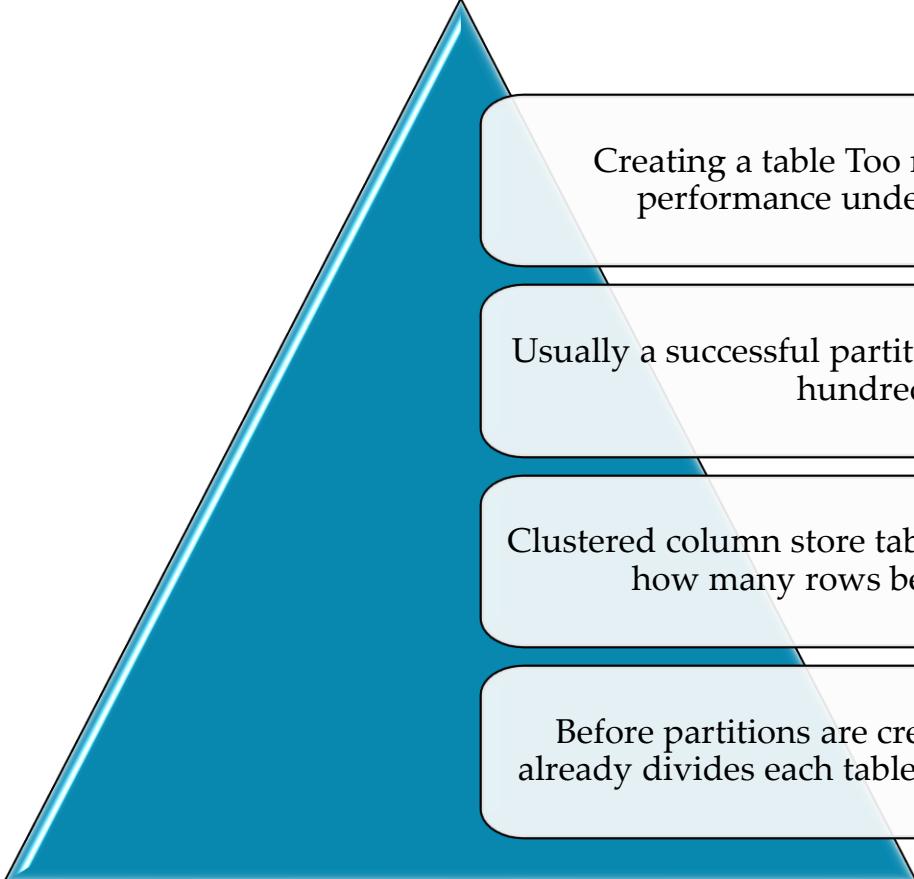
# Why Partitioning?

Easy load or  
unload data

Easy  
maintenance,  
rebuids or  
reorganizes

Improve Query  
Performance

# Partitions best practices



Creating a table Too many partitions can hurt performance under some circumstances

Usually a successful partitioning scheme has 10 or a few hundred partitions

Clustered column store tables, it is important to consider how many rows belong to each partition

Before partitions are created, SQL Data warehouse already divides each table into 60 distributed databases



A highly granular partitioning scheme can work in SQL Server but hurt performance in Azure SQL Data Warehouse.

# Example

60 Distributions



365 Partitions



21900 Data Buckets

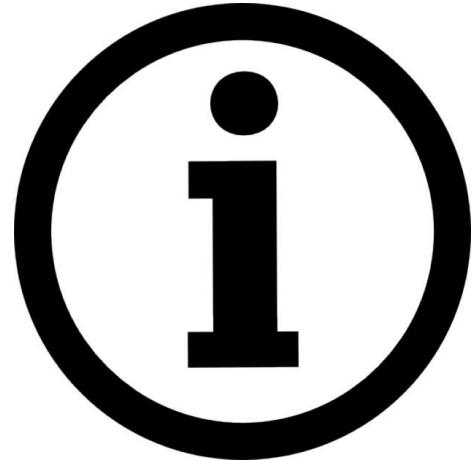
21900 Data Buckets



Ideal Segment  
Size (1M Rows)



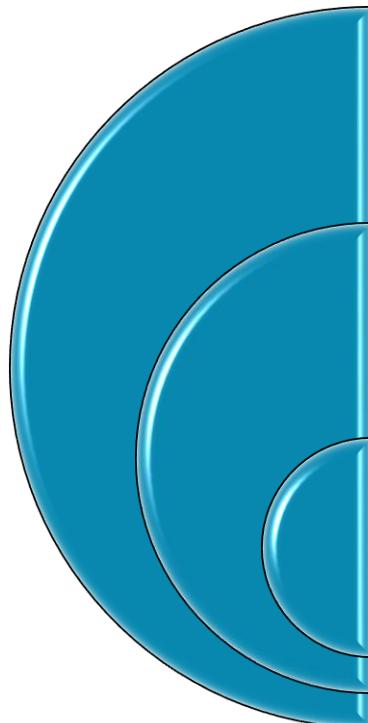
21 900 000 000 Rows



Lower Granularity (week, month)  
can perform better depending on  
how much data you have.

How do we apply these principles to a  
Dimensional model?

# Fact Tables

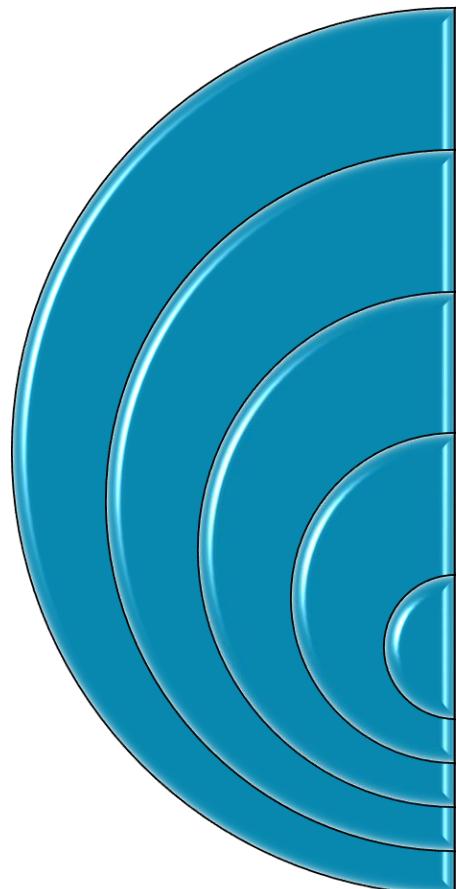


Large ones are better as Columnstores

Distributed through Hash key as much as possible as long as it is even

Partitioned only if the table is large enough to fill up each segment

# Dimension Tables



Can be Hash distributed or Round-Robin if there is no clear candidate join key

Columnstore for large dimensions

Heap or Clustered Index for small dimensions

Add secondary indexes for alternate join columns

Partitioning not recommended

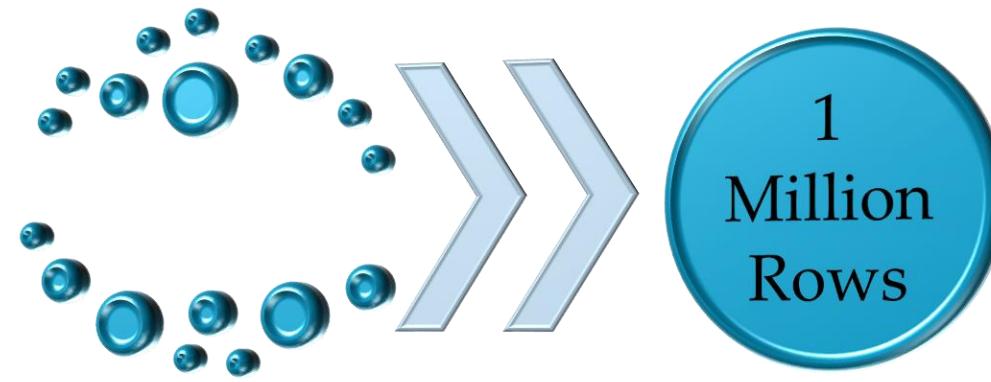
# Best Practices for Data Laod

# Data Warehouse Readers

Your DWUs have a direct impact on how fast you can load data in parallel

# Optimize Insert Batch Size

Avoid trickle insert pattern. Ideal batch size is 1 million or more direct or in a file

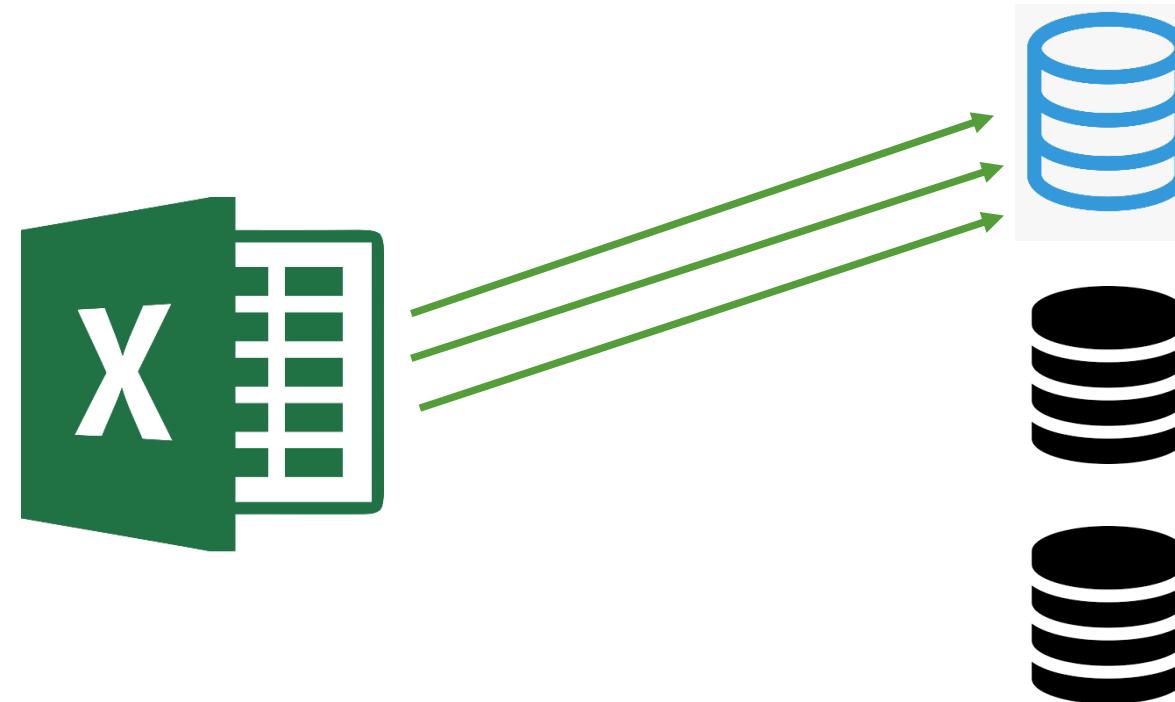


Combine rows to  
make it a batch of  
1 million

Ideal batch size

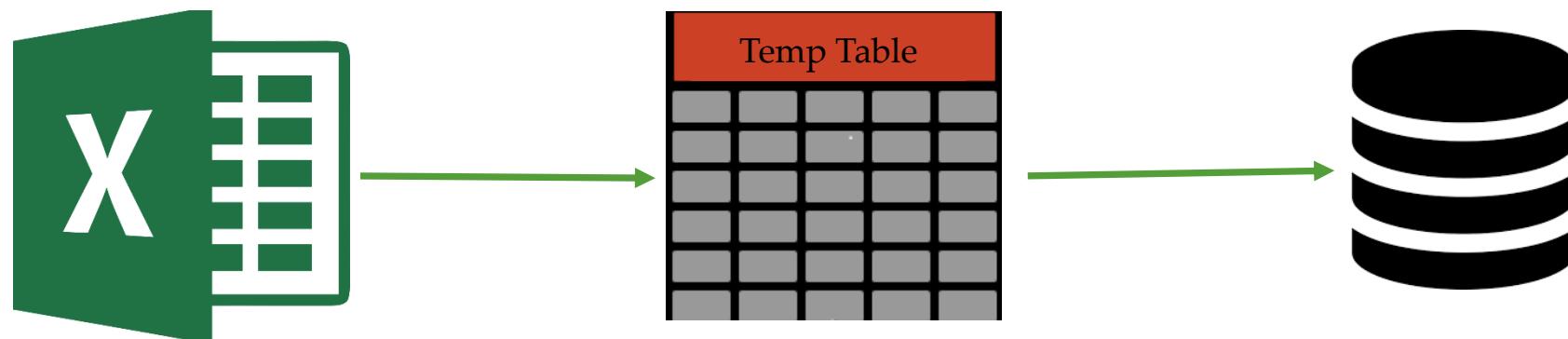
# Avoid ordered data

Data ordered by distribution key can introduce hot spots that slow down the load operation



# Using temporary tables

Stage and transform on a Temp Heap table before moving to permanent storage.



# CREATE TABLE AS

```
CREATE TABLE #tmp_fct
WITH
(
DISTRIBUTION = ROUND_ROBIN
)
AS
SELECT *
FROM
[dbo].[FactInternetSales];
```

- Fully Parallel operation
- It is minimally logged
- It can change: distribution, table type, partitioning

# Loading Methods

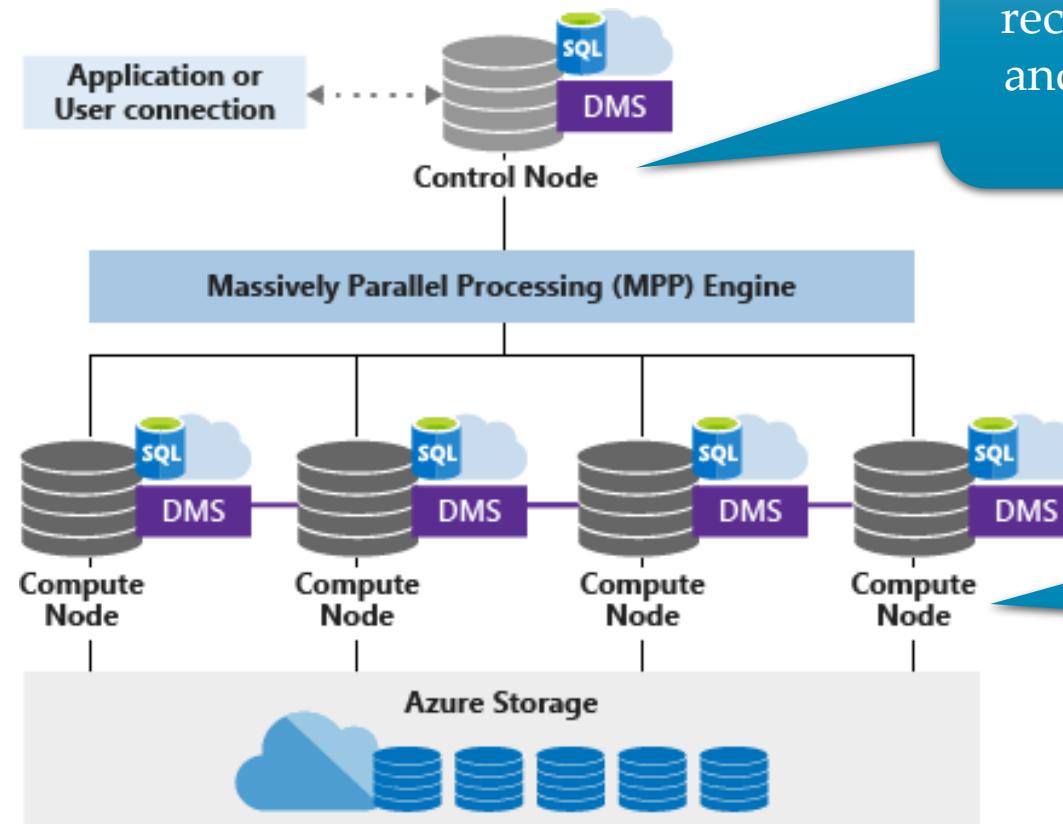
## Single client loading methods

- SSIS
- Azure Data Factory
- BCP
- Can add some parallel capabilities but are bottlenecked at the control node

## Parallel readers loading methods

- PolyBase
- Reads from Azure blob Storage and loads the contents into Azure SQL DW
- Bypasses the Control Node and loads directly into the Compute Nodes

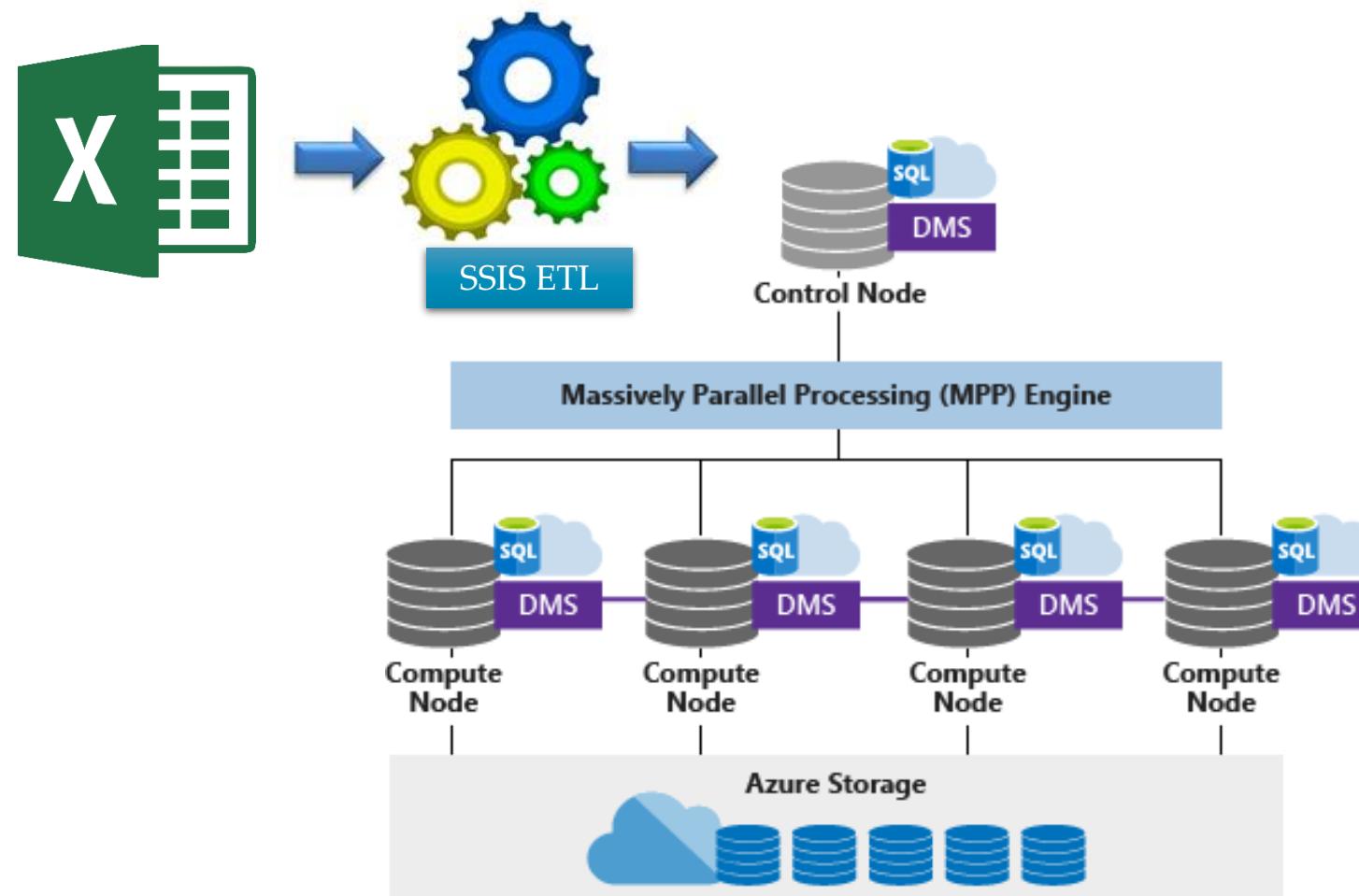
# Control Node



The **Control** node receives connections and orchestrates the queries

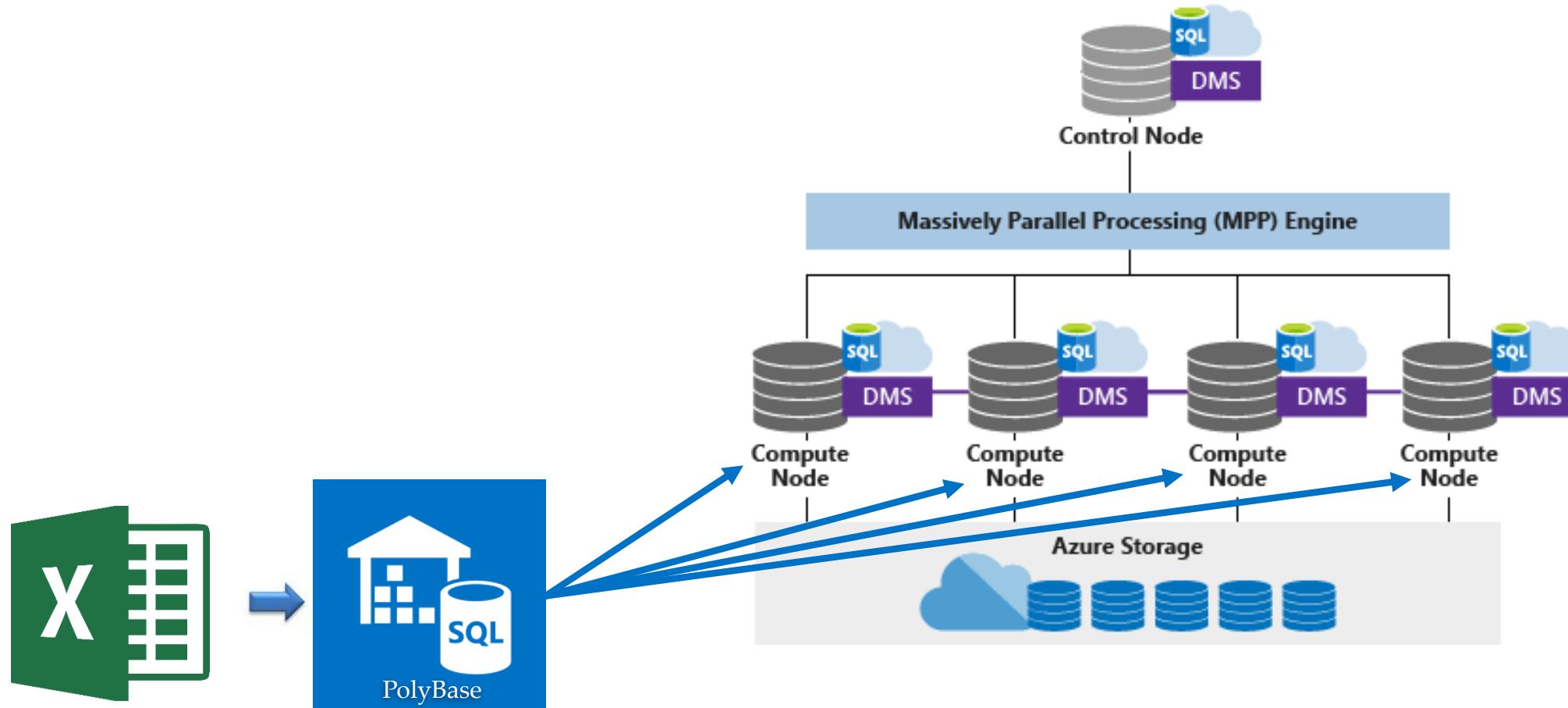
The **Compute** nodes do processing on the data and scale with the DWUs.

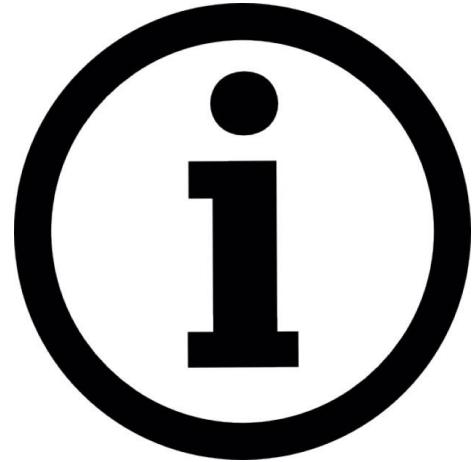
# Loading with SSIS



Source: Microsoft

# Loading with PolyBase



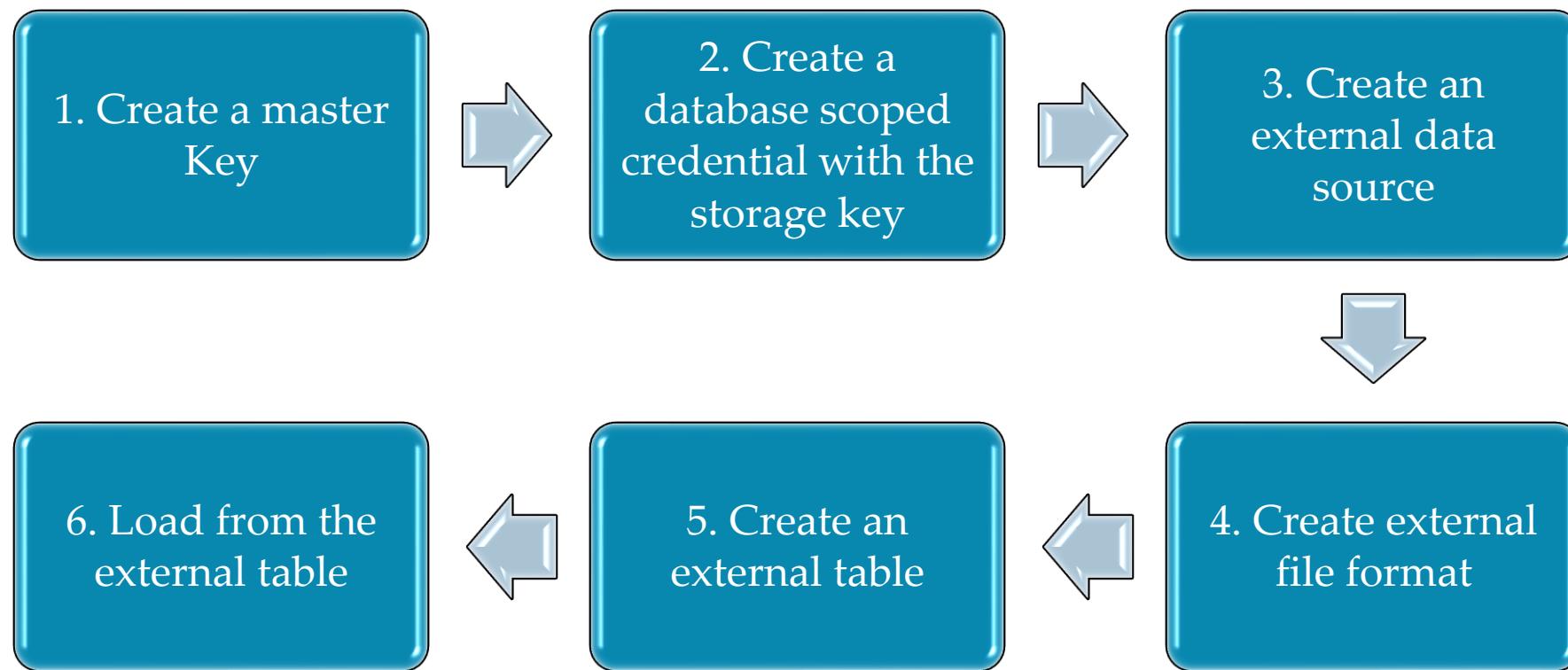


PolyBase can load data from UTF-8 delimited text files and popular Hadoop file formats like RC File, ORC, and Parquet



Multiple readers will not work against a  
compressed file

# PolyBase Setup



# DEMO

Load Dimension table from on-premises to Azure Synapse  
Data warehouse using SSIS

# Demo – PolyBase

1. Export table to flat file
2. Create blob storage account
3. Upload flat file to blob storage
4. Run PolyBase 6 steps process
5. Monitor and confirm successful migration
6. Confirm 60 distributions in destination table



Azure  
Synapse  
Analytics