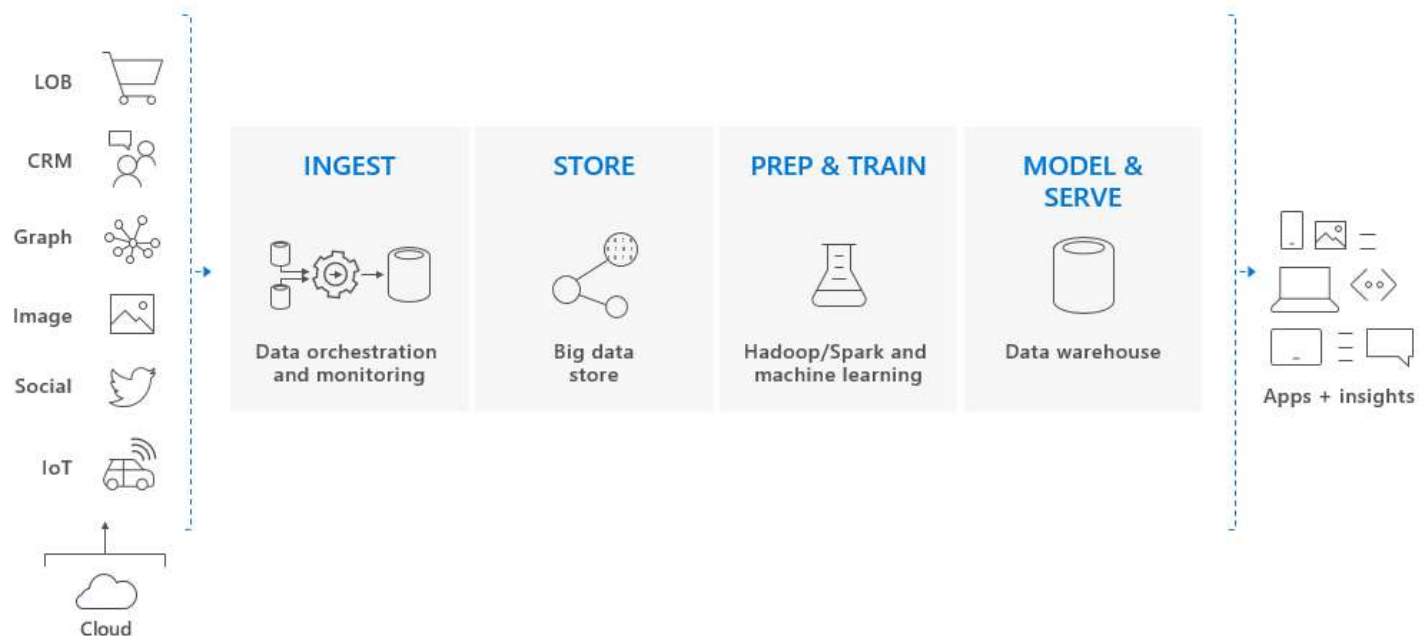# Caloudi 教育訓練課程

# Azure Synapse Analytics Workshop

# Agenda

- Data Analytics Evolution

- Data Lake Evolution

- Azure Synapse Analytics

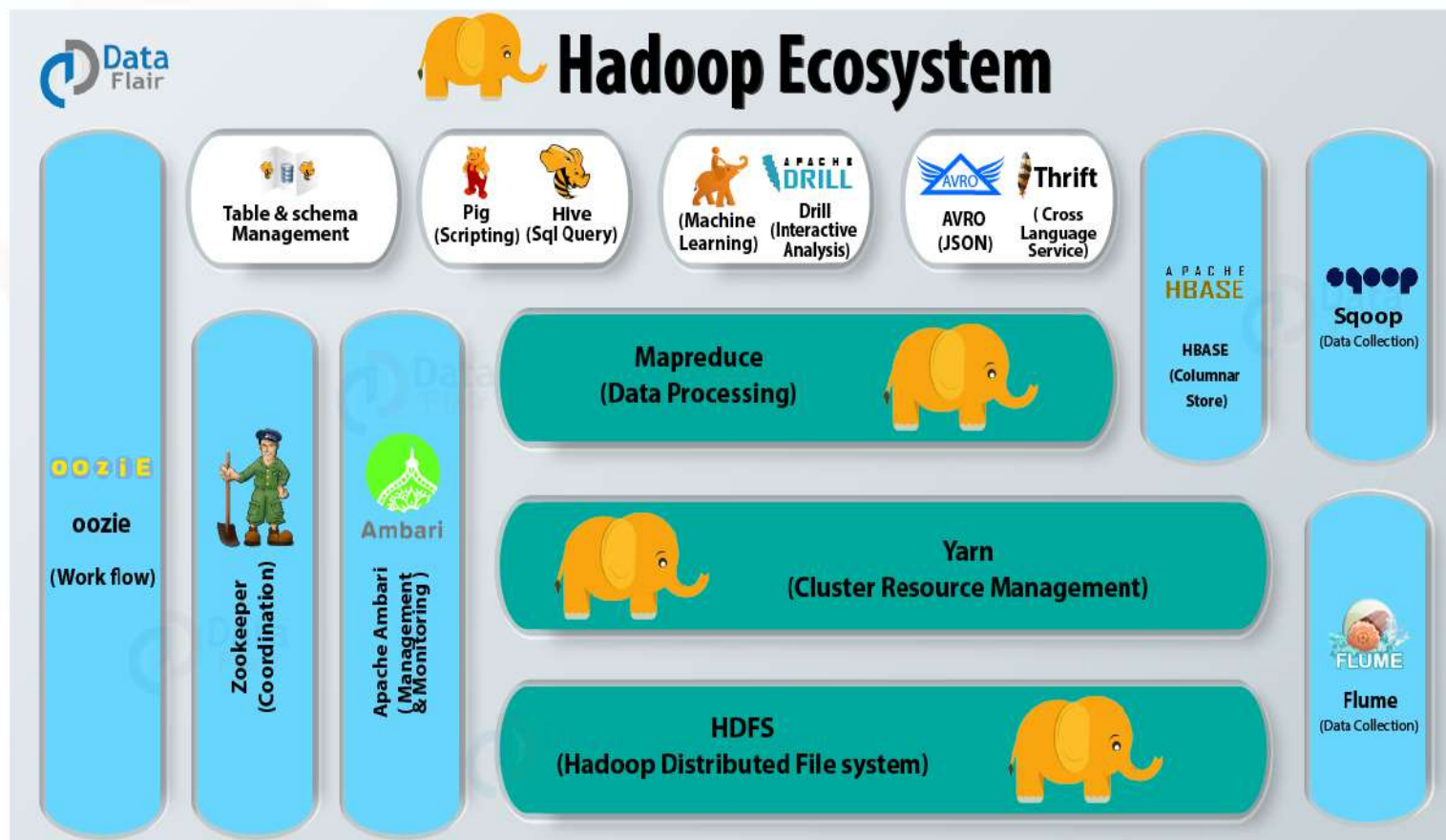- SQL DW / Dedicated SQL Pool

- Implementation Issues

- References

# Data Analytics Evolution

# Data Warehouse Solution



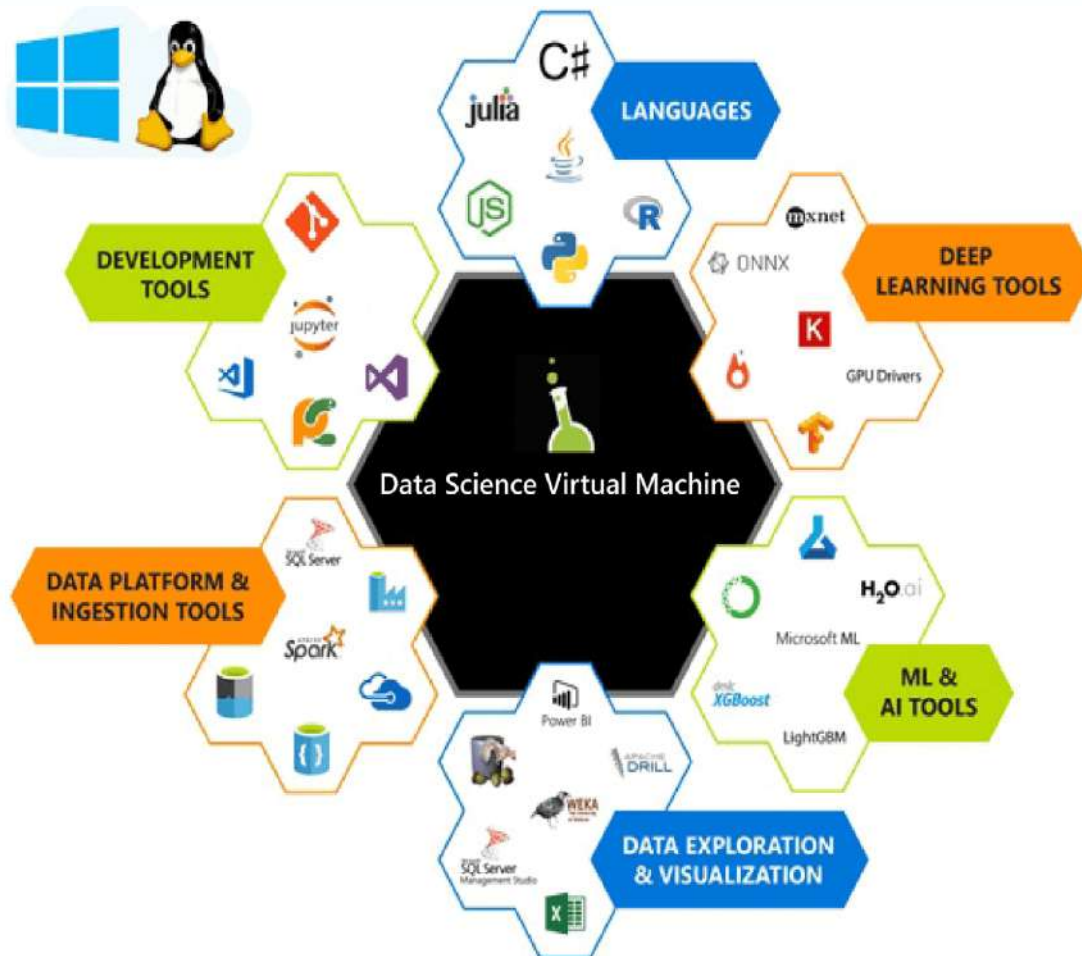參考資料：What is dedicated SQL pool (formerly SQL DW) in Azure Synapse Analytics?
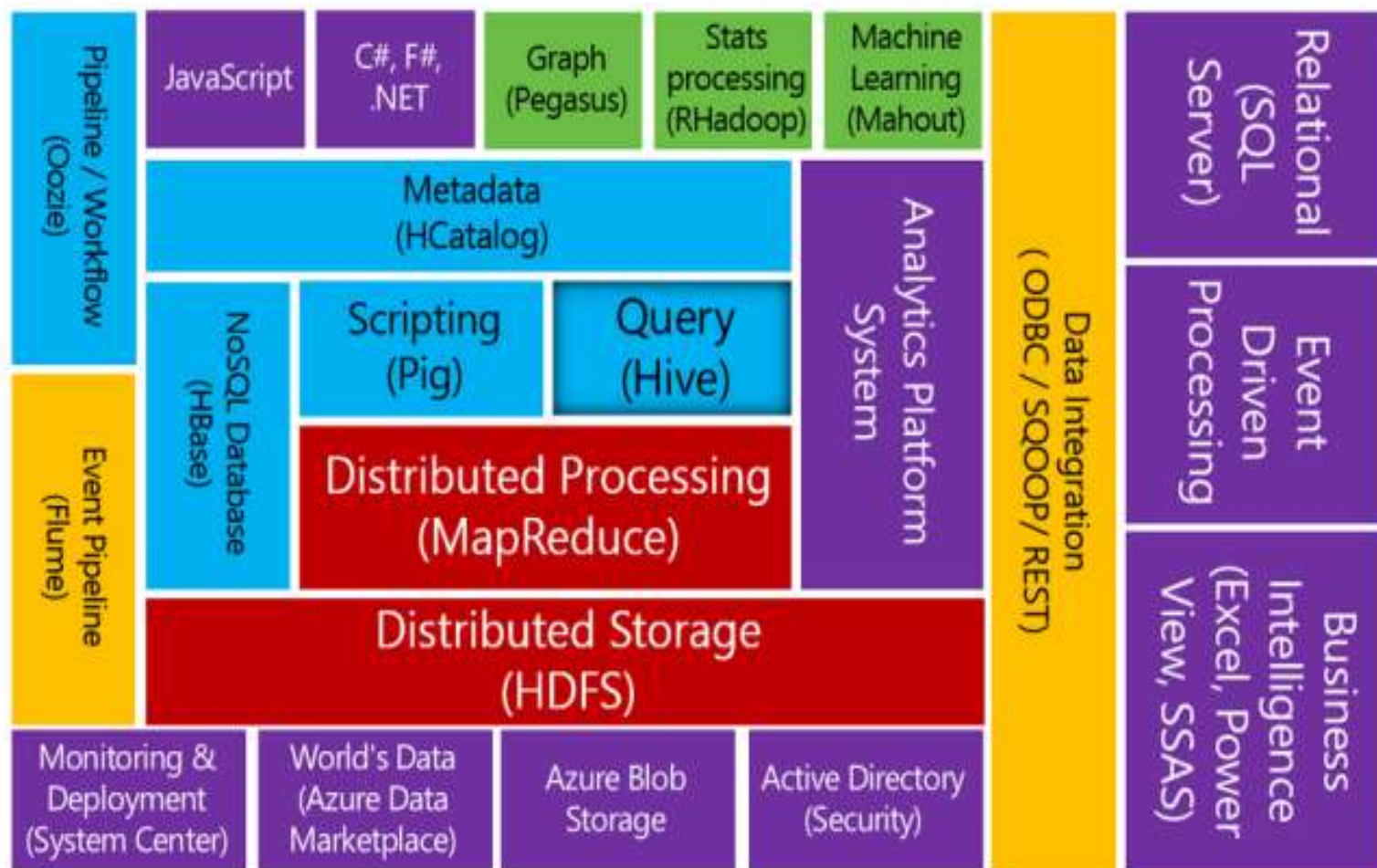
# Apache Hadoop Ecosystem

- 根據 Google 在 2003 年發表的 Google File System 論文而開發
- 主要包含 Hadoop Distributed File System (HDFS) 與 MapReduce Engine



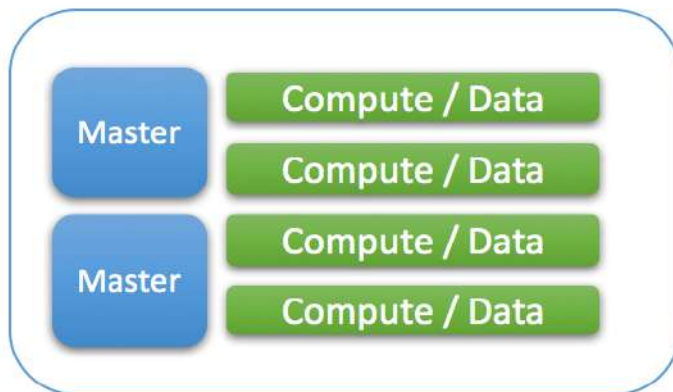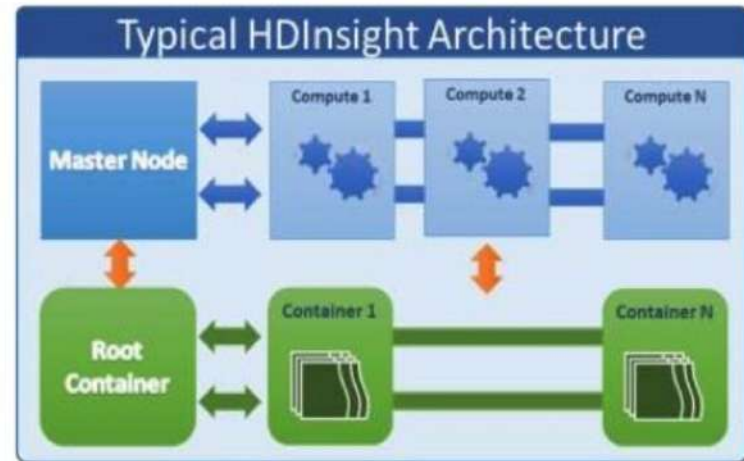參考資料：Hadoop Ecosystem – 15 Must Know Hadoop Components
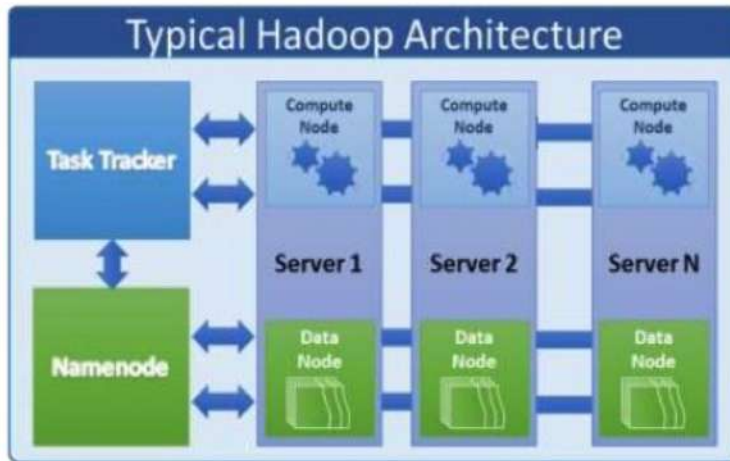
# Azure Data Science Virtual Machine



參考資料：Data Science Virtual Machines

# Azure HDInsight Ecosystem



参考資料：[Big Data] Hadoop/HDInsight : concepts généraux

# Apache Hadoop vs. Azure HDInsight

# Azure Storage Type for Apache Hadoop

`wasb`：

- Windows Azure Storage Blob
- HDFS File System Driver over Azure Blob Storage



參考資料：What is HDInsight?

# Apache Spark

- 2009 年開發
- 2010 年 Open Source
- 2013 年捐給 Apache Software Foundation

# Apache Spark on Azure HDInsight



參考資料：大象跳舞系列之Spark on HDInsight (1)

# Azure Storage Type for Apache Spark

`wasb`：

- Windows Azure Storage Blob
- HDFS File System Driver over Azure Blob Storage
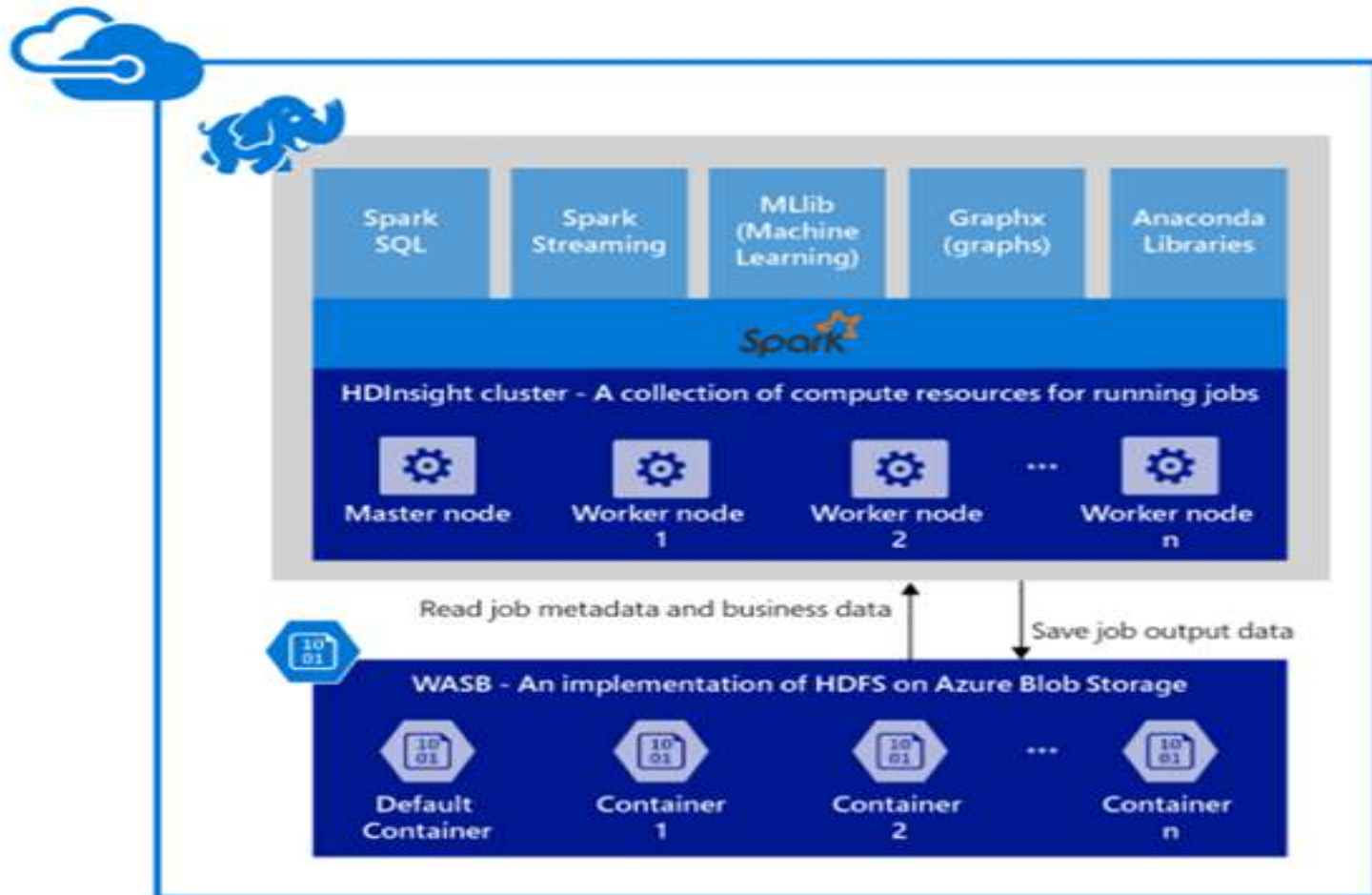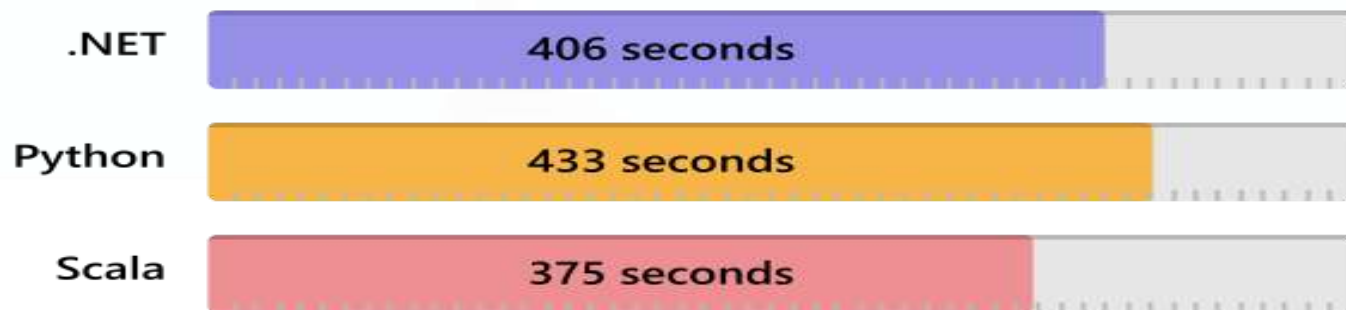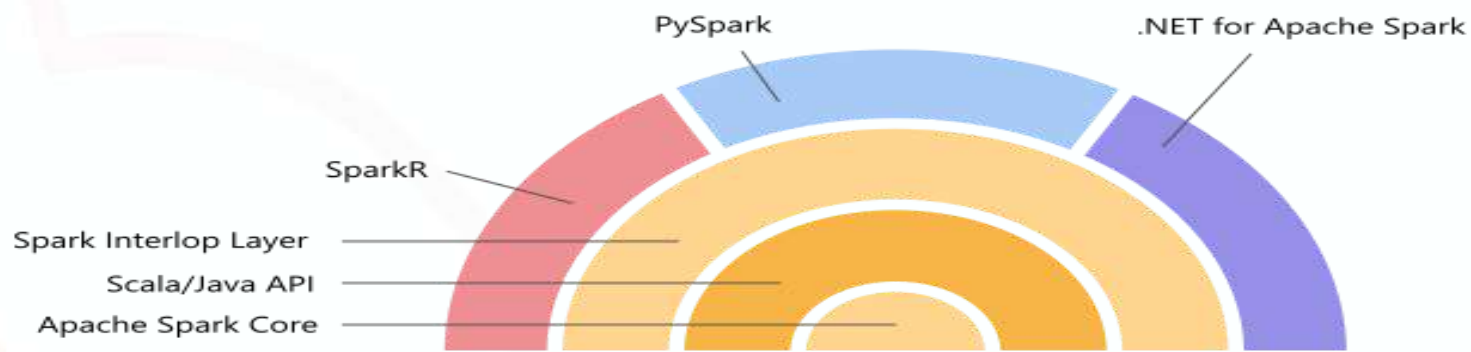
`abfs`：

- Azure Blob File System
- HDFS File System Driver Compatible with Azure Data Lake Storage Gen2

| Store Type | File System | Speed | Transient | Use Cases |
|---|---|---|---|---|
| Azure Blob Storage | wasb://url/ | Standard | Yes | Transient cluster |
| Azure Blob Storage (secure) | wasbs://url/ | Standard | Yes | Transient cluster |
| Azure Data Lake Storage Gen 2 | abfs://url/ | Faster | Yes | Transient cluster |
| Azure Data Lake Storage Gen 1 | adl://url/ | Faster | Yes | Transient cluster |
| Local HDFS | hdfs://url/ | Fastest | No | Interactive 24/7 cluster |

參考資料：Data storage optimization for Apache Spark

# .NET for Apache Spark

# Azure HDInsight and Datawarehouse



参考資料：What is Azure HDInsight?

# Databricks

- 由 Apache Spark 原始開發者建立
- 同時也開發 Delta Lake、MLflow、與 Koalas 等知名 Open Source 專案
- 2017 年 11 月推出與 Microsoft Azure 整合的 Azure Databricks 平台



參考資料：Announcing the Launch of Databricks SQL
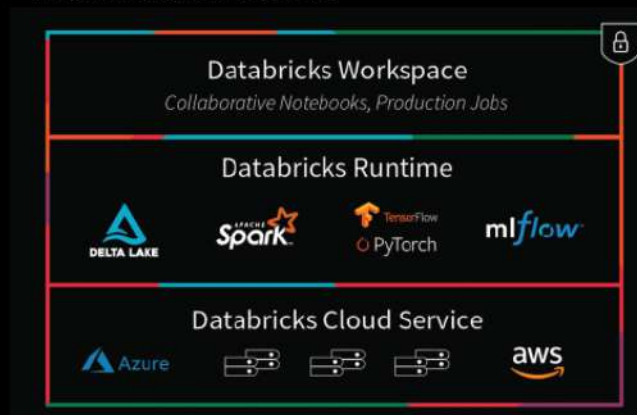
# Apache Spark vs. Databricks



參考資料：Comparing Apache Spark™ and Databricks

# Spark Comparisons

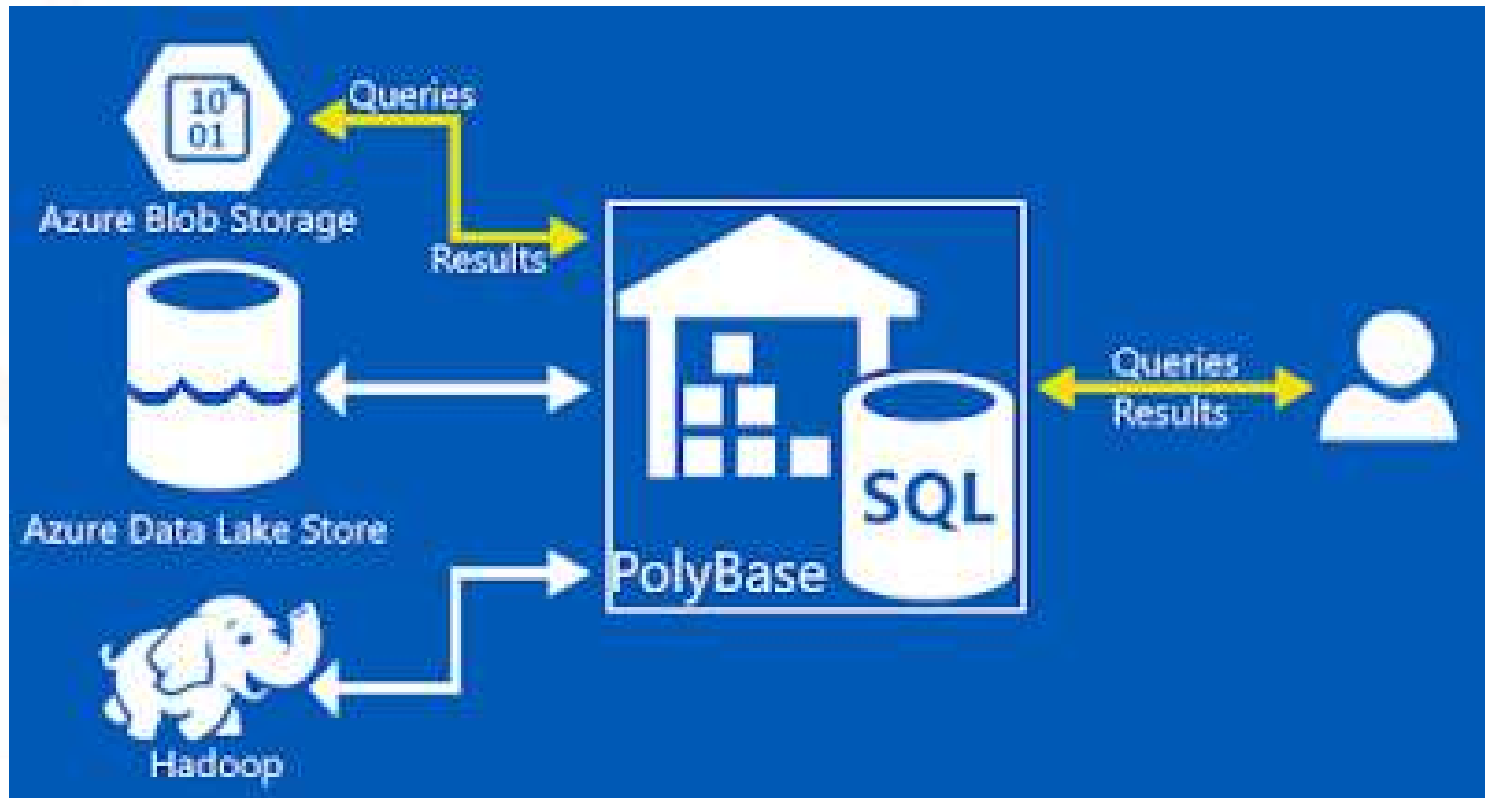Apache Spark vs. Azure HDInsight vs. Azure Databricks vs. Synapse Spark：

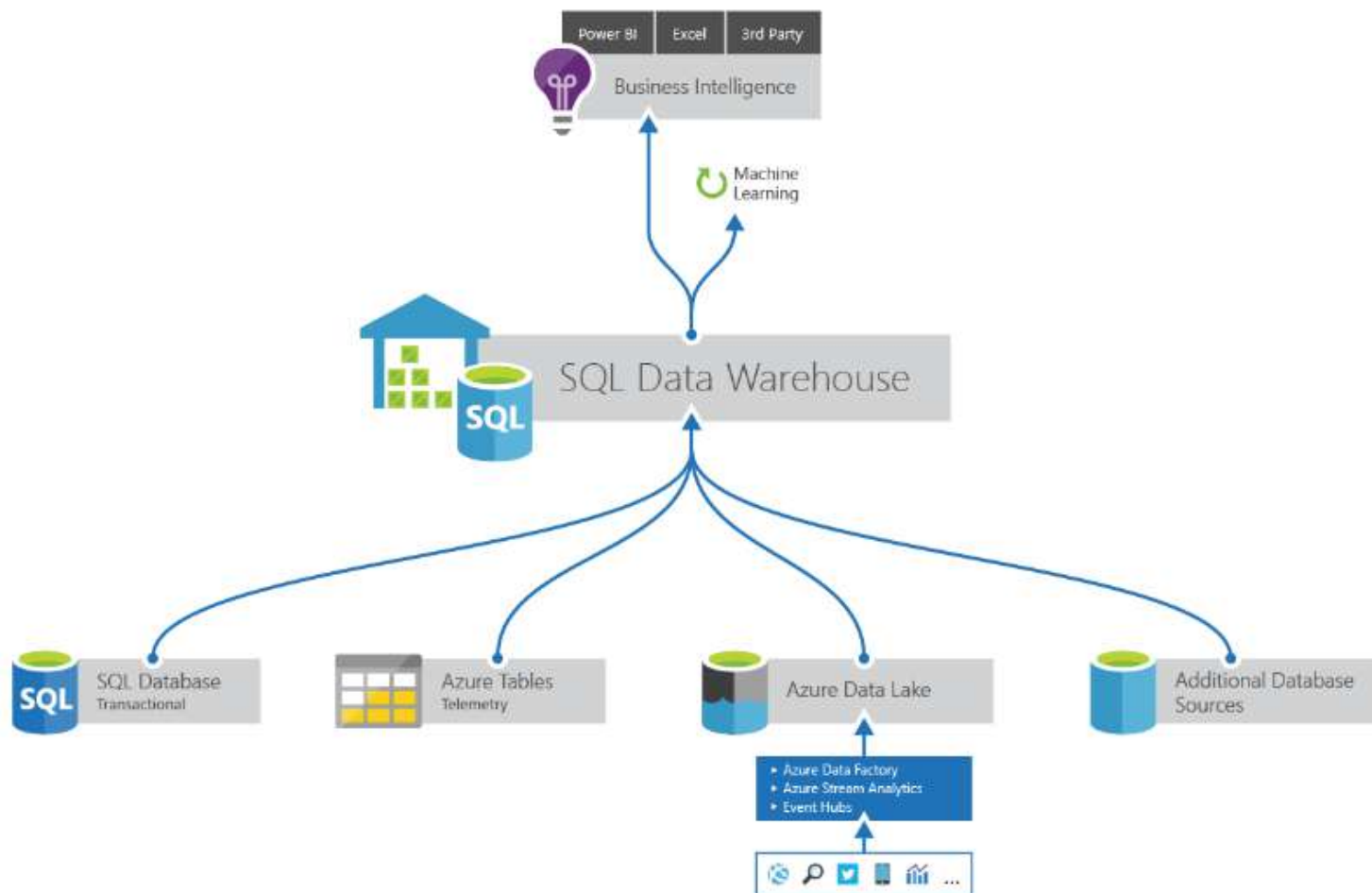| | Apache Spark | HDInsight | Azure Databricks | Synapse Spark |
|---|---|---|---|---|
| WHAT | Is an Open Source memory optimized system for managing big data workloads | Microsoft implementation of Open Source Spark managed within the realms of Azure | A managed Spark as a Service solution | Embedded Spark capability within Azure Synapse Analytics |
| WHEN | When you want to benefits of spark for big data processing and/or data science work without the Service Level Agreements of a provider | When you want to benefits of OSS spark with the Service Level Agreement of a provider | Provides end to end data engineering and data science solution and management platform | Enables organizations without existing Spark implementations to fire up a Spark cluster to meet data engineering needs without the overheads of the other Spark platforms listed |
| WHO | Open Source Professionals | Open Source Professionals wanting SLA's and Microsoft Data Platform experts | Data Engineers and Data Scientists working on big data projects every day | Data Engineers, Data Scientists, Data Platform experts and Data Analysts |
| WHY | To overcome the limitations of SMP systems imposed on big data workloads | To take advantage of the OSS Big Data Analytics platform with SLA's in place to ensure business continuity | It provides the ability to create and manage an end to end big data/data science project using one platform | It provides the ability to scale efficiently with spark clusters within a one stop shop DataWarehousing platform of Synapse. |

# Azure Databricks and Data Warehouse



參考資料：How to connect Databricks to Azure Data Lake?

# Microsoft SQL Server and PolyBase



參考資料：Use PolyBase to read Blob Storage in Azure SQL DW

# Azure SQL Data Warehouse (SQL DW)



參考資料：Azure SQL Data Warehouse

# Data Lake Evolution

# Database / Data Warehouse / Data Lake

Database：

- 比較 Operational 的 Data
- Data 可以新增刪除修改

Data Warehouse：

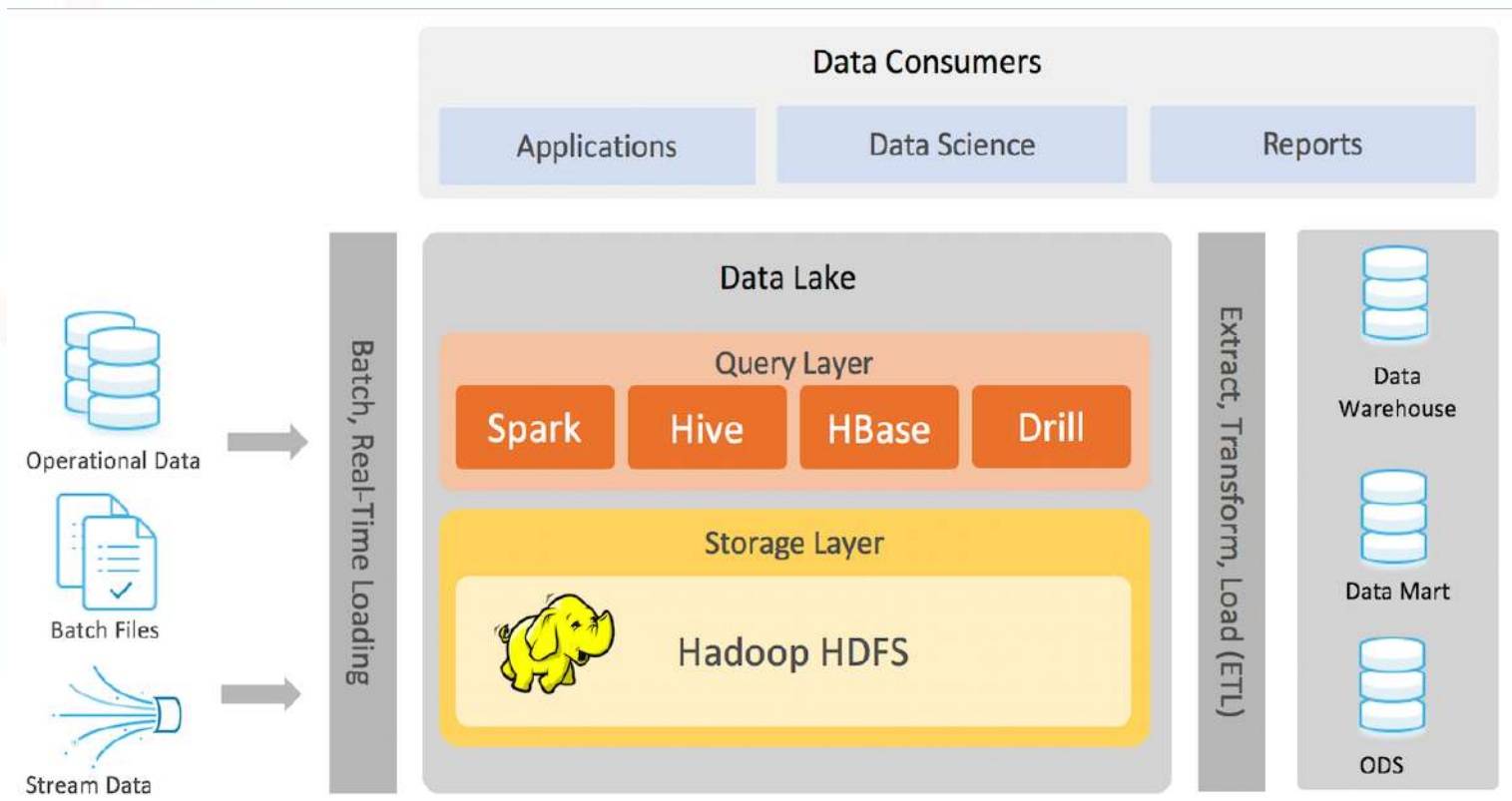- 比較 Analytical 的 Data，Quality 比較好的 Data
- Data 原則上是處理過的，不會刪除，不會修改，只會增加
- 通常是 Database 的 Data 累積一段時間之後，再轉到 Data Warehouse
- Subject-Oriented、Integrated、Time-Variant、與 Non-Volatile
- 是一種結構化的、可以分析查詢的 Transactional Data Copy
- 經常用來產生 Report 或 Dashboard

Data Lake：

- Repositories for Big Data in Its Native Form
- 儲存各種 Un-Structured Data / Semi-Structured / Structured Data 的大型 Data Storage
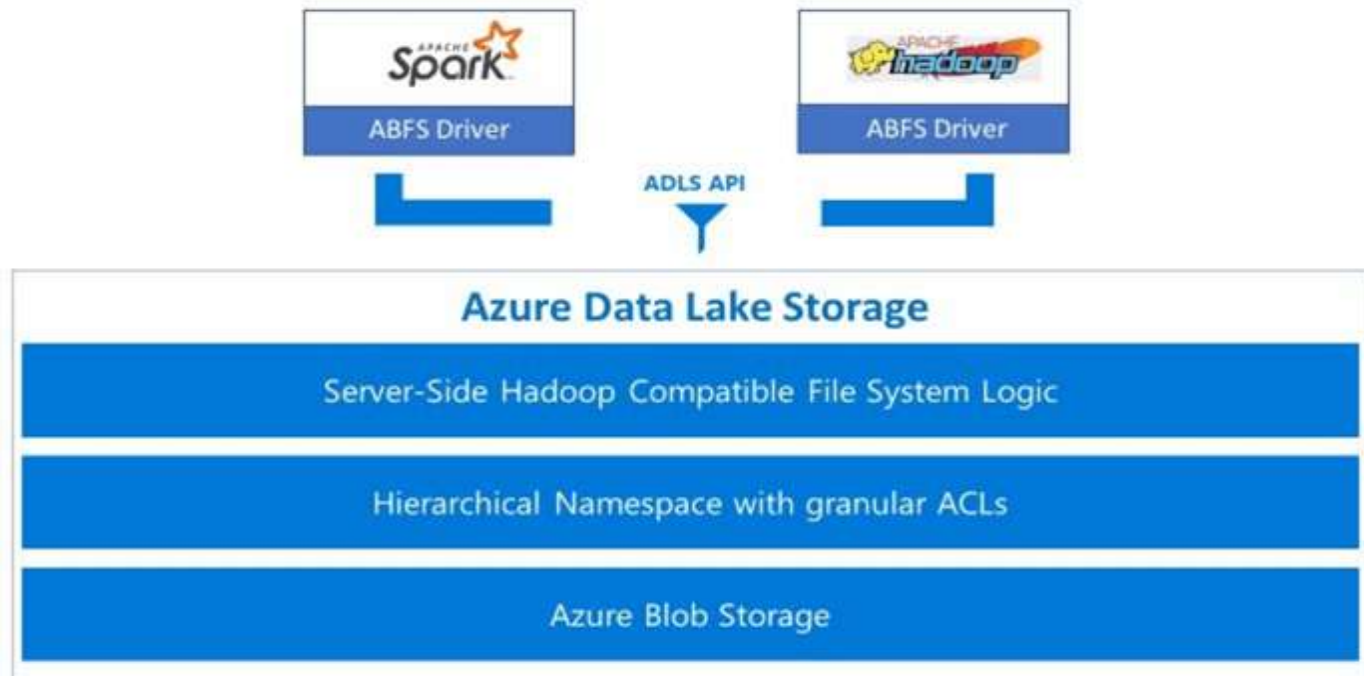- 不預先定義 Data 儲存的目的，或者說儲存就是它的目的
- 會比較在乎成本效益

參考資料：資料倉儲的類型

# Hadoop Data Lake Architecture

- Hadoop 是常見的一種 Data Lake 實現方式



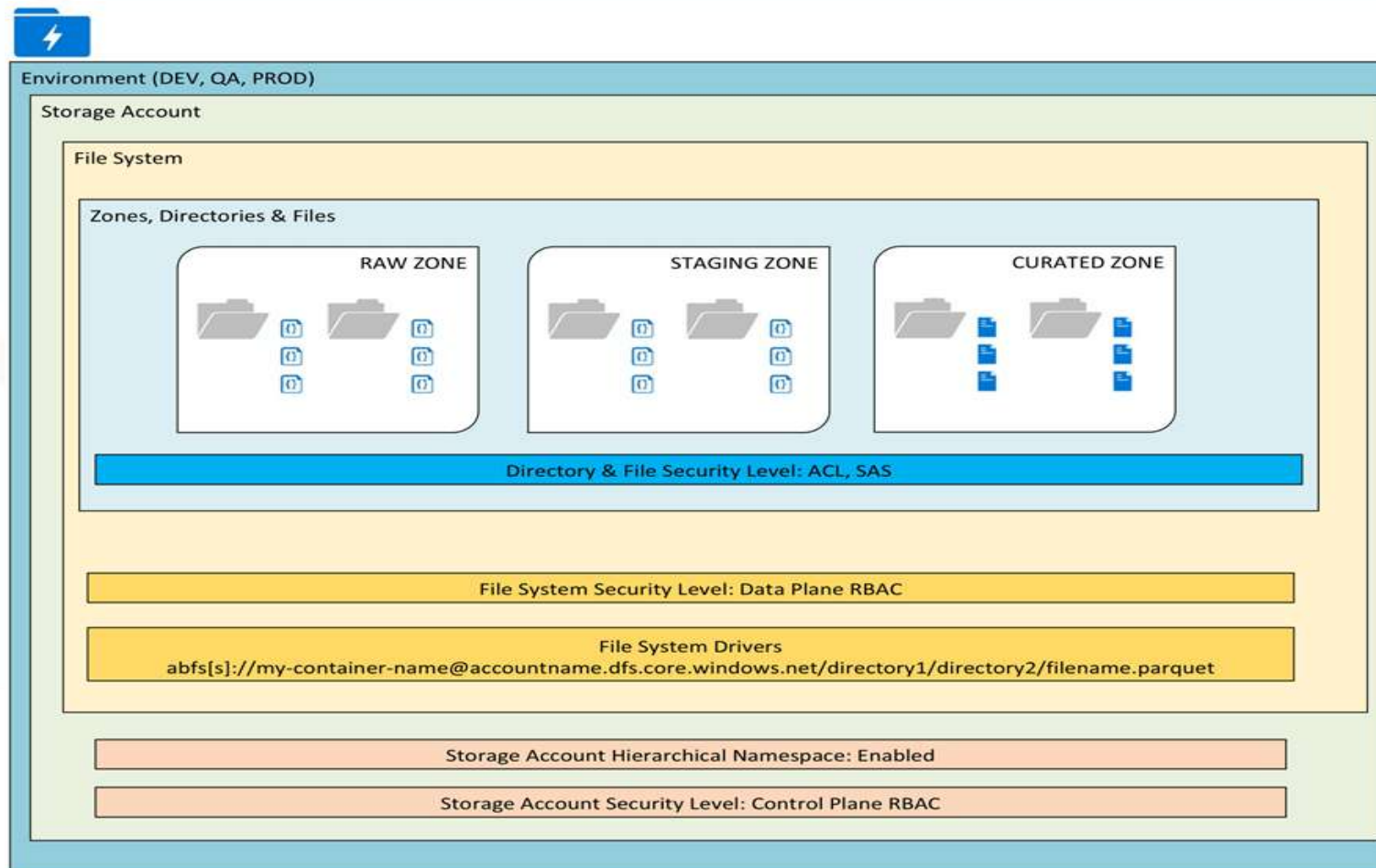參考資料：Webinar: Data Lake Advances in the Age of Operational ML/AI

# Azure Data Lake Storage Architecture

- File System Support
- Encryption of Data in Transit
- Encryption of Data at Rest
- Storage Account Firewall

- Compatible with HDFS
- Virtual Network Integration
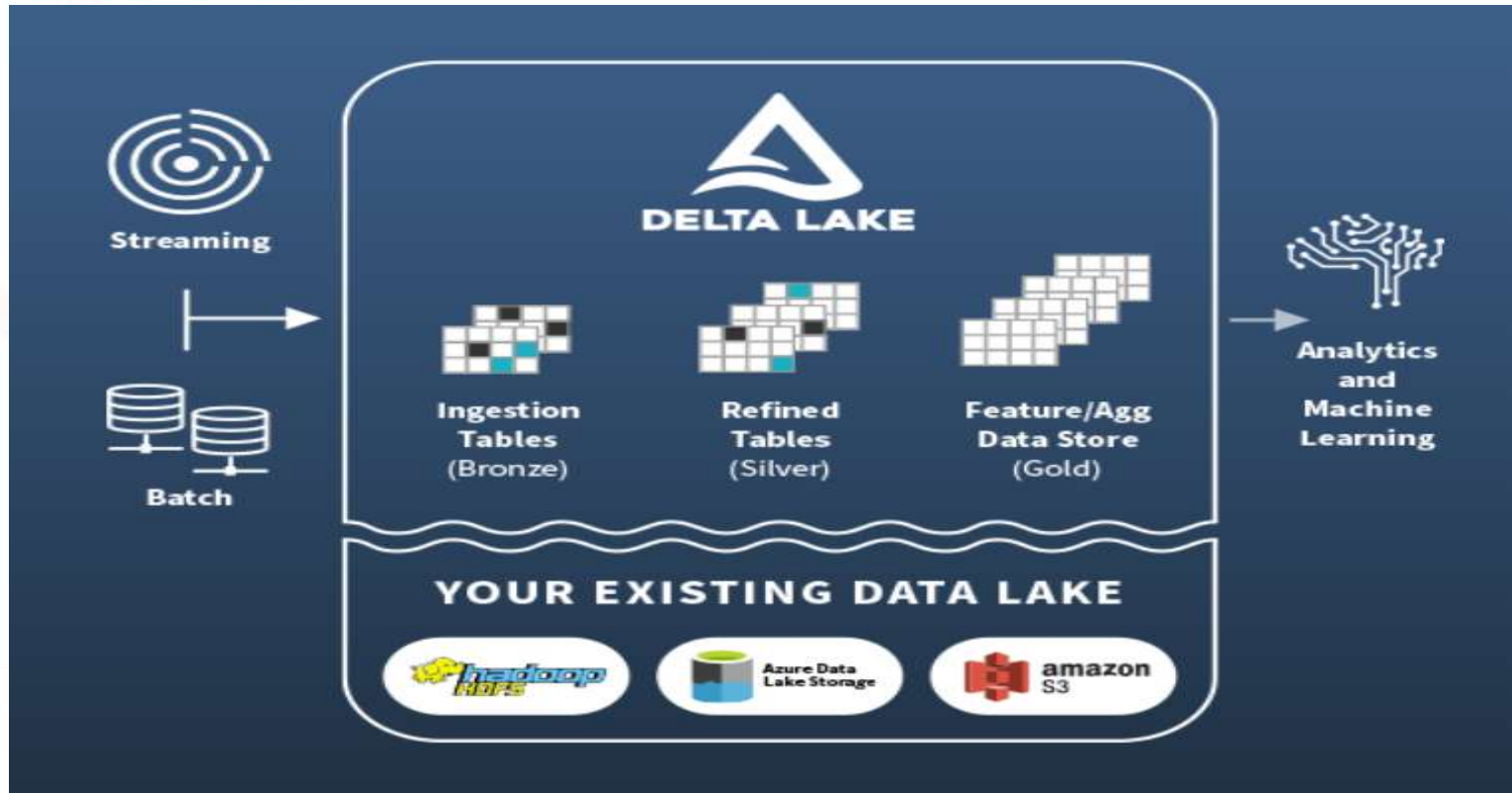- Role-Based Access Security
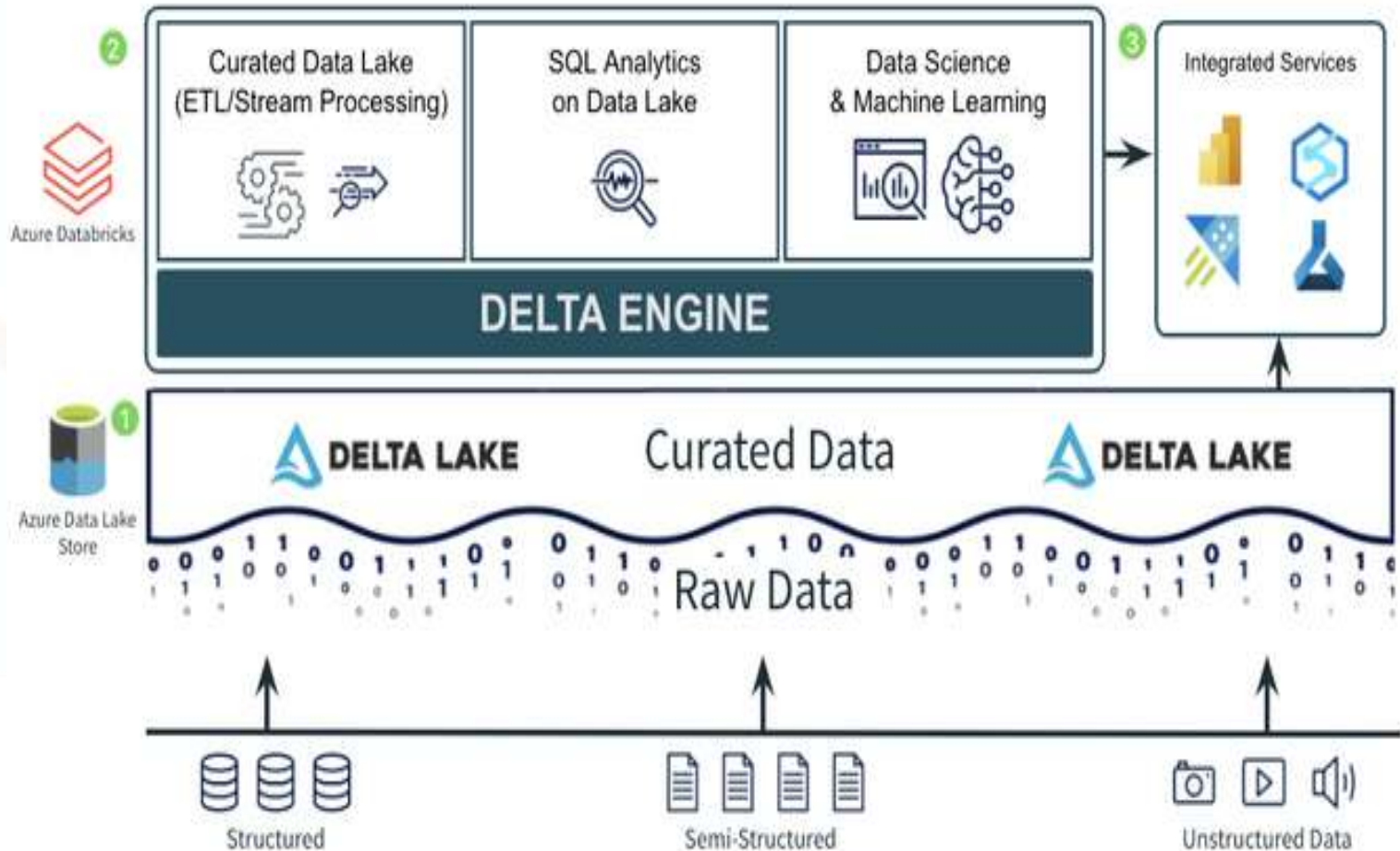- Hierarchical Namespace



參考資料：What's new in Azure Data Lake Storage Gen2

# Azure Data Lake Storage Layered Design



參考資料：Designing an Azure Data Lake Store Gen2

# Databricks Delta Data Lake Architecture

- Databricks 的 Open Source Data Lake 實現方式



參考資料：Data Lake or Warehouse? Databricks Offers a Third Way

# Data Lake + Warehouse = Lakehouse



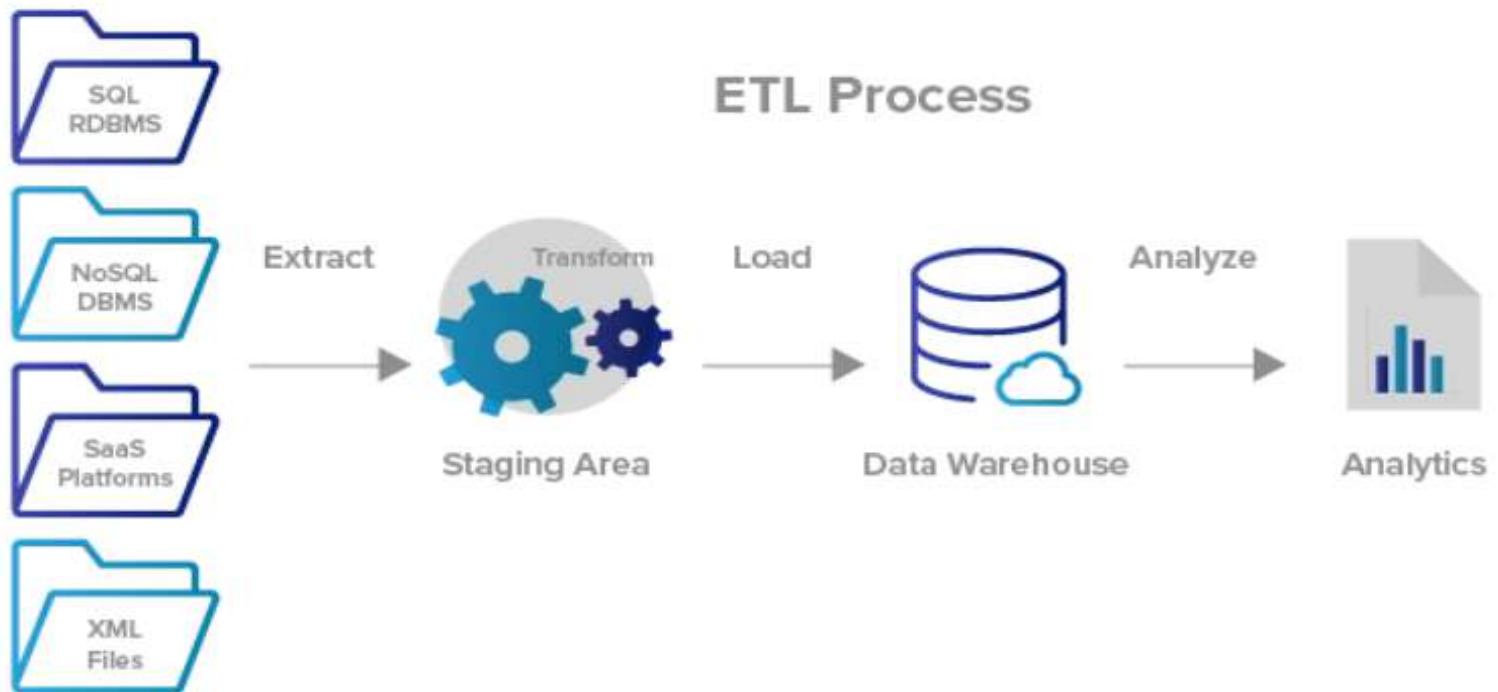參考資料：Simplify Your Lakehouse Architecture with Azure Databricks, Delta Lake, and Azure Data Lake Storage
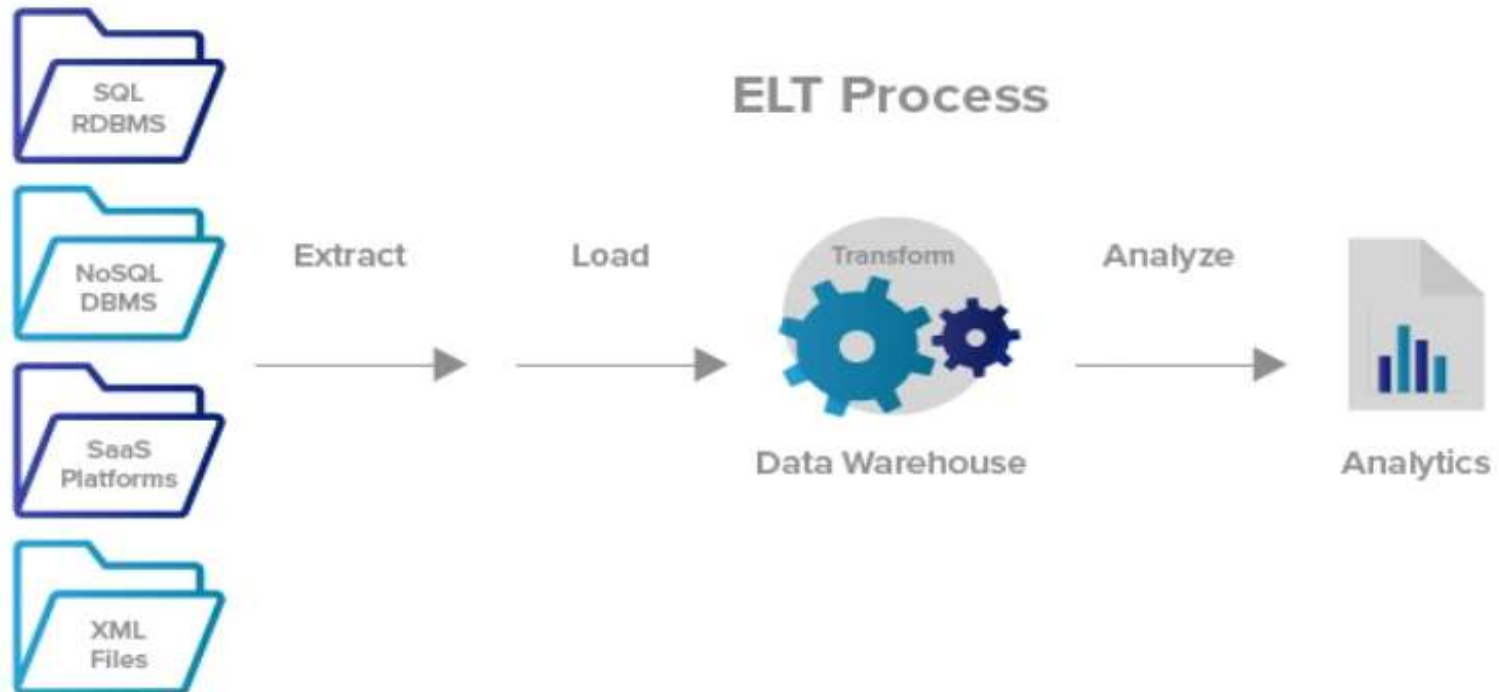
# Extract - Transform - Load (ETL)

- 習慣上會先把 Data Extract 到所謂的 Staging Area，進行 Transform，再 Load 到 Data Warehouse 裡頭
- 通常處理 On-Premises 的 Relational Structured Data
- 比較不支援 Data Lake 的作法



參考資料：ETL vs ELT: 5 Critical Differences

# Extract - Load - Transform (ELT)

- Data 直接 Extract 之後就 Load 到 Data Warehouse，然後在 Data Warehouse 裡面進行 Transform
- 比較貼近在 Cloud 同時處理 Structured 與 Unstructured Data 的想法
- 比較支援 Data Lake 的作法
- 也比較適合處理 Big Data 進行分析



參考資料：ETL vs ELT: 5 Critical Differences

# Azure Synapse Analytics

# Azure Synapse Analytics

Limitless Analytics Service with Unmatched Time to Insight.

- SQL DW：Data Warehouse
- Apache Spark / Azure Databricks：Big Data Analytics
- Azure Synapse：Data Ware House + Big Data Analytics



參考資料：微軟推出Azure Synapse Analytics以串連分析系統、資料倉儲

# Azure Synapse Studio

- Web-Based Single Hub of Azure Synapse Analytics
- Ingest / Explore / Prepare / Train / Analyze / Visualize Data
- 不需要知道 Data Schema 就可以進行 Data Exploration
- 透過 Knowledge Center 提供無止盡的 Learning Material / Sample Code / Template / Azure Open Dataset / Pipeline
- 內建 Pipeline Template 可以整合各種 Data Source
- 可以建立 Data Pipeline 進行 ETL
- 內建 AI / ML 功能，不需要另外整合 Azure ML 或是 Cognitive Services
- 以前的 SQL DW 變成現在 Built-in 的 Serverless SQL Pool
- 可以根據專案需求自己建立一或多個 Dedicated SQL Pool
- SQL Pool 以 Columnar Format 儲存 Data，搭配 `CLUSTERED COLUMNSTORE INDEX` 會有更好的 Performance
- SQL Pool 直接可以使用 `PREDICT` Function 在 Query 中進行 Prediction
- Apache Spark Pool 可以使用 C# / Python / Scala / Spark SQL 語言，支援 Spark ML 與 Spart Streaming
- 整合 AutoML，而且可以直接使用 Serverless Apache Spark Pool
- 可以存取 Azure Machine Learning Model Registry 直接進行 Model Scoring，不需要將 Data 搬進搬出
- 透過 Cognitive Services 可以 Enrich Data
- 支援與 Jupyter 相同的 Notebook 格式
- 整合 Power BI，可以直接使用 Power BI 的 Dataset / Report / Dashboard 功能

# Why Azure Synapse Analytics

Implementing an end-to-end analytics solution in Azure costs up to 13 percent less compared to AWS, up to 49 percent less compared to Google.

### The Four Offerings of the Vendor Stacks for Dedicated Compute

| | Vendor Offering | Pricing Used |
|---|---|---|
| Azure | Azure Synapse Analytics Workspace | Pay as you go ($1.20/hour per 100 DWU)[2] |
| aws | Amazon Redshift RA3 | 1-year commitment all-upfront ($8.61 effective hourly)[3] |
| Google | Google BigQuery | Annual slot commitment ($1,700 per 100 slots)[4] |
| Snowflake | Snowflake Computing | Enterprise+ ($4.00 per hour per credit)[5] |

### The Four Offerings of the Vendor Stacks for Storage

| | Vendor Offering | Pricing Used |
|---|---|---|
| Azure | Azure Synapse Analytics SQL Pool | $0.023 per GB-month[6] |
| aws | Amazon Redshift Managed Storage | $0.024 per GB-month[7] |
| Google | Google BigQuery Storage | $0.023 per GB-month (uncompressed)[8] |
| Snowflake | Snowflake Computing Storage | $0.04 per GB-month[9] |

参考資料：Cloud Analytics Platform Total Cost of Ownership
参考資料：Analytics in Azure is up to 14x faster and costs 94% less than other cloud providers. Why go anywhere else?
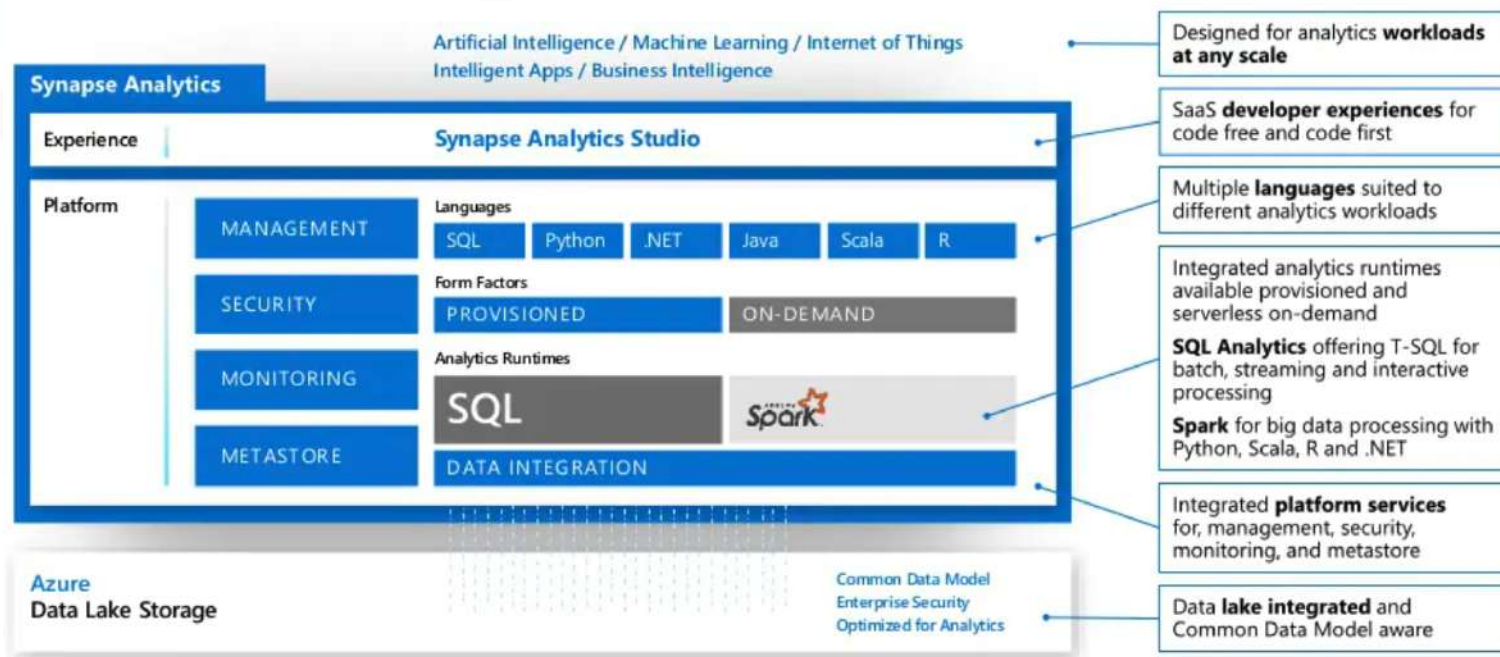
# Before vs. After Azure Synapse Analytics



参考資料：Microsoft Ignite postmortem: Cutting through the complexity

# Before vs. After Azure Synapse Analytics



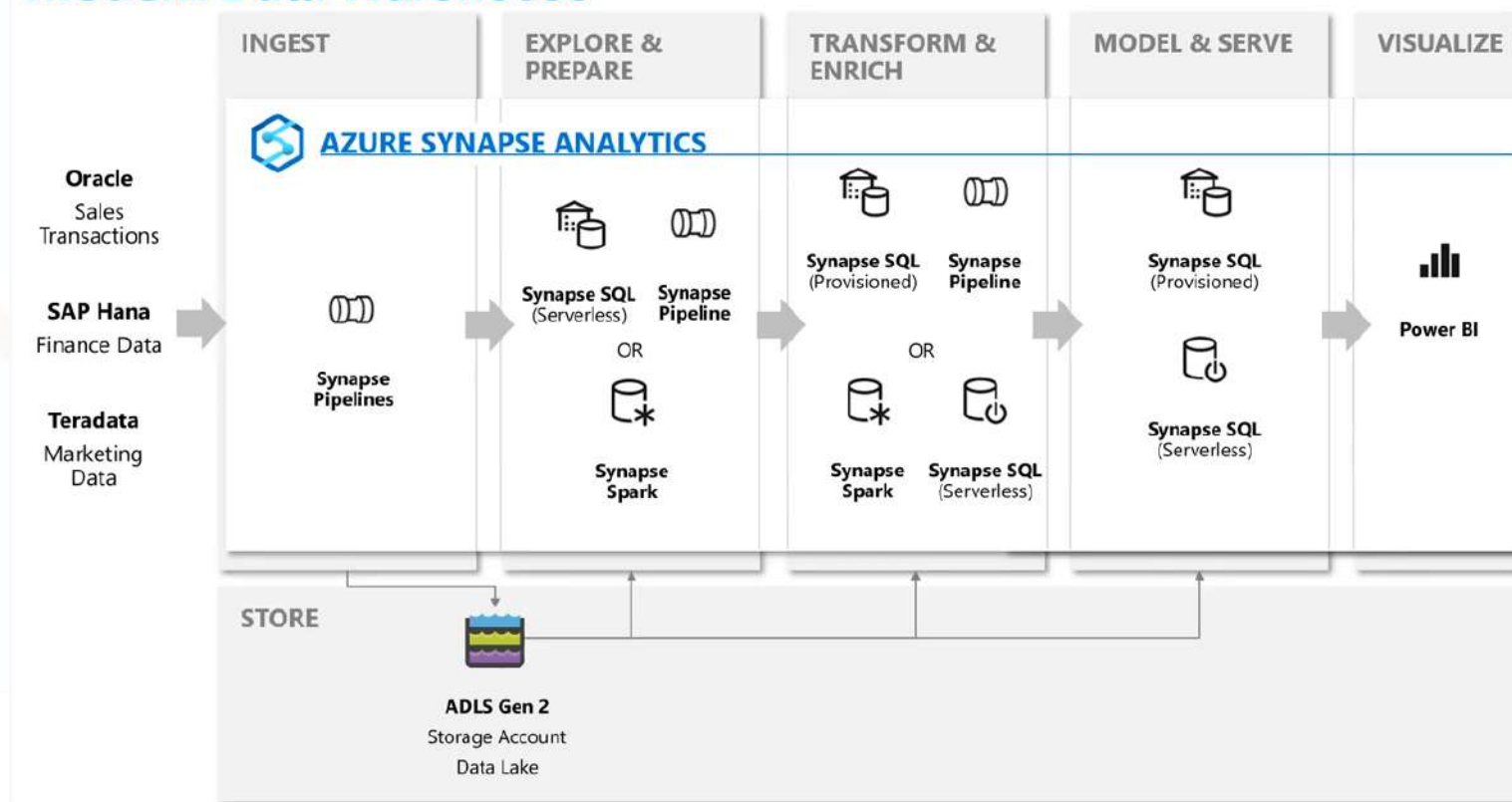参考資料：Microsoft Ignite postmortem: Cutting through the complexity

# Azure Synapse Analytics Architecture

- Data Storage 集中存放在 Azure Data Lake Storage Gen2 之上
- ETL/ELT 提供比 Azure Data Factory 更簡單好用的 Synapse Pipeline
- Data Analysis 提供 SQL-Based Runtime 與 Apache Spark Runtime
- Security 直接整合 Azure Active Directory
- 透過 Linked Service 架構整合各種 Data Source 與 Cloud Service
- 所有一切功能都整合在 Synapse Studio 之內



參考資料：Azure Synapse Analytics Overview (r2)

# Before vs. After (Detailed Edition)



參考資料：Azure Synapse Analytics For Beginner

# Data Warehouse/Analytics Solution

- Data Source：File / File Server / Database Server / Web / IoT
- Ingest：Synapse Pipeline / SQL BULK COPY
- Store：Azure Data Lake Storage Gen2
- Prepare & Train：Apache Spark Pool
- Model & Serve：Serverless SQL Pool (Metadata) / Dedicated SQL Pool (Data)
- Apps + Insights：Power BI / Tableau / Azure Service / O365 Service / Others



參考資料：What is dedicated SQL pool (formerly SQL DW) in Azure Synapse Analytics?

# Explore Data

- 收集：Azure Stream Analytics
- 儲存：Azure Data Lake Storage
- 預覽：Synapse Studio



參考資料：Azure Synapse Analytics end-to-end: From your database and sensors to a Power BI dashboard

# Ingest Data

- SQL BULK COPY



- Synapse Pipeline



參考資料：Azure Synapse Analytics end-to-end: From your database and sensors to a Power BI dashboard

# Knowledge Center Gallery

# Query Data

# Different Options of Compute / Query

Azure Synapse Analytics 整合了：
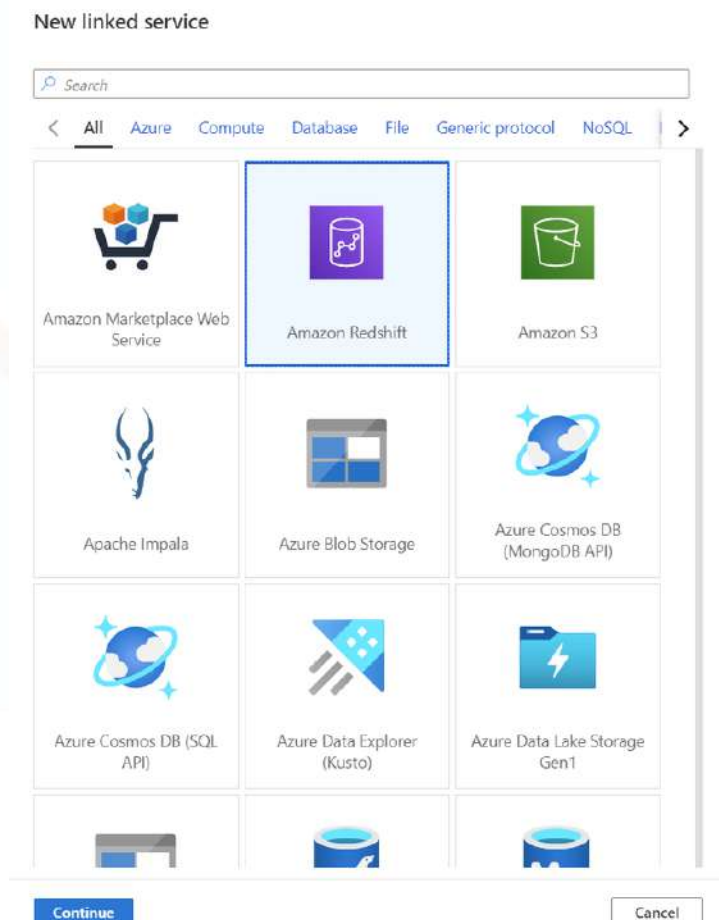
- Data Lake
- Data Warehouse
- Big Data Analytics

Azure Synapse Analytics 提供三種 Query Service：

- Dedicated SQL Pool (原來的 SQL DW，底下的 Provisoned SQL)
- Serverless On-Demand SQL Pool (支援 Unstructured Data)
- Serverless Apache Spark Pool (Distributed In-Memory Computing)

| | Relational Data | ADLS Gen2 | Spark Table | Cosmos DB |
|---|---|---|---|---|
| Provisioned SQL | Y | Y (external table) | Soon (parquet format) | X |
| On-demand SQL | X | Y | Y (parquet format) | Y (Synapse link) |
| Spark | Y | Y | Y | Y (Synapse link) |

SQL pool supported file formats in ADLS Gen2 are parquet, csv, json (Spark supports many more formats)

參考資料：Azure Synapse Analytics For Beginner

# Linked Service



参考資料：Azure Synapse Analytics end-to-end: From your database and sensors to a Power BI dashboard

# Analyze & Visualize Data



参考資料：Azure Synapse Analytics end-to-end: From your database and sensors to a Power BI dashboard
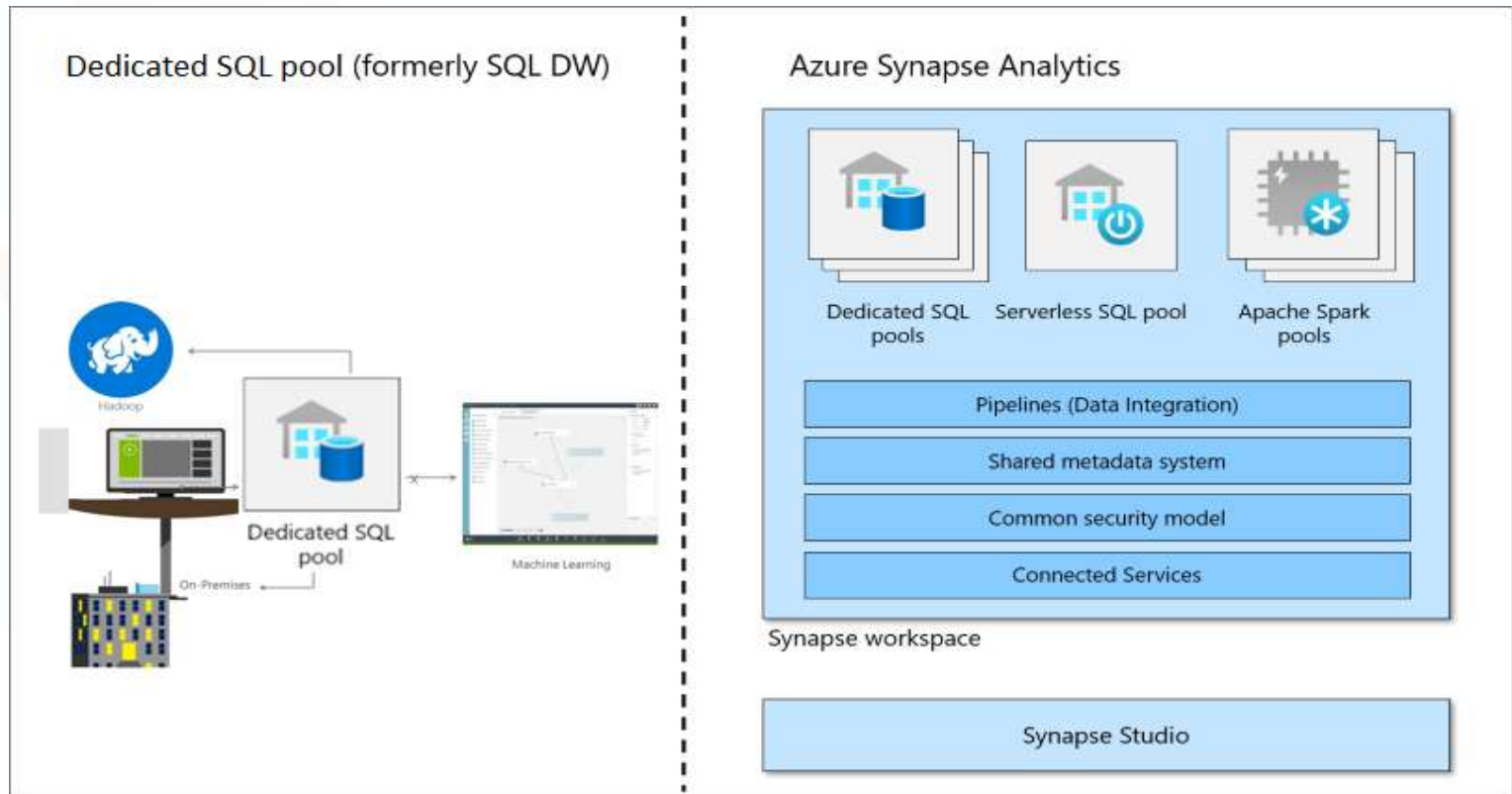
# SQL DW / Dedicated SQL Pool

# SQL DW vs. Azure Synapse Analytics

- 之前叫 SQL Data Warehouse
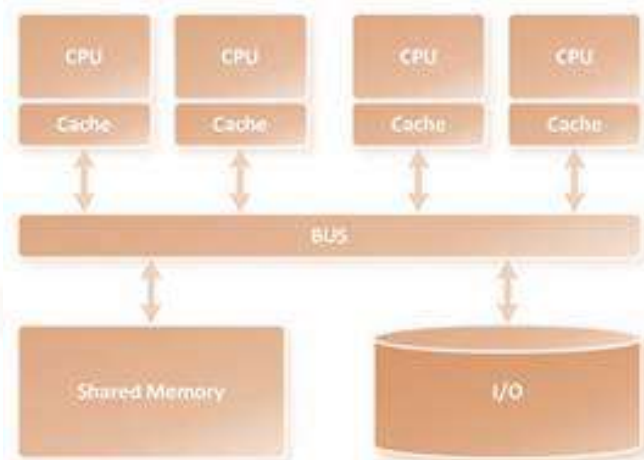- 現在是 Azure Synapse Analytics 裡面的一個 Dedicated SQL Pool
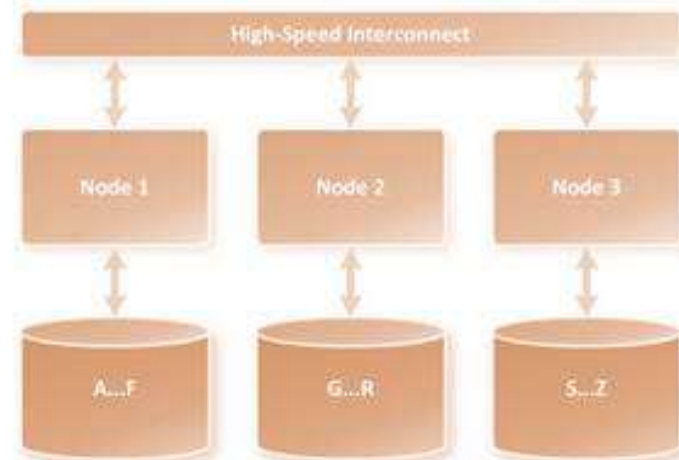


參考資料：
Dedicated SQL pool (formerly SQL DW) architecture in Azure Synapse Analytics
Azure SQL Data Warehouse

# Shared-Memory vs. Shared-Nothing

- SQL DW 採用 Massively Parallel Processing 架構

# Row Store vs. Column Store



參考資料：Understanding New Column Store Index of SQL Server 2012

# Separated Compute/Storage

- 由 1 個 Control Node (Master) + 1 或 多個 Compute Node (Slave) 組成
- Compute Node 數目由 Data Warehouse Unit (DWU) 決定
- 支援 PolyBase T-SQL Query Language
- 透過 Distributed Query Engine 進行 High-Performance Analytics



參考資料：Dedicated SQL pool (formerly SQL DW) architecture in Azure Synapse Analytics

# DWU vs. SLO vs. Performance Tier

Data Warehouse Unit (DWU)：

- 是 CPU、Memory、與 I/O 的組合
- Abstract Normalized Measure of Compute Resources and Performance
- 調整 Service Level Objective 就會改變 Data Warehouse Unit
- 改變 Data Warehouse Unit 不會影響 Storage Cost

Service Level Objective (SLO)：

- Scalability Setting of Cost and Performance Level
- Gen1 的 SLO 以 Data Warehouse Unit (DWU) 量測
- Gen2 的 SLO 以 Compute Data Warehouse Unit (cDWU) 量測

Performance Tier：

- 分為 Gen1 (`DW100-DW6000`) 與 Gen2 (`DW100c-DW30000c`) 兩種
- Gen2 提供 Local Disk-Based Cache，所以 Performance 更好
- 之前使用 Gen1 部分原因是 Gen2 只支援比較高的 DWU (Gen2 現在支援 `DW100c`)
- SQL-DW DWU x 7.5 = SQL Server Required DTU
- SQL-DW cDWU x 9 = SQL Server Required DTU
- SQL Server DTU 上限預設 `54,000` (`DW6000c`)，更高就要開 Support Ticket
- SQL DB Contributor 與 SQL Server Contributor 這兩種 Role 才能調整

參考資料：
Capacity limits for dedicated SQL pool in Azure Synapse Analytics
Azure SQL Data Warehouse Gen2 now supports lower compute tiers

# Service Level vs. Compute Nodes

| Performance Level | Compute Nodes | Distributions / Compute Node | Memory / Data Warehouse (GB) |
|---|---|---|---|
| DW100c | 1 | 60 | 60 |
| DW200c | 1 | 60 | 120 |
| DW300c | 1 | 60 | 180 |
| DW400c | 1 | 60 | 240 |
| DW500c | 1 | 60 | 300 |
| DW1000c | 2 | 30 | 600 |
| DW1500c | 3 | 20 | 900 |
| DW2000c | 4 | 15 | 1200 |
| DW2500c | 5 | 12 | 1500 |
| DW3000c | 6 | 10 | 1800 |
| DW5000c | 10 | 6 | 3000 |
| DW6000c | 12 | 5 | 3600 |
| DW7500c | 15 | 4 | 4500 |
| DW10000c | 20 | 3 | 6000 |
| DW15000c | 30 | 2 | 9000 |
| DW30000c | 60 | 1 | 18000 |

# Dedicated SQL Pool Pricing

| Service Level | DWU | Pay As You Go (隨用隨付) | 1 Year RC (節省約 37%) | 3 Year RC (節省約 65%) |
|---|---|---|---|---|
| DW100c | 100 | NT$54.399/hour | NT$34.2709/hour | NT$19.0391/hour |
| DW200c | 200 | NT$108.797/hour | NT$68.5417/hour | NT$38.0781/hour |
| DW300c | 300 | NT$163.195/hour | NT$102.8125/hour | NT$57.1172/hour |
| DW400c | 400 | NT$217.593/hour | NT$137.0833/hour | NT$76.1562/hour |
| DW500c | 500 | NT$271.991/hour | NT$171.3541/hour | NT$95.1952/hour |
| DW1000c | 1000 | NT$543.982/hour | NT$342.7081/hour | NT$190.3904/hour |
| DW1500c | 1500 | NT$815.972/hour | NT$514.0621/hour | NT$285.5856/hour |
| DW2000c | 2000 | NT$1,087.963/hour | NT$685.4161/hour | NT$380.7808/hour |
| DW2500c | 2500 | NT$1,359.953/hour | NT$856.7702/hour | NT$475.9759/hour |
| DW3000c | 3000 | NT$1,631.944/hour | NT$1,028.1242/hour | NT$571.1711/hour |
| DW5000c | 5000 | NT$2,719.906/hour | NT$1,713.5403/hour | NT$951.9518/hour |
| DW6000c | 6000 | NT$3,263.887/hour | NT$2,056.2483/hour | NT$1,142.3422/hour |
| DW7500c | 7500 | NT$4,079.858/hour | NT$2,570.3104/hour | NT$1,427.9277/hour |
| DW10000c | 10000 | NT$5,439.811/hour | NT$3,427.0805/hour | NT$1,903.9036/hour |
| DW15000c | 15000 | NT$8,159.716/hour | NT$5,140.6207/hour | NT$2,855.8553/hour |
| DW30000c | 30000 | NT$16,319.431/hour | NT$10,281.2413/hour | NT$5,711.7106/hour |

參考資料：Azure Synapse Analytics定價

# Apache Spark Pool Pricing

以分計費：

| Type | Price | Free Quantity |
|---|---|---|
| Memory Optimized | NT$4.646 / vCore-Hour | **2021-07-31 前** 120 Free vCore-Hours / Month |

# Data Exploration/Warehousing Pricing

Serverless / Dedicated 分別計費：

| Type | Price | Free Quantity |
|---|---|---|
| Serverless | NT$202.866 / TB Data Processed | **2021-07-31 前** 10 TB of Free Queries / Month |
| Dedicated | 與 Dedicated SQL Pool 相同計費 | |

參考資料：Azure Synapse Analytics定價

# SLO/DWU Verification (PowerShell)

在 PowerShell 或是 Azure Cloud Shell 手動安裝 `Az.Synapse` Module：

```
Install-Module -Name Az.Synapse

Untrusted repository
You are installing the modules from an untrusted repository... install the modules from 'PSGallery'?
[Y] Yes  [A] Yes to All  [N] No  [L] No to All  [S] Suspend  [?] Help (default is "N"): A
```

建立後檢視：

```
PS /home/kc_su> Get-AzSynapseSqlPool -WorkspaceName synapseworkspacem00 -Name sqlpoolm00

ResourceGroupName     :
WorkspaceName         : synapseworkspacem00
SqlPoolName           : sqlpoolm00
Sku                   : DW100c
MaxSizeBytes          : 263882790666240
Collation             : SQL_Latin1_General_CP1_CI_AS
SourceDatabaseId      :
RecoverableDatabaseId :
ProvisioningState     : Succeeded
Status                : Online
RestorePointInTime    : 1/1/0001 12:00:00 AM
CreateMode            :
CreationDate          : 7/15/2021 3:37:08 PM
Tags                  : {}
TagsTable             :
Location              : southeastasia
Id                    : /subscriptions/.../workspaces/synapseworkspacem00/sqlPools/sqlpoolm00
Type                  : Microsoft.Synapse/workspaces/sqlPools
```

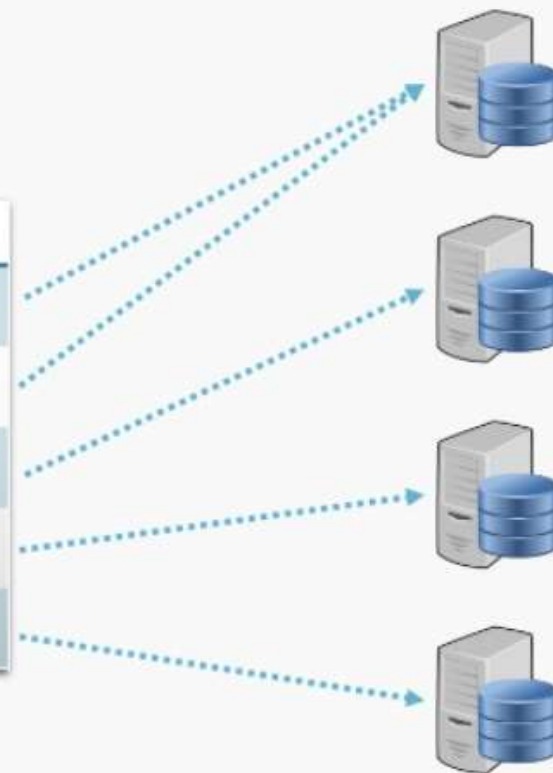# Data Distrubution (Sharding Pattern)

SQL-DW 載入資料的時候，最多會分成 60 個 Distribution：

- Hash

    - Distributed Using Hash Algorithm
    - Equal Values Hashed to the Same Distribution
    - Optimal for Joins and Aggregations on Large Fact Tables

- Round Robin

    - Distributed Evenly But Randomly
    - Not Requiring Knowledge about Data or Queries
    - Optimal for Large Tables without Good Hash Columns or Varied Queries
    - Fast Performance as a Staging Table for Loads

- Replicate

    - All Data Present on Every Node
    - Simplifying Query Plans and Reducing Data Movement
    - Best for Small Lookup Tables

參考資料：Design and Manage Azure SQL Data Warehouse

# Hash



Hashing by CustomerID

參考資料：Azure SQL Data Warehouse

# Round Robin



| RecordNo | CustomerID | InvoiceDate |
|---|---|---|
| 1 | 1000 | 2017-04-21 |
| 2 | 1000 | 2017-04-22 |
| 3 | 2000 | 2017-04-22 |
| 4 | 3000 | 2017-04-22 |
| 5 | 4000 | 2017-04-22 |

Rows distributed to all nodes

參考資料：Azure SQL Data Warehouse

# Replicate



All table rows are copied to each compute node

Table

Compute Nodes

Replicated table

參考資料：Dedicated SQL pool (formerly SQL DW) architecture in Azure Synapse Analytics

# Compute Node vs. Data Distribution



參考資料：Getting Started with Azure SQL Data Warehouse - Part 2

# Implementation Issues

# Permission Error

- 爲了存取 Azure Data Lake Storage Gen2 Storage Account，需要 Storage Blob Data Contributor 權限

# Problem and Answer

```
Error: org.apache.hadoop.hive.ql.metadata.HiveException: MetaException(
message:
   Got exception: org.apache.hadoop.fs.azurebfs.contracts.exceptions.AbfsRestOperationException
   Operation failed: "This request is not authorized to perform this operation using this permission.",
   403, HEAD,
   https://synapsem00.dfs.core.windows.net/data/synapse/workspaces/synapsem00/warehouse
       ?upn=false&action=getStatus&timeout=90);
```

ℹ️ We will automatically grant the workspace identity data access to the
specified Data Lake Storage Gen2 account, using the Storage Blob Data
Contributor role. To enable other users to use this storage account after
you create your workspace, perform these tasks:

- Assign other users to the **Contributor** role on workspace
- Assign other users the appropriate Synapse RBAC roles using
  Synapse Studio
- Assign yourself and other users to the **Storage Blob Data
  Contributor** role on the storage account

Learn more

# Subscription Owner vs. Non-Owner

底下這個選項只有在目前的 User 擁有 Subscription 的 Owner Role 時才會出現：

Assign myself the **Storage Blob Data Contributor** role on the Data Lake Storage Gen2 account to interactively query it in the workspace.

# Storage Account -> Workspace

# Storage Account -> User

# PolyBase OPENROWSET

- 可以存取 Azure Storage 裡面的 File Data
- 可以讀入 File Data，以 Relational Structure (A Set of Rows) 傳回

```
SELECT *
FROM OPENROWSET(
    BULK 'https://pandemicdatalake.blob.core.windows.net/.../covid-19/ecdc_cases/latest/ecdc_cases.csv',
    FORMAT = 'CSV',
    PARSER_VERSION = '2.0',
    HEADER_ROW = TRUE
) as [r]
```

```
SELECT TOP 1 *
FROM OPENROWSET(
    BULK 'https://azureopendatastorage.blob.core.windows.net/.../us_population/year=20*/*.parquet',
    FORMAT='PARQUET'
) WITH (
    [stateName] VARCHAR (50),
    [population] BIGINT
) AS [r]
```

```
SELECT TOP 1 *
FROM OPENROWSET(
    BULK 'https://azureopendatastorage.blob.core.windows.net/.../us_population/year=20*/*.parquet',
    FORMAT = 'DELTA'
) AS [r]
```

參考資料：How to use OPENROWSET using serverless SQL pool in Azure Synapse Analytics

# Serverless SQL Pool vs. Collation

- Serverless SQL Pool 可以將 UTF-8 Data 當成 `VARCHAR` Field 讀入，但是要注意 Collation 的影響
- 如果 `VARCHAR` Field 沒有指定 Collation，就會將 UTF-8 Character 強迫轉換為 Plain `CARCHAR` Character，可能會造成轉換錯誤
- 這類錯誤會發生在 `OPENROWSET` 沒有 `WITH` 子句，或是 `OPENROWSET` 與 External Table 傳回 `VARCHAR` Field 時沒有指定 Collation
- 這個議題對 `NVARCHAR` Field 不適用，因為 `NVARCHAR` Field 的轉換跟 Collation 無關，但是轉換成 `NVARCHAR` Field 會有一點 Performance 上的問題
- 改了 Database Colation 之後，既有的 External Table 必須重新建立才行

解決方式就是在 `CREATE DATABASE` 時指定 UTF-8 Collation：

```
CREATE DATABASE Database名稱 COLLATE Latin1_General_100_BIN2_UTF8
```

或是在 `OPENROWSET` 傳回的每個 `VARCHAR` Field 指定 UTF-8 Collation：

```sql
SELECT TOP 10 *
FROM OPENROWSET(
    BULK 'https://azureopendatastorage.blob.core.windows.net/.../us_population/year=20*/*.parquet',
    FORMAT='PARQUET'
) WITH (
    [stateName] VARCHAR (50) COLLATE Latin1_General_100_BIN2_UTF8 ,
    [population] BIGINT
) AS [r]
```

參考資料：
定序與 Unicode 支援
Always use UTF-8 collations to read UTF-8 text in serverless SQL pool

# Azure Data Factory vs. Spatial Data Types

- Azure Data Factory 處理 SQL Server 的 Spatial Data Type 常常有問題
- Microsoft Contoso BI Demo Dataset 的 `DimStore` Table 的 `GeoLocation` 與 `Geometry` 兩個 Field 會匯入失敗
- 解決方式就是先刪掉這兩個 Field

▲

0

▼

🕓

Do data factories support the Geography/Geometry data type?

No, Azure Data Factory does not support spatial types at this time. When selecting a table to sync via the copy wizard, if the table has any spatial columns, you will receive an error:

Error when processing request: Column: Location,The data type is not supported. activityId: [...]

Or if you select multiple tables, one of which has a spatial column, you will get the error:

Some tables contain unsupported data type or Object type: [dbo].[Table]. Please use Custom Query to exclude them.

參考資料：Azure Data Factory Geography or Geometry Data types

# Tumbling Window

A series of fixed-sized, non-overlapping and contiguous time intervals.



Tell me the count of tweets per time zone every 10 seconds

A 10-second Tumbling Window

```
SELECT TimeZone, COUNT(*) AS Count
FROM TwitterStream TIMESTAMP BY CreatedAt
GROUP BY TimeZone, TumblingWindow(second,10)
```

參考資料：Tumbling Window (Azure Stream Analytics)

# Power BI .pbids File

- `.pbids` 其實是 JSON 格式的檔案
- 所以也可以透過 Power BI 正常的 Get Data 程序連上 Data Source

```json
{
  "version": "0.1",
  "connections": [
    {
      "details": {
        "protocol": "tds",
        "address": {
          "server": "synapseworkspacem00.sql.azuresynapse.net",
          "database": "sqlpoolm00"
        }
      },
      "storageMode": "DirectQuery"
    }
  ]
}
```

# Power BI Get Data (01)

- 按下 Get data 圖示

# Power BI Get Data (02)

- 選取 Azure Synapse Analytics (SQL DW)，按下 Connect 按鈕

# Power BI Get Data (03)

- Server 輸入 Dedicated SQL endpoint
  `synapseworkspace員工編號.sql.azuresynapse.net`
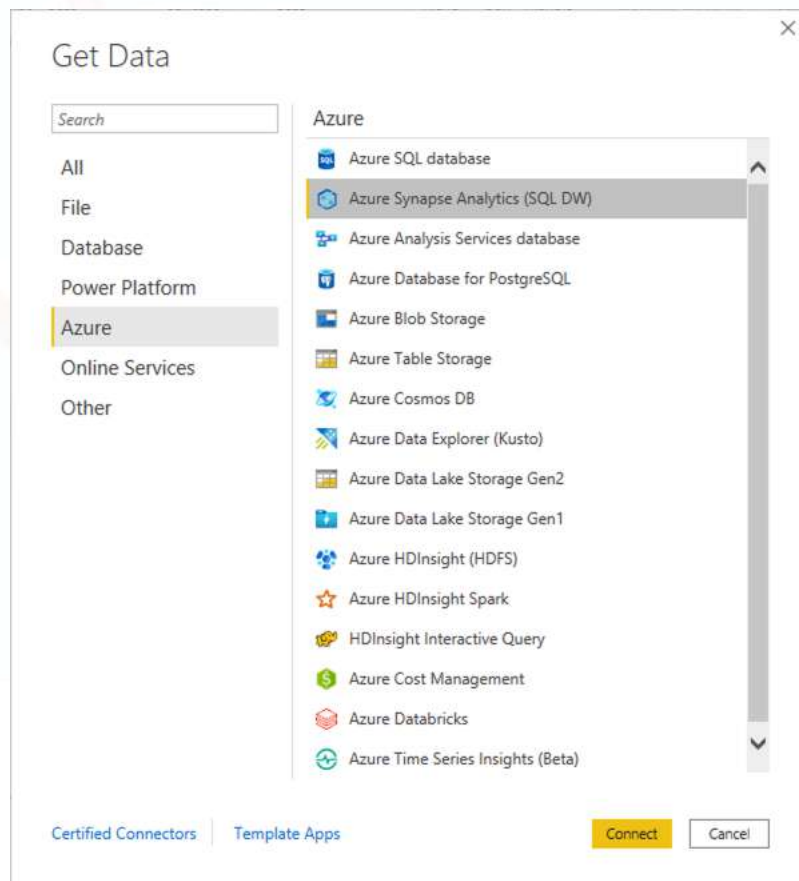- Database 輸入 Dedicated SQL Pool `sqlpoolm00`
- Data Connectivity mode 選取 Import
- 然後按下 OK 按鈕

SQL Server database                                              ×

Server ⓘ

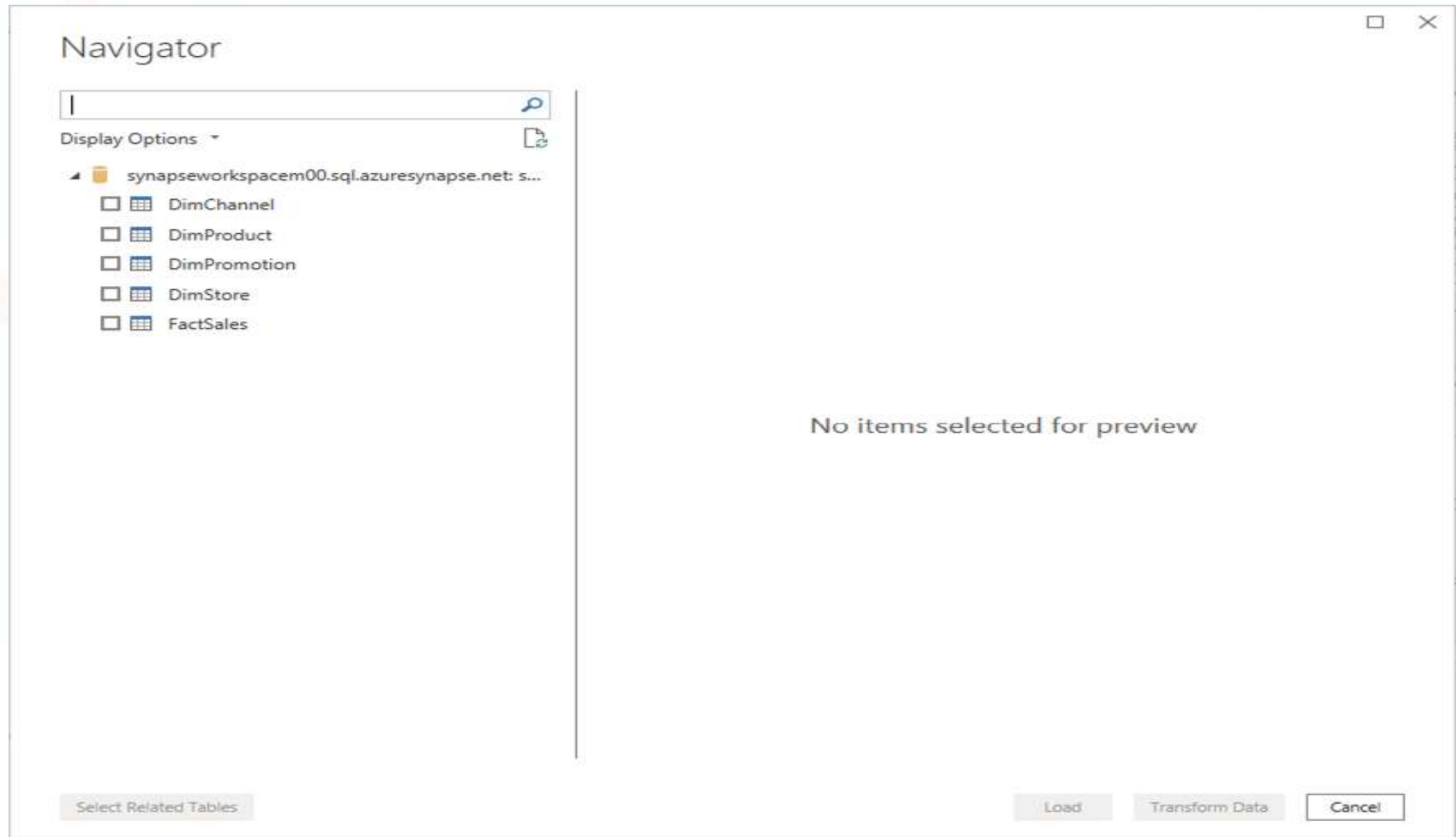synapseworkspacem00.sql.azuresynapse.ne

Database (optional)

sqlpoolm00

Data Connectivity mode ⓘ

⦿ Import

○ DirectQuery

▷ Advanced options

OK          Cancel

# Power BI Get Data (04)

- 之後就跟按下 `.pbids` 檔案的畫面一樣

# Power BI Import vs. DirectQuery

Import 與 DirectQuery：

- 所有 Data Source 都支援 Import，少部分 Data Source 才支援 DirectQuery
- Import 會真的把 Data 匯入，所以一開始 Load Data 比較慢
- DirectQuery 只是建立 Connection，所以一開始 Load Data 很快就結束
- 重點是：DirectQuery 不支援 `LOOKUPVALUE()` Function

## Remarks

- If there is a relationship between the result and search tables, in most cases, using RELATED function instead of LOOKUPVALUE is more efficient and provides better performance.

- The **search_value** and **alternateResult** parameters are evaluated before the function iterates through the rows of the search table.

- This function is not supported for use in DirectQuery mode when used in calculated columns or row-level security (RLS) rules.

參考資料：
LOOKUPVALUE
Use DirectQuery in Power BI Desktop

# References

# Reference Materials

- Azure Synapse Analytics

- What is dedicated SQL pool (formerly SQL DW) in Azure Synapse Analytics?

- Transact-SQL statements

- Az.Synapse

- Get Started with Azure Synapse Analytics

- Load data from Azure Data Lake Storage into dedicated SQL pools in Azure Synapse Analytics

- Load Contoso retail data into dedicated SQL pools in Azure Synapse Analytics