

## Slides

05 October 2021 10:05

dp-900-  
lecture-

# Core Data-Related Azure Services

**Azure Storage Accounts**

An umbrella service for various storage Types eg. Table, Files, Blob

**Azure Blob Storage**

Data which is stored as objects instead of files. Object storage is distributed storage (spanning multiple machines) for unstructured data

**Azure Tables**

A key/value NoSQL data store Intended for simpler projects

**Azure Files**

A managed file-shared NFS or SMB

**Azure Storage Explorer**

An application used to explore data within Azure Storage Accounts

**Azure Synapse Analytics**

Data warehouse and unified ai

**CosmoDB**

A fully-managed NoSQL database Can host various NoSQL engines Tables, Document, Key/Value, C

**Azure Data Lake Store (Gen2)**

A centralized data repository for Blob Storage designed for vast

**Azure Data Analytics**

Big Data as a Service (BDaaS) Write U-SQL to return data from

**Azure Data Box**

Import or export TB of data via Into Azure datacenters

# Core Data-Related Azure Services

**SQL Server on Azure Virtual Machines**

When migrating via a lift-and-shift and you want to bring your existing license and need OS access and control of VM

**Microsoft Office 365 Share**

A shared file-system for org The company owns all the f grain role-based access-con

**SQL Managed Instances**

A managed MS SQL server but with broad adaptability when migrating to Azure

**Azure Databricks**

A third-party provider part specializing in Apache Spark ELT jobs, as well as ML and

**Azure SQL**

Fully-managed MS SQL databases.

**Microsoft Power BI**

A Business Intelligence tool dashboards and interactive business decisions

**Azure Databases for <Open-Source>**

managed relational databases on Azure eg. MySQL, PostgreSQL, MariaDB

**HDInsights**

A full-managed Hadoop System many open-source Big Data

**Azure Cache for REDIS**



In-Memory data-store for returning data  
extremely vast but is also extremely volatile

data transformations for Str

## Core Data-Related Azure Services



### Azure Data Studio

An IDE that looks very much like Visual Studio Code  
But designed around data related tasks. Cross-platform  
Similar to SSIS but broader data workloads



### Azure Data Factory

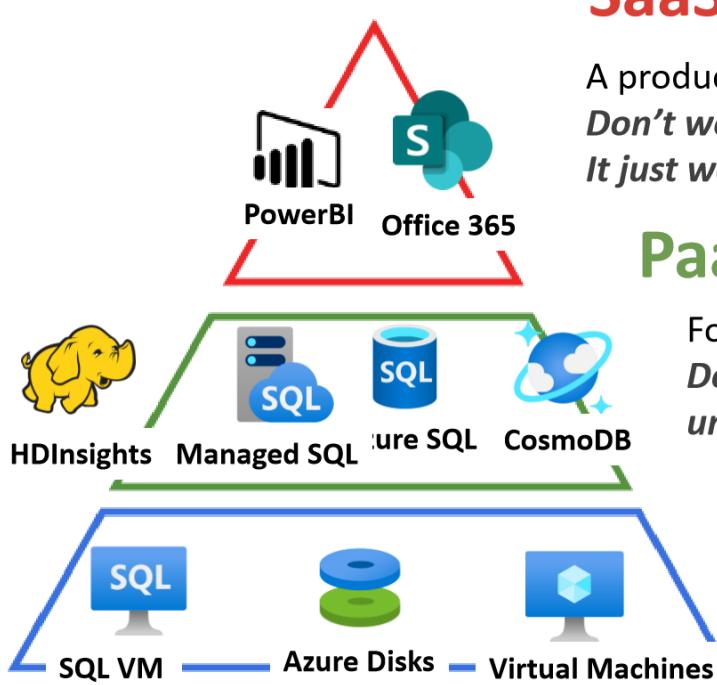
A managed ETL/ELT pipeline builder  
Easily build transformation pipelines via a web-interface



### SQL Server Integration Services (SSIS)

A stand-alone Windows app to prepare data for  
SQL workloads via transformation pipelines

## Types of Cloud Computing



### SaaS Software as a Service For Custom

A product that is run and managed by the service  
*Don't worry about how the service is maintained  
It just works and remains available.*

### PaaS Platform as a Service For Devel

Focus on the deployment and management  
*Don't worry about provisioning, configuring, understanding the hardware or OS.*

### IaaS Infrastructure as a Service

The basic building blocks for cloud IT.  
networking features, computers and  
*Don't worry about IT staff, data cent*

# Azure Data Roles

**Database Administrator:** configures and maintains a databases eg. Azure Data services or SQL server

## Responsibilities

- Database management
- Manage security, granting user access
- Backups
- Monitors Performance

## Common Tools

- Azure Data Studio
- SQL Server Management Studio
- Azure Portal
- Azure CLI

**Data Engineer:** Design and implement data tasks related to the transfer and storage of big data

## Responsibilities

- Database pipelines and process
- Data ingestion storage
- Prepare data for analytics.
- Prepare data for analytical processing

## Common Tools

- Azure Synapse Studio
- SQL
- Azure CLI

**Data Analyst:** Analyzes business data to reveal important information

## Responsibilities

- Provides insights into the data
- Visual reporting
- Modeling data for analysis
- Combines data for visualization and analysis

## Common Tools

- Power BI Desktop
- Power BI Portal
- Power BI services
- Power BI report builder

# Database Administrator – Common Tools

configures and maintains a databases

## Azure Data Studio Looks like VSC



Connect to Azure SQL, Azure SQL data warehouse, Postgres SQL and SQL Server (big data)

- Various libraries and extensions along with automation tools.
- Graphical interface for managing on-premises and cloud-based data services.
- runs on Windows, macOS, Linux
- Possibly a replacement for SSMS (still lacks some features of SSMS)

## SQL Server Management Studio (SSMS)



- Automation tooling for running SQL commands or common database
- Graphical interface for managing on-premises and cloud-based data services
- **Runs on Windows**
- More mature than Azure Data Studio



## Azure Portal and CLI

- Manage SQL database configurations. eg create, deleting, resizing, number of cores
- Manage and provision other Azure Data Services
- Automate the creating, updating or modifying resources via Azure Resource Manager

# Data Engineering – Common Tools

Design and implement data tasks related to the transfer and storage of big data



## Azure Synapse Studio

azure portal integrated to manage azure synapse, data ingestion (Azure data factory), management of azure synapse assets (SQL Pools/Spark Pool)



## Knowledge SQL

Create databases., tables, views, etc



## Azure CLI

Support operations SQL cmd to connect to Microsoft server Azure SQL data and run a talk queries and commands



## HDInsights

Streaming data via Apache Kafka or Apache Spark  
Applying ELT jobs via HIVE, PIG, Apache Spark



## Azure Databricks

Using Apache Spark to create ELT or streaming jobs to data ware houses or data lakes

# Data Analyst – Common Tools

Analyzes business data to reveal important information



## Power BI Desktop

A stand alone application for data visualization  
You can do data modelling  
Connect to many data sources  
Create interactive reports



## Power BI Portal/Power BI Service

A web ui for creating interactive dashboards



## Power BI Report builder

Create paginated reports (printable reports)

# Data Overview

**Data – units of information****Data Documents** – types of abstract groupings of data**Data Sets** – unstructured logical grouping of data**Data Structures** – structured data**Data Types** – how single units of data are intended to be used**Batch and Streaming Data**

How do we move our data?

**Relational and Non Relational**

How do access, query and store data?

**Data Modelling**

How do we prepare and model data?

**Schemas and Schemata**

How do we structure data?

**Data Integrity and Data Consistency**

How do we trust our data?

**Normalized and Denormalized**

How do we trade quality for performance?

# Introduction to Data

**What is data?**

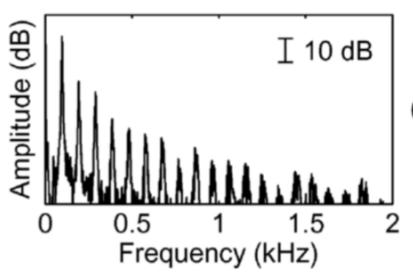
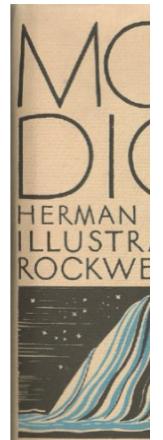
Data is **units of information** that could be in the form of **numbers, text or machine code, images, videos, audio or physical (handwriting)**



```

00000060 03 01 A1 50 53 00 3A 00 C3 00 36 00 32 25 00 00 ...PS...6.2%
00000070 FF FF 83 00 50 0E 00 56 00 41 00 0A 00 4A 26 00 00 ...P.V.A.J6...
00000080 03 00 01 50 0E 00 56 00 41 00 0A 00 4A 26 00 00 ...&A.p.l.y...
00000090 FF FF 80 00 26 00 41 00 70 00 70 00 6C 00 79 00 ...t.o.a.l.l...
000000A0 20 00 74 00 6F 00 20 00 61 00 6C 00 6C 00 00 00 ...P...
000000B0 00 00 00 00 00 00 00 00 00 00 00 01 00 01 50 ...}2....P
000000C0 7E 00 7D 00 32 00 0E 00 01 00 00 FF FF 80 00 ...}2...
000000D0 4F 00 4B 00 00 00 00 00 00 00 00 00 00 00 00 00 ...O.K...
000000E0 00 00 01 50 B4 00 70 00 3E 00 0E 00 02 00 00 ...}2...
000000F0 FF FF 80 00 43 00 61 00 6E 00 60 00 65 00 6C 00 ...C.a.n.c.e.l...
00000100 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ...}2...
00000110 EA 00 7D 00 32 00 0E 00 09 00 00 FF FF 80 00 ...}2...
00000120 26 00 48 00 65 00 6C 00 70 00 00 00 00 00 00 00 ...&H.e.l.p...
00000130 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ...P...
00000140 3B 00 0E 00 2F 25 00 00 FF FF 81 00 00 00 00 00 ...%...
00000150 00 00 00 00 00 00 00 00 00 02 50 00 00 30 00 ...P.0...
00000160 1E 00 08 00 EE 25 00 00 FF FF 82 00 46 00 69 00 ...%...F.i...
00000170 6C 00 65 00 20 00 54 00 79 00 70 00 65 00 00 00 ...l.e.T.y.p...
00000180 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ...P...
00000190 54 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ...T.O...%
000001a0 50 00 61 00 72 00 73 00 69 00 6E 00 67 00 20 00 ...P.a.r.s.i.n.g...
000001b0 52 00 75 00 6C 00 65 00 73 00 00 00 00 00 00 00 00 ...R.u.l.e.s...
000001c0 00 00 00 00 00 00 00 00 07 00 00 50 00 06 07 00 ...W...
000001d0 41 03 34 00 FD 25 00 00 FF FF 80 00 00 00 00 00 ...

```



$$\int_a^b f'(x) dx = f(b) - f(a)$$

$$\frac{d}{dx} \int_a^x f(t) dt = f(x)$$

# Data Documents

**What are data documents?**A data document defines the **collective form in which data exists**

A data document defines the **context** in which data exists.

Common types of data documents:

- **Datasets** — a logical grouping of data
- Databases — structured data that can be quickly accessed and searched
- Datastores — unstructured or semi-structured data housing data
- Data warehouses — structured or semi-structured data for creating reports and analysis
- Notebooks — data that is arranged in pages, designed for easy consumption



MNIST Dataset  
Dataset



Azure SQL  
Database



Azure Data Lake  
Datastore



Azure Synapse Analytics  
Data warehouse



Jupyter  
Notebook

## Data Sets

### What is a dataset?

A data set is a **logical grouping of units of data** that generally are closely related and share the same data structure.

There are **publicly available data** sets that are used in the **learning of statistics, data analytics, machine learning**

#### MNIST database

Images of **handwritten digits** used to test classification, clustering, and image processing algorithms.

Commonly used when learning how to build computer vision ML models to translate handwriting into digital text

0 0 0 0 0 0 0 0 0 0  
1 1 1 1 1 1 1 1 1 1  
2 2 2 2 2 2 2 2 2 2  
3 3 3 3 3 3 3 3 3 3  
4 4 4 4 4 4 4 4 4 4  
5 5 5 5 5 5 5 5 5 5  
6 6 6 6 6 6 6 6 6 6  
7 7 7 7 7 7 7 7 7 7  
8 8 8 8 8 8 8 8 8 8  
9 9 9 9 9 9 9 9 9 9

## More Data Sets....

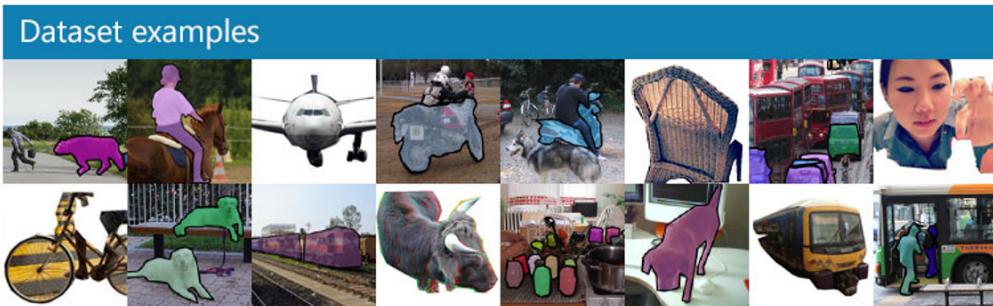
#### Common Objects In Context (COCO) dataset

A dataset which contains many common images using a JSON file (coco format) that identify objects or

This dataset features:

- Object segmentation
- Recognition
- Supervision

segments within an image.



- superpixel segmentation
- 329K images
- 0.5 million objects
- 79 object categories
- 90 stuff categories
- 4 captions per image
- 249,000 people

## More Data Sets...

### IMDB Reviews Dataset

A **movie review** dataset with 25,000 highly polar movie reviews for training, and 25,000 for testing

Could be useful for determine customer sentiment analysis.

★ 8/10

**A fun and enjoyable surprise.**  
cedrickroberts 9 September 2004

I am an Eddie Murphy fan, but I did not go to the theater to see this movie so much, that when I finally saw the movie, and was not repulsed by it, I began to relax and enjoy it. I think that many of the critics were expecting a Cop-type laugh fest, but that is not a comedy that's really an action-comedy. This is strictly a sci-fi comedy, as such, it does not take itself seriously. Neither should you if you decide to watch it. Just see if you enjoy it. :)

## More Data Sets...

### Free Music Archive (FMA)

A dataset for music analysis

- 106,573 tracks
- 163 genres

### LibriSpeech

A dataset of 1000 hours of English speech

There are many more datasets online

(some could be paid, or you need to extract the data via an API or scrape the data)

- Crunchbase
- Glassdoor Research
- Open Corporates
- FBI Uniform Crime Reporting
- Uppsala Conflict Data Program
- Dbpedia
- Google Trends
- DataHub – Stock Market
- Center of Disease Control (CDC)
- World Health Organization
- Statista Video Games data
- Data.gov.uk
- Open Data Canada
- NVC Taxi Trip Data
- Weather.gov
- AWS open registry
- Google Public Data
- Reddit Datasets
- USD Food Consumption
- *and many more...*

# Data Types

## What is a data type?

A data type is **a single unit of data** that tells a compiler or interpreter (computer **data is intended to be used**).



The **variety** of data types will **greatly vary based on the computer program**.

Let's us take the most common data types.

### Numeric Data Types: A data type involving mathematical **numbers**

- **Integer** – a whole number, (could be negative or positive) -100, 7, 11, 2190381
- **Float** – a number that has a decimal e.g. 1.5 , 0.0, -10.24, 9.432363535345



```
my_int = 1
my_float = 2.2
```

# Data Types

### Text Data Types: A data type that contains readable and non readable **letters**

- **Character** – a single letter, alphanumeric (A-Z), digit (0-9), blank space, punctuation, special character
- **String** – a sequence of characters eg. Words, sentences and paragraphs



```
my_char = 'a'
my_string = "We prefer to help ourselves"
```

### Composite: A data type that contains **cells of data** that can be **accessed via an index or a key**

- **Array** – a group of elements that contain the **same data type**, can be accessed via their index
  - **Hash (Dictionary)** – a group of elements where a key can be used retrieve a value
- Composites can be both data-types and data structures*



```
my_arr = ['live', 'long', 'and', 'prosper']
my_dict = { "Speed": 1, "Accuracy": 2 }
```

## Data Types

**Binary Data Type** – represented by a **bit or a series of bits (a byte)**, Which is either 0 (off) or 1 (on)

```
one_byte = int('11110000',
```

**Boolean Data Type** – A datatype that is either **True or False**

- Some languages represents a Boolean as
  - a bit as a Boolean eg. 0 (false) or 1 (true)
  - the first letter eg. t (true) or f (false)



```
my_bool = True
```

**Enumeration (Enum) Data Type** – a group of constant **(unchangeable)** variables eg. DIAMOND, SPADE, HEART, CLUBS

- *Can be a data type and/or a data structure, varies on the language*

↑  
All cast groups

```
class Shake(Enum):
    VANILLA = 7
    CHOCOLATE = 4
    COOKIES = 9
    MINT = 3

Shake.VANILLA
Shake.CHOCOLATE
Shake.COOKIES
Shake.MINT
```

## Schema

### What is a schema?

A schema (in terms of databases) is **a formal language which describes the structure of data** (blueprint)  
A schema can **define many different data structures** that serve different purposes for a database.

Different data structures (relational databases):

- |                 |                      |
|-----------------|----------------------|
| • Tables        | • Queues             |
| • Fields        | • Triggers           |
| • Relationships | • types,             |
| • Views         | • Sequences          |
| • Indexes       | • materialized views |
| • Packages      | • Synonyms           |
| • Procedures    | • database links     |

Data type

```
ActiveRecord::Schema.define(version: 20140212175753) do
  # These are extensions that must be enabled in order to support
  enable_extension "plpgsql"

  create_table "activities", force: true do |t|
    t.integer "company_id"
    t.integer "project_id"
    t.integer "target_id"
    t.string "target_type"
    t.integer "user_id"
    t.string "name"
    t.string "key"
    t.datetime "created_at", null: false
    t.datetime "updated_at", null: false
  end
end
```

- Functions
- Directories
- XML schemas

↑ use schema to  
create there

```
t.json "data"
end
add_index "activities", ["company_id"], name: "index_activities_
add_index "activities", ["project_id"], name: "index_activities_
add_index "activities", ["target_id"], name: "index_activities_o
t.integer "project_id"
end
# ...
```

A Ruby on Rails schema that defines the structure for

## Schemaless

### What is a schemaless?

Schemaless is when the **primary “cell” of database can accept many types**  
This allow developers to **forgo upfront data modelling**

Common schemaless databases are:

- Key/Value
- Document
- Columns
  - Wide Column
- Graph

Data type more flexible

## Query and Querying

SELECT \* FROM  
CrewMembers



### What is a query?

A query is a request for data results (reads) or to perform operations updating deleting data (writes).

A query can perform maintenance operations on the data and is not just working with the data that resides within the database.

### What is a data result?

Results are the data returned from a query

### What is querying

The act of performing a query

### What is query language?

A scripting or programming language designed to submit a request or actions to a database. Notable query languages:

- SQL

ID	Name	Title	StarDate
10	Picard	Captain	41133.7
11	Riker	Commander	41154.2
12	Data	Lt Commander	41114.3

13	Troi	Lt Commander	41412.1
14	Crusher	Command	41520.4

- GraphQL
- Kusto
- Xpath
- Gremlin

# Batch vs Stream Processing

## Batch Processing

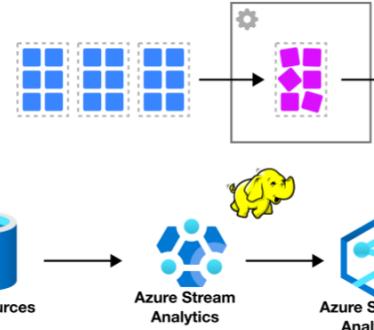
When you **send batches** (a collection) of data to be processed.

Batches are generally scheduled: eg. Every day at 1PM

Batches are **not real-time**

Batches processing is ideal for very large processing workloads

Batch processing is **more cost-effective**



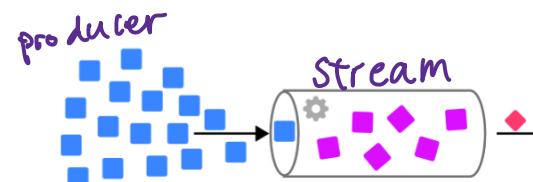
## Stream Processing

When you process data as soon as it arrives:

- **Producers** will send data to a stream and
- **Consumers** will pull from the stream

Stream processing is good for real-time analytics or  
real-time processing (streaming video)

**Much more expensive** than batch processing



# Relational Data

## Tables

A logical grouping of rows and columns. Think like a Excel spreadsheet

- Tabular data – data that makes use of table data structures

## Views

Views is a result set of a stored query on data **stored in memory** (a temporary or virtual table)

## Materialized Views

Material Views is a result set of stored query on data **stored on disk**

## Indexes

A copy of your data sorted by one or multiple columns for faster reads at the cost of storage

## Constraints

rules applied to writes, that can ensure data integrity: eg. don't allow duplicate records

## Triggers

a function that is trigger on specific database events

## Primary Key

one or multiple columns that uniquely identify a table in a row eg. **Id**

## Foreign key

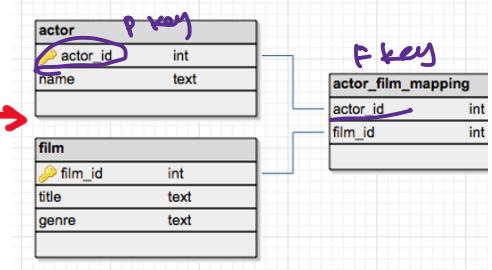
a column which holds the value of primary key from another key to establish a relationship

- A relationship is when two tables have a reference to one another to join data together

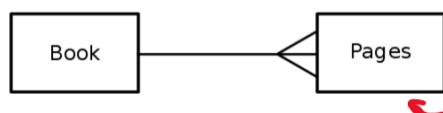
# Relational Data – Relationships

Relational databases **establish relationships to other tables** via **foreign keys** Referencing another table's primary key.

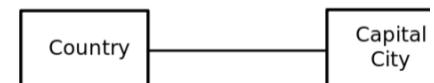
These are the relationships found in a relational data:



**One-to-one**  
A monkey has a banana



**One-to-many**  
A store has many customers



**Many-to-many**  
A project has many tasks and Tasks can belong to many projects

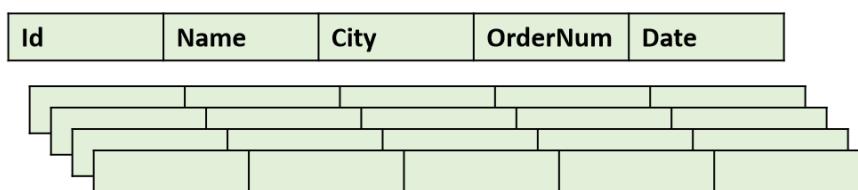
**Many to Many (via Join/Junction Table)**  
A student has many classes through enrollments  
A class has many students through enrollments

*A book has many authors through a library, vice versa*



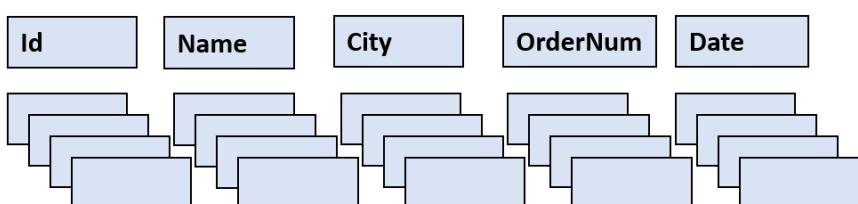
# Row-store vs Column-store

When looking at tabular data there are two ways data can be stored:  
Row-store or Column Store



## Row-store

- Data is organized in rows
- Traditional relational databases are
- Good for general purpose database
- Suited for Online transaction proces
- Great when needing all possible col**
- Not the best at analytics or massive



## Column-store

- Data is organized into columns
  - Faster at** aggregating values f
- NoSQL store or SQL-Like database
- Great for vast amount of data
- Suited for Online analytical proces
- Great when you only need a few c

# Indexes

A database index is a data structure that improves the speed of reads from the database by storing the same or partial redundant data organized in a more efficient logical order.

The logical ordering is commonly determined by one or more columns: sort key(s)

A common data structure of an index is a **Balanced Tree (B-Tree)**

**Creating an index in Postgres**

```
CREATE INDEX idx_address_phone
ON address(phone)
```

Table (phone book)      Index

Name	Number	ID
Andrew	626-9009	1
Carly	641-5324	2
Rishab	345-4121	3
Cindy	767-3423	4
Peter	623-2413	5
Lisa	767-5235	6
Otto	344-5353	7
Mona	345-6189	8
Zack	626-4421	9
Maya	767-7771	10
		.....

# Data Integrity and Data Corruption

Data integrity is the **maintenance and assurance of, data accuracy and consistency** over time.

Used at as **proxy term for data quality**, data validation is a pre-requisite for data integrity.

The goal of data integrity:  
**ensure data is recorded exactly as intended**

Data integrity is the **opposite** of Data corruption

**Data corruption** is the **act or state of data not being in the intended state** and will result in data loss or misinformation.

**Data corruption** occurs when unintended changes result when reading and writing:

- Unexpected hardware failure
- Human error when inputting, modifying data
- Malicious actors with intent of corrupting your data
- Unforeseen side effects for automated operations via computer code

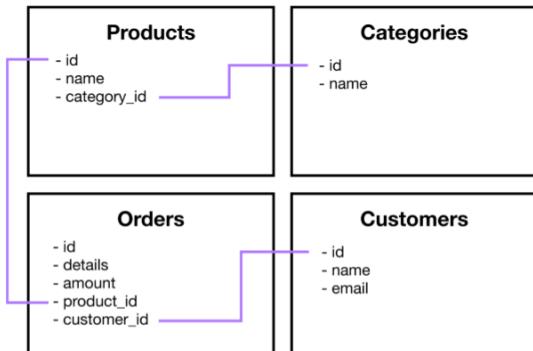
**How do we ensure data integrity?**

- Have a well defined and documented schema
- Logical constraints on your database
- Redundant and versions of your data to restore
- Human analysis of the data
- Hash functions to determine if changes have been tampered
- Principle of least-privilege, (limiting actions for specific user roles)

# Normalized vs Denormalized Data

**Normalized**

A schema design to store **non-redundant** and **consistent data**

**Denormalized**

A schema that **combines data** so that **accessing data (querying) is fast**

*combine*

Customer Orders
- id
- product_name
- product_code
- category_name
- customer_name
- cusomter_email
- order_id
- order_details
- order_amount

- Data Integrity is maintained
- Little to no redundant data
- Many tables
- Optimizes for storage of data
- **Data Integrity is not maintained**
- Redundant data is common
- Fewer tables
- Excessive data, storage is less opt

## Pivot Table

A pivot table is **a table of statistics that summarizes** the data of a more extensive table from a Database, Spreadsheet or Business intelligence (BI) tool

- Pivot tables are a technique in **data processing**
- They arrange and rearrange (or "pivot") statistics in order to **draw attention to useful information**
- This leads to **finding figures and facts quickly** making them integral to data analysis.



In **Microsoft Excel** its very easy to create Pivot Tables. Think of a pivot table as an interactive report where you can quickly aggregate (group) your data based on various factors eg.

- By Year Month, Week or Day
- Sum, Average, Min or Max

Sample sales data					
Date	Color	Region	Units	Sales	
3-Jan-16	Red	West	1	\$11.00	
13-Jan-16	Blue	South	8	\$96.00	
21-Jan-16	Green	West	2	\$26.00	
30-Jan-16	Blue	North	7	\$84.00	
7-Feb-16	Green	North	8	\$104.00	
13-Feb-16	Red	South	2	\$22.00	
21-Feb-16	Blue	East	5	\$60.00	
1-Mar-16	Green	West	2	\$26.00	
13-Mar-16	Blue	East	8	\$96.00	
23-Mar-16	Blue	North	7	\$84.00	
28-Mar-16	Green	West	2	\$26.00	
3-Apr-16	Blue	South	8	\$96.00	

<https://exceljet.net/>

"PivotTable" used to be a trademarked word owned by Microsoft

## Strongly Consistent vs Eventually Consistent

### What is data consistency?

When data being kept in two different place and **whether the data exactly match** or do not ma

When you have to have **duplicates** of your data in many places and need to **keep them up-to-date to be exact matching**, based on how data is transmitted and service levels cloud service providers will use these two terms:

### **Strongly Consistent**

Every time you request data (query) you can expect consistent data to be returned with x time (1 seconds)

We will never return to you old data. But you will have to wait at least 2 seconds for the query to return

### **Eventually Consistent**

When you request data you may get back within 2 seconds.

We are giving you whatever data is currently available. You may get new data or old data, but if you wait generally be up to date.

## Synchronous vs Asynchronous

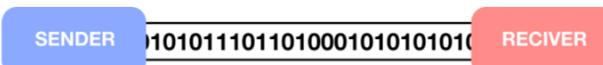
**Synchronous and Asynchronous** can refer to **mechanism for data transmission or data replication**

### **Synchronous**

continuous stream of data that is synchronized by a timer or clock (guarantee of time)

Can only access data once transfer is complete

- Guaranteed consistency of data return at time of access
- Slower access times



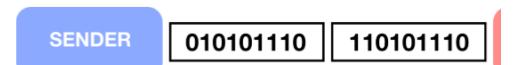
A company has a primary database, but they need to have a backup database in case their primary database fails. The company cannot lose any data, so everything must be in-sync. The database is not going to be accessed while it is standing by to act as replacement.

### **Asynchronous**

continuous stream of data separated by time (no guarantee of time)

Can access data anytime but may return old placeholder

- Faster access times, not guarantee of consistency



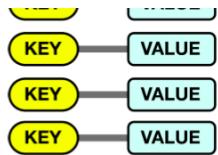
A company has a primary database, and a read-replica (copy of the database). An analytics person can create computationally intensive reports that do not impact the primary database. It does not matter if the data is not 1-to-1 at time of access.

## Non Relational Data

A non-relational database **stores data in a non-tabular form** and will be optimized for different kind

**Types of non-relational databases:**  
Key/Value

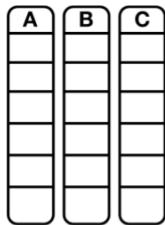
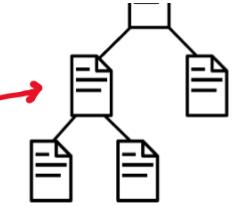


**Key/Value**

- Each value has a key
- Designed to scale
- Only simple lookups

**Document**

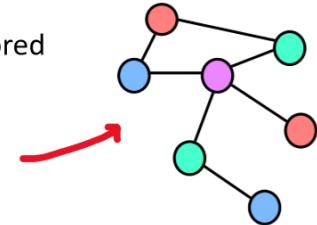
- Primary entity is a JSON-like data-structure called a document

**Columnar**

- Has a table-like structure but data is stored around columns instead of rows

**Graph**

- Data is represented with nodes and structures. Where relationships matter



Sometimes non-relational database can be both Key/Value and Document  
eg. Azure Cosmo DB or Amazon DynamoDB

## Wide-Column vs Columnar

Columnar data store will store each column separately **on disk**

Wide-column database is a type of columnar database that supports a column family together on disk, **not just a single column**

## Data Sources

### What is a data source?

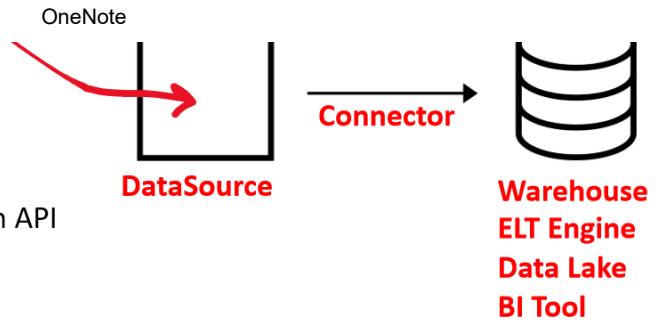
A data source is **where data originates from**.

An analytics tool may be connected to various data sources to create a visualization.

A **data source** could be a:

*Data f(x)*

- Data lake
- Data warehouse
- Datastore
- Database
- Data requested on demand from an API endpoint from a web-app
- Flat files (e.g. excel spreadsheet)



### Extracting data from data sources

A data tool like Business Intelligence (BI) software would establish a connection to multiple data sources. A BI will *extract* data which could be pulled at the time report, or could pull data on schedule, or data could be streamed. The mechanism for extracting data will vary per data source.

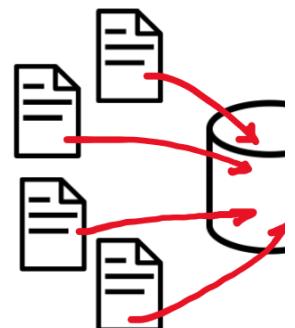
## What is Datastore?

A Datastore is a **repository for persistently storing and managing collections of unstructured or semi-structured data**

Data store is a very broad term, and interchangeable used with databases. But generally a data store indicates working **unstructured or semi-structured data**

A data store can be specialized in storing:

- Flat Files (documents, spreadsheets)
- Email (an email server)
- Databases (**complex data stores developed using formal design and modeling techniques**)
  - Relational Databases
  - NoSQL Databases
  - Object-Oriented Databases
- Data stores designed to be distributed across many machines
- Directory services



## What is a Database?

A database is a **data-store that stores semi-structured and structured data**.

A database is more **complex data stores** because it **requires using formal design and modeling techniques**

Databases can be generally categorized as either:

- **Relational databases**
  - Structured data that strongly represents **tabular data** (tables, rows and columns)
  - Row-oriented or Columnar-oriented

### • Non-relational databases

- Semi-structured that may or may not distantly resemble tabular data.

Databases have a rich set of functionality:

- specialized language to query (retrieve data)
- specialized modeling strategies to optimize retrieval for different use cases
- more fine tune control over the transformation of the data into useful data structures or reports

e.g. SQL

**SQL**

Database



Normally a databases infers someone is using a **relational row-oriented data store**

## What is Data Warehouse?

A **relational datastore** designed for **analytic workloads**, which is generally **column-oriented**

Companies will have **terabytes and millions of rows of data**, and they need a fast **way to be able to produce analytics reports**

Data warehouses generally perform **aggregation**

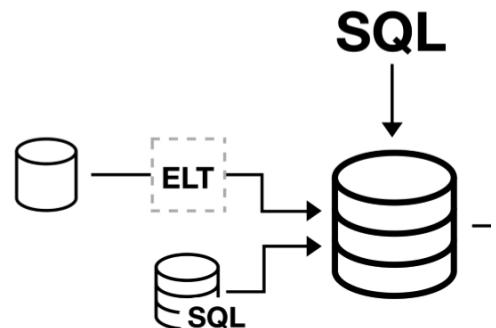
- aggregation is grouping data eg. find a total or average
- Data warehouses are optimized around columns since they need to quickly aggregate column data

Data warehouses are generally designed be **HOT**

- Hot means they can return queries very very fast even though they have vast amounts of data

Data warehouses are **infrequently accessed** meaning they aren't intended for **real-time reporting** but maybe once or twice a day or once a week to generate business and user reports.

A data warehouse needs to consume data from a relational databases on a regular basis.



generally are  
not for trav

## What is Data Mart?

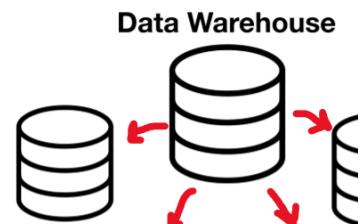
A Data mart **subset of a data warehouse**

A data mart will store under **100 GB** and has a **single business focus**

Data mart allows different teams or departments to have control over their own dataset for their specific use case

Data marts are generally designed be **read-only**

Smaller dataset, lower cost of query



Data marts also increase the frequency at which data can be accessed.

The cost to query the data is much lower and so queries can be performed multiple times a day or even hourly.



## What is a Data Lake?

A data lake is a centralized storage repository that holds a vast amount of raw data (big either a semi-structured or unstructured format).

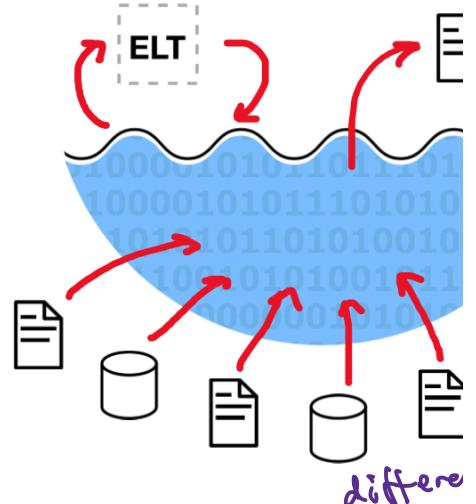
A data lake lets you store all your data without careful design or having to answer questions on the future use of the data. (*Hording for data scientist*)

A data lake is commonly accessed for data workloads such as:

- Visualizations (Business Intelligence)
- Real-time analytics
- Machine Learning
- On-premise data

Data lakes are great for data-scientists but its very hard to use data lake for BI reporting

If data lakes are not well-maintained they can become **data-swamps** (a mess of data)



## What is a Data Lakehouse?

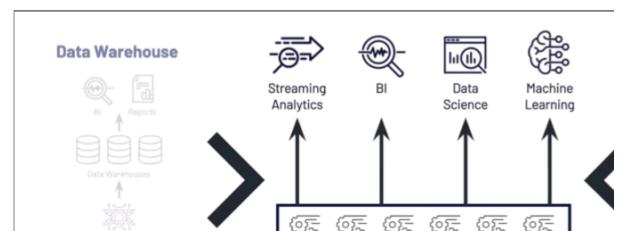
A Data Lakehouse combines the best elements of a **data lake** and a **data warehouse**.

Data Lakehouses compared to a Data warehouse can:

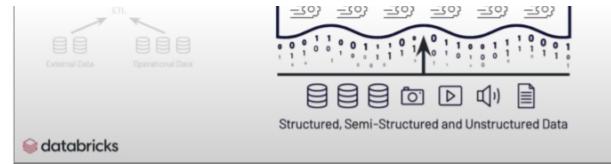
- support video, audio and text files
- Support data science and ML workloads
- have support for both streaming and ELT
- Work with many open-source formats
- Data will generally reside in a **data lake or blob stores**

Data Lakehouses compared to a Data lake:

- perform BI tasks very well
- much easier to setup and maintain
- has management features to avoid data swamp
- more performant than a data lake



An example of a Data Lakehouse  
is **Apache Delta Lake**



## Data Sources Comparisons

Datastore — a repository for data for **unstructured or semi-structured** data

Database — a repository for structured data, commonly refers to **Relational databases**

Data warehouse — a datastore designed for analytic workloads, generally **columnar store**

Data Lake — a centralized datastore for store vast amounts of data with distributed storage, generally **for data analytics**

Data Lakehouse — combines the best elements of a data lake and data warehouse

Data mart — **a subset of data warehouse**, it stores under 100 GB and has a single business focus

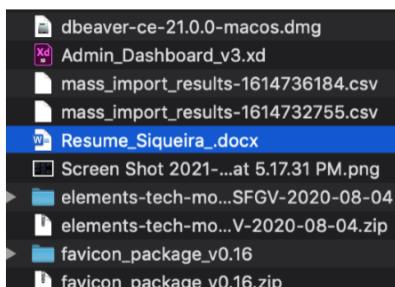
## Data Structures

### What is a Data Structure?

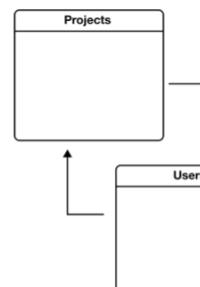
Data that is organized in a **specific storage format**, that enable easy access and modification.  
A data structure can store various data types.

Data can be abstractly described to have a degree of structure:

- **Unstructured** — a bunch of lose data that has no organization or possible meaning.
- **Semi-Structured** — data that can be browsed or searched (with limitations).
- **Structured** — data that can be easily browsed or searched



```
<?xml version="1.0"?>
<Root>
  <Customer>
    <Order OrderId="152313">
      <Name>Self-sealing stem bolts</Name>
      <Quantity>100</Quantity>
    </Order>
  </Customer>
</Root>
```



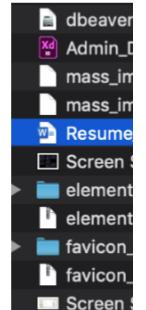
Unstructured

Semi-structured

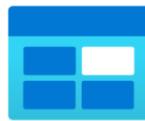
Struct

# Data Structures – Unstructured

**Unstructured data** is just a **bunch of loose data**, think of a junk folder on your computer with a bunch of random files, not optimized for search or analysis, or simply no relation between various data.



Microsoft and Azure services that store unstructured data

**Microsoft SharePoint**Shared documents for  
an organization**Azure Blob Storage**

Unstructured object datastore

**Azure Files**Mountable file system for  
storing unstructured files**Azure Data Lake Storage**For big data  
From many

# Data Structures – Semi-Structured

**Semi-structured data** is (\*no schema) data **has some form of relationship**, its easy to browse data to find related data, you can search data but there are limitations or when you search you will pay at a computative or operational cost

Concrete semi-structure data structures:

- XML, JSON, AVRO, PARQUET



Azure and other services that store semi-structured data

**Azure Tables**

A key/value data store

**Azure Cosmos DB**A document data store  
designed for global scale**MongoDB**Open-source document store  
NoSQL database**Apache Cassandra**Open-source  
NoSQL database

# Semi Structured Data structure

## What is semi-structured data?

Semi-structured data is data **that contains fields**.

The fields don't have to be the same in every entity.

You only define the fields that you need on a per-entity basis.

## Common semi-structured data structures:

### JavaScript Object Notation (JSON)

Format used in JavaScript notation; Store data in memory, read and write from files.

### Apache Optimized Row Columnar format (ORC)

organizes data into columns rather than rows (columnar store data structure).

### Apache Parquet

Another columnar data structure. A Parquet file contains **row groups**.

### Apache AVRO

row-based format. Each record contains a header that describes the structure of the data in the record.

# Semi Structured – JSON

**JSON** (JavaScript Object Notation) is a lightweight data-interchange format

- It is easy for humans to read and write.
- It is easy for machines to parse and generate
- It is based on a subset of the JavaScript

```
{
  "starships": {
    "enterprise": {
      "registry": [
        "NCC-1701",
        "NCC-1701-B",
        "NCC-1701-C",
        "NCC-1701-D"
      ]
    }
  }
}
```

JSON is built on two structures

1. **A collection of name/value pairs**
  - In other languages: realized as an *object*, *record*, *dictionary*, *hash table*, *keyed list*, or *associative array*
2. **An ordered list of values**
  - In other languages: realized as an *array*, *vector*, *list*, or *sequence*

JSON is a **text format** that is completely language independent

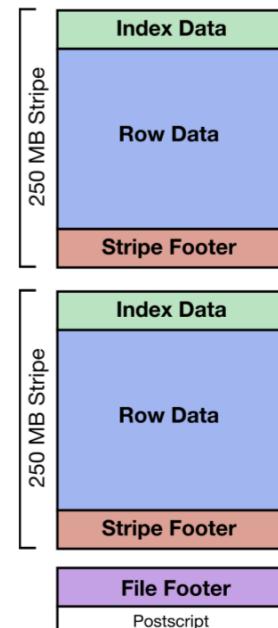
# Semi Structured – ORC

## Apache ORC (Optimized Row Columnar) a storage format of the Apache Hadoop ecosystem

- It is similar to RCFile and Parquet files and is the successor to RCFiles
- It was developed by Facebook to support columnar reads, predictive pushdown and lazy reads
- It is more storage efficient than RCFiles (taking up 75% less space)
- ORC only supports Hadoop's HIVE and PIG
- ORC performs better with HIVE than Parquet files
- ORC files are organized into **stripes of data**

### The Anatomy of an ORC file

- **File footer** stores auxiliary information
  - list of stripes in the file
  - number of rows per stripe
  - each column's data type.
  - column-level aggregates count, min, max, and sum
- **stripe footer** contains a directory of stream locations
- **Row data** is used in table scans
- **Index data** includes min and max values for each column and the row positions within each column
- The default stripe size is **250 MB**
- Large stripe sizes enable large, efficient reads from HDFS

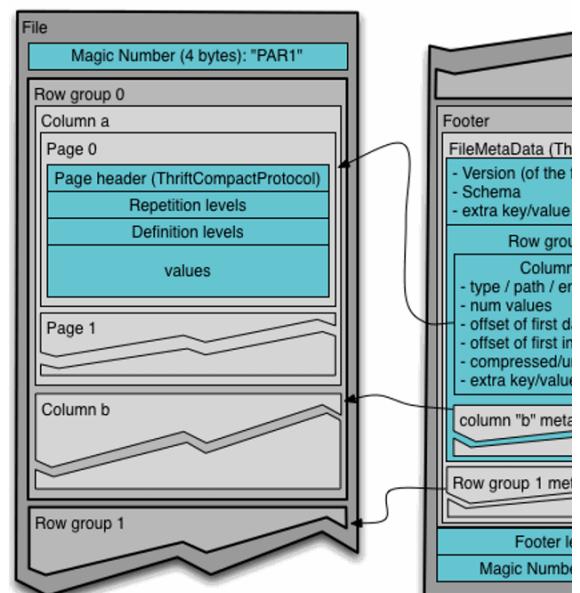


## Semi Structured – Parquet

Apache Parquet is a **columnar storage file format** available to **any project in the Hadoop ecosystem** (Hive, Hbase, MapReduce, Pig, Spark)

- Parquet is built to support very efficient compression and encoding scheme
- uses the record shredding and assembly algorithm

*fix*



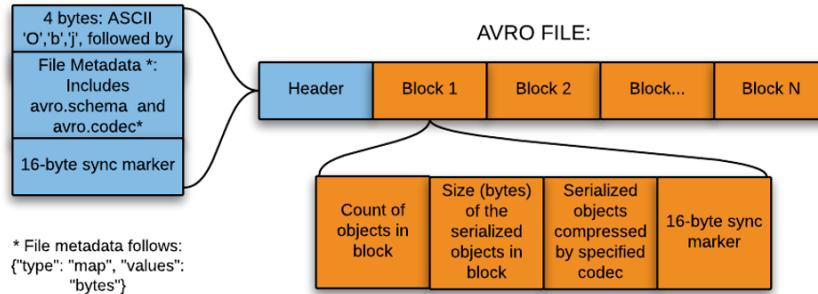
## Semi Structured – AVRO



Apache AVRO is a row-based format that provides:

- Rich data structures.
- A compact, fast, binary data format.
- A container file, to store persistent data.
- Remote procedure call (RPC).
- Simple integration with dynamic languages

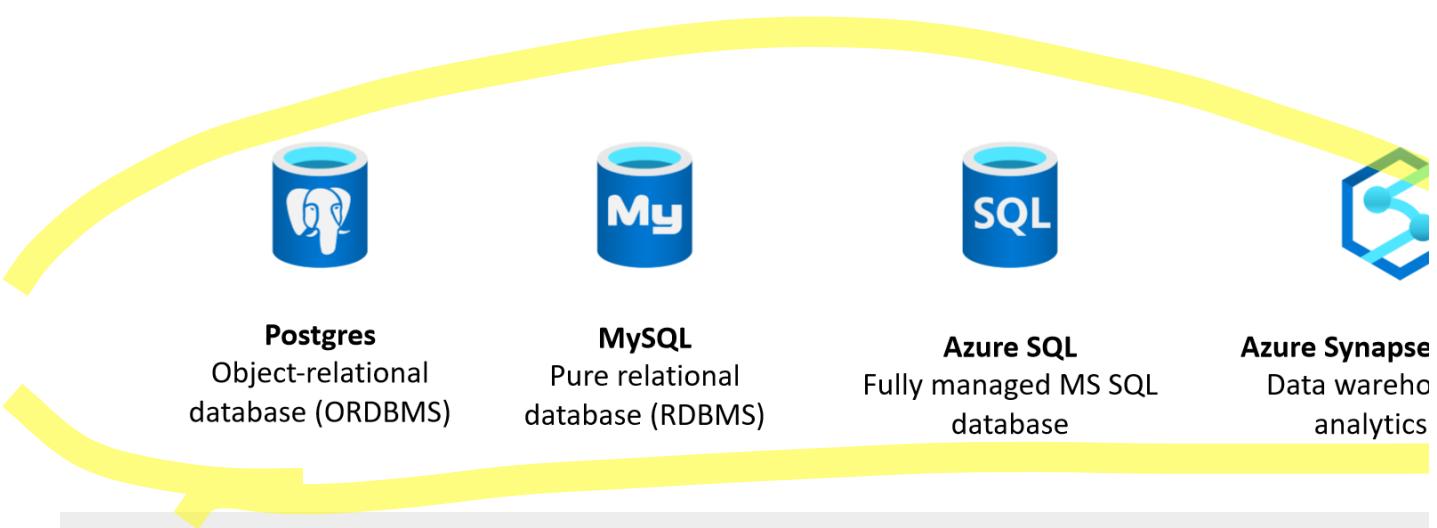
Avro provides functionality similar to systems such as Thrift, Protocol Buffers



## Data Structures – Structured

**Structured data** is (schema) data has a relationship, its easy to browse to find related data, its easy to search data. The most common structured data is tabular data (representing row and columns)

The most common data structure is **tabular data**



## What is Data Mining?

### What is Data Mining?

The **extraction of patterns and knowledge** from large amounts of data (**not the extraction**)

Cross-industry standard process for data mining (CRISP-DM) defines Data Mining into 6 phases:

### 1. Business understanding

What does the business need

### 2. Data understanding

What data do we have, and what data do we need?

### 3. Data preparation

How do we organize the data for modeling?

### 4. Modeling

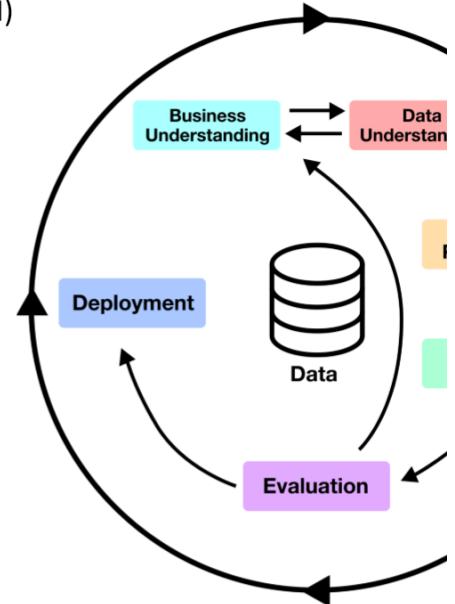
What modeling techniques should we apply?

### 5. Evaluation

Which data model best meets the business objective?

### 6. Deployment

How do people access the data?



# Data Mining Methods

**Data Mining methods** or techniques is a way to **find valid patterns and relationships**

### Classification

classify data in different classes

### Clustering

a division of information into groups of connected objects

### Regression

**identify and analyze the relationship** between variables because of the presence of the other factor

### Sequential

evaluating sequential data to discover sequential patterns

### Association Rules

discover a link between two or more items, finds a **hidden pattern in the data set**

These common **constraints** (math formulas) are used to determine significant and interesting links:

- Support — indication of how frequently the itemset appears in the dataset.
- Confidence — indication of how often the rule has been found to be true
- Lift — indication of importance compared to other items
- Conviction — indication of the strength of the rule from statistical independence

### Outer Detection

**outliers**

observation of data items in the data set, which do not match an expected pattern or expected beha

### Prediction

used a combination of other data mining techniques such as trends, clustering, classification to pred

# What is Data Wrangling?

### What is Data Wrangling?

The **process of transforming and mapping data from one "raw" data form into another**

intent of making it more appropriate and valuable for a variety of downstream purposes

Also known as data munging

There **5 core steps** behind data wrangling.

**1] Discovery**

understand what your data is about and keep in mind domain specific details about your data as you move through the other steps

**3] Cleaning**

remove outliers, change null values, remove duplicates, remove special characters, standardize formatting

**5] Validating**

authenticate the reliability, quality, and safety of the data

**2] Structuring**

you need to organize your content into that will be easier to work for our end

**4] Enriching**

appending or enhancing collected data context obtained from additional source

**6] Publishing**

place your data in a datastore so you can

# What is Data Modeling?

**What is a Data Model?**

an abstract model that **organizes elements of data and standardizes how they relate to one another** and to the properties of real-world entities eg. A data model could be a **relational database** that contains many tables.

A data model could be:

**Conceptual**

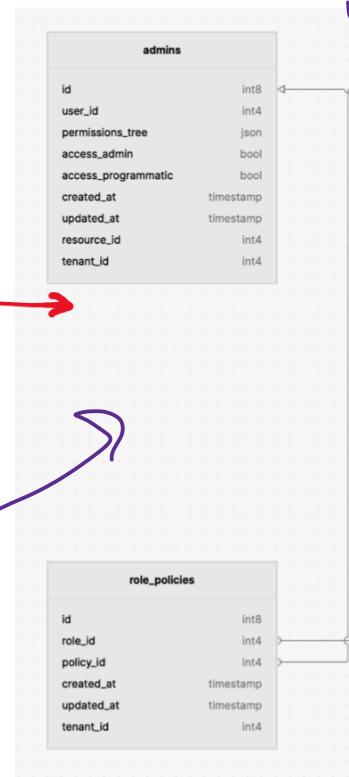
- How data is represented at the organization level abstractly without concretely defining how it works in software
- eg. People, Orders, Projects, Relationships

**Logical**

- How data is presented in software
- eg. tables and columns, object oriented classes

**Physical**

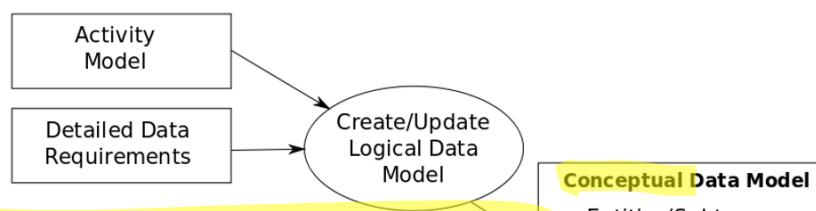
- How data is physically stored:
- eg. partitions, CPUs, tablespaces

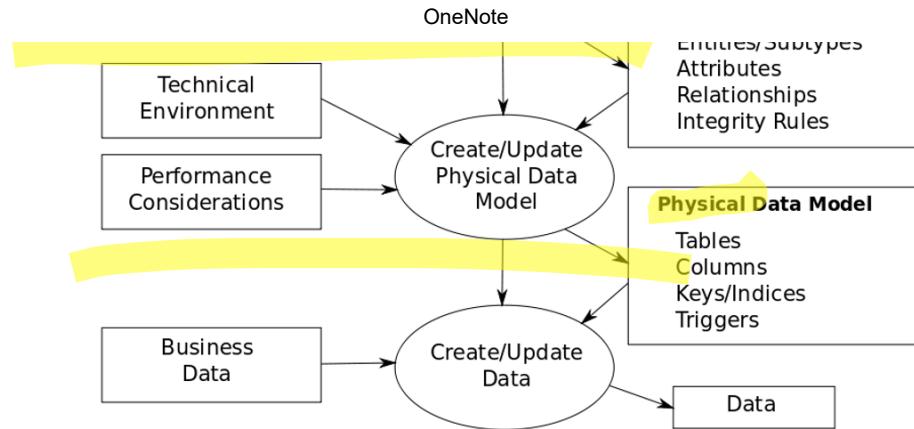


# What is Data Modeling?

**What is data modeling?**

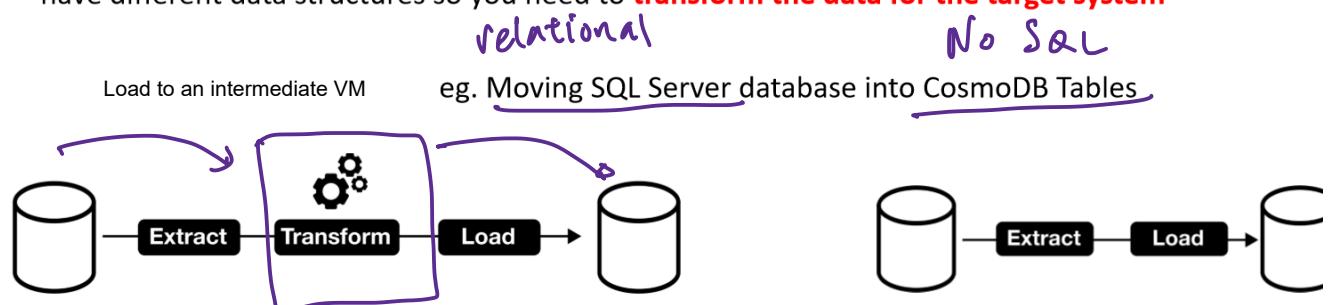
A process used to **define and analyze data requirements needed to support the business processes** scope of corresponding information systems in organizations





## ETL vs ELT

ETL and ELT is used when you want to **move data from one location to another**, where the data to have different data structures so you need to **transform the data for the target system**



## Extract, Transform and Load (ETL)

- loads data first into a staging server and then into the target system
  - used for on-premises, relational and structured data
  - used for a small amount of data
  - doesn't provide data lake support
  - easy to implement
  - Mostly supports relational data

## Extract, Load and Transform (ELT)

- loads data directly into the target system
  - used for scalable cloud structured and unstructured data
  - used for large amounts of data
  - provides data lake support.
  - requires specialized skills to implement a data pipeline
  - Support for unstructured data readily available

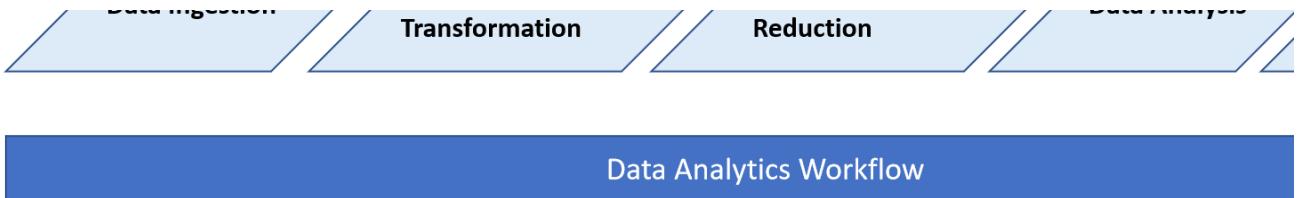
# Data Analytics

# What is Data Analytics?

Data analytics is concerned with examining, transforming, and arranging data so that you can extract and study useful information.

A data analyst commonly uses SQL, Business Intelligence (BI) tools and Spreadsheets





## Data Analytics Workflow

# Key Performance Indicators (KPIs)

**Key Performance Indicators (KPIs)** are type of performance measurement that a company or organization to determine **performance over time**



KPIs can **evaluate the success** of an organization or of a specific organization a

There are two categories of measurements for KPIs

### 1. Quantitative

- properties can be measured with a numerical result
- facts presented with a specific value
  - eg. monthly revenue, number of signups, number of reported defects

### 2. Qualitative

- properties that are observed and can generally not be measured with a numerical value
- numeric or textual value that represent of personal feelings, tastes, or opinions
  - eg. Customer sentiment

# Data Analytic Techniques

### Descriptive Analytics — **What happened?**

- Specialized metrics
  - Key Performance Indicators (KPI)
  - Return on Investment (ROI)
- Generating **sales and financial reports**
- Accurate, comprehensive , live-data and effective visualizations

### Diagnostic Analytics — **Why did it happen?**

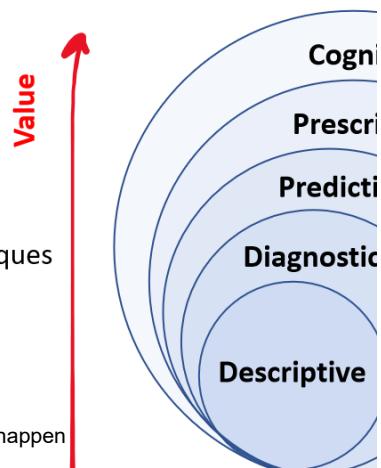
- Supplemental to descriptive analytics
- Drill down, or investigate descriptive metrics to determine root cause
- Find and isolate anomalies into its own datasets and apply statistical techniques

### Predictive Analytics — **What will happen?**

- Use historical data to predict trends or reoccurrence
- Statistical and Machine Learning techniques applied
  - eg. Neural Networks, Decision Trees, Regression, Classification

### Predictive Analytics — **What will happen?** Prescriptive analytics - how can we make it happen

- Goes a step further than predictive and uses Machine Learning by injecting



- goes a step further than predictive and uses machine learning by ingesting hybrid data to predict future scenarios that are exploitable

### Cognitive Analytics — What-if this happens?

- Using analytics to draw patterns to create what-if scenarios and what actions can be taken if those scenario become reality
- Uses ML, Natural Language Processing (NLP)
  - Eg call center conversation logs and product reviews

## Microsoft OneDrive



Microsoft OneDrive is **a storage and storage synchronization service** for files which are stored in the cloud. Similar products: DropBox, GoogleDrive, Box and \*Microsoft SharePoint

OneDrive is intended for personal storage for a single individual

You pay for different sizes of storage (5GB **free**, 100GB, 1TB, 6TB)

You do not worry about the underlying hardware's the, durability, resilience, fault tolerance, or availability.

Files can be shared easily to other users via:

- a shareable link
- or to a specific email that also has a OneDrive Account

Files are accessed via:

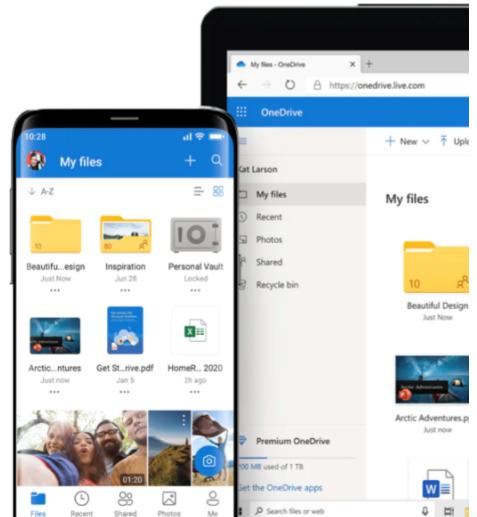
- A web-application (web interface)
- Via shared folders that hold a reference to the files stored in the cloud

Files can be synchronized

- a copy resides on a local computer hard drive is copied to the cloud
- A file residing in the cloud can be copied to a local computer hard drive
- Copying occurs automatically when files are changed
  - difference in files could result in conflicts a user must choose to which file to keep

Files can be versioned (you can recover older versions of a files)

- Older files may retain for 30 days and be automatically deleted



## Microsoft 365 – SharePoint



Microsoft 365 SharePoint is **a web-based collaborative platform** that integrates with Microsoft Office. Intended **for document management and shared storage**.

### SharePoint Sites

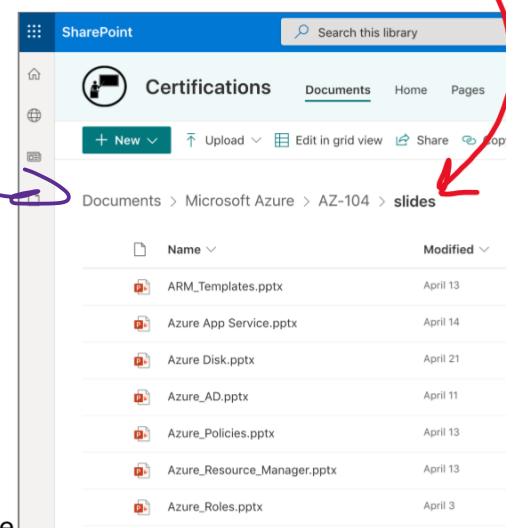
Data within SharePoint organized around Sites

A site a collaborative space for teams with the following components:

- **Document Library**
- Pages
- Web Parts
- And more....

### SharePoint Document Library

A document library is file storage and synchronization but **designed for teams**. It is very similar to OneDrive, but files are owned by the company and not an individual.



You can apply robust permissions to access files within or outside

your organization

A Site always has default Document Library called "documents"

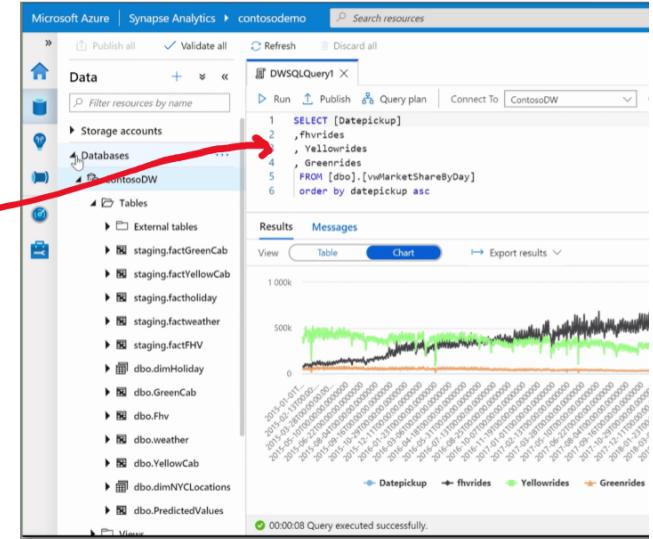


# Azure Synapse Analytics

Azure Synapse Analytics is a **data warehouse** and **unified analytics platform**

Build ETL/ELT processes:

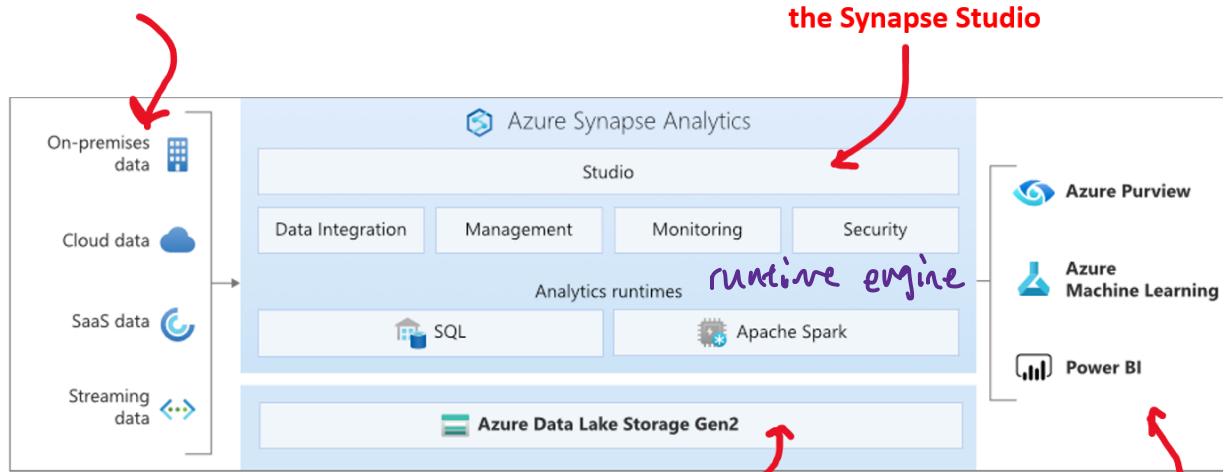
- in a code-free visual environment
- easily ingest data from more than 95 native connector
- Deeply integrated Apache Spark
- use **T-SQL** queries on both your data warehouse and Spark engines
- Supports multiple languages: T-SQL, Python, Scala, Spark SQL, and .Net
- Integrated with Artificial Intelligence (AI) and Business Intelligence tools (BI)
  - Azure Machine Learning
  - Azure Cognito Services
  - Microsoft Power BI



# Azure Synapse Analytics

You can ingest data from many data sources.

Azure Synapse Analytics Is interfaced via the Synapse Studio



The data is stored in Object Storage  
Via Data Lake Storage Gen 2

You can output to  
Azure Service



# Synapse SQL

Synapse SQL is a distributed version of T-SQL designed for data warehouse workloads.

- extends T-SQL to address streaming and machine learning scenarios
- use built-in streaming capabilities to land data from cloud data sources into SQL tables
- Integrate AI with SQL by using ML models to score data using the T-SQL PREDICT function
- and offers both **serverless** and **dedicated** resource models

For **unpredictable** workloads (unplanned or bursty) use the **always-available, serverless SQL endpoint**.

For **predictable** workloads

- create **dedicated SQL pools** to reserve processing power for data stored

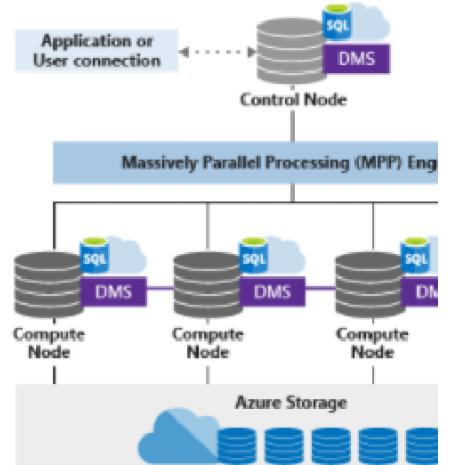


# Dedicated SQL Pool

Dedicated SQL pool is a query service over the data in your **data warehouse**

The unit of scale is an abstraction of compute power that is known as a **data warehouse unit (DWU)**.

Once your dedicated SQL pool is created, you can import big data with simple PolyBase T-SQL queries, and then use the power of the distributed query engine to run high-performance analytics



# Serverless SQL Pool



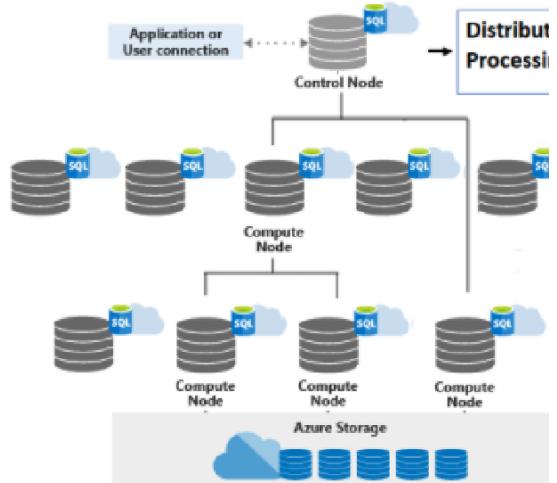
## Serverless SQL Pool

Serverless SQL pool is a query service over the data in your **data lake**

Scaling done automatically to accommodate query resource requirements.

As topology changes over time by adding, removing nodes or failovers

it adapts to changes and makes sure your query has enough resources and finishes successfully.



# Apache Spark for Synapse

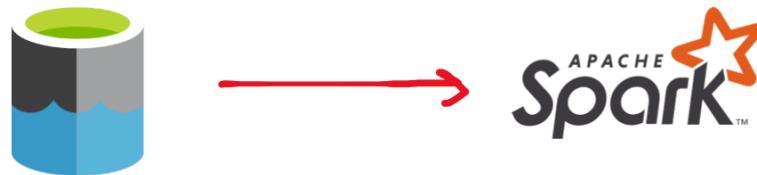
Azure Synapse can deeply and seamlessly integrates with Apache Spark



- ML models with SparkML algorithms and AzureML integration for Apache Spark
    - with built-in support for Linux Foundation Delta Lake.
  - Simplified resource model that frees you from having to worry about managing resources.
  - Fast Spark start-up and aggressive autoscaling.
  - Built-in support for .NET for Spark allowing you to reuse your C# expertise and code within a Spark application.

# Apache Spark with Data Lake

Azure Synapse removes the traditional technology barriers between using SQL and Spark together. You can seamlessly mix and match based on your needs and expertise.



- Tables defined on files in the data lake are seamlessly consumed by either Spark or SQL.
- SQL and Spark can directly explore and analyze Parquet, CSV, TSV, and JSON files in the data lake.
- Fast, scalable data loading between SQL and Spark databases

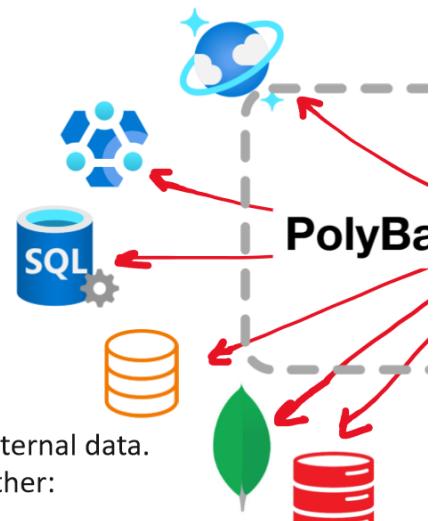
## PolyBase

PolyBase is a **data virtualization feature** for **SQL Server**

PolyBase enables your SQL Server instance to query data with T-SQL directly from:

- SQL Server
- Oracle
- Teradata
- MongoDB
- Hadoop clusters
- Cosmos DB

without separately installing client connection software.



PolyBase allows you to join data from a SQL Server instance with external data.

**Prior to PolyBase** to join data to external data sources you could either:

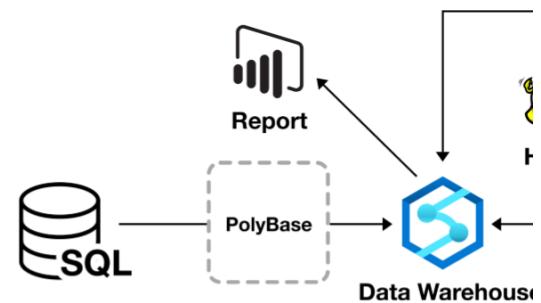
- Transfer half your data so that all the data was in one location.
- Query both sources of data, then write custom query logic to join and integrate the data at the client level.

## Azure Synapse Analytics – ELT

You can **perform ELT using Synapse SQL** in Azure Synapse Analytics.

The fastest and most scalable way to load data is through PolyBase external tables and the COPY statement

With PolyBase and the **COPY statement**, you can access external data stored in Azure Blob storage or Azure Data Lake Store via the T-SQL language.



The basic steps for implementing ELT are:

1. Extract the source data into text files.
2. Load the data into Azure Blob storage or Azure Data Lake Store.
3. Prepare the data for loading.
4. Load the data into staging tables with PolyBase or the COPY command.
5. Transform the data.
6. Insert the data into production tables.

## SQL Server Management Studio (SSMS)

SQL Server Management Studio (SSMS) is an IDE for managing any SQL infrastructure.

Access, configure, manage, administer, and develop all components of

- SQL Server
- Azure SQL Database
- Azure Synapse Analytics

### Object Explorer

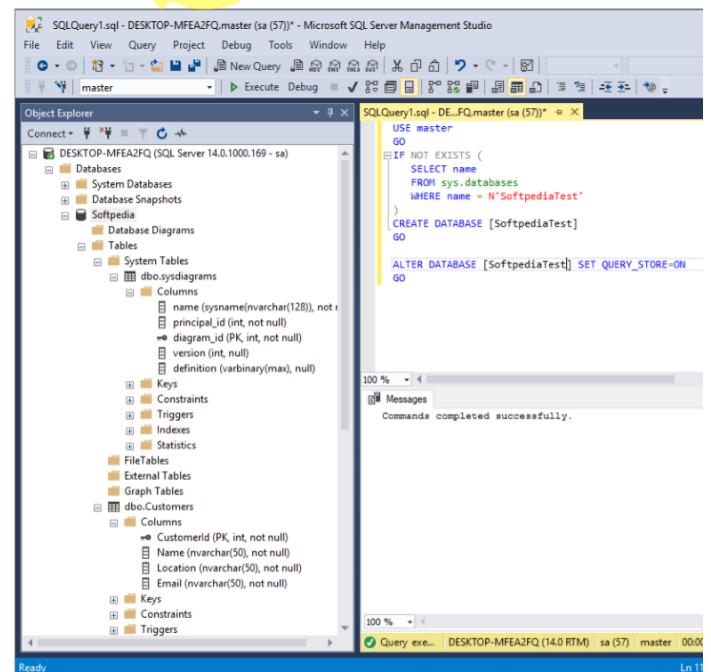
view and manage all of the objects in one or more instances of SQL Server

### Template Explorer

build and manage files of boilerplate text that you use to speed the development of queries and scripts

### Solution Explore (deprecated)

build projects used to manage administration items such as scripts and queries



## SQL Server Data Tools (SSDT)

VSC tool

SQL Server Data Tools (SSDT) transforms database development by introducing a ubiquitous, declarative language that spans all the phases of database development inside Visual Studio

T-SQL

use SSDT Transact-SQL to build, debug, maintain, and

refactor databases.

SSDT also provides a visual Table Designer for creating and editing tables in either database projects or connected database instances  
Be able to view control data related files  
Easy to publish to SQL Database or SQL Server

**SQL Server Object Explorer** in Visual Studio offers a view of your database objects similar to SQL Server Management Studio (SSMS)

- allows you to do light-duty database administration and design work
- easily create, edit, rename and delete tables, stored procedures, types, and function
- edit table data, compare schemas, or execute queries by using contextual menus right



## Azure Data Studio

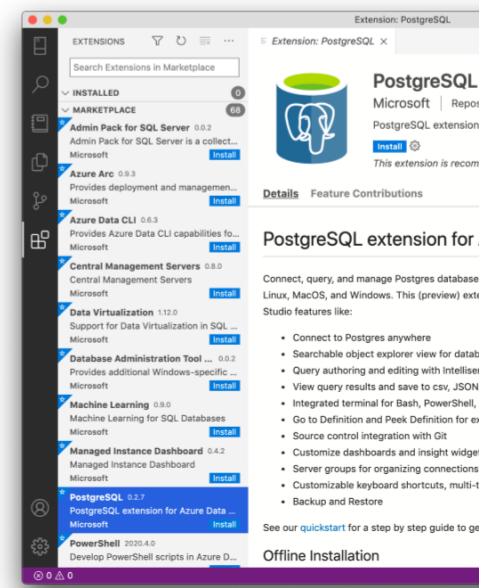


Azure Data Studio is a **cross-platform database tool** for data professionals using **on-premises and cloud data platforms** for Windows, macOS, and Linux.

Query, design, and manage your databases and data warehouses

Azure Data Studio offers:

- a modern editor experience with IntelliSense
  - **Very similar experience to Visual Studio Code.**
- code snippets
- source control integration
- integrated terminal
- built-in charting of query result sets
- customizable dashboards
- Jupyter Notebooks connected to your datasets
- A marketplace of free extensions
  - SQL Database Inspector (inspect data with just a few clicks)
  - Kusto (KQL) extension for Azure Data Studio
  - PostgreSQL extension for Azure Data Studio
  - And many more!



## Business Intelligence (BI)

Business intelligence (BI) is both a data-analysis strategy and **technology** for business information.

The most popular BI tools are:

- Tableau
- Microsoft Power BI
- Amazon QuickSight

Business intelligence (BI) helps organizations make data-driven decisions by (BI) combining:

- business analytics
- data mining
- data visualization
- data tools
- infrastructure

- best practices



# Microsoft Power BI

Power BI is a Business Intelligence tool for **visualization business data**

The screenshot shows the Microsoft Power BI Desktop interface. The ribbon menu includes File, Home, Insert, Modeling, View, and Help. The Home tab is selected. The main workspace displays several data visualizations: a treemap labeled "Sales Amount by Brand Name", a bar chart labeled "Units by Country and Sales Size", a line chart labeled "Units Sold by Year, Quarter and Manufacturer", and a scatter plot labeled "Sales Amount by Year, Month and Brand Name". The left sidebar shows "Key Influences" and "Top segments" with various filters applied. The right sidebar contains sections for "Visualizations", "Filters", "Drill through", "Cross-report", and "Keep all filters". A red arrow points from the text "A way to design an" to the top right corner of the Power BI Desktop window.

**Power BI Desktop**  
A way to design an

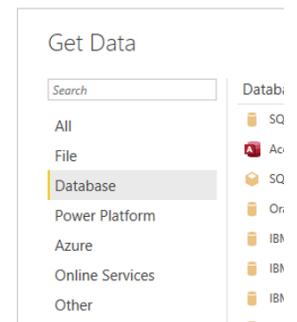
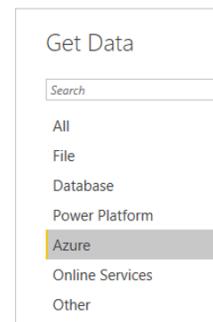
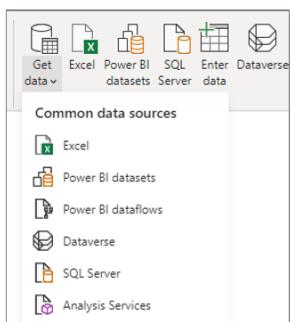
**Power BI Mobile**  
View reports on th

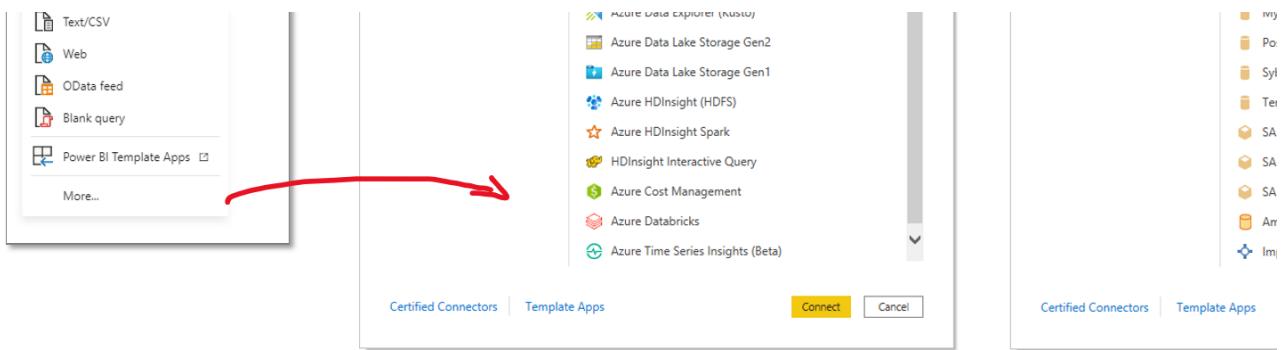
**Power BI Service**  
Access and modify  
cloud

**Power BI embedde**  
A way to embed Po  
components into y

# Microsoft Power BI Desktop - Data Sou

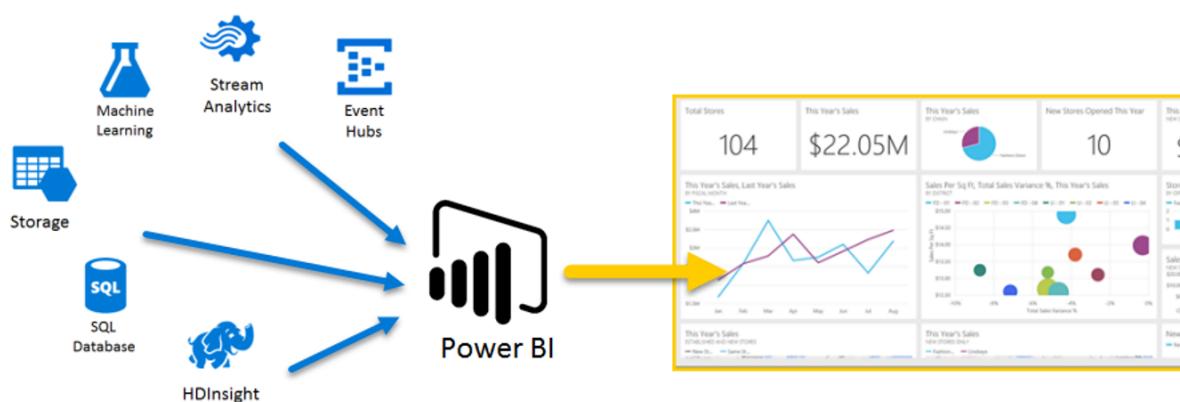
Microsoft Power BI can ingest data from **many data sources**.





## Azure and Power BI

Power BI can directly integrate with Azure Services:



## Power BI Desktop vs Power BI Service

### Power BI Desktop

downloaded as a free Windows application and installed on a local Windows computer

**Report designers** use the Desktop application to publish the Power BI reports to the Power BI Service

### Power BI Service

cloud-based service where users interact with the repository

**Users** in the Power BI Service can create reports and create visualizations based on the data model and they can share a report with co-workers

#### Power BI Desktop

- Many data sources
- Transforming

#### Both

- Reports
- Visualizations

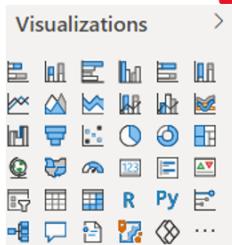
#### Power BI Service

- Some data sources
- Dashboards

- | OneNote  | Power BI   |   |
|--|--|---|
| <ul style="list-style-type: none"> <li>• Shaping and modeling</li> <li>• Measures</li> <li>• Calculated columns</li> <li>• Python</li> <li>• Themes</li> <li>• RLS creation</li> </ul> | <ul style="list-style-type: none"> <li>• Security</li> <li>• Filters</li> <li>• Bookmarks</li> <li>• Q&amp;A</li> <li>• R Visuals</li> </ul> | <ul style="list-style-type: none"> <li>• Apps and workspace</li> <li>• Sharing</li> <li>• Dataflow creation</li> <li>• Paginated reports</li> <li>• RLS management</li> <li>• Gateway connection</li> </ul> |

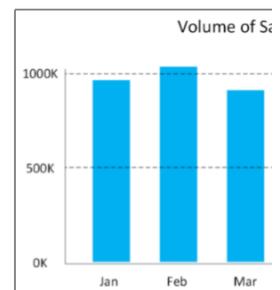
# Data Visualizations and Chart Types

Power BI has **many kinds of visualizations**. We'll cover the most common ones.



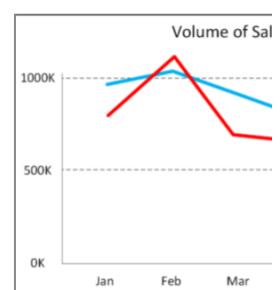
## Bar and column charts

See how a set of variables changes across different categories



## Line charts

overall shape of an entire series of values



# Data Visualizations and Chart Types

## Matrix

tabular structure that summarizes data

Quarter Year	Q1 Revenue	YTD Revenue	Q2 Revenue	YTD Revenue
2015	\$45,186	\$45,186	\$70,609	\$115,795
2016	\$52,154	\$52,154	\$73,542	\$125,696
2017	\$51,388	\$51,388	\$68,149	\$118,537
2018	\$48,281	\$48,281	\$66,853	\$115,134
2019	\$53,145	\$53,145	\$49,135	\$102,280

## Key influencers

the major contributors to a selected result or value.



## Treemap

charts of colored rectangles, with size representing the relative value of each item



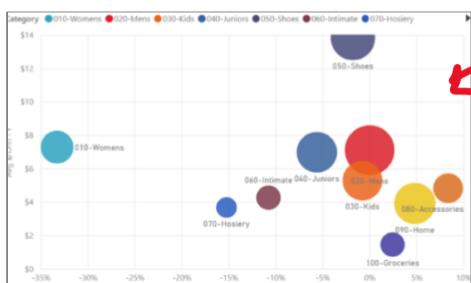
REPRESENTING THE RELATIVE VALUE OF EACH ITEM



# Data Visualizations and Chart Types

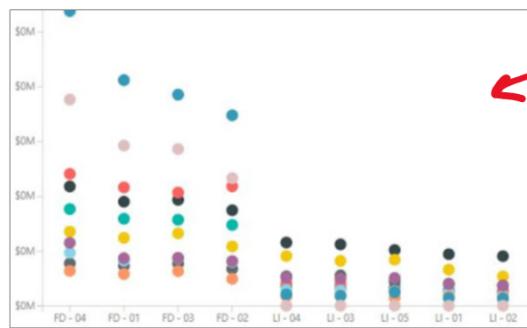
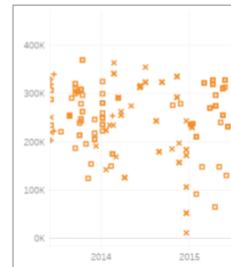
## Scatter

relationship between two numerical values (X and Y Axis)  
A bunch of dots on a graph



## Bubble Chart

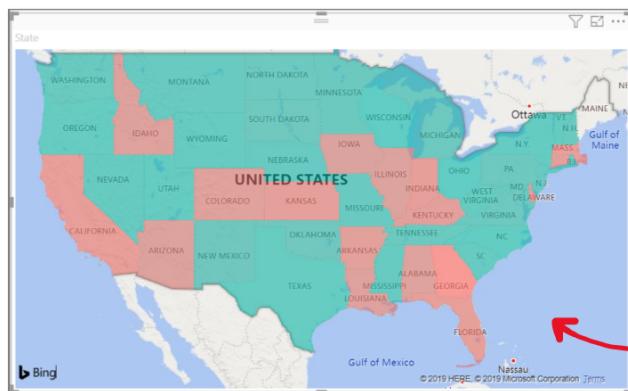
a scatter chart that replaces data points with bubbles, larger bubbles representing a third dimension



## Dot plot Chart

similar to a bubble chart and scatter chart, but can plot categorical data along the X-Axis.

# Data Visualizations and Chart Types



## Filled map

A geographic map where different areas can be filled in eg. States different colors or ranges of colors

# Power BI Embedded



Azure Power BI Embedded is a platform-as-a-service (PaaS) analytics service that allows you to quickly embed **visuals, reports and dashboards** into your application.

**For independent software vendors (ISV)**  
enables you to visualize application data, rather than building that service yourself.

**For developers**  
embed reports and dashboards into their application for their customers.

## To use Azure Power BI embedded

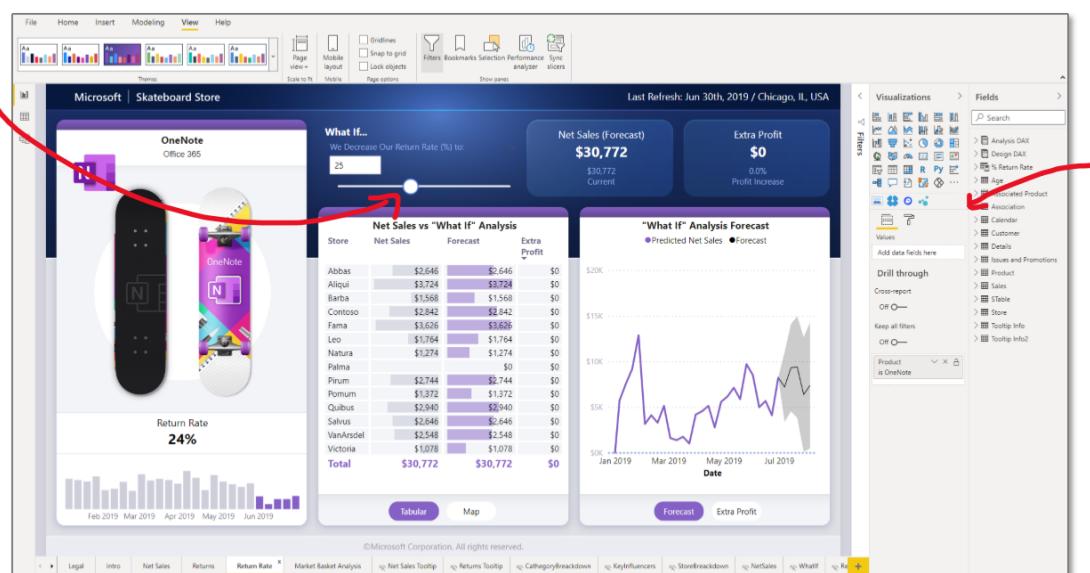
- You need a Power BI Pro user
- You need to create an App workspace
- You need to choose a capacity
  - Billing works via a capacity-based, hourly metered mode

# Power BI – Interactive Reports

One way you can make reports interactive is by **having knobs and controls** directly in the report

Reports can be highly **stylized**

Power BI allows you to generate **reports which are interactive**



A report can contain many **pages**

# Power BI – Interactive Reports

You can view the **underlying data** for an Interactive Report

(you can't do this with Dashboards)

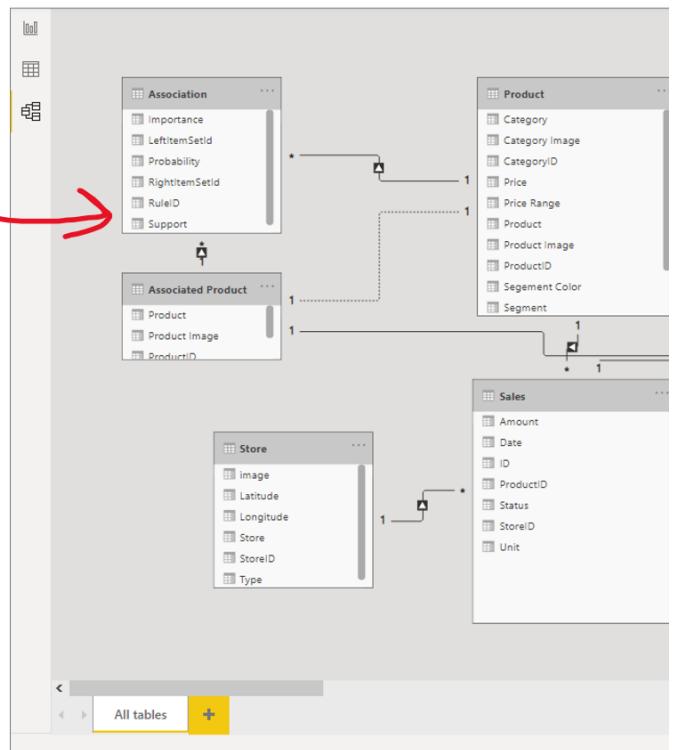
*PowerBI desktop*

Product	ProductID	Category	CategoryID	Segment	SegmentID	Produ
Access	1	Office 365	1	Red	7	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Excel	2	Office 365	1	Jade	4	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Exchange	3	Office 365	1	Cyan	1	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
OneNote	4	Office 365	1	Magenta	8	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Outlook	5	Office 365	1	Cyan	1	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
PowerPoint	6	Office 365	1	Orange	6	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Publisher	7	Office 365	1	Turquoise	3	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
SharePoint	8	Office 365	1	Turquoise	3	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Skype	9	Office 365	1	Cyan	1	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Visio	10	Office 365	1	Royal Blue	2	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Word	11	Office 365	1	Royal Blue	2	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
XBOX	12	XBOX	3	Green	4	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
OneDrive	13	Office 365	1	Blue	2	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Yammer	14	Office 365	1	Blue	2	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
XBOX ONE	15	XBOX	3	Green	4	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Power BI	16	Power Platform	3	Yellow	5	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Kaizala	17	Office 365	1	Cyan	1	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Planner	18	Office 365	1	Jade	4	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Forms	19	Office 365	1	Turquoise	3	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
PowerApps	20	Power Platform	2	Magenta	8	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Teams	21	Office 365	1	Purple	9	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Stream	22	Office 365	1	Red	7	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
To-Do	23	Office 365	1	Royal Blue	2	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>
Flow	24	Power Platform	2	Neon Blue	2	<a href="https://imagizer.imagesh">https://imagizer.imagesh</a>

You can view **and modify the underlying data models** for a report

(you can't do this with Dashboards)

*PowerBI desktop*



## Power BI Service

Power BI is **cloud-based service** where users view and interact with the rep and where they can create **Dashboards**

The screenshot shows the Power BI Home page. On the left is a sidebar with links like Home, Favorites, Recent, Create, Datasets, Goals, Apps, Shared with me, Discover, Learn, Workspaces, and My workspace. The main area displays "Good afternoon, Andrew" and a message to "Select a tile to find and share data-driven insights". Below this are "Data stories from the Power BI community" featuring three tiles: "THE DEFINITIVE 100 MOST USEFUL PRODUCTIVITY TIPS" by Alice\_Drummond, "Cancer statistics in the USA" by immatey, and "Ranking Sports by Degree of Difficulty for Key Skills" by Patrick Baumgartner. At the bottom, there's a "Getting started with Power BI" section with links to Power BI basics, Sample reports, and How to create reports.

You access Pow by visiting [app](#)

## Power BI – Dashboard Tiles

A **tile** is a **snapshot of your data**, pinned to the dashboard.



A tile can be created from a:

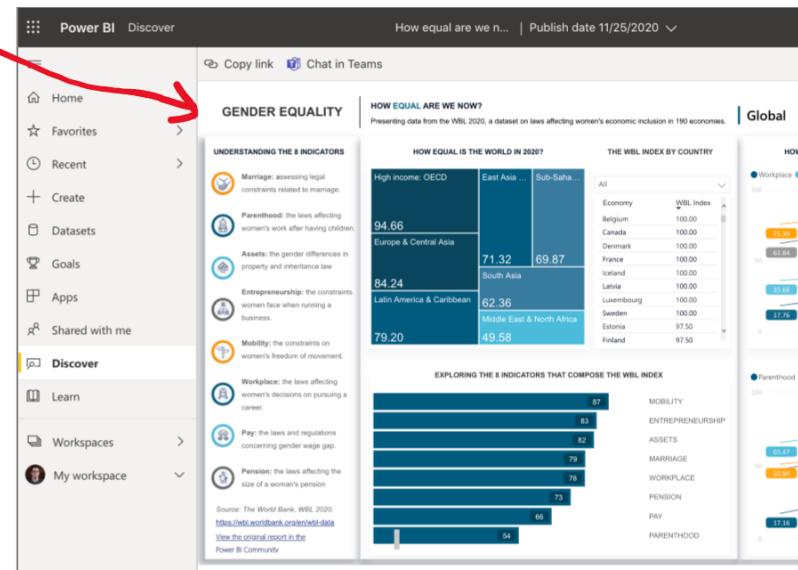
- **Report**
- **Dataset**
- Dashboard
- Q&A box
- Excel
- SQL Server Reporting Services (SSRS)
- and more

## Power BI Service – Dashboard

**Power BI dashboard** is a **single page**, often called a canvas, that tells a story through visualizations

The visualizations you see on the dashboard are called **tiles**.

You **pin** tiles to a dashboard from reports



## Power BI – Reports vs Dashboards

Capability	Dashboards	Reports
Pages	One page	One or more pages
Data sources	One or more reports and one or more datasets per dashboard	A single dataset per report
Filtering	Can't filter or slice	Many different ways to slice
Set alerts	Can create alerts to email you when the dashboard meets certain conditions	No
Feature	Can set one dashboard as your featured dashboard	Can't create a featured report
Can see underlying dataset tables and fields	No. Can export data but can't see the dataset tables and fields in the dashboard itself	Yes. Can see dataset table values that you have permission to view
Customization	No	Can filter, export, view reports, add bookmarks, generate reports, analyze in Excel, and more

## Paginated Reporting (RDL)

**Paginated Reports** are reports designed to fit into page format so they can be **printed** or **shared**. The data display of all data are tables which can span multiple pages.

### Report Definition Language (RDL)

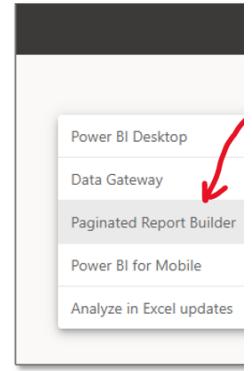
## An XML representation of a SQL Server Reporting Services report definition

A report definition contains data retrieval and layout information for a report.

Paginated Reports are just a visualization of an .rdl file.

# Power BI Report Builder

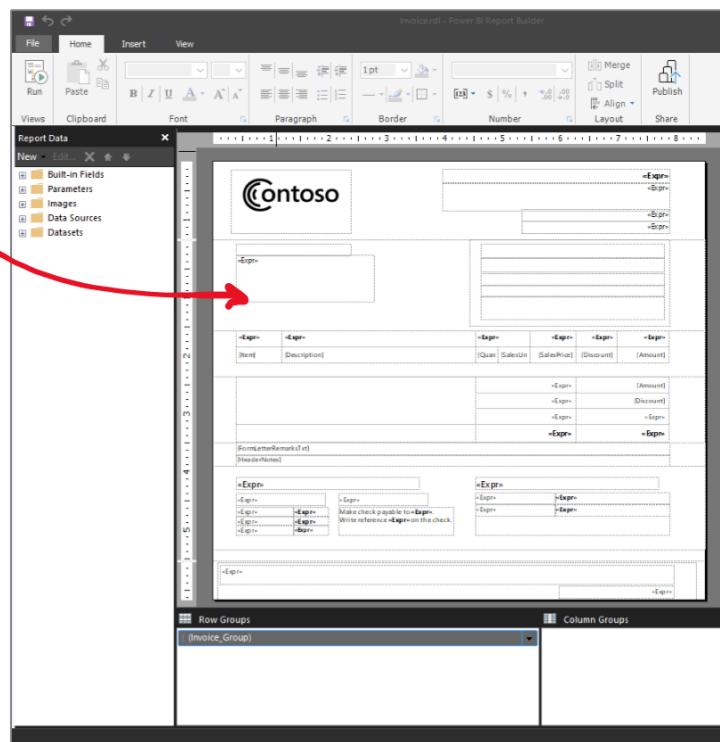
To design pixel-perfect paginated reports, you use Power BI Report Builder. It is a tool specifically for the creation of Paginated Reports



You can download from Power BI's

# Paginated Reporting (RDL)

This is what a **Paginated Report** looks like in Power BI Report Builder



# Structured Query Language

**Structured Query Language (SQL)** designed to access and maintain data for a relational database management system (RDBMS)

We use SQL to get to **insert, update, delete, view** data from our database's tables.

SQL can join many tables and include many functions

SQL can join many tables and include many functions  
to transform the final outputted result

The SQL Syntax was standardized (ISO 9075).

Relational databases will mostly adhere to this standard while adding in their own database specific features.

SQL is a highly transferable skill and we see SQL being used in Non-Relational databases to provide a popular and familiar querying tool.

```

AND questions.published = 1
) as published_questions_count
SELECT
domains.id,
domains.per,
domains.name,
(
  SELECT count(true)
  FROM question_tags
  INNER JOIN questions ON q
  WHERE
    question_tags.tag_id = ?
    AND questions.exam_id = ?
    AND questions.domain_id = ?
    AND questions.published = ?
) as published_questions_count
FROM domains
WHERE
  domains.exam_id = exam_sets.id
  AND domains.domain_id IS NOT NULL
ORDER BY domains.name ASC
{/array} as domains
FROM exam_sets
INNER JOIN tags ON tags.id = exam_sets.tag_id
WHERE
  exam_sets.exam_id = {exam_id}
ORDER BY
  exam_sets.position ASC

```

## SQL Syntax

### Clauses

- WHERE Clause
- ORDER BY clause
- HAVING Clause
- UPDATE Clause
- GROUP BY Clause

< MAYEB JUST SKIPT THIS S

### Expressions

Scalar or tabular data that is used along side a Predicate to provide conditional data to be returned

### Predicates

The condition of how to filter the data

### Queries

### Statements

## OLAP vs OLTP

### Online Transaction Processing (OLTP)

#### Database

A database is built to store current transactions and enable fast access to specific transactions for ongoing business processes e.g. SQL Servers

### Online Analytical Processing (OLAP)

#### Data Warehouse

A data warehouse is built to store large amounts of historical data and enable fast, complex queries across all the data



- Single Data Source
- **Short transactions** (small and simple queries)
  - with an emphasis on writes.
- Many transactions
- Latency sensitive
- Small payloads

Use Case: (General Purpose)

Adding Items to your shopping cart

- Multiple Data Sources
- **Long transactions** (long and complex)
  - with an emphasis on reads.
- Few transactions
- Throughput sensitive
- Large payloads

Use case: (Analytics)

Generating Reports

*small transactions*

## Open-Source Relational Databases



Created by **MySQL AB** then acquired by **Sun Microsystems** and then acquired by Oracle. MySQL is a **pure relational database (RDBMS)**.

It is a simpler database which makes it **easy to setup, use and maintain**.

Has multiple storage engines: **InnoDB** and **MyISAM**.

The **most popular** relational database.



MariaDB is a fork of MySQL by the original creators of **MySQL AB**.

After Oracle acquired MySQL there was concern that Oracle may change the open-source licensing or stop future MySQL being free to use.



PostgreSQL evolved from the Ingres project at the University of California.

Postgres is an **object-relational database (ORDBMS)**.

Just has a single storage engine.

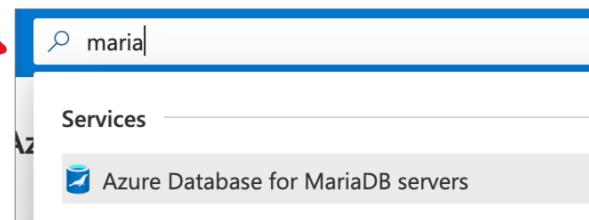
*Ingres engine*

The **most advanced** relational database.

- Full text search, table inheritance, triggers, rows and data types required.

## Azure Database for Maria DB, PostgreSQL and PostgreSQL

If you want to **deploy** an open-source database on Azure you need to search by its name.



# Azure DB Read Replicas

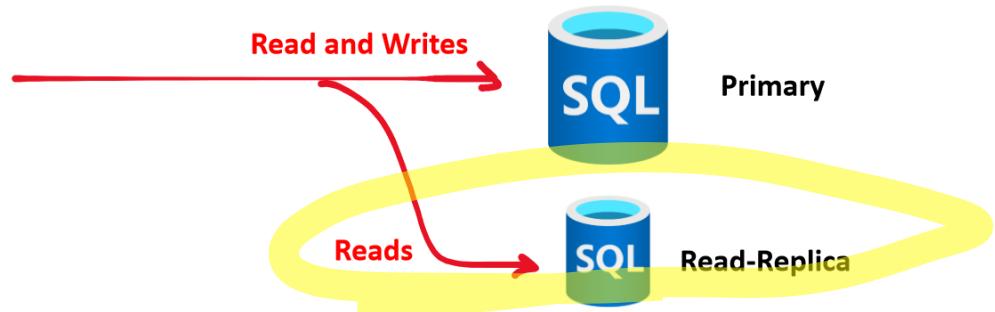
## What is a read replica?

A read replica is a copy of your database that is kept synced with your primary database. This additional database is used to improve read contention by offloading reads to a database dedicated to perform read operations.

Read Replicas are available for:

- Azure SQL Database
- Azure SQL Managed Instance

You can have multiple read-replicas for a dat



In simple use cases a Read Replica can act as an OLAP for a relational databases that is very si

# Citus on Azure



Citus is an **open-source Postgres extension** that transforms Postgres into a **distributed database**.

Distributed dat

Citus extends postgres to provide better support for:

- Database sharding (easy horizontal scaling)
- Realtime queries (great for real-time analytics dashboards)
- Multi-tenancy (great for SaaS company's)
- Time series workloads



**Hyperscale (Citus) server group**  
Best for ultra-high performance beyond 100GB.

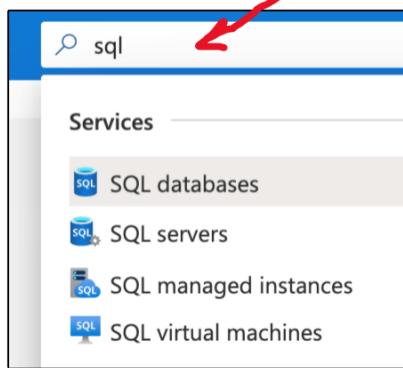
Azure Postgres for **HyperScale deploy option** is just a managed postgres cluster us the Citus extension

Ideal for multi-tenant application analytical workloads that need sub-second response. Supports both transactional workloads and hybrid transactional/analytical workloads.

[Create](#)[Learn more](#)

## Azure SQL Family

Azure has **multiple solutions** for relational databases



### SQL Server on Azure Virtual Machines

- When you need OS-level control and access
- When you need to **lift-and-shift** your workload
- When you have existing SQL licenses and you want to take advantage of the Azure Hybrid Benefit



### SQL Managed Instance

- When you have an **existing database** that you want to move to the cloud while maintaining broadest SQL server engine compatibility
- highly available, disaster recovery, automated backups
- Ideal for most migrations to the cloud



### Azure SQL Database

- Fully managed **SQL databases**
- designed to be fault-tolerant
- built-in disaster recovery
- highly available
- designed to scale



### SQL Servers

The underlying servers for Azure SQL Database

## Azure Elastic Pools



Azure SQL Database elastic pools are a simple, cost-effective solution for **managing and scaling multiple databases** that have **varying and unpredictable usage demands**

a f

Databases in an elastic pool are on:

- a single server
- share a set number of resources at a set price

Elastic pools in Azure SQL Database enable SaaS developers to optimize the price performance for a group of databases within a prescribed budget while delivering performance elasticity for each database.

# Connectivity Architecture – SQL Database and Azure Synapse

When a connection from a server to an Azure SQL database the client will connect to a **gateway** that listens on **port 1443**

Based on the **connection policy** the gateway will grant traffic and route access to the appropriate database

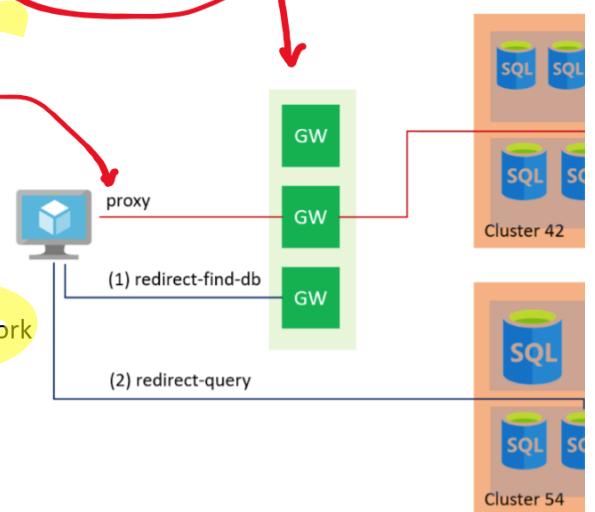
There are three kinds of **connection policies**:

## Proxy

connections are proxied through a gateway increased latency and reduced throughput intended for workloads connecting from **outside the Azure Network**

## Redirect (Recommended)

establishes a direct connection reduced latency and improved throughput intended for workloads connecting **inside the Azure Network**



## Default

When you create a database it will default to either Proxy or Redirect Based on workloads **inside or outside the Azure networking**

# MS SQL Database Authentication

During setup of your MS SQL databases must **select an authentication mode** .

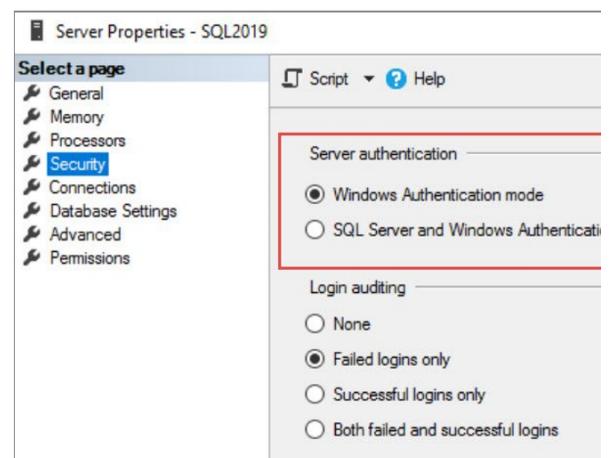
- **Windows Authentication mode**
  - enables Windows Authentication and disables SQL Server Authentication
- **Mixed mode**
  - enables both Windows Authentication and SQL Server Authentication

## Windows Authentication (recommended)

- Specific Windows user and group accounts are trusted to log in to SQL Server
- Very secure, and very easy to modify or revoke access

## SQL Server Authentication

- A username and password is set and stored on the primary database
- cannot use Kerberos security protocol
- **login password must be passed over the network at the time of the connection** (additional attack points)
- easier to connect to database from outside a domain or from a web-based interface



# Network Connectivity

For SQL Database you can choose network connectivity to have a **Public** or **Private** Endpoint

**Public Endpoints** are reachable outside the Azure Network over the internet

- You would use **Firewall rules** to ~~protect~~ ~~protect~~ the database

**Private Endpoints** are only reachable within the Azure Network (or a connection originating from

- You would use **Azure Private Link** to keep traffic within the Azure Network.

Connectivity method \* ⓘ

No access  
 Public endpoint  
 Private endpoint

Firewall rules

Setting 'Allow Azure services and resources to access this server' to Yes allows the Azure boundary, that may or may not be part of your subscription. [Learn more](#)

Allow Azure services and resources to access this server \*

No  Yes

Add current client IP address \*

No  Yes

Connectivity method \* ⓘ

No access  
 Public endpoint  
 Private endpoint

Private endpoints

Private endpoint connections are associated with a private IP address. The private endpoint connections for this server. Note that private endpoints and they provide access to all databases in the server. [Learn more](#)

+ Add private endpoint

Name \_\_\_\_\_ Subsc \_\_\_\_\_

Click on add to create private endpoint

# Azure Defender for SQL



Azure Defender for SQL is **a unified package for advanced SQL security** including **Vulnerability Assessment** and **Advanced Threat Protection**



Azure Defender is available for:

- Azure SQL Database
- Azure SQL Managed Instance
- Azure Synapse Analytics

What it does:

- discovering and classifying sensitive data
- surfacing and mitigating potential database vulnerabilities
- detecting anomalous activities

You can turn it on at anytime and you pay a monthly cost

## Azure Defender for SQL

Protect your data using Azure Defender for SQL, a unified security package including vulnerability assessment and advanced threat protection for your server. [Learn more](#)

Get started with a 30 day free trial period, and then 15 USD/server/month.

Enable Azure Defender for SQL \* ⓘ

- Start free trial  
 Not now

# Azure Database Firewalls Rules

## Diff from Azure firewall

Azure databases are protected by **server firewalls**

A server firewall is an internal firewall that resides on the database server

All connections are **rejected by default** to database

You can set server firewall rules within the Azure Portal

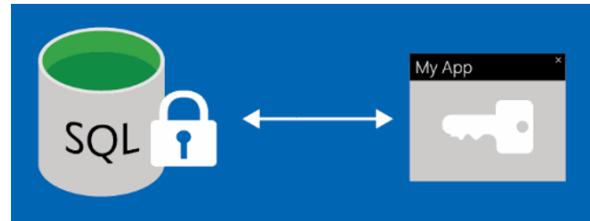
Rule name	Start IP	End IP
AllowAzure	0.0.0.0	0.0.0.0

You can set server firewall rules via T-SQL

```
EXECUTE sp_set_database_firewall_rule N'OnlyAllowServer', '0.0.0.0'
```

## Always Encrypted

**Always Encrypted** is a feature that **encrypts columns in an Azure SQL Database**



If you had a column for credit cards you would want to Always Encrypt

Always Encrypted uses two types of keys:

- column encryption keys — used to encrypt data in an encrypted column
- column master keys — a key-protecting key that encrypts one or more column encrypt

You apply Always Encrypted using T-SQL

## Role-Based-Access-Controls (RBAC)

**Role-Based-Access-Controls (RBAC)** is when you can apply roles to users to manage the fine-grade actions for specific Azure services

#### SQL DB Contributor

- manage SQL databases, but not access to them
- can't manage their security-related policies or their parent SQL servers

#### SQL Managed Instance Contributor

- manage SQL Managed Instances and required network configuration
- can't give access to others

#### SQL Security Manager

- manage the security-related policies of SQL servers and databases
- but not access to SQL servers

#### SQL Server Contributor

- manage SQL servers and databases
- but not access to them SQL servers

## Transparent Data Encryption (TDE)

**Transparent Data Encryption (TDE)** **encrypts data-at-rest** for Microsoft Data



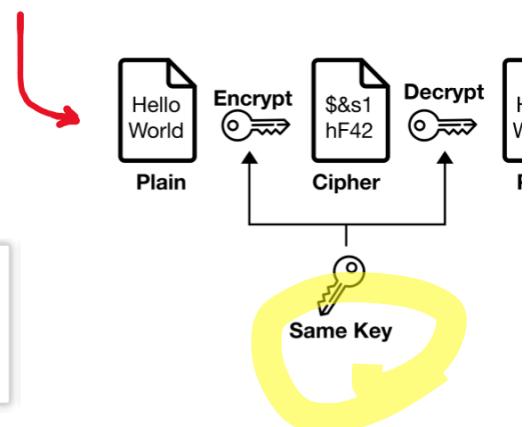
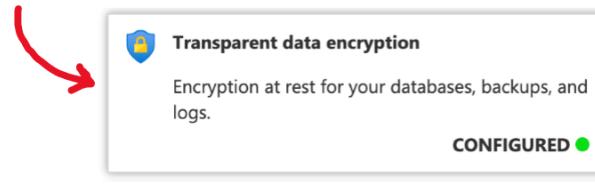
TDE can be applied to:

- SQL Server
- Azure SQL Database
- Azure Synapse Analytics

- TDE does real-time I/O encryption and decryption of data
- encryption uses a **database encryption key (DEK)**
- database boot record stores the key for availability during failover
- The DEK is a **symmetric key** (same cryptographic **keys** is used for the **encryption** of plaintext and the **decryption** of ciphertext)

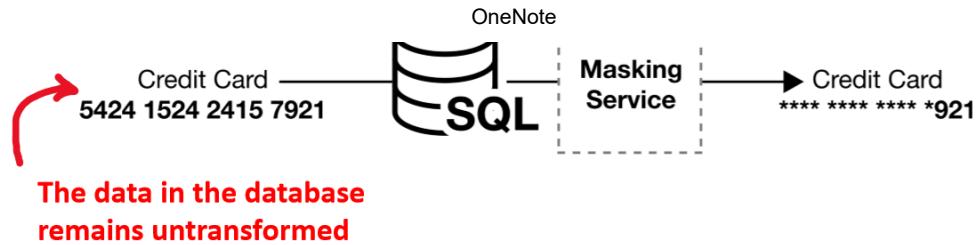
Steps to apply TDE to a database:

- Create Database Master Key
- Create a Certificate to support TDE
- Create Database Encryption Key
- Enable TDE on Database**



## Dynamic Data Masking

Data Masking is when a **request for data is transformed to mask sensitive information**



- { Dynamic Data Masking can be applied to:
- Azure SQL Database
  - Azure SQL Managed Instance
  - Azure Synapse Analytics

**Dynamic Data Masking**  
Limit sensitive data exposure by masking it to non-privileged users.

NOT CONFIGURED

You create a Dynamic Data **Masking Policy**:

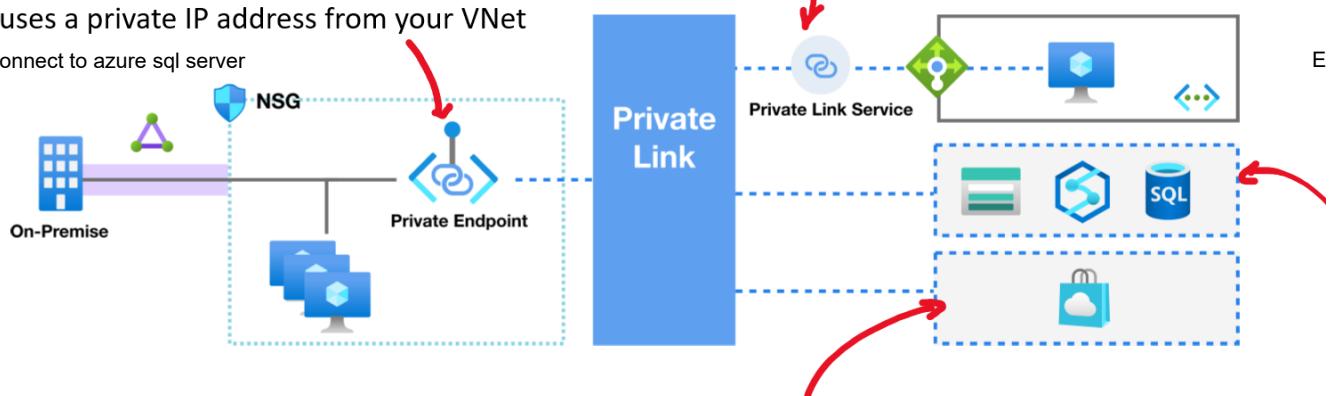
- SQL users excluded from masking — users who can get data unmasked
- Masking rules — what fields should be masked
- Masking functions — how to apply masking to fields

## Private Links

Azure Private Links allows you to establish secure connections between Azure resources so traffic **remains within the Azure Network**

**Private Link Endpoint** is an **Network Interface** that connects you privately and securely to a **service** powered by Azure Private Link. Private Endpoint uses a private IP address from your VNet

E.g. connect to azure sql server



**Private Link Service** allows you to connect your workload to Private Link. You need **Internal Load Balancer** and associated **IP Address Range**

Many Azure services by default support Private Link eg. Azure Sto

## T-SQL

**Transact-SQL** (T-SQL) is a set of programming extensions from Sybase and Microsoft that add several features to the Structured Query Language (SQL),

- T-SQL expands on the SQL standard to include:
- procedural programming
  - local variables

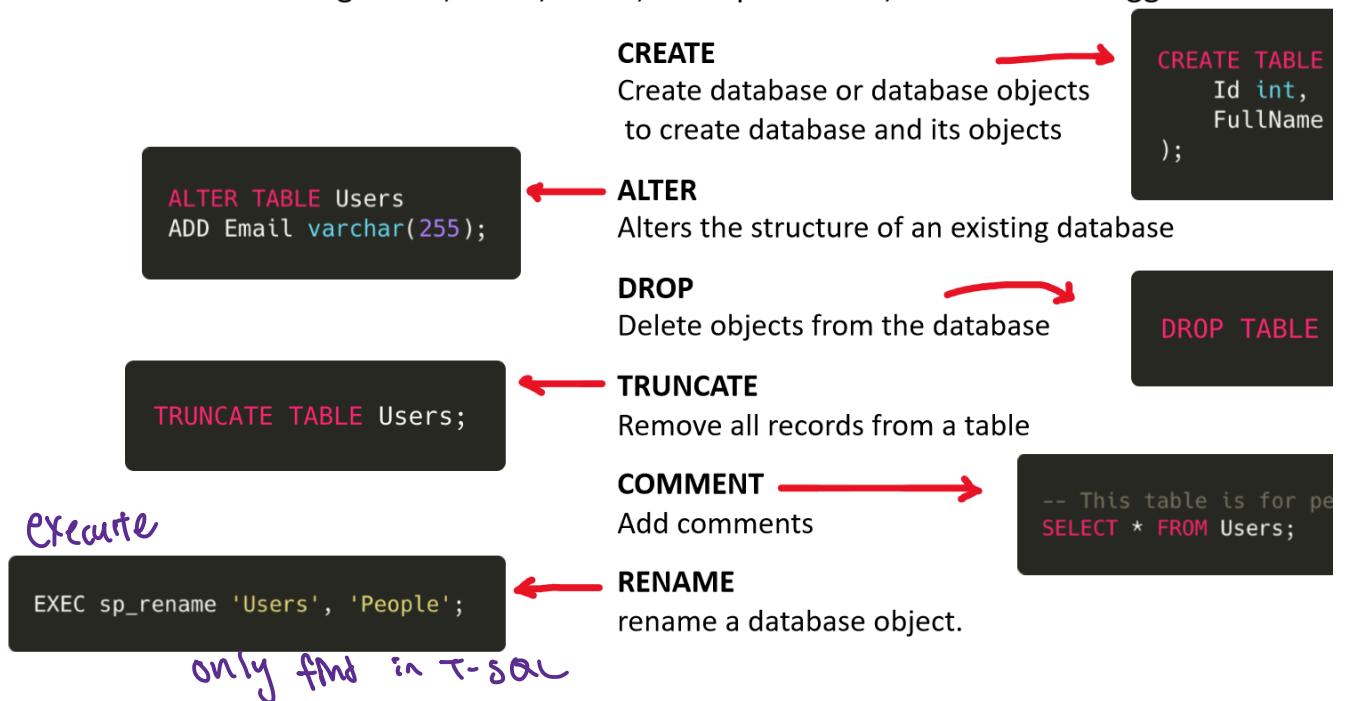
- For Microsoft SQL Server there are five groups:
- Data **Definition** Language (DDL)
  - used to define the database schema

- various support functions for string processing
- date processing
- Mathematics
- changes to the DELETE and UPDATE statements

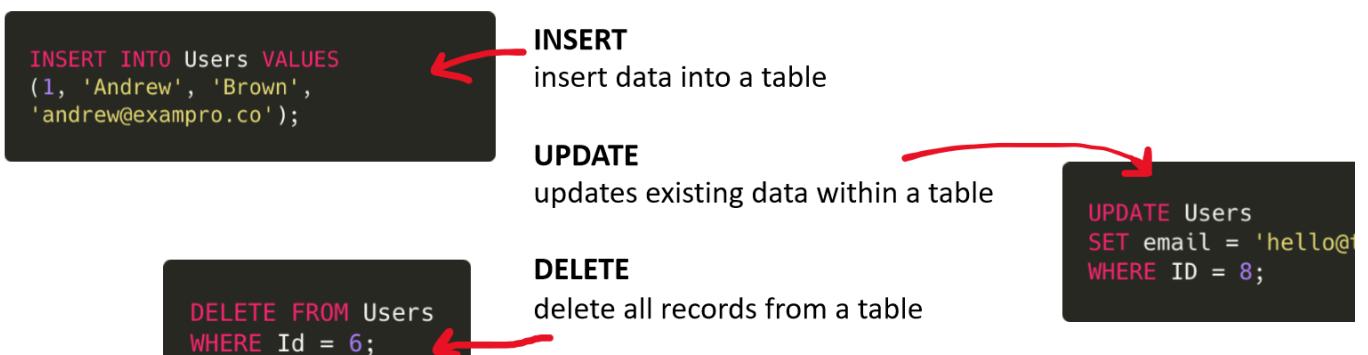
- Data **Query** Language (DQL)
  - used for performing queries on the data
- Data **Manipulation** Language (DML)
  - manipulation of data in the database
- Data **Control** Language (DCL)
  - rights, permissions and other controls
- **Transaction Control** Language (TCL)
  - transactions within the database

## Data Definition Language (DDL)

A Data **Definition** Language (DDL) is SQL syntax commands for **creating and modifying the database or database objects**  
eg: table, index, views, store procedure, function and triggers



## Data Manipulation Language (DML)



**MERGE – UPSERT**

To insert and update records at the same time

**CALL**

call a PL/SQL or Java subprogram.

**T-SQL Specif.**

CALL CalcDistanace 'To'

**LOCK TABLE**

concurrency control to ensure two people are not writing to the program at the same time

## Data Query Language (DQL)

**SELECT**

```
Users.Id,  
Users.FullName,  
Users.Email  
FROM Users;
```

**SELECT**

Select data from a table, And be able to specific exactly which columns to return

**SHOW**

Describe what a table looks like (what column it contains)

T-SQL

EXEC sp\_columns U

**EXPLAIN PLAN**

Returns the query plan for a Microsoft Azure Synapse Analytics SQL statement without running the statement

**HELP**

Reports information about a database object

EXEC sp\_help Users

## Data Control Language (DCL)

GRANT SELECT, INSERT, UPDATE  
DELETE ON employees TO andr

**GRANT**

allow users access privileges to database.

**REVOKE**

withdraw users access privileges given by using the GRANT command.

rights { }

```
REVOKE DELETE ON employees FROM
bayko;
```

# Transaction Control Language (TCL)

**Transaction Control Language** commands are used to manage transactions in the database.

↓ neat mu  
things

## COMMIT

set to permanently save any transaction into the database

## ROLLBACK

restores the database to last committed state

## SAVEPOINT

used to temporarily save a transaction so that you can rollback to that point whenever needed

## SET TRANSACTION

specify characteristics for the transaction.

# MS SQL Commands

## SQL Commands

### DDL

- CREATE
- ALTER
- DROP
- TRUNCATE
- COMMENT
- RENAME

### DML

- INSERT
- UPDATE
- DELETE
- MERGE
- UPSERT
- CALL
- LOCK TABLE

### DQL

- SELECT
- SHOW
- EXPLAIN PLAN
- HELP

### DCL

- GRANT
- REVOKE

Defining

Manipulating

Querying

Controlling

These are the two the SQL documents the exam will focus upon

## What is a Key / Value store?

A key/value stores a **unique key** alongside a value

Key	Value
Data	1010101000101011001010010101001
Worf	0110101100010101010101011100010
Ro Laren	001010100101011001010101010101010

Key values stor  
They generally  

- Relationship
- Indexes
- Aggregation
- **transaction**

Key	Value
Data	{species: android, rank: 'Lt commander' }
Worf	{species: klingon, rank: 'Lt commander' }
Ro Laren	{species: bajoran, affiliation: 'maquis'}

A simple key  
interpret thi  
dictionary (a  
or hash)

A key/value store can resemble  
tabular data, it does not have to  
have the consistent columns per  
row (hence its schemaless)

Key (Name)	Species	Rank	Affiliation
Data	android	Lt commander	
Worf	klingon	Lt commander	
Ro Laren	bajoran		maquis

Due to the  
they can s  
relational

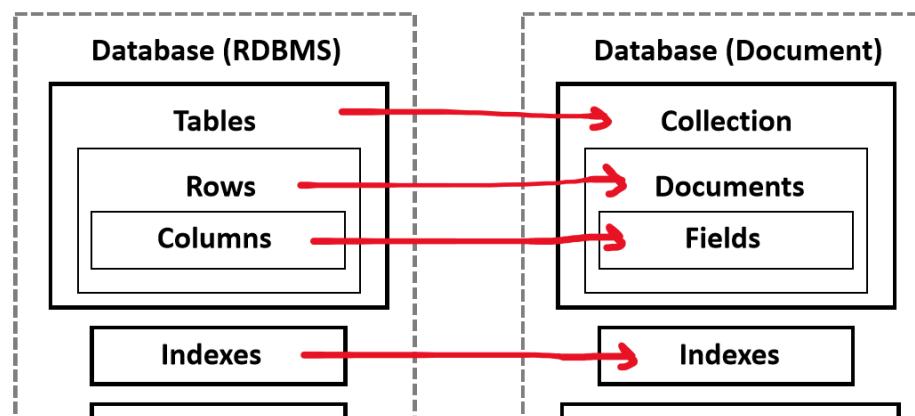
## What is a Document store?

A **document store** is a NOSQL database that stores **documents** as its primary data structure.

A document could be an XML but more commonly is JSON or JSON-Like

Document stores are sub-class of Key/Value stores

The components of a document store compared to Relational database





# MongoDB

BSON



**MongoDB** is an open-source document database **which stores JSON-like data**.  
The primary data structure for MongoDB is called **BSON**

## Binary JSON (BSON)

- BSON is a subset of JSON and so its data structure is very similar.
- BSON is designed to be **efficient both in storage space and scan-speed** compared to JSON
- BSON has more data-types than JSON:
  - Eg. Datetime, byte arrays, regular expressions, MD5 binary data, javascript code

```

BSON:
\x16\x00\x00\x00          // total document size
\x02                      // 0x02 = type String
hello\x00                  // field name
\x06\x00\x00\x00world\x00  // field value (size)
\x00                      // 0x00 = type E00 (null)
  
```

What it looks like to perform an → operation on a MongoDB databases

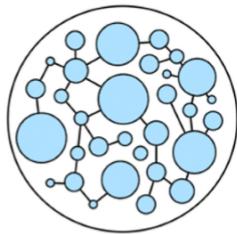
```

db.inventory.insertMany([
  { item: "journal", qty: 25, size: { h: 14, w: 21, uom: "cm" },
  { item: "notebook", qty: 50, size: { h: 8.5, w: 11, uom: "in" },
  { item: "paper", qty: 100, size: { h: 8.5, w: 11, uom: "in" },
  { item: "planner", qty: 75, size: { h: 22.85, w: 30, uom: "in" },
  { item: "postcard", qty: 45, size: { h: 10, w: 15.25, uom: "in" }
]);
  
```

# MongoDB

- MongoDB supports searches against:
  - fields
  - ranged queries
  - regular-expressions
- MongoDB supports **primary** and **secondary** indexes
- High availability can be obtained via replica sets (replica to offload reads or acts a stand-by in case of failure)
- MongoDB scales horizontally using sharding
- MongoDB can run over multiple servers via load balancing
- MongoDB can be used as a file system, called **GridFS**
  - with load balancing and data replication features over multiple machines for storing files.
- MongoDB provides three ways to perform aggregation (grouping data during a query)
  - aggregation pipeline
  - map-reduce
  - single-purpose aggregation
- MongoDB supports fixed-size collections called capped collections
- MongoDB claims to support multi-document ACID transactions

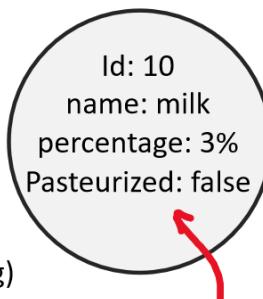
# What is a Graph database?



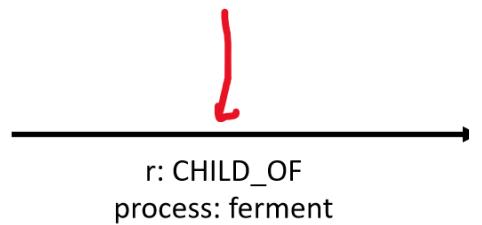
A graph database is a database composed of a data structure that uses vertices (nodes, dots) which form **relationship** to other vertices via edges (arcs, lines)

Use Cases for Graph Database:

- Fraud detection
- Real-time recommendation engines.
- Master data management (MDM)
- Network and IT operations.
- Identity and access management (IAM)
- Traceability in Manufacturing
- Contact Tracing
- Data Lineage for GDPR
- Customer 360-degree analysis (marketing)
- Product recommendations
- Social Media graphing
- Feature Engineering (ML)



The edge can contains **relational data** such as direction and data properties



The node will contain **data properties**

# What is a Apache TinkerPop?



Apache TinkerPop is a **graph computing framework** for databases (OLTP) and graph analytic systems (OLAP)

TinkerPop enables developers to use **an vendor-agnostic distributed framework** to traverse (query) many different graph systems:

Alibaba Graph Database  
**Amazon Neptune**  
 ArangoDB  
 Bitsy  
 Blazegraph  
**CosmosDB**  
 ChronoGraph  
 DSEGraph  
 GRAKN.AI

**Hadoop (Spark)**  
 HGraphDB  
 Huawei Graph Engine Service  
 HugeGraph  
 IBM Graph  
 JanusGraph  
 JanusGraph (Amazon)  
**Neo4j**  
 neo4j-gremlin-bolt

**OrientDB**  
 OverflowDB  
 Apache S2Graph  
 Sqlg  
 Stardog  
 TinkerGraph  
**Titan**  
 Titan (Amazon)  
 Titan (Tupl)  
 Unipop



TinkerPop includes a graph traversal language called **Gremlin**  
 Which is the single language that can be used for all these graph systems

## What is Gremlin?

## Graph Databases

Gremlin is the **graph traversal language** for Apache TinkerPop

```
g.V().has("name", "gremlin").as("a")
  .out("created").in("created")
    .where(neq("a"))
      .groupCount().by("title")
```

Gremlin is designed to the "Write once, run anywhere" (**WORA**)

Gremlin traversal can be evaluated as

- **real-time database query (OLTP)**
- or as a **batch analytics query. (OLAP)**



**Gremlin Host**  
means You can use  
programming lan

## Azure Tables

Azure Table storage is a **NoSQL key/value datastore** within Azure Storage Accounts

Azure Table stores **non-relational** structured data with a schemaless design

There are two ways to interact with Azure Tables:

- Azure Table Storage API
- **Microsoft Azure Storage Explorer**

The screenshot shows the Microsoft Azure Storage Explorer interface. The left pane, titled 'EXPLORER', lists storage accounts and containers. The right pane, titled 'mytable', shows the table structure with columns 'PartitionKey' and 'RowKey'. The table contains several entries, including 'Klingon' and 'Worf'.

## Azure Tables - Adding Entries

When you enter data you must provide a

- **Partition Key** — unique identifier for the partition within a given table
- **Row Key** — unique identifier for an entity within a given partition

*Data type*

Azure Tables supports:

- String
- Boolean
- Binary
- DateTime
- Double
- Guid
- Int32
- Int64

## Add Entity

Property Name	Type	Value
PartitionKey	String	Klingon
RowKey	String	Worf
Rank	String	Lt Commander
Age	Int32	39
<input type="text" value="Enter a name up to 255 characters in size"/> <div style="border: 1px solid #ccc; padding: 2px; margin-top: 5px;"> <span style="color: #0078d4;">✓</span> String  <span style="color: #0078d4;">✓</span> Boolean  Binary  DateTime  Double  Guid  Int32  Int64 </div>		
<input type="button" value="Add Property"/>		

## Azure Tables – Querying

You can **query** data along the Partition and Row Key.

You also further filter for any filters

The screenshot shows the Azure Table Query interface. At the top, there are buttons for Close Query, Import, Export, Add, Edit, Select All, Column Options, Delete, and Table. Below this is a toolbar with icons for filtering, sorting, and other operations. The main area displays a query builder with three clauses:

And/Or	Field	Type	Operator	Value
+	PartitionKey	String	=	Klingon
+	RowKey	String	=	Worf
+	Age	String	>	20
<a href="#">Add new clause</a>				

## Introduction to CosmoDB

Azure CosmoDB is a service for fully-managed **NoSQL databases** that are designed to scale and high performance

CosmoDB supports **different kinds** of NoSQL database engine which you interact with

- **Core SQL (document)** datastore
- Azure Cosmos DB API for **MongoDB (document)** datastore
- **Azure Table (key/value)** datastore
- **Gremlin (graph)** datastore
  - Based on **Apache TinkerPop**



NoSQL

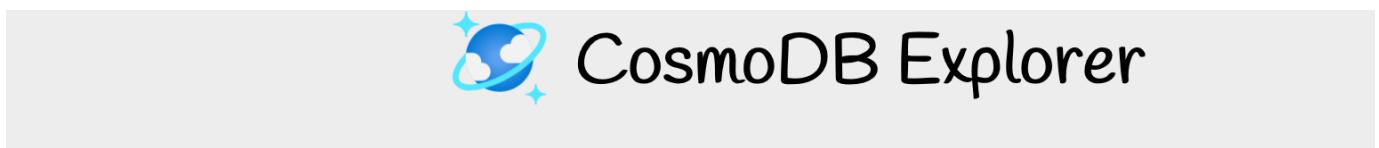
All of these NoSQL engines specific **capacity**:

- Provisioned throughput (pay for guarantee of capacity )
- Serverless (pay for what you use)

Capacity mode ⓘ

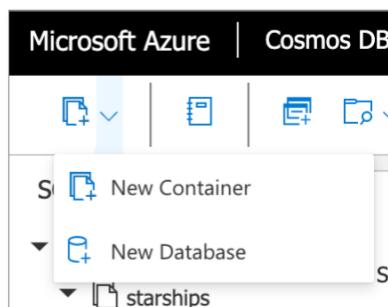
Provisioned throughg

[Learn more about capaci](#)



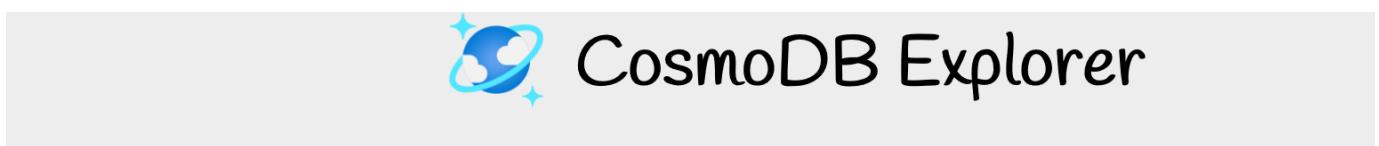
CosmoDB Explorer is a **web interface** to explore and interact with your Cosmo

[cosmos.azure.com](https://cosmos.azure.com)



id	_rid
"id": "NCC-1701-D", "name": "enterprise-d", "details": {   "captain": "Jean-Luc Picard" }	

Here is a **CosmoDB Core SQL**  
In CosmoDB Explorer adding a  
document to the database



When navigating within Azure to a ComsoDB Account under **Data Explorer** is the same as CosmoDB Explorer

The screenshot shows the Azure portal interface. On the left, a sidebar for the 'Starships' Cosmos DB account lists various management options: Overview, Activity log, Access control (IAM), Tags, Diagnose and solve problems, Quick start, Notifications, and Data Explorer. The 'Data Explorer' option is highlighted with a red arrow. To the right, the main content area displays a 'Welcome to Cosmos DB' message and a 'Start with Sample' button. Below that are sections for 'Common Tasks' (New Database) and 'Recents' (starships, Families). At the bottom, there are performance metrics: 0 requests/s, 0吞吐量 (throughput), and 0 errors.



## CosmoDB Explorer

Using CosmoDB Explorer  
With a **Graph Database**  
(Gremlin)

The screenshot shows the CosmoDB Explorer interface for a 'gremlin-cosmodb-exampro' database. The top navigation bar includes Microsoft Azure, Cosmos DB, and a dropdown for the database. The main area is titled 'GREMLIN API' and shows a tree view of a 'Starsystems' collection with a 'starsystems' sub-collection. A query 'g.V()' is run against the 'Graph' tab, and the results are displayed in a table with two rows: '0002087e-7bad-4b0e-932f...' and '13ca6de2-21cf-42d7-ad7d...'. To the right, a graph visualization shows a single blue vertex with the ID '13ca6de2-21cf-42...'. A sidebar on the right shows a 'Label' section with 'Bezad' selected and a 'system' property set to 2. There is also an 'Add Property' button. A red arrow points from the 'Using CosmoDB Explorer With a Graph Database (Gremlin)' text to the 'Graph' tab in the interface.



## CosmoDB Table API vs Azure Table

When comparing Account Storage Azure Table vs CosmoDB Table API

Feature	Azure Table Storage	Azure Cosmos DB Table API
Latency	Fast, but no upper bounds on latency.	Single-digit millisecond latency for read:
Throughput	Variable throughput model limit of 20,000 operations/s.	Guaranteed backed by SLAs. No upper limit on throughput

Global distribution

Single region with one optional readable

30+ regions

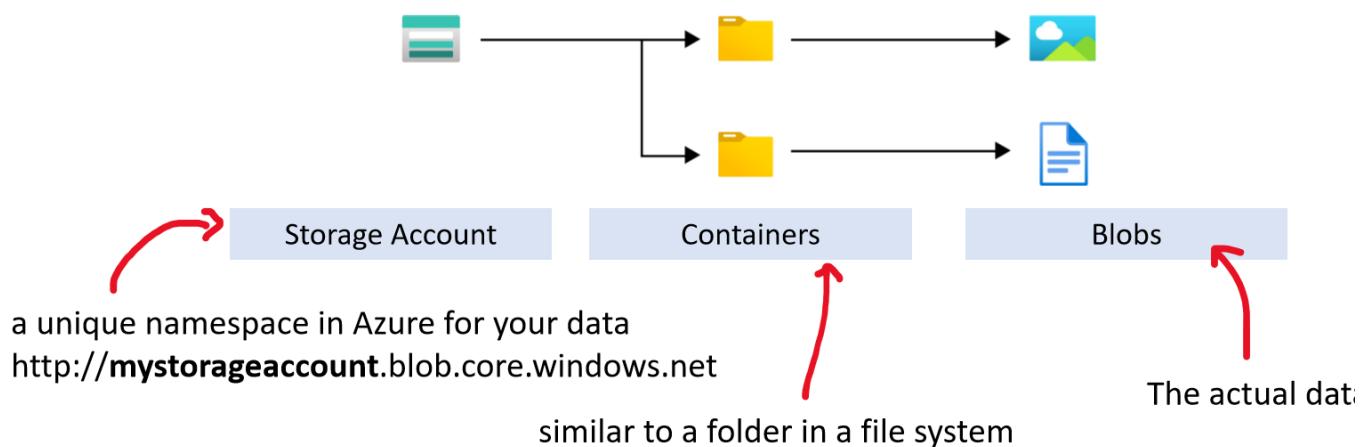
Indexing	Only primary index on PartitionKey and RowKey. Automatic and complete indexing on all management.	
Query	Query execution uses index for primary key, and scans otherwise.	Queries can take advantage of automation for fast query times.
Consistency	Strong within primary region Eventual within secondary region.	Five well-defined consistency levels
Pricing	Consumption-based	consumption-based or provisioned capacity
SLAs	99.99% availability	99.99% availability SLA (some conditions)



## Azure Blob

Blob storage is a **object-store** that is optimized for **storing massive amounts of unstructured** Unstructured data is data that doesn't adhere to a particular data model or definition, such as:

Azure Blobs are composed of the components:



## Azure Blob

Azure Storage supports **3 types** of blobs:



### 1. Block blobs

- store text and binary data
- made up of blocks of data that can be managed individually
- store up to about 4.75 TiB of data



### 2. Append blobs

- Optimized for append operations

- ideal for scenarios such as logging data from virtual machine



### 3. Page blobs

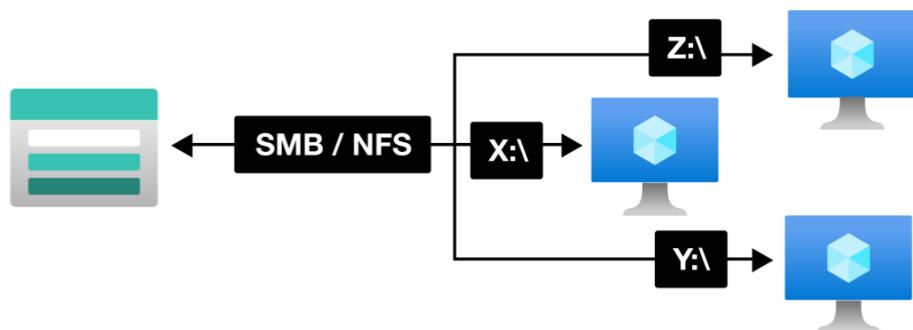
- store random access files up to 8 TB in size.
- store virtual hard drive (VHD) files and serve as disks for Azure virtual machines



## Azure Files

Azure Files is a fully managed **file share** in the cloud.

A file share is a **centralized server for storage** that allows **multiple connect**. It's like having one big shared drive that everyone (Virtual Machines) can work on at the same time.



To connect to the file share a **network protocol** is used:

- Server Message Block (SMB)
- Network File System (NFS)

When a connection is established the file share's filesystem will be accessible in the system as a new directory within your own directory tree. This is known as **mounting**.



## Azure Files – Use Cases

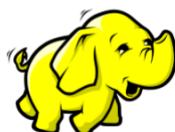
### Use Cases

- Completely **replace or supplement** on-premises file servers Network Attached Storage (NAS) devices
- **Lift-and-Shift** your on-premise storage to the cloud via Classic Lift or Hybrid Lift
  - “Lift-and-Shift” means when you move workloads without rearchitecting, e.g. importing local VHDs
  - Classic Lift — where both the application and its data are moved to Azure
  - Hybrid Lift — where the application data is moved to Azure Files, and the application continues to run on-premises
- **Simplify cloud development**
  - Shared application settings — Multiple VMs and developer workstations need to access the same configuration files
  - Diagnostic share — All VMs log to the file share, developers can mount and debug all logs in a central location
  - Dev/Test/Debug — Quickly share tools for developer needed for local environments
- **Containerization**
  - You can use Azure Files to persist volumes for stateful containers

VISIT USE AZURE FILES INSTEAD OF SETTING UP YOUR OWN FILE SHARE SERVER!

- **Shared Access** — Already setup to work with standard networking protocols SMB and NFS
- **Fully managed** — Its kept up to date with security patches, designed to scale
- **Scripting and Tooling** — You can automate the management and creation of file shared with Azure
- **Resiliency** — Built to be durable and always working

## What is Apache Hadoop?



Hadoop is an open-source framework for **distributed processing of large data sets**

Hadoop allows you to distribute:

- large dataset across many servers eg *HDFS*
- computing queries across many servers eg. *MapReduce*

These computer servers do not need specialized hardware and can run on standard hardware

The Apache Hadoop framework has the following:

- Hadoop Common — collection of common utilities and libraries that support other Hadoop modules
- Hadoop Distributed File System (HDFS) — a resilient and redundant file storage distributed on clusters
- Hadoop MapReduce — writes apps that can process multi-terabyte data in-parallel on large clusters
- Hbase — a distributed, scalable, big data store
- YARN — manages resources, nodes, containers and performs scheduling
- HIVE — used for generating reports using an **SQL language**
- PIG — A high-level **scripting language** to write complex data transformations

Hadoop can integrate with many other open-source projects via **Hadoop compatibility**

## What is Apache Kafka?



Apache Kafka is an **open-source streaming platform** to create **high-performance pipelines**, streaming analytics, data integration, and mission-critical applications

Kafka was originally developed by LinkedIn, and open-sourced in 2011

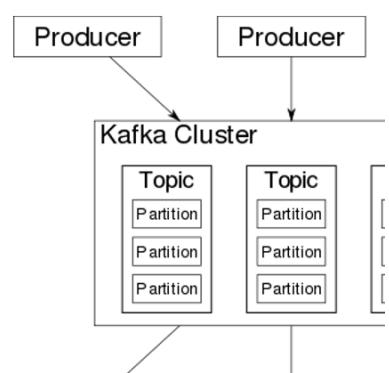


Kafka was written in Scala and Java.  
to use Kafka you need to write **Java** code.

In Kafka data is stored in partitions on a Kafka Cluster which can span multiple machines (distributed computing)

**Producers** publish messages in a key and value format using the Kafka Producer API

**Consumers** can listen for messages and consume them using the Kafka Consumer API



Messages are organized into **Topics**. Producers will push messages to topics and consumers will listen on topics.



# Introduction to Azure HDInsights

Azure HDInsight is **managed service to run popular open-source analytics s**



HDInsight supports the following frameworks:

- Apache Hadoop
- Apache Spark
- Apache Kafka
- Apache Storm
- Apache Hive
- Apache HBase
- Low Latency Analytical Processing (LLAP)
- R

HDInsight has broad range of scenarios

- Extract, Transform, and Load (ETL)
- Data Warehousing
- Machine Learning
- Internet of Things (IoT)

## Azure HDInsights – Apache Ambari

Apache Ambari is an **open-source Hadoop management web-portal** for **provisioning, managing, and monitoring** Apache Hadoop clusters

When you **create an HDInsights Cluster** you will get a Cluster Dashboard (Apache Ambari dash

# Apache Spark



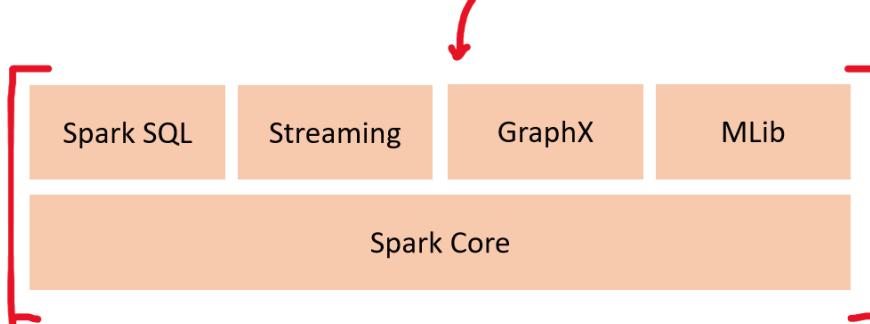
Apache Spark is an open-source **unified analytics engine** for **big data and machine learning**

Spark lets you run workloads **much faster than Hadoop**:

- 100x faster in memory
- 10x faster than disk

Which is why Spark is described as **lightning fast**

Apache Spark is a collection of libraries that work well together to form an **analytics ecosystem**



## Spark Core

The underlying engine and API. The AI programming languages:  
R, SQL, Python, Scala, Java

## Spark SQL

Introduces a data structure called a DataFrame used with DSL to work with structured data

## Spark Streaming

Allows Spark to ingest data from many sources like HDFS, Flume, Kafka, Twitter, Kinesis

## GraphX

distributed graph-processing framework

## Machine Learning Library (MLib)

a distributed machine-learning framework for machine learning and statistical algorithms

# Apache Spark – RDD API

To interact with Apache Spark:

Resilient Distributed Dataset (RDD) is a domain specific language (DSL) to execute various parallel operations on an Apache Spark cluster.

## Common RDD API functions

map  
flatMap  
mapPartitions  
filter  
distinct  
reduce  
count  
first  
take  
countByValue  
sortBy  
groupBy  
fold

union  
add  
subtract  
intersection  
saveAsTextFile  
saveAsHadoopFile  
saveAsPickleFile  
min  
max  
mean  
status  
parallelize

## Example of RDD API

```
text_file = sc.textFile("hdfs://...")  
counts = text_file.flatMap(lambda line: line.  
    .map(lambda word: (word, 1)) \  
    .reduceByKey(lambda a, b: a + b)  
counts.saveAsTextFile("hdfs://...")
```

# Databricks Platform



Databricks is a software company specializing in **providing fully managed Apache Spark**. The company founders were the creators of Apache Spark, Delta Lake and MLFlow.

Databricks has two offerings:

**Databricks Platform** — Databricks cloud-based Spark platform with an ease-to-

- Launch fully managed Spark clusters
- Launch notebooks to write code and interact with Spark
- Create workspaces to collaborate with team members
- Role Base Access Controls
- Create jobs for ELT or data analysis tasks that run immediately or on a schedule
- Create MLFlow Workflows
- **Available on all main cloud service providers eg. AWS, Azure, GCP**

*free*

**Databricks Community Edition** — free version of Databricks Platform for education

- Create a free micro-cluster that terminates after 2 hours when idle
- No workspace, jobs or RBAC

*Role Base Access (RBAC)*

## Introduction to Azure Databricks

Azure Databricks is a **partnership between Microsoft and Databricks** to offer **Databricks Platform within the Azure Portal** running on Azure compute serv

Azure Databricks offers two environments:

### Azure Databricks Workspace

- The DataBrick Platform with integrations to **Azure data-related services** for building big data pipelines
  - Batching: Azure Data Factory
  - Streaming: Apache Kafka Event Hub, or IoT Hub
  - Storage: Azure Blob Storage or Azure Data Lake Storage

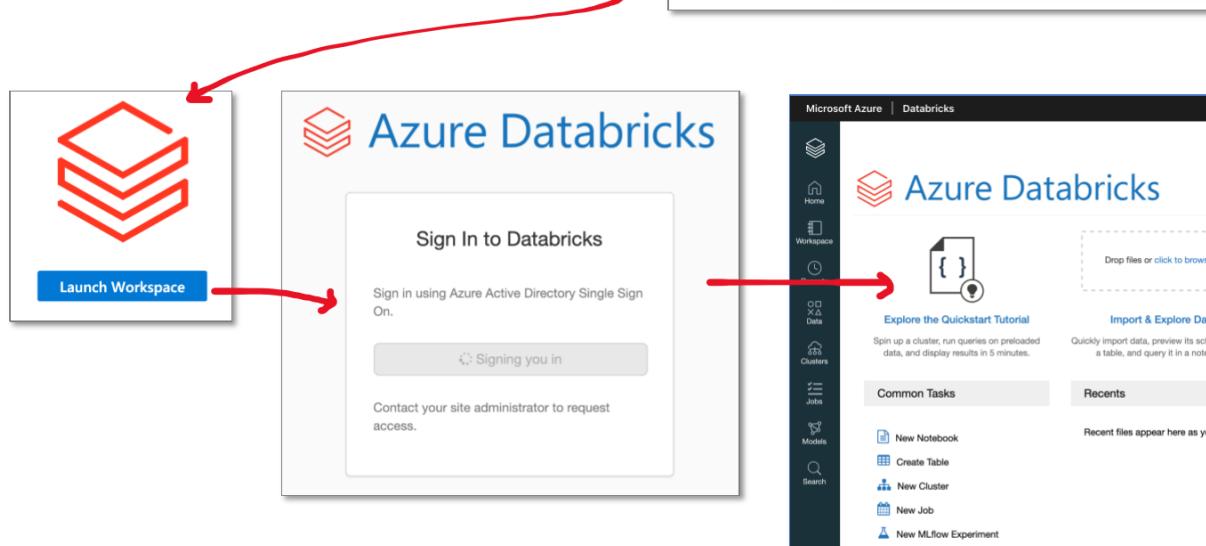
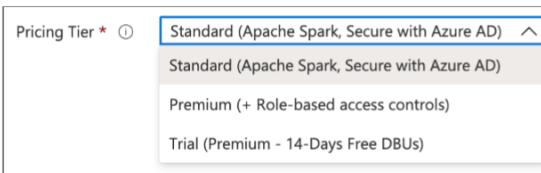
### Azure Databricks SQL Analytics

- run SQL queries on your data lake
- create multiple visualization types to explore query results
- build and share dashboards

## Azure Databricks Workspaces

Create a Databricks Workspace by:

- Creating a workspace and choosing a plan
- Launching the workspace
- SSO to the workspace
- Start using Databricks Platform

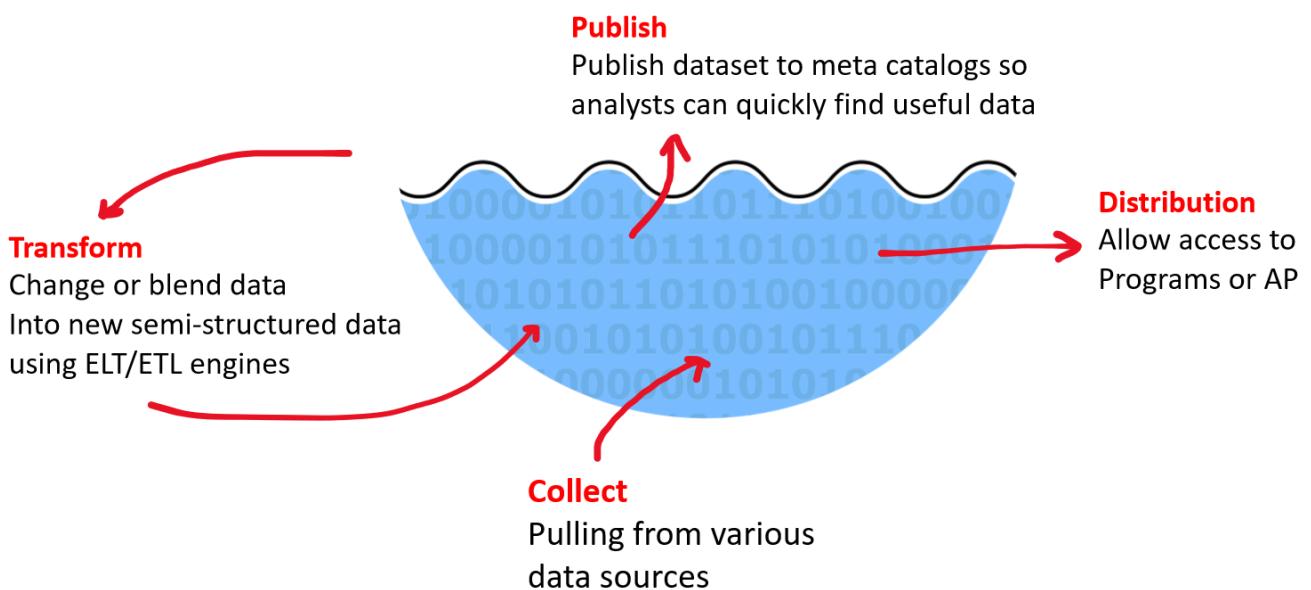


## Introduction to Data Lakes

A data lake is a **centralized data repository for unstructured and semi-structured data**

A Data Lake is intended to store vast amounts of data

Data lakes generally use **object (blobs) or files** as its storage medium.

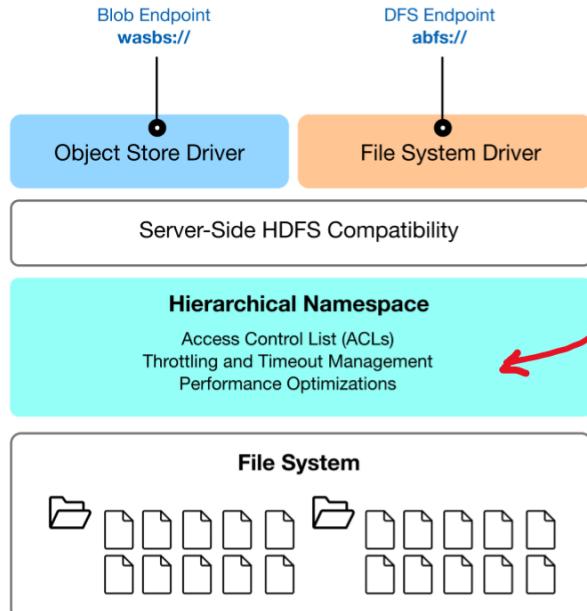


## Introduction to Azure Data Lake

### Azure Data Lake Storage Gen1 (*deprecated*)

The first version of Data Lake Storage and will be retired in 2024. New users should use Gen2.

[https://imperiallondon-my.sharepoint.com/personal/hz6215\\_ic\\_ac\\_uk/\\_layouts/15/Doc.aspx?sourcedoc={3644fae9-16d3-4f66-bda7-a4303813c...](https://imperiallondon-my.sharepoint.com/personal/hz6215_ic_ac_uk/_layouts/15/Doc.aspx?sourcedoc={3644fae9-16d3-4f66-bda7-a4303813c...) 70/75



## Azure Data Lake Storage Gen2

Data Lake Storage is Azure Blob storage which has been extended to support big data analytics workloads

- Designed to handle **petabytes of data** and **hundreds of gigabits of throughput**
- In order to efficiently access data, Data Lake Storage uses a **hierarchical namespace** to Azure Blob Storage

# Azure Data Lake Storage Gen2

Azure Data Lake Storage can be accessed two ways:

### Azure Blob File System over SSL (abfss)

An Hadoop filesystem driver that is compatible with Azure Data Lake Storage Gen2

The abbfss driver will have a URL that will allow you to access and write to the underlying Data Lake Storage

```
abfs[s]://<file_system>@<account_name>.dfs.core.windows.net/<path>/<file>
```

### Windows Azure Storage Blob over SSL (wasbs)

<MAYBE JUST SKIP THIS SLIDE>

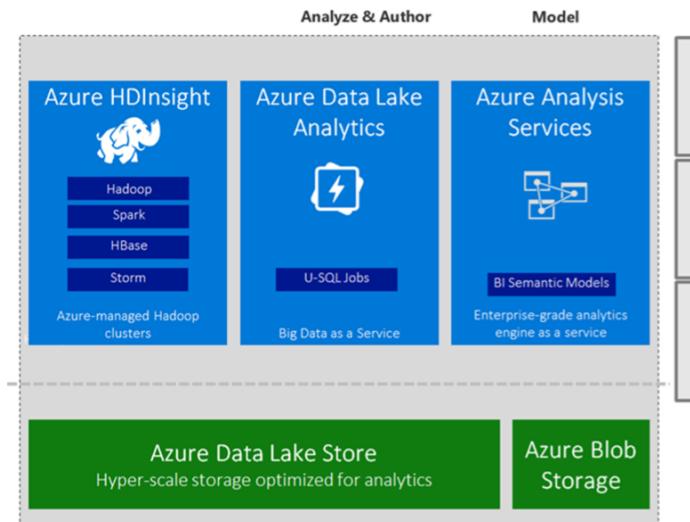
# Azure Data Lake Analytics

Azure Data Lake Analytics is an **on-demand analytics job service** that simplifies

Instead of deploying, configuring, and tuning hardware...

You write queries (via U-SQL) to transform your data and extract valuable insights

Exporting approximately 2.8 billion rows of TPC-DS store sales data (~500 GB) into a CSV file took less than 7 minutes and importing the full 1 TB set of source data into Azure Analysis Services by using the Azure Data Lake connector took less than 6 hour



## Azure Data Lake Analytics – U-SQL

**U-SQL** is a structured query language included within Data Lake Analytics to **perform queries** on your data lake.

**U-SQL** can query and combine data from a variety of data sources, including:

- Azure Data Lake Storage
- Azure Blob Storage
- Azure SQL DB
- Azure SQL Data Warehouse,
- SQL Server instances running in Azure VMs



You can install **Azure Data Lake Tools** for Visual Studio to perform U-SQL jobs on your Azure Data Lake

```

DECLARE @in string = "/Samples/DataLakeAnalytics/SearchLog.csv";
DECLARE @out string = "/output/results/searchlog.parquet";

@searchlog =
    EXTRACT UserId      int,
             Start        DateTi
                         string,
             Region       string,
             Query        string,
             Duration     int?,
             Urls         string,
             ClickedUrls string
    FROM @in
    USING Extractors.Tsv();

@rs1 =
    SELECT Start, Region, Duration
    FROM @searchlog
    WHERE Region == "en-gb";

@rs1 =
    SELECT Start, Region, Duration
    FROM @rs1
    WHERE Start >= DateTime.Parse("2017-01-01")
        AND Start <= DateTime.Parse("2017-01-02");

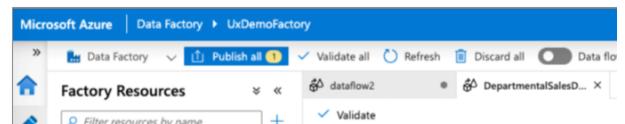
OUTPUT @rs1
TO @out
  
```



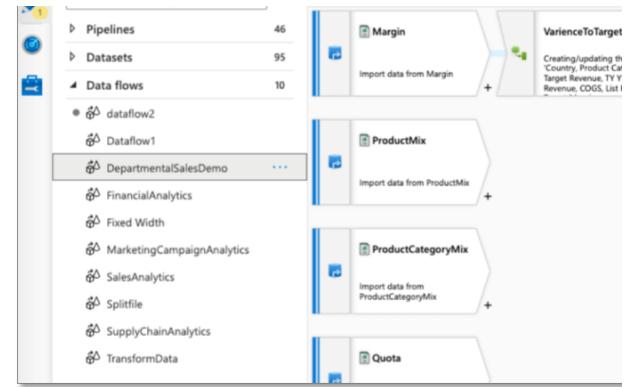
## Introduction to Azure Data Factory

Azure Data Factory is a managed service for **ETL, ELT and data integration**. Create **data-driven workflows** for orchestrating **data movement** and **transform**

- Create Pipelines to schedule data-driven workflows



- build complex ETL processes that transform data visually with data flows
- using compute services such as Azure HDInsight Hadoop, Azure Databricks, and Azure SQL Database
- publish your transformed data to data stores such as Azure Synapse Analytics
- raw data can be organized into meaningful data stores and data lakes



## Azure Data Factory – Core Components

### Pipelines

a logical grouping of activities that performs a unit of work

### Activities

A processing step in a pipeline

### Datasets

data structures within the data store

### Linked services

define the connection information for data sources to connect to Data Factory



### Data Flows

logic to determine how data moves through a pipeline or is transformed

### Integration Runtimes (RI)

compute infrastructure used by Azure Data Factory



### Control flow

orchestration of pipeline activities that includes chaining activities in a sequence, branching

## Microsoft SQL Server Integration Services



**Microsoft SQL Server Integration Services (SSIS) is a platform for enterprise-level data integration and data transformation.**

You can perform the following tasks with SSIS

- Copy files
- Download files
- Loading data into data warehouses
- Cleansing data
- Mining data

SSIS can be used to automate SQL Server database

SSIS can be used as an integration runtime in Azure

SSIS has...

- built-in tasks and transformations

~~Managing Data~~

- managing SQL Server objects
- ~~Managing SQL~~ Server data

Perform ELT with variety of sources:

- XML
- Flat files
- Relational data sources

## OneNote

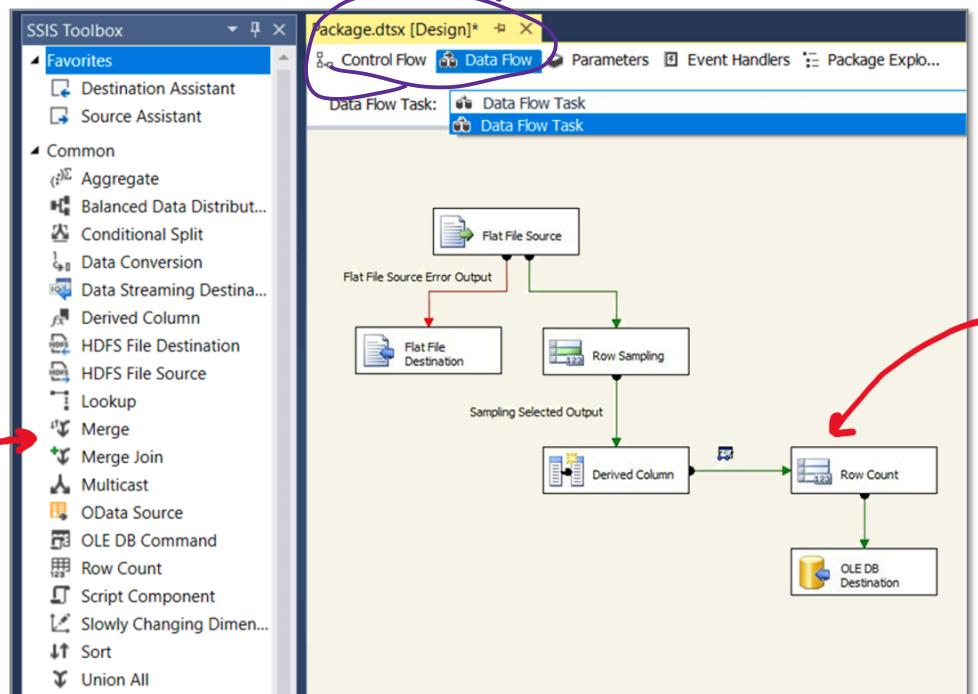
- ~~MANAGE TASKS AND TRANSFORMATIONS~~
- graphical tools for building packages
- Integration Services Catalog database
  - where you store, run, and manage packages

Use **Graphical Integration Services tools** for to integrate and transform data without having to write code

**SSIS Designer** is a graphical tool that you can use to maintain Integration Services packages.

# Microsoft SQL Server Integration Serv

*similar to Data factory*



SSIS allows you to drag out data transformations

