

# Pre-requisites

13 February 2023 10:04

CPU i5/i6,  
RAM - at least 12 GB,  
Disk Storage - 500GB

Windows 10, 11 Professional / Enterprise  
SQL Server 2017 / 2019 Developer / Enterprise  
MS Office  
PowerBI  
Alteryx

username	email	Password	invitation_link
db01@mml.local	db01@mml.local	uavdxa5yg	<a href="https://laas.makemylabs.in/wsm/KPMG-Databrick-skip20230213170942?invitation_id=35acb6a3-341b-4330-b87a-8a810472703c&amp;context=True">https://laas.makemylabs.in/wsm/KPMG-Databrick-skip20230213170942?invitation_id=35acb6a3-341b-4330-b87a-8a810472703c&amp;context=True</a>

# Database Fundamentals and SQL Server

13 February 2023 10:07

## Database types & Models

1. **Flat file database** - kind of text database where each line of the plain text file holds only a single record (e.g. MS access)
2. **Hierarchical database** - based on hierarchical data model, it's viewed as a collection of tables, data is designed into a tree like structure where each record consists of one parent record and many child record. (e.g. IBM DB2 - IBM information Management system (IMS), Windows Registry, XML data storage.
3. **Network Model database** - it can consists of parent segments and this segment can be grouped together as levels but there always exists a logical association between the segments belonging to any level.
4. **Relational database** - consists of tables and columns, rows.
5. **Object-oriented database** - information can be represented in the form of object - oriented programming inclined towards the objects like e.g. multimedia records in a relational database can be definable data object.
6. **Distributed database** - consists of two or more files located in different sites / location.
7. **NoSQL database** - non-relational db which has support for unstructured, semi-structured data as well as it can include dynamic schema, flexible data model for faster data retrieval  
e.g. Mongo db, Cassandra, Azure Cosmos db, Couch db etc.
8. **Graph database** - node-entity (rows in the table), attributes (relationship/columns) . e.g. Neo4j, Azure Cosmos db Graph API etc.

## ACID properties in RDBMS

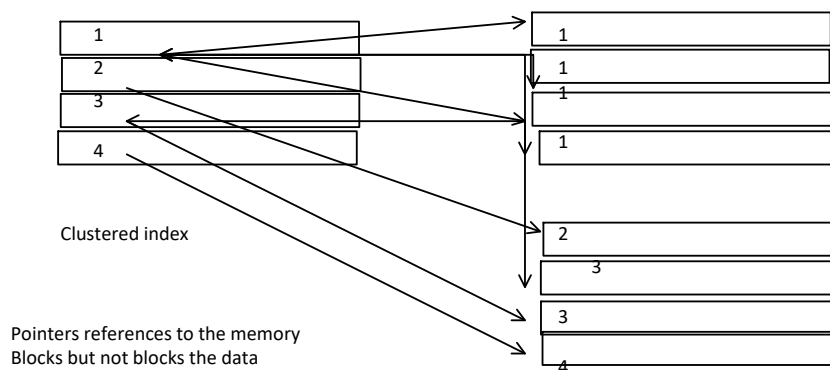
A = Atomicity	The entire transaction should take place at once or doesn't happen at all
C = Consistency	The database must be consistent before and after the transaction
I = Isolation	Multiple transactions can occur independently without interference
D = Durability	The changes of a successful transaction occurs even if a system failure occurs.

## Clustered & Non-clustered Index

1. In a table, there can be one clustered index, can be multiple non-clustered index.
2. Clustered index is much more faster, non-clustered index is slower.
3. Clustered index requires less memory for operations, non-clustered indexes requires more memory for operations.
4. In clustered index, index is the main data. In case of non-clustered index, index is the copy of the data.
5. Clustered index store pointers to blocks, not the data. Non-clustered index store both the value and a pointer to actual row which holds the data.
6. Primary Keys of the table by default is considered as a clustered index. Composite key used with unique keys of the table defines the non-clustered index.
7. A clustered index is a type of index in which table records are physically recorded to match the index. A non-clustered index is a special type of index in which the logical order of the index doesn't match physical stored order of the rows on disk.
8. Clustered index size is larger, non-clustered index size is smaller.

## Features of Indexes

1. A index can speed up the data retrieval and query execution very quickly by optimization
2. Indexes can be created or dropped with no effect on the data
3. When an index is created, it includes a column containing a wide range of values.



Primary Key	Unique Key
A table can have only one primary key	A table can have more than one unique key unlike the primary key
A primary key can't accept null values	Unique key constraint can accept null values for a column
	Unique key constraints are also referenced by the foreign key of another table, it can be used when developer wants to enforce a unique constraints on a column or group of columns which is not a primary key
Primary key has the support of auto-increment values.	A unique key does not support auto-increment value
We can't change or delete values stored in primary key	We can change the unique key values

#### Surrogate Keys

Surrogate keys are called synthetic primary keys which are generated when a new record is inserted into the table automatically by the database which can be declared as the primary key of the table.

#### Features of surrogate key

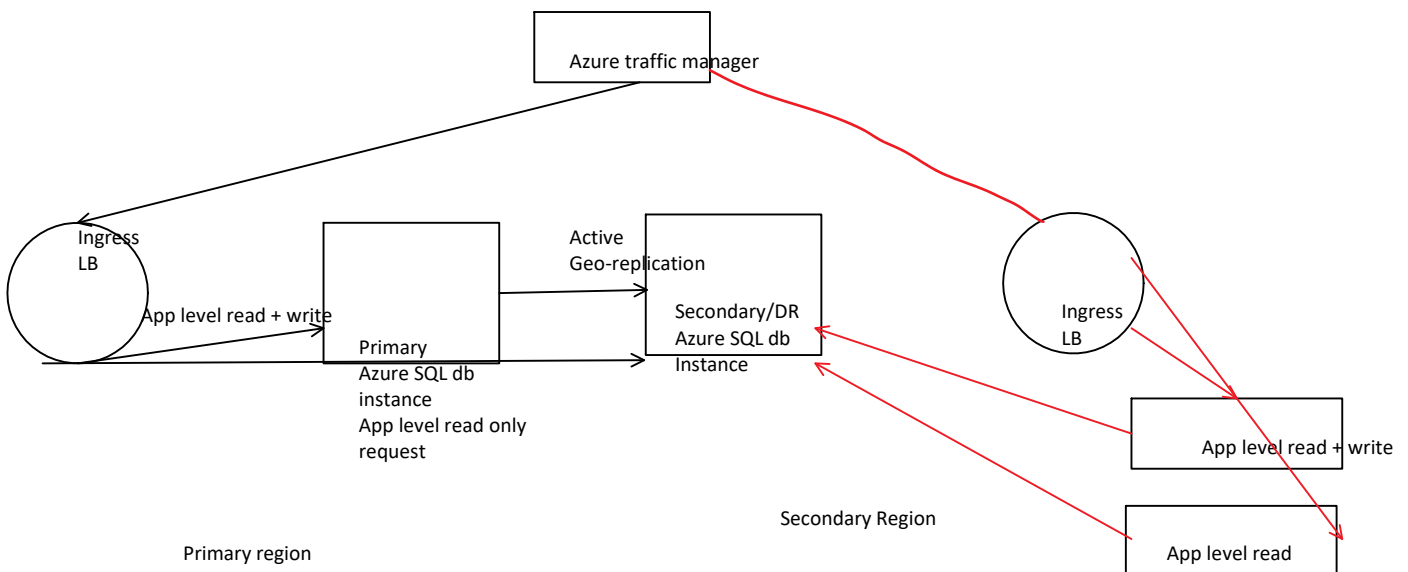
1. It's a sequential number outside the database which's made available to the user and application or it just acts as an object which's present in the database but not visible to the user/application.
2. It's automatically generated by the system
3. It holds an anonymous integer
4. It contains unique values for all records of the table
5. The value can never be modified by the user or application
6. Surrogate keys are called the factless key which is added just for the case of identifying unique values and contains no relevant fact which is useful for the table.

Surr_no	Reg_no	Name	Marks
1	21101	Mark	50
2	32281	Henry	70
3	43353	Alan	60
4	CS101	Maria	80
5	CS201	John	67

# Azure SQL

14 February 2023 09:31

SQL Server on Azure VM (IaaS)	Azure SQL Managed Instance (PaaS)	Azure SQL database (PaaS)
Azure SQL VM is preferred for the scenario of simple lift-shift or rehost of existing full SQL server database engine with core administration services like e.g. service broker, SQL Server Agent, SQL mirrors, DTC, .NET and CLR data types & functions, SQL server BI support.	Lift - shift of database migration including core SQL db features like service broker, SQL server agent, distributed transactions and .net/CLR data types etc. without any management overhead.	Purely managed SQL database offering with no licensing requirement, no underlying administration of server, SQL database or network is required.  End to end database migration is feasible with a just few clicks into the core sql development features.
License is required. Either through BYOL with AHB or pay-as-you-go model, license has to be procured.	No license is required, entire SQL on-premise feature can be availed with the managed platform service.	No license is required. Only SQL db dev specific features are available. There's a limitation of db sizes (~100TB)
SQL Server BI (SSRS, SSAS, SSIS), SQL Server broker, .net framework runtime, CLR integration, distributed transactions, database mail etc. all features are supported	.net framework, functions, distributed transactions, ACID , automated backups, HA are also supported along with no administration required.	No support for .net framework runtime, distributed transactions, CLR related functions, stored procedures, windows runtime etc. are not supported
SLA - 99.95% Automated backups and business continuity through FCI and availability groups. Migration can be implemented through Azure migrate (Azure site recovery).	SLA - 99.99% automated backups, point-in time restore, geo-replication, high availability, automated patching etc. features are available.	SLA - 99.99% , automated backups, active geo-replication, high availability, and disaster recovery, replication on both availability zone and regional level.



## Migration to Azure SQL database

We can migrate SQL Server database running on-premises or to :

- SQL Server on Azure VM
- SQL Server Managed instance
- Azure SQL db

We can migrate the following db types -

1. SQL Server running on-prem or SQL server running on VM (vmware/hyper-v etc.)
2. SQL server running on AWS EC2, Google compute engine
3. SQL server on AWS RDS
4. Cloud SQL for SQL server - GCP

1. Azure Migrate - Discovery and assess single database or at scale from different env
2. Azure SQL migration extension for Azure Data Studio - migrate to Single databases at scale, it can run in both online and offline modes.
3. Import-export service / BACPAC - migrate individual LOB apps databases , suited for smaller databases and doesn't require a separate migration service or tool.
4. Bulk Copy - migrate / transform data from source SQL server db to Azure SQL db. There's a downtime for exporting data at source and importing at the target.
5. Azure Data Factory (ADF) - source SQL server db to target Azure SQL db. Cost is important consideration and is based on factors like pipeline triggers, activity runs and duration of data movements.
6. SQL data sync - synchronize data between source and target databases. Suitable to run continuous sync between Azure SQL db and on-premise / Azure SQL db. It can have higher performance impact depending on the workload.

#### Migration Steps

1. Discover -
2. Assessment
3. Migrate
4. Cutover
5. Optimize

#### A) Data Migration Assistant (DMA)

#### Scenario to choose for Azure SQL Managed Instance

Pools are well suited for a large number of databases with specific utilization patterns, for a particular database, the pattern is characterized by low average utilization with infrequent utilization spikes.

Conversely, multiple databases with persistent medium-high utilization shouldn't be placed in the same elastic pool.

The more databases, we can add into a pool, the greater the savings become. Depending on the app utilization pattern, it's possible, to see savings as few as two AWS S3 db.

Overutilization and underutilization of DTU usage can be overcome through Azure SQL elastic db. A elastic pool allows these unused DTUs to be shared across multiple databases. A pool reduces the DTUs needed and the overall cost.

The best size for a pool depends on the aggregate resources required for all databases in the pool

1. Maximum compute resources utilized by all databases in the pool. Compute resources are indexed by either eDTUs or vCores depending on the purchasing model.
2. Maximum storage bytes utilized by all databases in the pool.

#### Business Continuity for Azure SQL Elastic db

1. Point-in time restore - point-in-time restore uses automatic database backups to recover a database in a pool to a specific point in time.
2. Geo-restore - Geo-restore provides the default recovery option when a database is unavailable because of a incident in the region where the db is hosted.
3. Active geo-replication - For applications, which has more aggressive recovery requirements, than geo-restore can offer, we can configure active geo-replication or auto-failover group.

#### Azure SQL Managed Instance

##### vCore based purchasing model

A vCore represents a logical CPU and provides the option to choose the physical characteristics of the hardware (the no of cores, the memory, the storage sizes). The vCore based purchasing model gives us the flexibility, control, transparency of individual resource consumption and a straight forward way to translate on-prem workload requirements to the Azure platform.

vCore based purchasing model depends on -

- Service tier
- Hardware configuration
- Compute resources (the no of vcores and amount of memory)
- Reserved database storage
- Actual backup storage

Benefits of vCore based purchasing model used by Azure SQL managed instance

- Control over hardware configuration to better match the compute and memory requirements of the workload
- Pricing discounts for AHB and reserved instance
- Greater transparency in the hardware details which empowers compute, helping facilitate planning for migrations from on-prem deployments
- Higher scaling granularity with multiple compute sizes available.

Backup Storage

- Point in time restore - The storage consumption depends on the rate of the change of database and retention period configured for backups. We can configure a separate retention period of each database between 0 to 35 days for SQL managed instance. A backup storage amount is equal to the configured max data size is provided at no extra charge.
- Long term retention (LTR) - customers have the option to configure the long term retention of full backups for upto 10 years,

Feature	General Purpose	Business Critical
Scenario	Most standard business workloads. Offers budget-oriented, balanced and scalable compute and storage options	Offers business applications with highest resilience to failures by using several isolated replicas, and provides the highest I/O performance.
Read-only Replicas	0	1
HA replica	One replica is available	Three HA replica is available & one read-scale replica
Read-only replicas with failover groups enabled	One additional read-only replica and two total readable replicas, which includes the primary replica	Two additional read-only replicas, three total read-only replicas, four total readable replicas which includes the primary replica.

SQL Server on Azure VM (HADR configurations)

A Windows Server Failover Cluster is used for high availability and disaster recovery (HADR) with SQL Server on Azure VM.

Best Practices

- Deploy the SQL Server VM to multiple subnets to avoid the dependency on the Azure LB or a distributed network to route traffic to HADR solution.
  - Change the cluster to less aggressive parameters to avoid unexpected outages from transient network failure to Azure platform maintenance.
  - Place the SQL Server VM in a AS or different Azs.
  - Use a single NIC per cluster node
  - Configure the cluster quorum voting to use 3 or more odd numbers of votes.
- a) Cloud witness - it's ideal for deployments in multiple sites, multiple zones and multiple regions. Use a cloud witness for disk quorum whenever possible unless using a shared-storage cluster solution.
- b) Disk witness - it's the most resilient quorum option and is preferred for any cluster which uses Azure shared disks (like shared SCSI, iSCSI or fiber SAN). A clustered shared volume can be used for disk witness.
- c) Fileshare witness - is suitable when the disk witness and cloud witness are unavailable

SQL database auditing

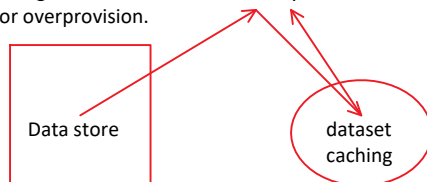
- Retain - an audit trail of selected events. We can define categories of database actions to be audited.
- Report - on database activities. We can use pre-configured reports and a dashboard to get started quickly with activity and event reporting.
- Analyse - reports, we can find suspicious events, unusual activity and trends.

#### Server level vs database level audit policy

- An audit policy can be defined for a specific database or as a default server policy in Azure SQL db.
- A server policy applies to all existing and newly created databases on the server.
- If server auditing is enabled, it always applies to the database. The database will be audited, regardless of the database auditing settings.
- When auditing policy is defined at the database level to a log analytics workspace or an event hub subscription destination, these operations will not keep the source database level auditing policy like
  - Database copy
  - Point-in time restore
  - Geo-replication

#### Performance improvement measures for Azure SQL db

1. Caching can improve app performance - database queries can be stored over the cache and return the queries much faster than the db when the app requests them. Not only this approach, helps for significant reduction for latency but also reduces the load on the database, lowering the need for overprovision.



- Determine whether the item is held in cache
- If the item is not available, currently in cache, read them from the data store
- Store a copy of the item in the cache

2. Caching are better than db at handling high throughput of requests - enabling the app to handle more simultaneous users.
3. Caches are typically most popular beneficial for read-heavy workloads where the same data is being accessed again and again. Caching pricing, inventory, session state or financial data are some of the examples.
4. Resolve the server index fragmentation with automatic tuning - keep the fast query performance is paramount for app with rdbms, one of the common cause for degraded performance is index fragmentation. With the indices, The SQL server can quickly locate the row with the data which user is requesting for, the first step is to identify the degree of fragmentation.
5. Resolve the fragmentation - there're three primary options for automatic tuning with Azure SQL db
  - a) CREATE INDEX - create new indices which can improve the performance
  - b) DROP INDEX - Drops redundant and unused indices (>90 days)
  - c) FORCE LAST GOOD PLAN - identifies queries using the last known good execution plan

#### MAXDOP configuration

- a) In Azure SQL database (both for single and elastic pool db), the default MAXDOP setting for each new single database and elastic pool database is 8.
- b) For Azure SQL managed instance, the max degree of parallelism instance option will be set to 8 by default.

This default prevents unnecessary resource utilization, while still allowing the database engine to execute queries faster using parallel threads.

- Azure SQL db level, MAXDOP can be controlled at the db level using MAXDOP database-scoped configuration.
- For Azure SQL managed instance, customers can also set the server 'max degree of parallelism' configuration option & can control MAXDOP at the resource governor workload group level.
- For all of Azure SQL deployment options, MAXDOP can additionally be controlled at the individual query level by using OPTION (MAXDOP) query hint where it actually overrides MAXDOP configurations set in the database or instance scope.

# Azure Data Factory

14 February 2023 09:31

ETL is the process for integrating and loading of the data for computation and analysis. It's also the primary method to process data for traditional data warehousing and BI applications.

### Benefits

1. Extract the data from the legacy systems
2. Cleanse the data to improve data quality and establish consistency
3. Load the data into the target database.

Azure Data Factory is a managed data orchestration and integration platform which helps to build complex Extract, Transform and Load (ETL) or ELT projects with data integration features.

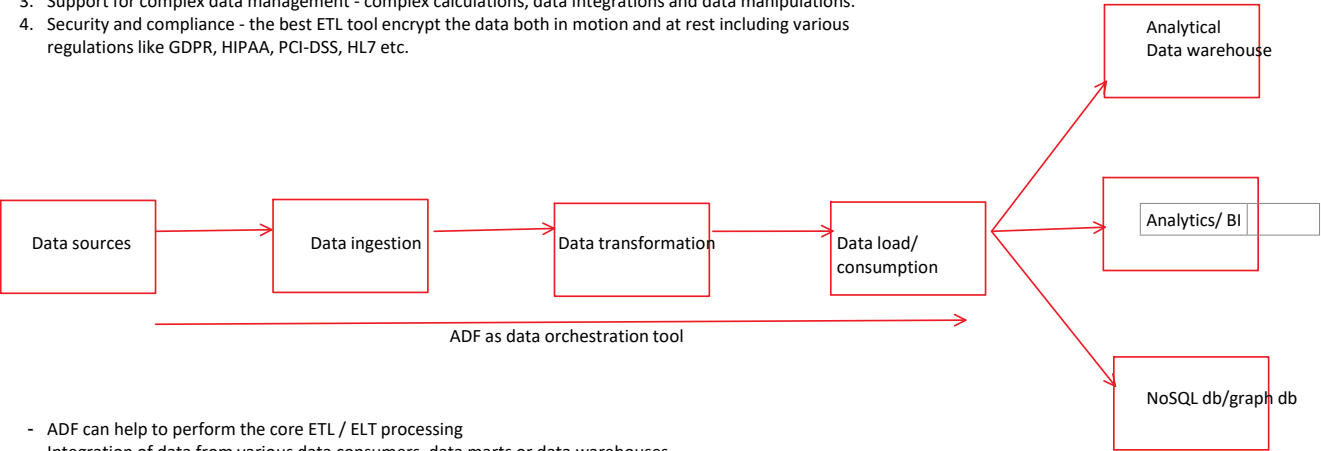
Data Orchestration is the practice of acquiring, cleaning and matching, enriching and making data accessible across the technology teams. The effective data orchestration captures data from several sources and unifies it within a centralized system, making it organized and ready to use for getting insights from data.

In ETL perspective, orchestration means automated management, co-ordination and management of complex data pipelines. ETL tools run on a schedule.

ETL (Extract, Transform & Load)	ELT (Extract, Load and Transform)
ETL can perform the end to end data transformation after loading the raw data into Transformation layer and captures the insights from the data.	ELT copies or exports the data from the data source locations, but instead of loading it to a staging area for transformation, it loads the raw data to the target data store to be transformed as required.
The ETL process involves more structured datasets and data has to be rational in nature, should have proper keys (PK, FK) to integrate the relationships between multiple tables	ELT is useful for high-volume, unstructured datasets as loading can occur directly from the data source.  ELT is more ideal for Big Data Analytics pipeline because it can clean, parse the unstructured, semi-structured data and can load the data directly into the native format. It's ideal for big data use cases.
In on-premise, ETL data pipeline is more common	On Cloud, Azure ELT pipeline is more popular

### Azure Data Factory (ETL/ ELT tool)

1. Comprehensive automation support - leading ETL tool (ADF) can automate the entire data flow, from data sources to the target data warehouse. Many tools recommend rules for extracting, transforming and loading of the data.
2. Visual drag-drop support - the functionality can be used for specifying rules and data flows.
3. Support for complex data management - complex calculations, data integrations and data manipulations.
4. Security and compliance - the best ETL tool encrypt the data both in motion and at rest including various regulations like GDPR, HIPAA, PCI-DSS, HL7 etc.



- ADF can help to perform the core ETL / ELT processing
- Integration of data from various data consumers, data marts or data warehouses
- Allows us to create the data driven workflows through ingestion, transformation and consumption
- Code free ETL tool, it helps to integration of data through ingestion, transformation and consumption through the data flows, control flows and scheduling of the ADF pipeline.

### ETL process

1. Extract

Raw data gets copied into / exported from source location to a staging area. Data management can extract the data from the variety of data sources that can be structured or unstructured.

- SQL / noSQL
- CRM/ERP/MDM
- Flat files



- Web pages

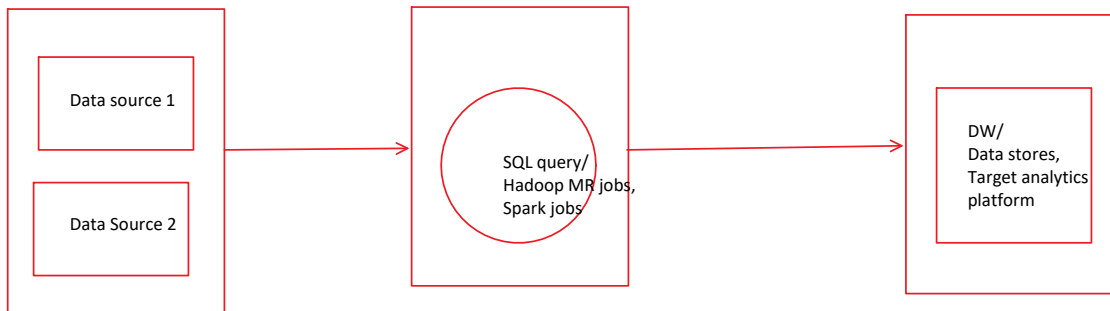
## 2. Transformation

In the staging area, the raw data undergoes the data processing, data has to be transformed and consolidated for its intended analytical use cases -

- Filtering, cleansing and de-duplicating, validating and authenticating the data
- Performing calculations, translations, summarizations of raw data
- Conducting audits to ensure data quality and compliance
- Removing, encrypting or protecting data governed by regulators
- Formatting the data into tables, joined to match the schema of the target data warehouse.

## 3. Load

In the last phase, the transformed data is moved from the staging area to a target data warehouse. Which involves loading of the data, followed by the periodic loading of incremental data changes and full refreshes to erase and replace the data in the data warehouse.



- The data transformation can take place usually involves various operations, such as filtering, sorting, aggregating, joining of the data, cleaning, deduplicating and validating the data.
- As part of offering, these three ETL processes can execute in parallel to save time. e.g. while the data is being extracted, and a loading processing can begin working on the prepared data. Rather than just waiting for the entire extraction process to complete.

## Core Components of Azure Data Factory

1. **Linked Service** - Creates a linking connection between the data source and the ADF pipeline.

We must create the linked service to link up the data source with the data factory.

e.g. db connection string which defines the connection information required for the service to connect to the external resource.

2. **Dataset** - It's the named view of the data which simply points to or references the data which required to be used in the ADF activities as inputs and outputs.  
e.g. SQL server db, files

3. **Activities** - The activities refers to the actions, jobs / tasks can be performed on the data. Activities also can produce a dataset & the datasets consume the activities.

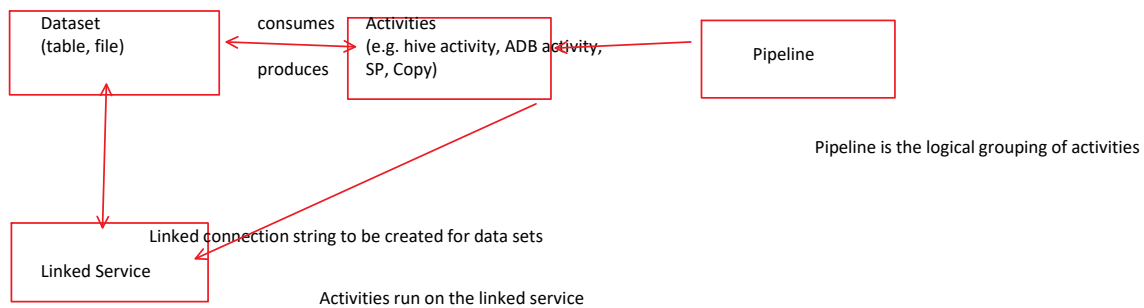
while copying/moving the data from Azure Blob storage to Azure SQL db, the storage and the linked service connection string need to be created.

- Activity can help us to copy the data
- Deduplicating the data
- Formatting the data removing all nulls, NaNs
- Applying normalization, remove redundancies
- Adding column headers and formatting

4. **Pipeline** - It's logical grouping of activities, which can be monitored, managed and scheduled. The activities in a pipeline define actions to perform on the data.

e.g. an ADF pipeline could contain the set of activities which can ingest and clean log data, then transform the mapping data flow to analyse the log data.

- The pipeline allows to manage the activities as a set instead of each one individually. We can deploy and schedule the pipeline instead of activities independently.
- The activities in a pipeline defines the action to perform on the data.



#### Hands-on Lab 01

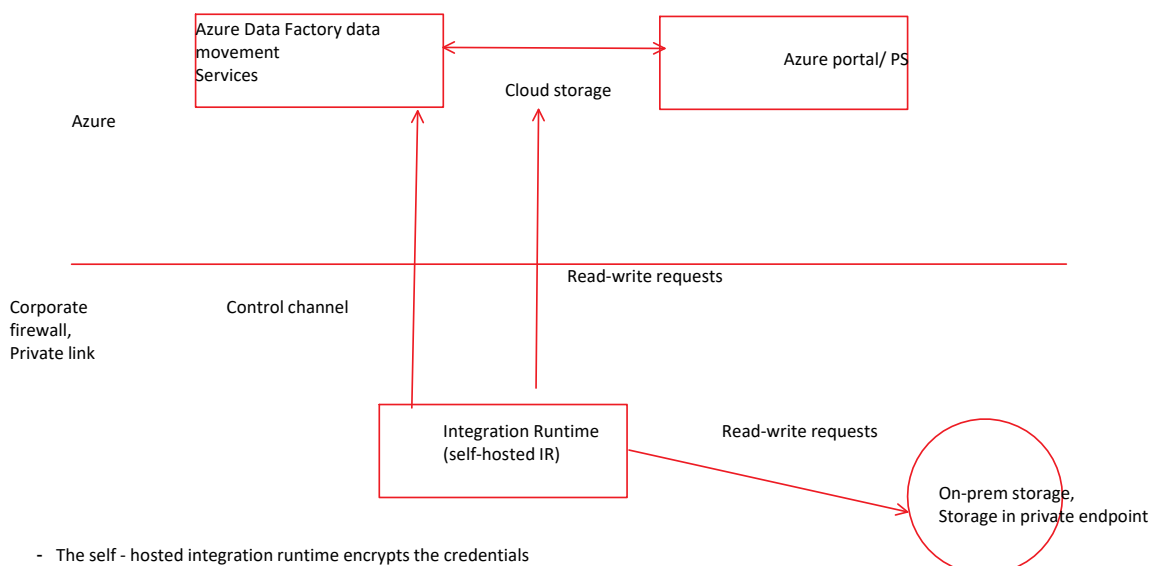
1. Create Azure Data factory
2. Explore the Data Factory Studio
3. Explore the ADF activities, pipelines and datasets, linked services
4. Monitoring and management of ADF

#### Hands-on Lab 02

1. Transferring data from Azure SQL db to Azure Blob storage using Copy activity and Copy data tool.
  - a) Provision Azure SQL db with sample db
  - b) Provision Azure blob storage
  - c) Use the copy data tool to copy data from Azure SQL db to Azure blob storage/ADLS gen2
  - d) Create the ADF pipeline
  - e) Monitor the pipeline

The Integration Runtime (IR) is the compute infrastructure used by ADF to provide the data integration capabilities across the different environments.

- Data Flow - Execute a data flow in a managed Azure compute environment.
- Data movement - Copy data across the data stores in a public or private networks (for both on-prem or Vnets). This service provides support for built-in connectors, format conversion, column mapping and performant / scalable data transfer.
- Activity dispatch - Dispatch and monitor transformation activities running on a variety of compute services like ADB, Azure HDI, ML Studio, Azure SQL db and SQL server.
- SSIS package extension - natively execute SSIS packages in a managed Azure compute environment.
- So, the integration runtime provides the bridge between the activities and linked services. It specifies the compute environment where the activity is to be performed in the closest region to the target data store or compute service to maximize the performance while allowing flexibility to meet security and compliance requirements.
- 
- 
- a) Azure - Data Flow, Data movement, Activity Dispatch (public and private endpoints)
- b) Self-hosted - Data movement, Activity dispatch (public and private endpoints)
- c) Azure-SSIS - SSIS package extension (public and private endpoint)



availability, the credentials are further synchronized across other nodes.

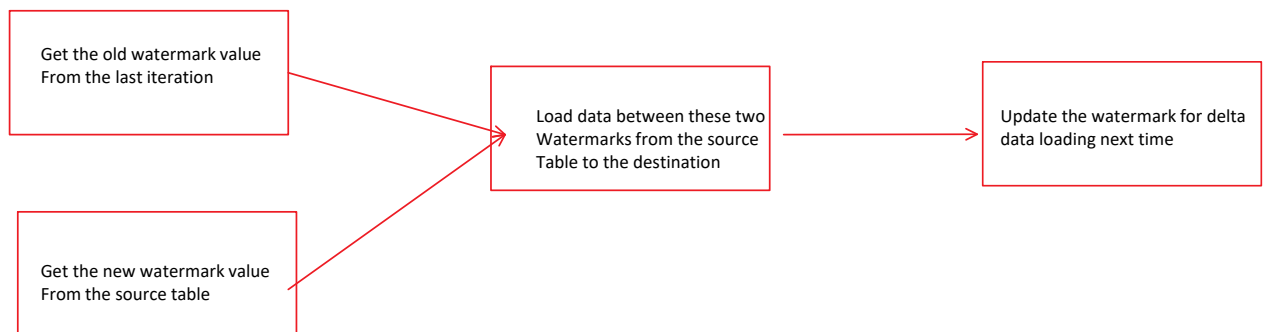
- The self hosted integration runtime can directly communicate over cloud based storage service like Azure blob storage over a secure HTTPS channel.

### Hands-on Lab 03

1. Create SQL Server db (on-premise)
2. Create a Blob storage
3. Provision over self - hosted integration runtime on ADF
4. Create the pipeline

Incrementally (or delta) loading of data after an initial full data load is a widely used context.

#### 1. Delta data loading from database/data store by using watermark

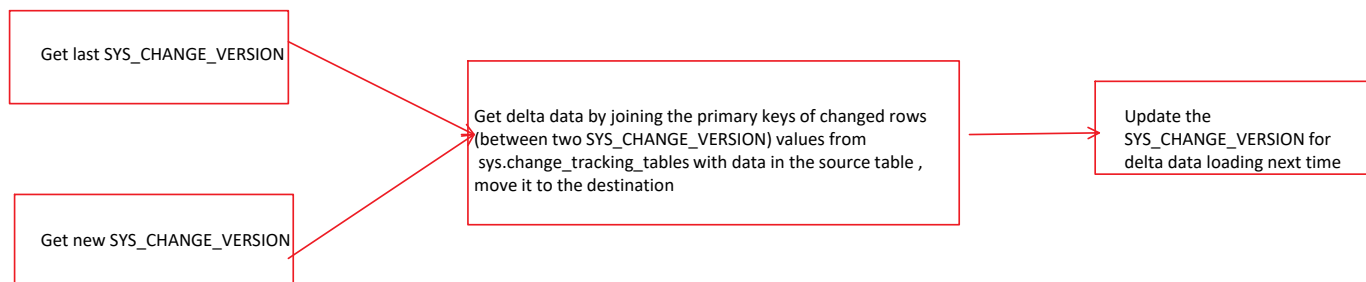


We can define a watermark in the source database. A watermark is a column which has the last updated timestamp or an incrementing key. The delta loading solution loads the changed data between the old watermark and a new watermark.

- Create two lookup activities - first lookup activity is to retrieve the last watermark value. Use the second Lookup activity to retrieve the new watermark value. These watermark values are passed to the Copy activity.
- Create a Copy activity which copies the rows from the source data store with the value of the watermark column > old watermark value and < the new watermark value. It also copies the the delta data from the source data store to Blob storage as a new file.
- Create a StoredProcedure activity which updates the watermark value for the pipeline which runs the next time.

#### 2. Delta table data loading from SQL db using The Change Tracking Technology

Change tracking technology is the lightweight solution in SQL/ Azure SQL db which provides efficient change tracking mechanism for apps. It enables an application to easily identify the data which was inserted, updated or deleted.



1. Initial loading of historical data (run once)
2. Incremental loading of delta data on a schedule
  - Get old and new SYS\_CHANGE\_VERSION values
  - Load the delta data by joining the primary keys of the changed rows (between two SYS\_CHANGE\_VERSION values) from sys.change\_tracking\_tables with data in the source table

- then move the delta table to destination.
- Update the sys\_change\_version for the delta loading next time.

#### Incremental Load

- Create two lookup activities to get old and new SYS\_CHANGE\_VERSION from Azure SQL db/ SQL db to pass it to copy activity.
- Create one copy activity to copy the inserted / updated/ deleted data between the two SYS\_CHANGE\_VERSION values from Azure SQL db to Blob storage
- Create one Stored Procedure activity to update the value of SYS\_CHANGE\_VERSION for next pipeline run.

#### 3. Loading of new & changed files by using LastModifiedDate & by using time partitioned folder or file name

ADF will scan all of the files from the data source store, can apply the filters over by their last modified data and to copy only the new and updated file since the last time to the destination store.

#### 4. Loading of new & changed files by using time partitioned folder or file name

Copy of new files/ folders can be possible through ADF, when they're partitioned with timeslice information as part of file or folder name

# Azure SQL Data warehouse

14 February 2023 09:31

# Azure Data Lake

14 February 2023 09:31

# Azure Databricks

14 February 2023 09:31

# Data Analytics

14 February 2023 09:32



