

# Pre-requisites

12 December 2022 19:53

<b>Lab Setup Requirement</b>	<b>Hardware</b>  CPU - Intel Core i3/i5/i7 processor, RAM - at least 8 GB HDD- 512 GB / 1 TB, OS - Windows 10 /8.1/11, MS Word/excel, PowerBI desktop (optional)
------------------------------	--

<b>Pre-requisites</b>
<b>Software -</b> 1. SQL Server 2016, 2017 or 2019 Enterprise/Developer edition, 2. SQL Server Management Studio/Azure Data Studio, 3. Visual Studio 2019/2022/VS code, (IDE) 4. A Valid Azure Subscription, Azure CLI, Storage explorer, Microsoft Azure Subscription, Git tools and GitHub account, Azure PowerShell, PowerShell ISE, Note: Linux environment VMs (Apache hadoop/big data) (Linux OS) Tool: Putty.exe Azure based Linux servers, Vmware player/ virtual box AWS Tools for VS code, GCP Tools for VS code.

Azure Subscription (trial)

- 12 months of free service + 30 days of 200 USD credit
- enterprise subscription

# Database Fundamentals and SQL Server BI

12 December 2022 19:54

Application metadata

MSSQL - 1433  
MySQL - 3306



1. User related attributes (which users, port, MSSQL - 1433, encryption, protocol (TCP))
2. SqlConnection
3. ADOConnection (ADO.Net)

1960 (IBM) - developed the integrated Management system which is based on hierarchical database model.

1970 (IBM) - the relational database model was developed by E.F Codd.

1980 (IBM) - developed the Structured Query Language (SQL). It is declared as the standard language for the queries by ISO (International Standard Organization) and ANSI (American National Standards Institute).

OOPs principles

- Abstraction
- Encapsulation
- Inheritance
- Polymorphism



```

Type Car = {
  Types: 'small' | 'medium' | 'heavy'
  Color: 'white' | 'black' | 'red'
  Full_efficiency: 'efficient' | 'non-efficient'
  Price: 'chaper' | 'moderate' | 'expensive'
}
  
```

Data Abstraction is the root principle of design.



1. Cohesion - represents a relationship within the module where a single class, has its well-defined purpose.

- It refers to what the class (a module) can do.

a) Low cohesion - low cohesion means a class which can perform a great variety of actions - broad, involves multiple operations.

Contracts - interface, abstract class, types

What: Data Modelling concepts using the contracts and concretions. (objects)

Benefits: Simply design API, simply data design, should Separate how the data/ abstract can be designed Declarative vs imperative

```

Staff Class
checkEmail()
sendEmail()
emailValidate()
PrintLetter()
  
```

b) High cohesion - means the class is to be focused on what it should be implementing. It includes only the methods related to the intension of the class.

Dependency Inversion - forms the pillar of object oriented programming to define as couple the software modules loosely so that one form of objects/ class while changing, doesn't affect others.

Polymorphism - Compile time / Static polymorphism (method overloading, operator overloading)  
Runtime / dynamic polymorphism (function overriding)

```

Staff Class
salary
emailaddr
setSalary(newsalary)
getSalary()
setEmailAddr(newEmail)
getEmailAddr()
  
```

2. Coupling - refers to the principle of how related or dependent two classes / modules are toward each other. For low coupled classes, changing something major in one class should not affect other.

## Coupling definition

In OOD, the coupling refers to the degree of the direct knowledge that one element/ object has of another element/object. How often the changes in class A forces the changes in class B.

1. Tight coupling - means when two classes often change together, when class A cant be changed directly because it's going to make changes in class B.

```
Class Subject {
Topic t = new Topic();    // subject class is tightly coupled the topic class
Public void letsRead()
{
t.understand();
}
}
```

```
Class Topic {
public void understand()
{
System.out.println("Tight coupling scenario");
}
}
```

2. Loose coupling - when two classes are just aware of each other , like class A and class B. Also, class B is expose through its interface, then class A and class B are loosely coupled.

e.g. Microservices application

```
Public interface Topic
{
void understand();
}
```

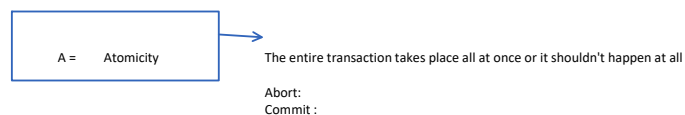
```
Class Topic1 implements Topic {
Public void understand()
{
System.out.println("This is loose coupling example");
}
}
Class Topic2 implements Topic {
Public void understand()
{
System.out.println("This is another loose coupling class");
}
}
Public class Subject{
Public static void main(String[] a)
{
Topic a = new Topic1();
t.understand();
}
}
```

e.g. microservices application design

Empaddr and emp classes are separate classes where change of a simple address of a emp entity / object does not affect to the emp class and its related methods.

High coupling scenario - it would make difficult to change or maintain the code, since classes are closely knit together , making a change could require an entire system revamp.

Best practice - good software design should always has **high cohesion and low coupling**.



## Different Types of Databases

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, MS Excel)
2. **Hierarchical database** -- based on hierarchical data model, where the data is viewed as a collection of tables, data is designed into a tree like strudture 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** - can consists of multiple parent segments and this segments can be grouped together as levels but there's always exists a logtcal association between the segments belonging to any level.
4. **Relational database** - consists of set of tables with columns and rows
5. **Object-oriented database** - information can be represented in the form of object-oriented programming. Inclined towards the objects e.g. multimedia records in a relational database can be definable data object.  
Mongo db is a object-oriented database.
6. **Distributed Database** - Consists of two or more files located in different sites/ location.  
e.g. SQL Server mirror databases,

7. **NoSQL databases** - non-relational db has support for unstructured, semi-structured data and it can include dynamic schema, flexible data model for faster data retrieval  
e.g. Mongo db, Cassandra, Couch db, Azure Cosmos db
8. **Graph database** - node - entity (rows in the table), attributes (relationship/columns)

e.g. Neo4j, Azure Cosmos db graph API



#### Advantages

- As the database is based on the hierarchical data models, the relationship between various layers are logically simple, it has a very simple hierarchical database structure.
- It has the data sharing facility since all data are held in a common database data and sharing of data becomes practically feasible
- It offers data security and integrity since it's based on parent-child relationship.

#### Disadvantages:

- Design is simple but implementation is complex
- This model also lacks flexibility as the changes to the new tables or segments often leads to very complex system management tasks.
- It has no standards as the implementation of the model does not provide any specific standard and limited to the relationship s which does not conform to 1:N format.



#### Advantages

- This model is simple and easy to design compared to hierarchical data model
- This model is capable of handling multiple types of data easily and we can access data easily, we can implement 1:1, 1:M, M:N relationships.

#### Disadvantages

- The schema of the network database mode is quite complex and records are maintained through pointers.
- The design of the network database mode is not user-friendly
- The model does not have any scope of automated query optimization

#### Data Dictionary or metadata in DBMS

A data dictionary is an integral part of a database. It holds the information about the database and the data which it stores named as metadata.

#### Characteristics of data dictionary

- A metadata is defined as the data about the data
- It's the self-describing in nature of databases
- It holds the information about each data elements in the database
- Such as names, types, ranges of values, access authorization, include the application details / program which uses the data
- Metadata is being used by programmers to develop the application, queries to manage and manipulate the data.

#### Types of Data Dictionary

1. Active Data Dictionary -- self updating
  - a) managed automatically by the data management software
  - b) It's always consistent with the current structure of the database (e.g. RDBMS)
2. Passive Data Dictionary - mainly for documentation purpose
  - Managed by user of the system and is modified manually by the user. (e.g. Dataedo by IBM IMS)

#### Active Data Dictionary



Data considered in DBMS for data Dictionary

#### Schema

- Tables
- Columns
- Constraints
- Foreign keys
- Indexes
- Sequences

#### Benefits of data dictionary

1. Improves the data quality
2. Spot the data anomalies
3. Implement transparency and collaboration
4. Get access to the good data
5. Involve regulatory compliance
6. Enables fast and accurate data analysis

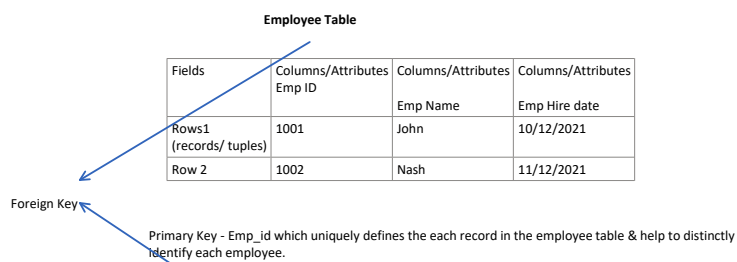
#### Programs in SQL

- Views
- Stored Procedures
- User defined Functions (UDFs)
- Triggers

#### Storage

- Size of tables and indexes stored inside the db
- Number of rows in table

OLTP - Online Transaction Processing (SQL server database engine / MSSQL, mysql, oracle)  
OLAP - Online Analytical Processing (SQL server datawarehouse)



#### Employee\_Address table

Fields	Attributes (emp_id)	Emp_primary_address	Emp_secondary_address
Row1	1001		
Row2	1002		

Primary Key = emp\_id

#### Primary Key definition

A primary key is a column or a set of columns in a table whose values uniquely identifies a row in the table. A relational database is designed to enforce the uniqueness of primary keys by allowing only one row with a given primary key value in a table.

#### Foreign Key definition

A foreign Key (FK) is a column or combination of columns which is used to establish and enforce a relationship between the data in two tables.

- A Foreign key is a column whose value corresponds to the values of the primary key in another table
- A foreign key is a column or a set of column in table whose values corresponds to the values of the primary key in another table.
- In order to add a row with a given foreign key value, there must exist a row in the related table with the same primary key value.
- Emp\_id the foreign key column of employee\_address table whose value corresponds to values of the primary key (emp\_id) of employee table

#### Integrity Constraints in SQL

- Constraints are the set of rules which enforce on the data to be entered into the database table. Basically, constraints are used to restrict the type of data which can be inserted into a database table.

#### Feature of Integrity Constraints

- The integrity constraints are one of the protocols which should be followed by the table's data columns. The constraints are generally established to restrict multiple types of information which can be entered into a table to ensure the integrity of data.
- Achieving the complete configuration of the constraints ensures that the data in the db is accurate and reliable to be consumed in the application.
- We can apply the integrity constraints at the column level or table level.
- The table level integrity constraints are applied on the entire table
- While the column level integrity constraints apply to the only one column.

Constraints can be defined in two ways:

- Column level: The constraints can be defined immediately after the column definition with the CREATE TABLE statement. So, these are called as the column-level constraints.
- Table level: The constraints can be defined after all the columns specified, with the ALTER TABLE statement. This is referred as the table level constraints.

SQL Server has the support for six types of constraints

- 1. Primary Key constraint** - The **primary key** is a set of one or more fields / columns of a table which helps to identify the records in the db table. It can not accept any null or duplicate values.

#### Features of Primary Key Constraint

- The primary key must have distinct values which means one of the columns of the table should have a unique value from the other.
- By default with the primary key, null values are not allowed in a primary key column of a table.
- Any table may have only a single primary key which can be made up of multiple fields.

Primary key constraint defined at the column level:

```
Create table table_name
(
Column1 datatype [CONSTRAINT constraints_name] PRIMARY KEY,
Column2 datatype
);
```

- 2. Unique Key Constraints** - A unique key is a set of one or more columns/fields which uniquely identifies each row / record in a table.

#### Features of Unique Key Constraint

- A table can have more than one unique key unlike the primary key
- Unique key constraint can also accept null values for a column
- Unique constraints are also referenced by the foreign key of another table. It can be used when developer wants to enforce unique constraints on a column and a group of columns which is not a primary key.

#### Students table - 1

Roll_no	Student_Name	Batch	Phone_no	Citizen_ID	Country_of_origin
001	Alan	01	212-334-2234	AB01	US
002	Matt	02	202-440-2335	CD04	Canada
003	Hans	03	304-223-2337		

Roll\_no is a primary key

**Citizen\_ID** , Country\_of\_origin are the unique key for the Students table

#### Student table -2

Roll_no	Student_Name	Batch	Phone_no	Citizen_ID	Country_of_origin
001	Alan	01	212-334-2234	AB01	US
002	Matt	02	202-440-2335	CD04	US
003	Hans	03	304-223-2337		Holland

Primary Key - Roll\_no

Unique Key - Citizen\_ID

#### Difference between Primary Key and Unique Key

Features	Primary Key	Unique Key
Basic	Used to serve as a unique identifier for each row in a table	Uniquely identifies a row which is not a primary key
NULL value acceptance	Cannot accept any NULL values	Can accept NULL values
Number of Keys can be defined in the table	Only one primary Key in a table	More than one unique key
Auto-increment	A Primary Key has the support for auto-increment value.	A unique Key does not support the auto-increment value
Modification	We cannot change or delete values stored in primary key	We can change the unique key values

IDENTITY (2,1)

Seed - initialization value (2)

Increment - increment value (3, 4, 5)

**Definition** - An index is a schema object. It's used by the db server to speed up the data retrieval or rows/records by using a pointer. It can reduce disk (I/O) by using a rapid path access method to locate data quickly.

#### Features of Indexes

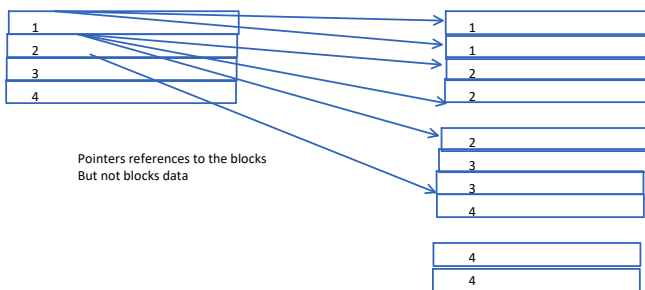
1. An Index can help to speed up the select queries and where clauses.
2. Indexes can be created or dropped with no effect on the data.
3. When an index is created, it includes a column contains a wide range of values.

Create index index\_name on table\_name column\_name;

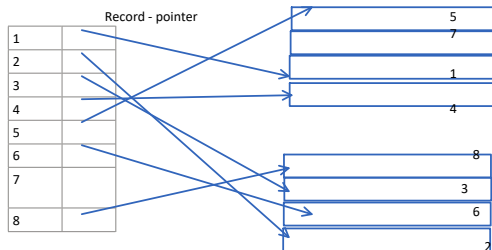
Table - table\_name  
Column - column\_name

Create index index\_name on table\_name (col1, col2...)

- a) **Clustered Index** - Clustered index is a kind of index which should satisfy the conditions
- The data or file, which moving to the secondary memory should be in sequential or sorted order,
  - There should be a key value, means it cant have any repeated value.
  - When applied clustered index on a table, it should perform sorting in that table only.
  - We can create only one clustered index in a table like primary key
  - Clustered index is as same as the dictionary where the data is arranged by alphabetical order.
  - In clustered index, index contains pointer to block but not direct data gets blocked.
  - We can have one clustered index on multiple columns and this kind of index is called composite clustered index.



- b) **Non-Clustered index** - index contains the corresponding pointer to the data



#### Difference between Clustered and Non-clustered index

Clustered index	Non-Clustered Index
Faster	Is slower
Required less memory for operations	Required more memory for operations
In clustered index, index is the main data	Index is the copy of the data
A table can have only one clustered index	A table can have multiple non-clustered index
Clustered index store pointers to block not data	Non-clustered index store both the value and a pointer to actual row which holds the data
It has the ability to store data on disk	Non-clustered index does not have inherent ability to store data on disk
Primary Keys of the table by default is considered as clustered index	Composite key / used with Unique key of the table defines the non-clustered index
A clustered index is 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 logical order of the index doesn't match physical stored order of the rows on disk
Clustered index size is larger	Non-clustered index size is smaller.

#### Surrogate Keys

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

#### Features of Surrogate Key

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

Student\_reg\_A

Reg_no	Name	% obtained
210101	Harry	80

210102	Matt	70
210103	Chris	84
210104	Lee	85

Student\_reg\_B

Reg_no	Name	% obtained
CS101	Maria	70
CS102	Simon	60
CS203	Sam	90

Surr_no	Reg_no	Name	% obtained
1	210101	Harry	80
2	210102	Matt	70
3	210103	Chris	84
4	210104	Lee	85
5	CS101	Maria	70
6	CS102	Simon	60
7	CS203	Sam	90

Item1
Item2
Item3
Item4
Item5
Item6
Item7
Item8

OFFSET 3 ROWS

FETCH 5 ROWS

#### Benefits

1. There's no direct information related with the table, so the changes are only based on the requirements of the application
2. Performance enhanced as the value of the key is relatively smaller
3. The key value is guaranteed to contain unique information
4. As it holds smaller constant values, the integration of the table easier
5. Enables to execute fast queries.



- a) Key Attribute - key attribute is used to represent the main characteristics of an entity. It represents the primary key.



- b) Composite attribute - composed of many other attributes altogether is called as composite attribute. The composite attribute is multiple simple attributes joined together.

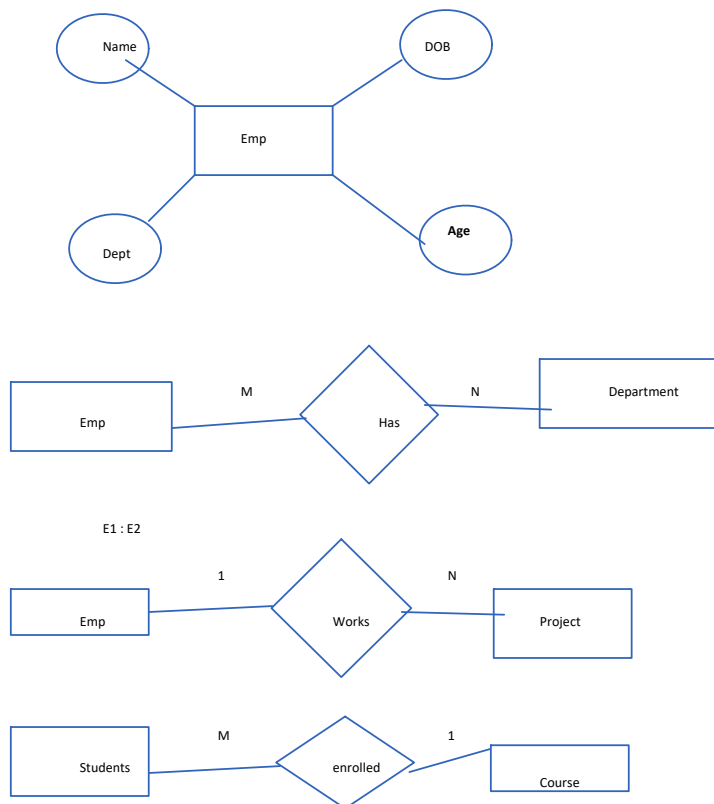


- c) Multi-valued attribute - this attribute can have more than one value. These attributes are called as multivalued attributed.



- d) Derived attribute - An attribute can be derived from other attribute,





#### Data types in SQL Server

##### Numeric data type in SQL server

Numeric	Storage space		
Tinyint	1 byte		
smallint	2 byte		
int	4 byte		
bigint	8 bytes		
Decimal(p,s)	5-17 byte		
Numeric(p,s)			
smallmoney	4 byte		
money	8 byte		
real	4 byte		
Float(n)	4-8 bytes		

P = precision (total number of digits can be stored both to the left and right of the decimal)  
 S = scale (max no of digits to the right of the decimal)

Decimal(8,3) allow the storage of total 8 digits and 3 of the digits to the right of the decimal or values.

##### Character data type

Char(n)	1 byte per character defined upto n
Varchar(n)	1 byte per character stored upto max 8000 bytes
text	1 byte per character upto 2 GB
Nchar(n)	2 bytes per character
Nvarchar(n)	2 bytes per character upto max 4000 bytes
ntext	2 bytes per character stored upto 2 gb

Nvarchar(max)

##### Date & time data type

datetime	0.00333 sec	8 bytes
Datetime2	100 nanosec	6-8 bytes
Date	1 day	3 bytes
time	00:00:00 to 23:59:59	3-5 bytes

##### Binary Data type

binary	Fixed with binary data	Upto 8000 bytes
Varbinary	Variable length binary data	Upto 8000 bytes

#### Column Property in SQL Server

##### IDENTITY property in SQL server

Identity property on columns can be defined to have a numeric data type. When the IDENTITY property is defined, SQL server manages the values in the columns on behalf of user.

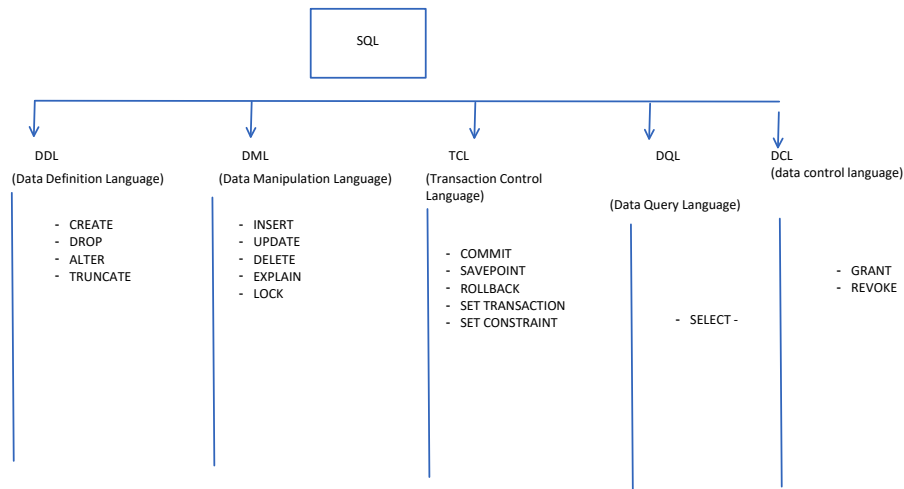
IDENTITY consists of two parameters

- Seed - seed parameter defines the first number which will be assigned when the data is inserted into the table
- Increment - defines the number which will be added to the previous value for every subsequent row inserted into the table.

IDENTITY(1,2) - it'll start the value with 1 and will increment by 2 everytime when a new row is inserted into the table.

Nullability constraint in SQL server

Nullability is the most common property assigned to a column, whenever the column is defined as NOT NULL, we're required to assign a value to the column, Whereas, when a column is defined as NULL, you don't have to assign a value to the column.



TCL

The transactions group a set of tasks into a single execution unit. Each transaction begins with a specific task and ends when all the tasks in the group successfully complete.

If any of the tasks fail, the transactions fail. Therefore, a transaction has only two results,

- Success
- Failure

MARS in SQL Server

The multiple active result sets is a feature with SQL server to allow the execution of multiple batches on a single connection. When MARS is enabled for use in SQL Server, each command object used adds a session to the connection.

- Applications can have multiple default result sets open and can interleave reading from them.
- MARS enables the execution of multiple db connection requests within a single connection. It allows batches to run. The database applications could not maintain multiple active statements in a connection.
- Applications can execute others statements (INSERT, UPDATE, DELETE, stored procedure calls) while the default result sets are open.

#### Benefit of schema

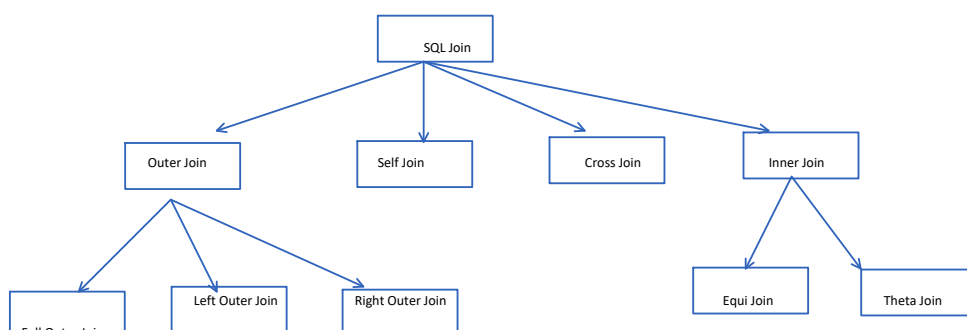
A database schema is considered the blueprint of a database object which describes how the data may relate to other tables or other data models.

The default schema for sql server is dbo.

But it's always recommended to create a customized schema to define how the different types of data (product, customer, customerAddress, employee) have been related to other tables / data models.

#### SQL Server Joins

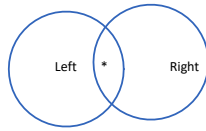
A SQL Join is a special form of generating a meaningful data by combining multiple tables relate to each other using a 'Key'. Typically, relational tables, must be designed with a unique column and this column is used to create relationships with one or more other tables. When you require a result-set that includes related rows from multiple tables, then the SQL joins are required on the column.



SQL Inner Join - The most simple and common form of join which is default join type.

- An Inner join combines two tables based on the criteria in the IN clause while also eliminating any rows from both tables which do not meet the criteria.

Cross - Join



Self join



Joining the table itself, is referred as self join

Normalization

- Benefits
1. Normalization is the process of organizing data in the database
  2. It's used to eliminate undesirable characteristics from a set of tables. Used to eliminate the errors/anomalies from insertion, update or delete.
  3. Normalization divides the larger table into smaller tables linking them with relationships.
  4. The normal form is used to reduce redundancy from the database table.

First Normal Form

- A relation will be 1NF if it contains an atomic value
- It can include a single value only in an attribute/field
- An attribute of a table cannot hold multiple values. It must hold single-valued attribute.

Second Normal Form

- All of the tables should be 1NF
- In the second normal form, all non-key attributes are fully functional dependent on the primary key

A B C - three variables

A → B A is dependent on B

B → C B is dependent on C

A → C (A is dependent on C), (transitive dependency)

Third Normal Form

- A relation will be called in 3NF, if it is already in 2NF, and not contains transitive dependency
- 3NF is used to reduce the data duplication.
- It's also used to achieve data integrity
- If there's no transitive dependency exists for non-prime attributes for a table, then the relation/table must be in third normal form (3NF).

A relation is in third normal form if it holds at least one of the conditions for dependency X → Y

1. X is super key
2. Y is prime attribute/ primary key, each element of Y is dependent on to some part of candidate/super key.

# SQL Data Warehouse

12 December 2022 19:54

OLAP - Online Analytical Processing platform (Data Warehouse)

1. Datawarehouse definition
2. Concepts on Data Marts
3. What are the Data Lakes? Why we should design a Data Lake?
4. Examples of data warehouse
5. Examples of Data Lake & Data Mart
6. Data Warehouse Architecture
7. Tables design in Datawarehouse
8. Star schema design in SQL Server db through SSMS (SQL Server Management Studio)
9. - Fact table
  - Dimension table
10. Define the concepts on Data Modelling
  - Star Schema (Demo/lab)
  - Snowflake Schema
11. Definition of data integration
12. OLAP (Online Analytical Processing) , benefits, use cases
13. Difference between OLAP( SQL Data warehouse) and OLTP (SQL database) database
14. Data Mining concepts

Tools -

1. SQL Server Management Studio (SSMS)

A Data-warehouse is a process of collecting and managing data from varied sources to provide a meaning business insights.

e-Commerce use case

Customer sample dataset

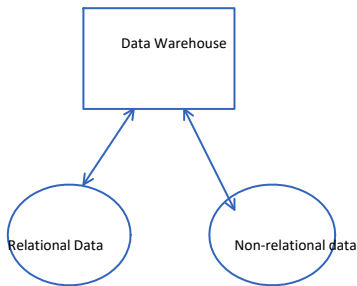
ID	Name	Location	Items Purchased	Time of purchase	IP address	Items kept in Cart
1	Matt	NJ	Laptop, Monitor	2022-9-6 17:00:00	120.34.23.5	Printer
2	Joe	NY	Headset	2021-12-01	192.168.1.1	Book

Problem Statement

- a) Define which customers have made maximum purchases over the last three months
- b) Define which customers have made minimum purchases (sales) over the last three months
- c) Define which locations customers visited maximum to the ecommerce website?
- d) Define how to make a recommendation of a product to customer which he/she can purchase next?
- e) Define which of items has been kept in shopping cart of the customers but never made purchase?

- Sales Analytics
- Sentiment Analytics
- Weblog Analytics
- Recommendation Engine

Retail domain  
Healthcare domain  
Banking and Financial domain  
Manufacturing domain



To ingest the relational / non-relational dataset into the SQL Data warehouse, the following analytical steps has to be performed.

Non-relational data - json, xml format

Real-time devices where data is extracted (IoT devices, weather sensors etc.)

**ETL - Extract - Transform - Load model in Data Analytics**

- a) Data Extraction

Tasks-

- Data Cleaning (remove all duplicate data)
- Data parsing (sort the data as per the analytical requirement)
- Identify the structuredness/ formatting of the data (into proper format, in terms of rows/columns design)
- Impose the data into the tabular structure (table format with specific attributes and records)

- b) Data Transformation

- Sql queries to start the analytical references (joins, functions, indexes, complex queries, views/stored procedures)
- Total no of customers who made the max purchase for products
- Total no of customers who made the min purchase of products
- What're the location(city/country) where customers visited max to the ecommerce website

- c) Data Load

The processed data can be stored into various data warehouses for future analytical purposes. Load the data into the SQL datawarehouse. Design further the schema of the dataset with Star and snowflake schema.

```
<xml>

<node1>
<node2> Name </node2>
<node3> Address </node3>
```

A XML data for Customer

file: Customer.xml (Extensible Markup Language)

```
<Customer>
  <Name> Matt </Name>
  <address> New York </address>
  <phone> 202-122-1222 </phone>
</Customer>
```

Customer.json (javascript object notation)

```
Customer:
  "name": "matt"
  "address": "New York"
  "phone": "212-122-1222"
```

file: Customer.json

#### Requirements of Data Warehouse

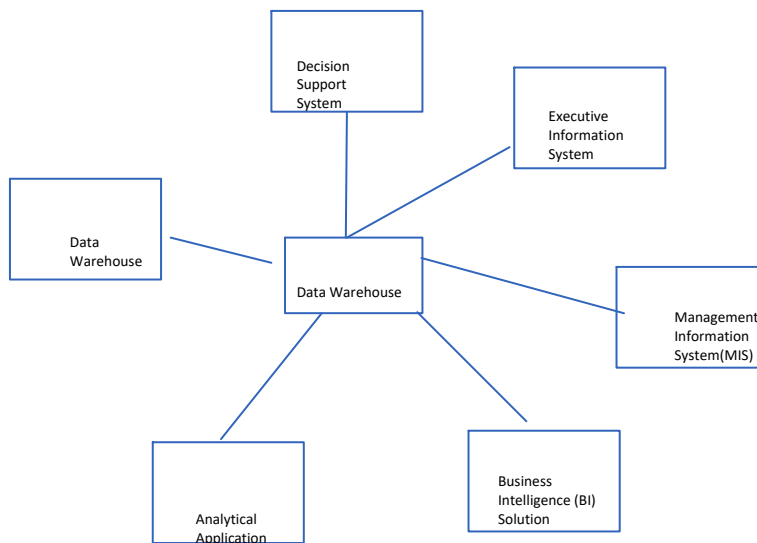
1. A Data warehouse (DW) is process of collecting and managing data from various sources to provide meaningful business insights.
2. A Data warehouse is typically used to connect and analyse business data from heterogenous sources. (data sources- traditional file system data, .csv, excel data, json/xml data, social media data, ERP/CRM dataset)
3. The data warehouse is the core of the Business Intelligence (BI) system which is built for data analysis and reporting.

#### Features of Data warehouse

- The data warehouse is maintained separately from the organization's operational database.
- We can call a data warehouse in some other names as well.

#### Examples of Data warehouse

- Informatica
- Vertica
- SQL Server Analysis Services (SSAS)
- Oracle Data warehouse
- Azure SQL Datawarehouse
- AWS Redshift
- GCP BigQuery



#### - How the data warehouse works

A data warehouse works as a central repository where the information arrives from one or more data sources. Data flows into the data warehouse from the transactional system and other relational databases.

Data may be

1. Structured
2. Semi-structured
3. Unstructured

- The data is processed, transformed and ingested so that users can access the processed data in the DW through the BI tools, SQL client tools and excel sheets.
- A DW also merges data coming from different sources into one comprehensive database.

#### - Types of Data Warehouse

1. Enterprise Data Warehouse (EDW)

Enterprise data warehouse (EDW) is a centralized warehouse. It provides decision support service across the enterprise. It offers a unified approach for organizing and representing data. It also provides the ability to classify the data according to the subject and gives access according to the division of the data.

e.g. SAP successfactor (employee organizational data)

2. Operational Data Store

Operational data store (ODS) is nothing but the data store required when neither the data warehouse nor the OLTP systems support organizations reporting requirements. In ODS, Data warehouse is refreshed in real time. It is widely preferred for routine activities like storing employee activities.

e.g. Salesforce CRM

### 3. Data Mart

A data mart is a subset of data warehouse. It specifically designed for a particular line of business such as Sales, Finance, Accounting, HR etc.

A data mart can collect data directly from various different data sources.

### Limitation of DBMS

1. Data volume , there 's a limitation. In Azure SQL db, maximum volume/size of the db can be 5 TB.
2. Traditional DBMS cant fulfill the requirements of Big Data (>100 TB, 500 TB, 100 PB, 500 PB)
3. A data warehouse s separate from DBMS, it stores a huge amount of data which is typically collected from multiple heterogenous sources like files, DBMS etc.
4. The goal is to produce statistical results which can help in decision making

e.g. a college might see the exams results, student performance analysis results.

### Need for Data warehouse

1. An ordinary RDBMS can store only in MB/GB amount of data and too specific purpose.
2. For storing in TB size, the data storage should be used Data warehouse,
3. A transactional database (e.g. RDBMS - SQL server db, mysql) doesn't offer to analytics.
4. To perform effective analytics, an organization needs to keep a centralized data warehouse to study its business by organizing , and using its historic data for taking strategies, decisions and analyse the trends.

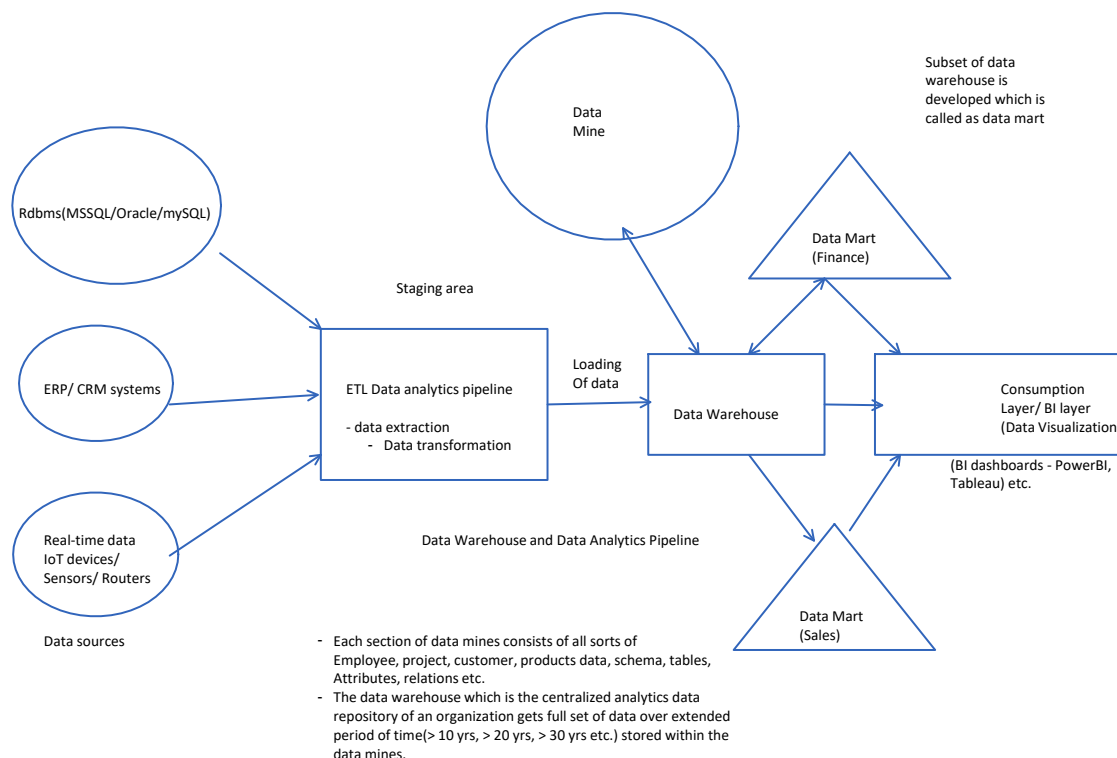
### Difference Between RDBMS and Data Warehouse

Features	RDBMS (OLTP)	Data warehouse (OLAP)
Purpose	A common transactional database on operational or transactional processing. Each operation is an indivisible transaction.	A data warehouse is based on analytical processing of data.  We can build data pipeline through ETL jobs. (Extract - Transform - Load)
analysis	The RDBMS system stores the current and upto-date data which is been used for daily operations.	A Data warehouse is integrated generally at the organization level by combining data from different databases.  Example - A data warehouse integrates data from one or more databases. So that analysis can be done to get the results such as top sales of the month of an ecommerce website.
	A database is generally application specific.  e.g. more data includes for app specific detail  Customer database contains various tables customer, customerAddress, customerPurchase , customerOrders etc.	A data warehouse is integrated generally at the organization level , by combining data from various databases.  e.g. A data warehouse integrates data from one or more multiple databases, so that the analysis can be implemented to get the results.  "top product purchased by customer over last six months" in case of retail/ecommerce domain.
	Construction of the database is not so expensive.	Construction of data warehouse can be extensive.

### Examples of Data warehousing

1. Social Media Websites - The social networking sites - LinkedIn, Twitter are based on large data sets.
2. Banking - Credit card holders information, Spending/purchase pattern of customers, stored with the centralized DW of the banks.
3. Governments - data warehouse can be used to store and analyse tax payments which used to detect the tax liabilities.

e-commerce, Retail, healthcare, marketing.



### Definition of Data Mining

The Data Mining refers to the knowledge mining from data, knowledge extraction, data / pattern analysis, data insights analysis. It's basically the process carried out for the extraction of useful information from the bulk of data or data warehouses.

The data mining is the result of extraction of the data patterns and knowledge after the data is processed through ETL jobs.

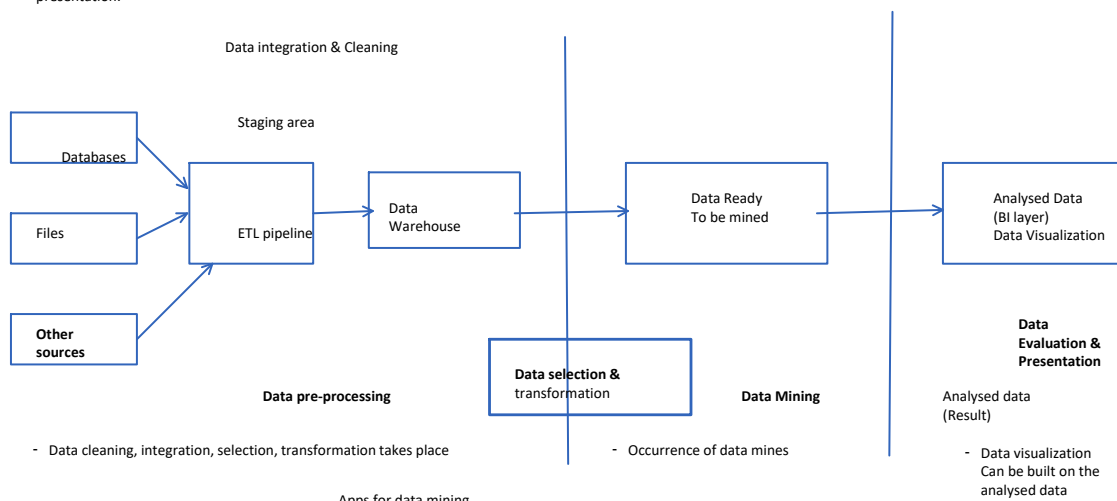
Technical data mining is the computation process of analysing data from different perspectives, dimensions, angles, categorizing/ summarizing it into meaningful information.



Data Mining can be applied to any type of data

- Data warehouses
- Transactional databases
- Relational databases
- Multimedia databases
- WWW (web apps)

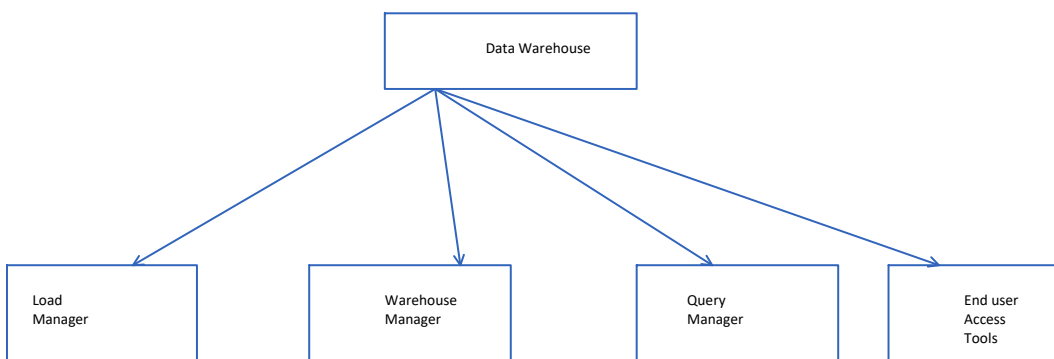
Data mining is a while process involving into data pre-processing, data mining and data evaluation and presentation.



Apps for data mining

- Financial Analysis
- Biological analysis
- Scientific analysis
- Fraud detection analysis
- Research analysis

- Components of Data Warehouse



- Load manager is called the front-end Component.

- It performs operations associated With the management of data into the data warehouse

- Query Manager is the backend

- Kinds of tools can be used

- Load manager is called the front-end Component.
- It performs all of the operations associated with the extraction and loading of the data into the data warehouse.

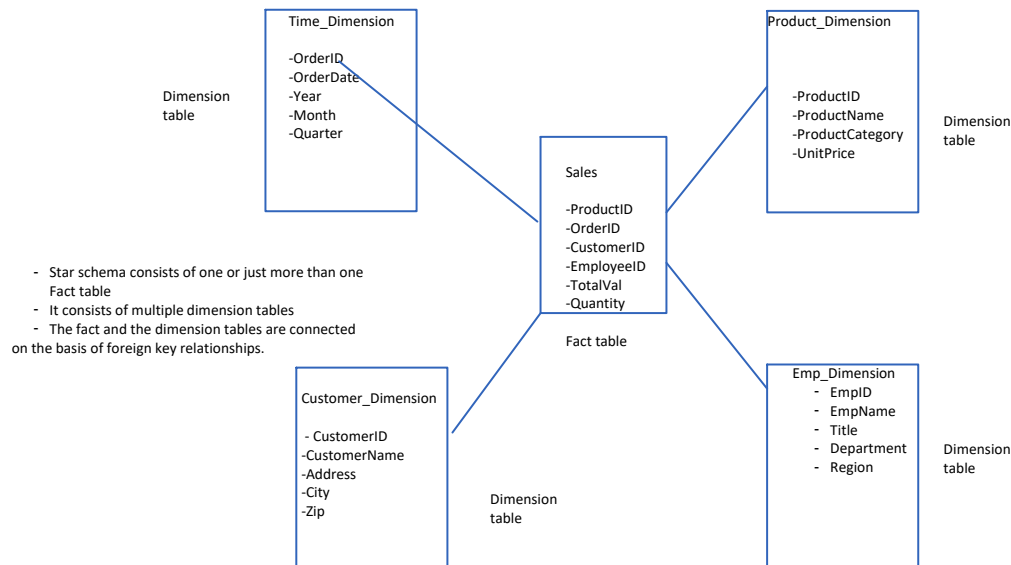
- It performs operations associated With the management of data into the data warehouse
- It performs the analysis of the data to ensure consistency , creation Of index, views,
  - Generation of normalization, aggregation, transformation & merging Of all source and backing up data.

- Query Manager is the backend Component of SQL DW.
  - Performs all the operations related to the management Of user queries.
  - Operations of the DW components are direct Queries to the appropriate tables for scheduling of the query execution.

- Kinds of tools can be used
  - Data Reporting tools (SSRS, Excel, PowerBI, Tableau etc.)
  - Query tools (SQL client tools)
  - App development tools (Visual studio)
  - OLAP tools / data mining tools (DB2, informatica)

#### Star Schema in Data Warehouse modelling

Star schema is the fundamental schema among the data mart and warehouse. It's simplest. This schema is widely used to develop and build a data warehouse and dimensional data marts. It includes one or more fact tables indexing any number of dimension tables.



- Star schema consists of one or just more than one Fact table
- It consists of multiple dimension tables
- The fact and the dimension tables are connected on the basis of foreign key relationships.

#### Features of Star Schema on data modelling

- In star schema, the business process data, which holds the quantitative data about the a business which is distributed in fact and dimension tables.
- The fact table consists of more number of columns / attributes
- The dimension are smaller in size compared to fact tables, contains more number of rows/records.

#### Advantages of Star schema

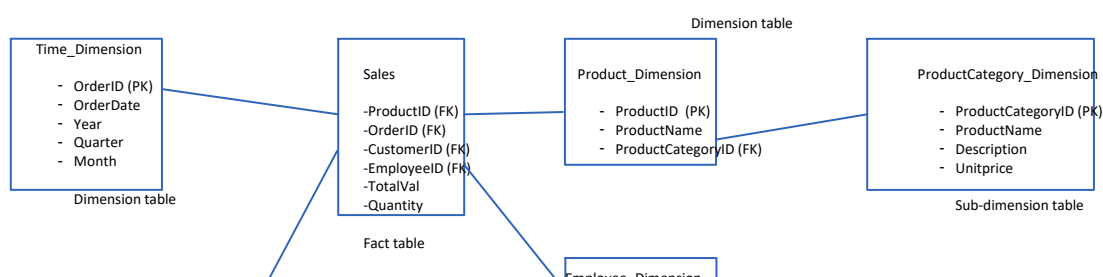
1. Simpler Queries - Join logic of star schema is quite simple and in comparison other SQL joining logic can be seen to fetch data from a transactional schema. The tables are highly normalized.
2. Simplified Business Reporting logic - In comparison to a transactional schema, which is highly normalized, the star schema makes simpler common business reporting logic like as reporting and ad-hoc analysis on data warehouse.
3. Build OLAP systems - Star schema is widely used in data warehouses (OLAP systems). In fact, the Star schema can help to deliver the major Relational OLAP model which can use star schema as a source of data.

#### Snowflake Schema

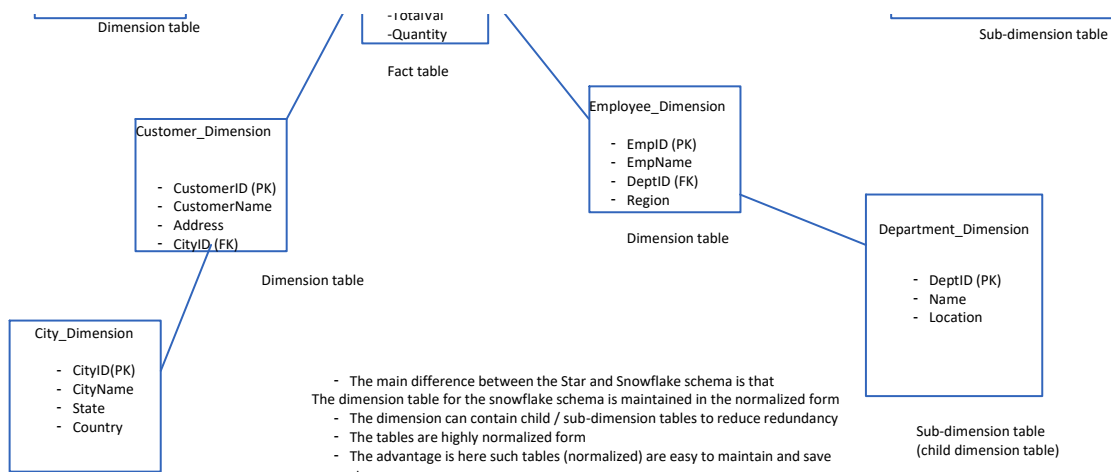
The snowflake schema is a variant of the star schema. The centralized fact table is connected to multiple dimensions. In the snowflake schema, dimensions are presented in a normalized form in multiple tables.

The snowflake structure is materialized when the dimensions of a star schema are detailed and highly structured, also having several levels of relationship, the child tables have multiple parent tables.

The snowflake schema effect affects only the dimension tables and does not affect the fact tables.







- The main difference between the Star and Snowflake schema is that the dimension table for the snowflake schema is maintained in the normalized form
- The dimension can contain child / sub-dimension tables to reduce redundancy
- The tables are highly normalized form
- The advantage is here such tables (normalized) are easy to maintain and save storage space.
- It also means more joins will be required to execute the query.

Sub-dimension table  
(Child dimension table)

#### Characteristics of Snowflake schema

- The tables are highly normalized compared to star schema, they take less disk storage space .
- It's easy to implement dimension which is added to the schema
- They are multiple dimension tables, so performance is sometimes reduced, querying is bit more complex.
- Structured data which reduces the complexity around the problem of data integrity
- It uses small disk space because the data is highly structured.

#### Difference between Star and Snowflake schema

Star Schema	Snowflake schema
In Star schema, the fact and the dimension tables are contained.	While in snowflake schema, The fact tables, the dimension tables as the sub-dimension tables are contained.
Star schema is top-down approach	While it's a bottom-up approach
Star schema uses more space	While, it uses less space because all of the dimension tables are normalized upto one or more sub-dimension tables.
It takes less time for the execution of queries.	While for snowflake schema, it takes more time for the execution of queries.
In star schema, normalization is not used	In snowflake schema, both normalization and denormalization is used
It's quite simple to design	While, design is complex compared to star schema
The query execution & complexity for star schema is quite low.	Query execution and complexity of snowflake schema is higher than star schema
It has less number of foreign keys because of less number of dimension tables	While, snowflake schema has more of foreign keys because of higher number of dimension tables
Star schema has high data redundancy	While, it has less data redundancy because all of the dimension tables are normalized to sub-dimension tables.

#### Difference between Fact and Dimension table

Fact table	Dimension table
Fact table contains the measuring of the attributes Of a dimension table.	Dimension table contains the attributes on that truth table which calculates the metric.
In fact table, there's less attributes than dimension table. It means fact table contains more number of Rows/ records.	While in dimension table, there's more attributes than the fact table. More number of columns are available in the dimension tables.
In fact table, there' more records than dimension Table,	There's less records, more columns than fact table
Fact table forms a vertical table	While, dimension table forms a horizontal table.
The attribute format of fact table is in numerical format and text format.	While, the number of dimension is more than fact table in a schema.
It's used mainly for analysis purpose and decision Making systems	While the main task of the dimension table is to store the information about a business and its process.

#### Hands-on Lab

##### Create Star Schema with SQL Server Management Studio

##### Tasks

1. Create a dedicated db for Star Schema
2. Create Database diagram support for the db
3. Create new database diagram
4. Create the data types for Fact and dimension tables (Under Database -> Programmability -> Types -> User-Defined-data-type)
  1. dbo.udt\_surrogate\_key (int, not null)
  2. dbo.udt\_business\_key ( nvarchar(20), not null)
  3. dbo.udt\_dollar\_amount (money, not null)
  4. dbo.udt\_quantity(int, not null)

5. Create Dimension tables
6. Create Fact tables
7. Create Build out relationships between Fact and dimension tables.

1. Click on empty space & select "Select All"
2. Right click on any table & click "Table view" & choose the option "Name only"
3. Right click on the any of the table and click on option "Authorize selected tables"
4. Right Click on the empty space & click on "arrange tables"

#### Dimension tables

1.dim_date	primary_key- (date_key)
2. dim_individual_customer	PK - (individual_customer_key)
3. dim_territory	PK - (territory_key)
4. dim_store	PK - (store_key)
5. dim_employee	PK - (employee_key)

#### Dimension Data Modelling

Dimension data modelling is comprised of fact and dimension tables. The main objective of the dimension data modelling is to improve the data retrieval so it is optimized for SELECT operation.

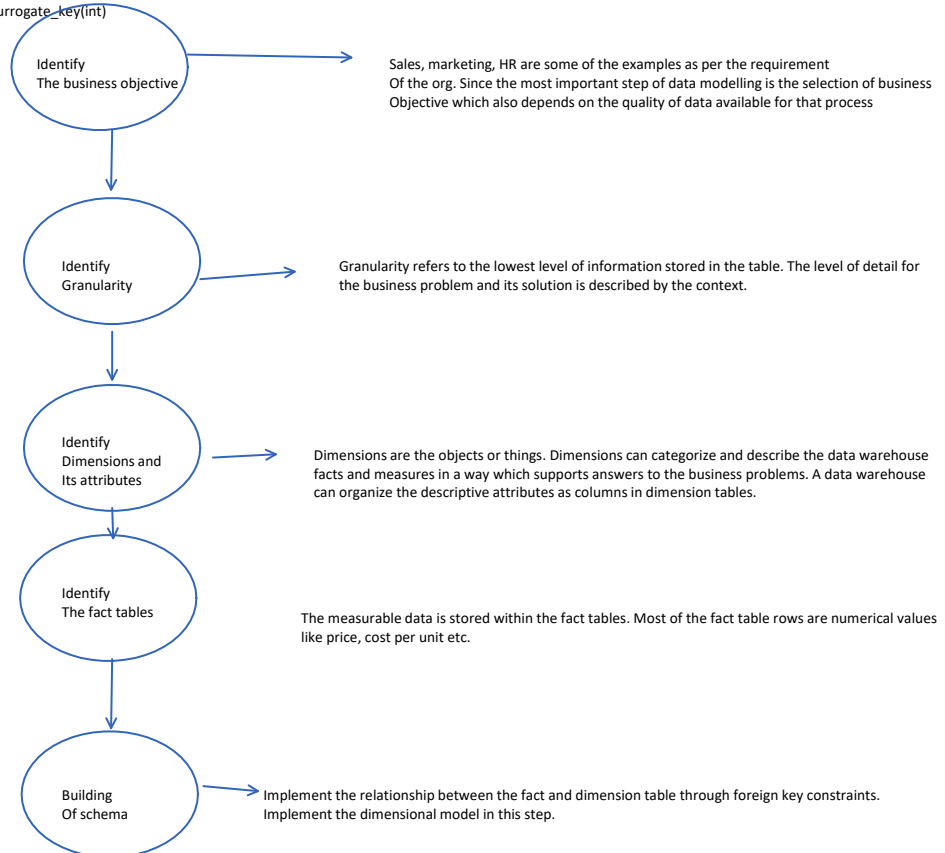
Advantage - we can store the data in such a way so it's becomes easier to retrieve the data once stored in the data warehouse.

Dimensional model is the data model used by many OLAP systems.

Steps for creating the dimension data model

Fact table (fact\_sales\_order)

date_key	udt_surrogate_key(int)
territory_key	udt_surrogate_key(int)
employee_key	udt_surrogate_key(int)
store_key	udt_surrogate_key(int)
individual_customer_key	udt_surrogate_key(int)



#### Measures

Measures in Data warehouse are the set of aggregates which we can calculate, such as sum of total orders placed by the customer.

#### Measures example

Qualitative - productID

Quantative - like the price of the product

The source of all OLAP structures is a table within a data source. OLAP has two terms to refer to the

source tables.

- Fact table - is used to define the measures
- Dimension table - stores the data used to define dimensions

In data warehouse, the dimensions are pieces of data which allows us to understand and index measures in the data models. Dimensions are either characteristics of a measure or pieces of data which helps to contextualize the fact.

- Dimensions example

-

- Product which was sold online
- Product color
- Product name
- Name of the customer purchased the product
- Name of the employee sold the product
- Store location
- Warehouse location

- **Dimensions**
- **Attributes**
- **Measures**
- **Hierarchies**

1. Dimensions define the basic business attributes like we want to analyse. It refers to customers, products, and time.
2. The attributes within a dimension define the columns within a dimension which are used for analysis like CustomerName, ProductName, City, PortalCode and OrderDate.
3. Measures are the set of aggregates which we can calculate such as sum of total orders, number of orders.
4. Hierarchies allows us to define a navigation structure within a dimension such as
  - Customers within cities within states within countries
  - Calendar months within the calendar quarter within the calendar years.
  - Fiscal months within the fiscal quarter within fiscal year

Cube - A cube contains of the analysis objects you define with the highest level of security that can be assigned. Within SQL Server Analysis Services, we can multiple cubes for the users within the SQL Server analysis services database.

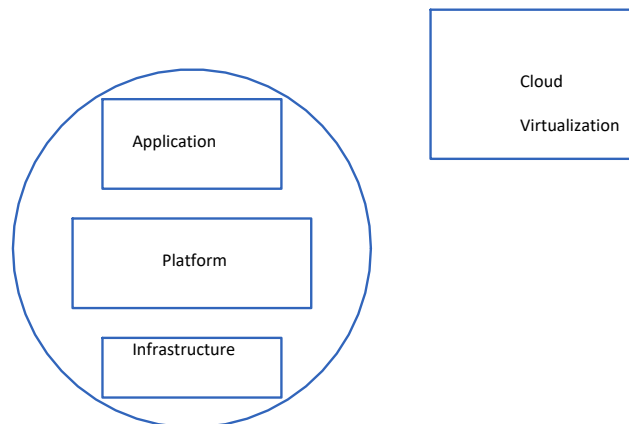
# Cloud Fundamentals

12 December 2022 19:54

1. Cost efficiency
2. Security
3. Flexibility
4. Mobility
5. Insights
6. Increased collaboration
7. Quality Control
8. Disaster Recovery (infrastructure, platform, application, database, network)
9. Loss Prevention (any data deleted should be able to recover it)
10. Automated Software Updates
11. Sustainability

By definition, the Cloud computing refers to the internet based computing. Cloud computing refers to the components and sub-components required

- Front-end (thin client (web apps), thick client (desktop apps))
- Backend (servers, storage)
- Cloud based network (internet, intranet & intercloud)



## Benefits of Cloud Computing

1. Cloud is easy to access for everyone. All the resources servers, disk, network, IP address, storage, databases, data warehouses etc. all accessible to everyone globally through internet.
2. Developing new applications, services, storage, databases are easier to build onto cloud.
3. It's useful to migrate the existing applications, databases, storage to the cloud as per supportability of the cloud vendor.
4. Software on demand.
5. Cost effectiveness - **Pay as you go**. (pay by hour)
6. Analysis of data
7. Streaming of high quality OTT resources (audio/video streaming)

Cloud Migration - Capex (capital expenditure) -> Opex (operational expenditure)  
(on -prem env) (cloud env)

## Features of Cloud

- a) Scalability - increase the number of instances or sizes of server, storage or databases on cloud.
  1. - **horizontal scalability ( add more server instances)**
    - 1.1 Scale out - add more server instances
    - 1.2 Scale in - reduce the number of server instances
  2. **Vertical scalability (add more compute capability - more CPU core, memory, disk storage etc.)**
    - 1.1 scale up - add more compute capacity to the server instances
    - 1.2 scale down - reduce the compute capacity to the server instances

As per the Cloud design principles, **horizontal scalability (scale out) is recommended to go for in business.**

- b) **SLA - Service Level Agreements** - 99.99% of the availability will be there all the time of the Azure SQL db instance.
  - there will be only 0.001sec of downtime per sec/week/month basis will be there.
- Associated with each and every cloud provider (Amazon web services, Microsoft Azure, Google Cloud platform, IBM cloud, Oracle Cloud)
- Cloud provider can compensate for any downtime for any missed SLA / guaranteed uptime.
- This is the part of cloud based business model for the cloud provider to their customer

## Different types of Cloud

Public Cloud  
Private Cloud  
Hybrid Cloud

## Cloud Models

IaaS (Infrastructure as Service)  
PaaS (Platform as Service)  
SaaS (Software as Service)

## Different types of Cloud Computing Platform

	Public Cloud	Private Cloud	Hybrid Cloud
Definitions	Public cloud is open to all to store and access information via the internet, using pay-as-you-go model. In Public cloud, computing resources are managed and operated by the Cloud Service provider (i.e. Microsoft Azure, Amazon Web Services, Google Cloud platform etc.)	Private cloud is known the organization's internal cloud. The applications deployed over the private cloud are accessible only within the private network.	It's the combination of both public and private cloud.  Hybrid cloud = public cloud + private cloud
Features	Public cloud is secure with the cloud controls (policy, Certificates, keys/passwords)	Private cloud is the most secure cloud platform since the applications are being deployed into the organization's own server and datacenters.	Hybrid cloud is partially secure because the services which are running on the public cloud can be accessed by anyone. While the services which are running on a private cloud can be accessed only by the org's users.
Benefits	Public cloud can be owned at a lower cost than the private and hybrid cloud	Private Cloud provides high level of security and privacy to the users	Hybrid cloud is suitable for organizations which require more security than the public cloud
Benefits	Public cloud is maintained by the cloud service provider (CSP) and users do not need to worry about the maintenance of the cloud	Private cloud offers better performance with improved speed and space capacity.	Hybrid cloud helps to deliver new products and services more quickly
Benefits	Public cloud is easier to integrate, offers a better flexibility approach to the customers, public cloud is delivered through internet, apps, databases, virtual machines are accessible over the internet.	Each organization has their own private cloud where they have full control over the cloud because it's managed by the organization itself.	Hybrid cloud offers flexible resources because public cloud and secure resources because of private cloud.
Cons	Public cloud is less secure because resources are shared publicly	Private cloud is accessible only within the organization level, the area of operations in private cloud is limited.	Security feature is not good as private cloud.  Managing a hybrid cloud is complex because it's difficult to manage more than one type of deployment model.
cons	Performance of the public cloud resources Depends upon the high-speed network linked to the cloud provider.  The client no control to the data underlying to the infrastructure	Private cloud is not suitable for organizations which has a high user base and organizations that do not have the pre-built infra, sufficient developers and can manage the cloud platform	In the Hybrid Cloud, the reliability of the services depends on the cloud providers
examples	Microsoft Azure, AWS, GCP, IBM bluemix, Oracle Cloud	HP data centers, Azure Stack, Vmware private cloud, Ubuntu canonical private cloud	Openstack, Azure hybrid cloud



## Different types of Cloud Models

	Cloud Model	Benefits	Examples
IaaS (Infrastructure as Service)	Rent for respective infrastructure required for business.  End users are responsible to manage the underlying infrastructure  - Compute - Cpu	Less complex Less development cycle	Azure Virtual Machine Azure Virtual Network Azure Disk

	<ul style="list-style-type: none"> <li>- Memory</li> <li>- Disk storage</li> <li>- Containers</li> <li>- network</li> </ul>		
PaaS (Platform as a Service)	<p>Managed services (Compute, storage, database, data warehouse, analysis services)</p> <p>-- end users are not responsible to manage infrastructure in PaaS</p> <p>To manage</p> <ul style="list-style-type: none"> <li>- Compute</li> <li>- Disk storage</li> <li>- Network</li> <li>- OS</li> <li>- OS update/patch management</li> <li>- Security of the infra managed by the cloud providers</li> </ul> <p>Only responsible to manage their application on Azure</p>	<p>Less burden on infra provisioning</p> <p>Developers can focus more into their application or business logic</p> <p>Development language agnostic (.net, java, springboot, react/angular, python, ruby etc.)</p>	<p>Azure App Service (Web apps)</p> <p>Azure Storage</p> <p>Azure SQL database</p> <p>Azure Hadoop Services</p> <p>Azure Analysis Services</p> <p>Azure Data Lake Services etc.</p>
Software as Service (SaaS)	<p>Managed services,</p> <p>End users/customers are not responsible for the application development and deployment.</p>	<p>Almost zero burden on the end user in terms of app development.</p> <p>No burden on app scalability, reliability, disaster recovery, app backup or restore</p>	<p>Office 365</p> <p>Gsuite gmail, Google drive</p> <p>Salesforce</p>

#### Hands-on Lab

#### Creation of Azure Storage Account

1. Login to Azure portal (<https://portal.azure.com>)

#### Azure IaaS, PaaS & SaaS

#### Azure IaaS

IaaS - Infrastructure as Service (IaaS) is a type of cloud computing service which offers compute, storage and networking resources on demand on a pay-as-you-go basis.

#### Features

1. Migrating the organization's infrastructure to an IaaS solution helps to reduce the maintenance of on-premise data centers. It helps to save money on hardware costs, and gain real-time business insights.
2. IaaS solutions give us the flexibility to scale the IT resources up and down with demand.
3. They can help to quickly provision new applications and increase the reliability of the underlying infrastructure.
4. IaaS lets us to bypass the cost and complexity of buying and managing physical servers and datacenter infrastructure.
5. A cloud computing service provider like Microsoft Azure manages the entire infrastructure, when you can purchase, install, configure and manage your own software - operating systems, middleware and applications.

#### Advantages of Azure IaaS -

1. **Reduce the capital expenditure and optimizes the cost** - IaaS can eliminate the cost of configuration and managing a physical datacenter, which is a cost-effective choice of migrating to the cloud. The pay-as-you-go subscription model used by IaaS providers helps to reduce hardware costs and maintenance. It enables the IT team to focus on the business.
2. **Increases scale and performance of IT workloads** - IaaS helps to scale globally and deliver the applications globally.
3. **Increases stability, reliability and supportability** - With IaaS, there's no need to maintain and upgrade software and hardware or troubleshoot equipment issues. With the appropriate SLA, portal service, we can get 99.99% of uptime of cloud infra resources (VM, Azure Virtual Network etc.)
4. **Enhance security** - appropriate SLA (99.99%) in Microsoft Azure offers better security for the applications and data we can achieve.

Azure IaaS components - Azure Virtual Machine, Azure Virtual Network, Azure Disk Storage etc.

#### Azure PaaS (Platform as a Service)

Platform as a service (PaaS) in Azure is a complete development and deployment environment in the

cloud. With resources, it enables us to deliver the applications from single cloud based apps to cloud enabled enterprise applications.

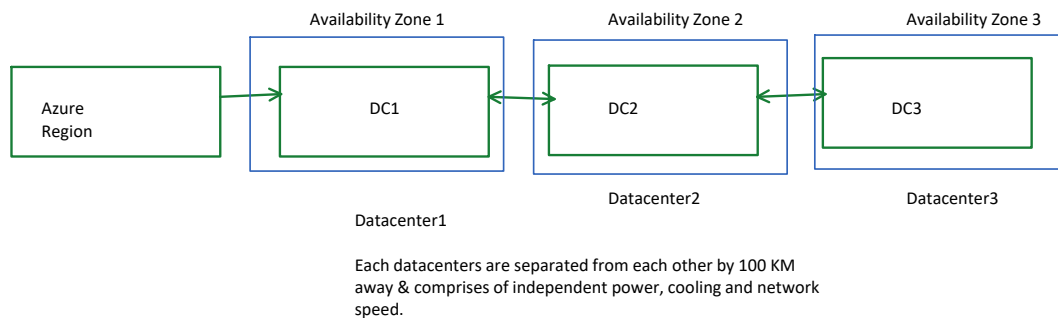
When we purchase the resources from a cloud service provider (CSP), on a pay-as-you-go model, can access them over a secure internet connection.

In PaaS, we focus only on the application development and deployment. Azure PaaS is designed to support the complete web application lifecycle - building(development), testing, deploying, managing and updating.

PaaS also allows us to avoid the expense and complexity of buying and managing software licenses.

#### Advantages of Azure PaaS

1. Cut the coding time - PaaS development tools can cut the time it takes to code new apps with pre-coded app components built into the platform - workflow, directory services or security features etc.
2. App development on Azure PaaS - PaaS can give the development team new capabilities without requiring new staff & having new skills.
3. Develop for web , mobile apps and develop new apps cross - platform - develop and deploy applications (web and mobile) apps much faster & deploy on Azure PaaS services.



All of these physical buildings (datacenters) are comprises of Azure region.

e.g. Central India (Pune) , South India (Chennai), Southeast Asia, East Asia

**Resource :** A resource in Azure is a manageable item which is available through Azure. Virtual Machines, Storage accounts, web apps, databases and virtual networks are examples of resources. Resource groups, subscriptions, management groups are also examples of resources.

**Resource Group:** A logical container which holds the related resources for an Azure solution. The resource group includes those resources which we can manage as a group. We can decide which resource belongs in a resource group based on what is required as per org policy.



#### Features of Azure resource group

- A resource group can be created at a specific Azure region

- Resource groups can be created, deleted, managed within the Azure portal.
- Azure resource groups can be used for administrative actions to manage the Azure resources.

Resource provider - A service which supplies the Azure resources.

All Azure Virtual Machine are part of "Microsoft.Compute" resource provider

All storage account resources are part of "Microsoft.Storage" resource provider



# Azure Fundamentals

12 December 2022 19:55

## Blob storage types -

Block Blob	Mainly used for storing of any objects data, like flat files, audio/video, media streaming, batch files, excel/csv/xml/json files
Page blob	Contains the virtual hard disk images
Append blob	Add new blob, made up blocks like block blobs. Append blobs are optimized for append operations.

## Azure Storage types

Blob storage	Store any kind of object based storage
Queue storage	Store randomized messages
Table storage	Store the data files which are not RDBMS compliant
File Share	Store any on-premise windows/linux compatible file systems
Disk storage	Store any OS compatible disk drives (VHD, VHDX)

# Basics of PowerShell Scripting

12 December 2022 19:55

Features	Programming Language	Scripting Language
Mode of execution	Programming language need to transform the high level Programming code into machine languages(byte code) through a compiler.	Compilation step is not required for the scripting language because they use an interpreter
	Compiler used in the programming language compiles the whole bunch of code at a time	Interpreter executes code line by line
Basic Difference	Objects are consisting of data and code,	Scripting languages are of functions & actions. The scripts are designed for a specific runtime environments.
Interpretation mechanism	Objects Oriented (OOPS) employ a compilation step.	Scripting language executes scripts inside the another parent language while OOPS languages like C, C++ are executed without a parent program.
Platform	OOPS can be executed in various platform (cross-platform)	Scripting languages are specific to the platform
Translation	In OOPS, a complete chunk of code is compiled and converted into machine level byte code	Scripting language are not self-executable and need to be encapsulated in a host to execute the scripts.
examples	Java, C#, Go	PowerShell, Shell scripting, Python, javascript, PHP

## Features of PowerShell

1. PowerShell Remoting - PowerShell remoting allows scripts and cmdlets to be invoked on a remote machine.
2. Background jobs - it helps to invoke script or the pipeline asynchronously, we can run the jobs either on the local machine or multiple remotely operated machines.
3. Transactions - enable the cmdlets and allows the developers to perform transactions
4. Network file transfer - PowerShell offers native support for prioritized, asynchronous, throttled, transfer of files between the machines using the background intelligent transfer (BITS) technology
5. Evening - This command helps us to listen, forwarding, and acting on management and system events.

## PowerShell cmdlets

A cmdlet is called Command let, it's a lightweight command used in the Windows base PowerShell environment. PowerShell invokes these cmdlets in the command prompt. We can create and invoke cmdlets command using PowerShell APIs.

Cmdlets	Command
Cmdlets in PowerShell are .NET framework class objects, so It can't be executed separately.	

Cmdlets can be constructed from a few lines of code	
Parsing, output formatting, and error presentation are not handled by the cmdlets	
Cmdlets are record based so that it processes a single object at a time	
Get - to get action Start- to run action/execution Out - to output something Stop - to stop something Set - to define something New - to create something	Get-Help

#### Concepts on PowerShell

Cmdlets	Cmdlets are the built-in commands written in .net language like C#, VB. It allows the developers to extend the set of cmdlets by loading and write powershell scripts
Functions	Functions are commands which's written in the PowerShell language. It can be developed without using any other IDE like VS or VS code
Scripts	Scripts are the text files on disk with .ps1 extension
Applications	Applications are existing windows programs
What if	Tells the cmdlet not to execute. But to define what would happen if the cmdlet is being executed
Debug	Instructs the cmdlet to provide the debugging information
ErrorAction	Instructs the cmdlet to perform a specific action when an error occurs.
ErrorVariable	It specifies the variable which holds the error information
OutVariable	Tells the cmdlets to use a specific variable to hold the output information

#### PowerShell datatypes

Boolean	True or false condition
Char	An 8-bit unsigned whole number from 0 to 255
Date	A 16 bit unsigned number from 0 to 65535
Decimal	A 128 bit decimal value
Double	A double precision 64 bit floating point number.
Integer	A 32 bit signed whole number
Long	A 64 bit signed whole number
Object	Description
Short	A 16 bit unsigned number
Single	A single precision 32 bit floating point number
String	A grouping of characters which is called text

# Intro to Big Data

12 December 2022 19:55

## Traditional Constraints & Reasons for evolving in Big Data

1. **Data volume** aspects - KB, GB, TB, PB, ZB (10 to the power 9-11)
2. **Data Variety** aspects - unstructured, semi-structured, quasi-structured format , structured

RDBMS , weblog, clickstream, BI -> consumption layer (semi-structured data, unstructured data or quasi-structured data)

3. **Data Velocity** - IoT devices, sensors (latitude, longitude, temperature), fast rate the data should be processed and transformed.
4. **Data Veracity** - measurement of different aspects of data. Twitter feed data, Clickstream from various ecommerce websites, weblog data comes from various web servers, crime data, weather data, social events data

## Big Data Processing

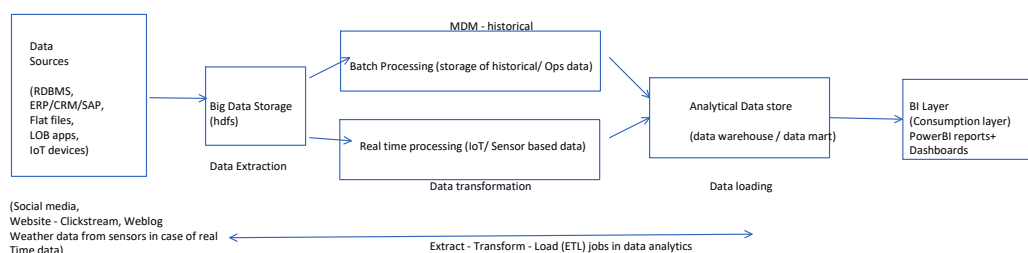
- Batch Data processing
- Real time streaming data processing

e.g. Social Sentiment Analytics (Twitter feed data)

## Core Concepts on Big Data framework

- Data computation happens where the data is being stored (big data storage systems).
- Data is never moved to anywhere for computation

## Big Data Analytics pipeline architecture



## Use cases of Big Data

Big Data involves the data produced by different devices and applications.

- Social Media apps - Social media like Twitter, LinkedIn which holds the versatile volume of data, consists of more petabytes to zetabytes level of data leads to the way of using Big Data platform
- Stock Exchange data - It holds the information related buyers and sellers decisions based on stocks made by different companies,
- Power Grid data - Power grid data which holds the information consumed by a particular node with the respect to the electrical base station of substation.
- Transform data - data includes model, capacity, distance and available of the vehicle
- Search engine data - search engines retrieve data from different databases

## Big Data examples according to the data variety

- a) **Structured data** - Relational data
- b) **Semi-Structured data** - XML, Json data
- c) **Quasi-Structured data** - Json, sensors data, flat file data
- d) **Unstructured data** - Word, PDF, text, media logs, weblogs

## Features of Big Data

1. Big Data Analytics is the use of advanced analytics technique against very large, diverse data sets include structured, semi-structured and unstructured data from different sources and in different sizes from terabytes to zettabytes.
2. Big Data is a term applied to data sets whose size or type is beyond the ability of traditional relational databases to capture, manage and process the data with low latency.
3. Big Data has characteristics - high volume, high velocity and high variety
4. Artificial Intelligence (AI), mobile, social and IoT are driving data complexity through which new forms and sources of data are considered.
5. Example - Big Data comes from Sensors, devices, audio/video, networks, log files, transactional applications, web and social media data. Much of it generates in real time and at a very large scale.

## Characteristics of Big Data

1. Big Data analytics is the advanced analytics techniques against very large, diverse data sets which include structured, semi-structured and unstructured data from different sources and in different sizes from terabytes to petabytes.
2. Big Data can be defined as data sets whose size or type is beyond the ability to traditional relational databases to capture, manage and process the data with low latency.
3. The characteristics of Big Data include high volume, high velocity and high variety.
4. Sources of data are becoming very complex than those for traditional data because they are driven by AI, mobile apps, social media data and IoT device data.
5. With Big Data analytics, we can ultimately get the real time data insights and leads the way to decision making, data modelling and prediction of data for enhanced business intelligence.
6. The Core Big Data frameworks like open source - Apache Hadoop, Apache Spark and the entire hadoop ecosystem tools which are cost-effective, flexible data processing, and storage tools designed to handle the volume of data being generated today.

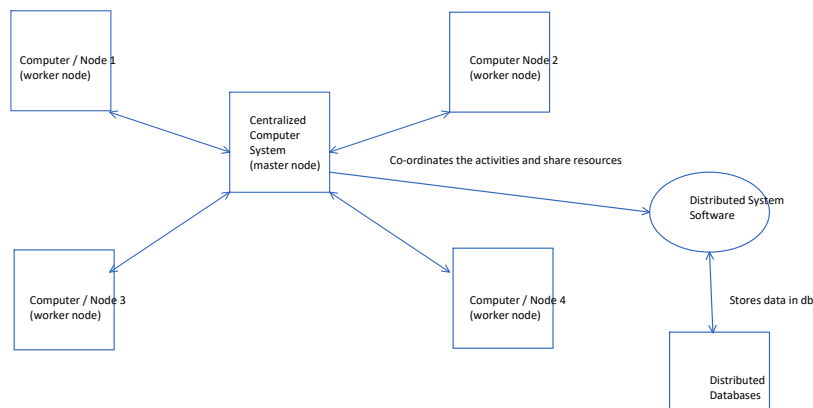
## Benefits of Big Data

1. Social media data helps to gather about the marketing campaigns, product promotions are analytics evolutions
2. Enterprises can use market basket analysis, customer churn analytics, predictive analytics, helps plan their new products, new campaigns & revenue model
3. Using Big Data in healthcare industry, the health records of the patients can be analysed in more details and can provide better insights
4. Banking industry, Big Data can provide more and more insights on customer purchase pattern, credit card usage, demographics which can lead to better sales profit margin analysis and recommend new products to customer.

## Distributed System

It's a collection of autonomous computer systems which are typically separated but are connected by a centralized computer network which is equipped with distributed system software.

The autonomous computers will communicate among other system by sharing resources and files and performing the tasks assigned to them.



### Characteristic of Distributed System

1. Scalable
2. Concurrency
3. Fault tolerance
4. Transparency
5. Heterogeneity

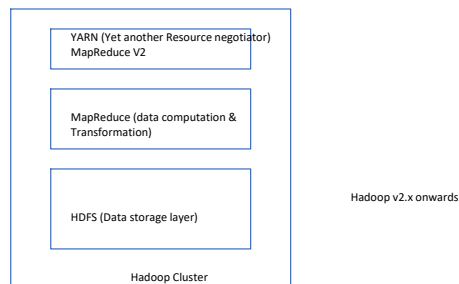
## Apache Hadoop

Apache Hadoop is the distributed, fault tolerant open source framework dedicated for Big Data analytics platform.

Apache Hadoop is the official Big Data framework/ tool released from ASF (Apache Software Foundation).

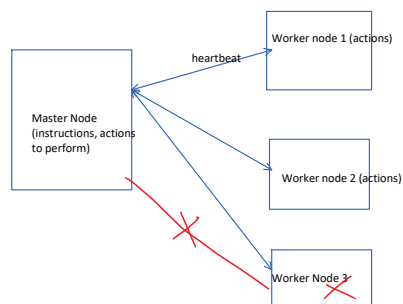
### Feature of Apache Hadoop

1. Apache Hadoop consists of two layers - Storage layer (HDFS - hadoop distributed file system)
2. MapReduce -> core data computation layer



### Rest of other Hadoop ecosystem components -

1. Apache Hive - Data warehouse and analytics tool in hadoop framework
2. Apache Pig - data querying tool with scripting language (ETL jobs)
3. Apache Kafka - real time data processing framework built on top of hadoop



- Apache hadoop is a distributed, fault tolerant, highly available Big data ecosystem cluster which works with master node & worker node design.

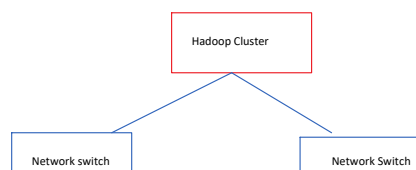
- The master node is responsible for sending the instructions and job details to the worker node(s).

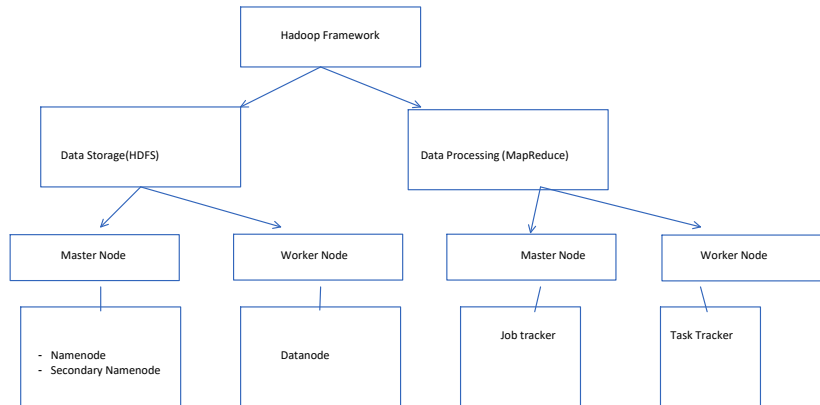
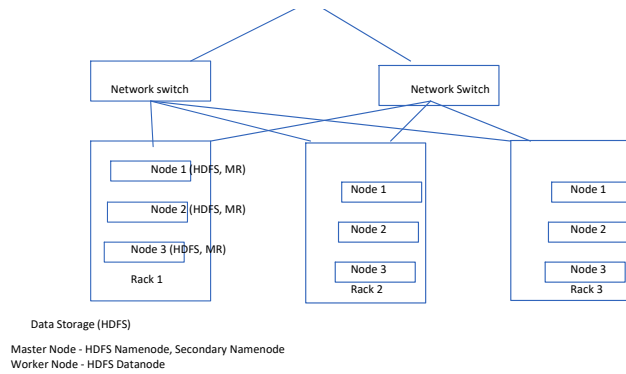
- The worker nodes processes the jobs and produces the output as per the instructions provided the master node

- Any point of time if any worker node failed to respond to the master on time, that worker node gets discarded, new worker nodes are being assigned as per the horizontal scaling feature.

- A collection of nodes is called cluster in hadoop

- A node is a point of intersection/ connection within a network, i.e. server.





#### HDFS Operations

1. HDFS has master and worker node architecture.
2. An HDFS cluster consists of a single Namenode, a master node that manages the file system namespace and regulates access to the files by clients.
3. In addition, there're number of datanodes, usually, one per node in the cluster which manages storage attached to the nodes that they run on.
4. HDFS exposes a file system namespace and allows user data to be stored in files.
5. Internally, a data file is being splitted into one or more blocks and these blocks are stored in a set of datanodes.
6. The namenode executes file system namespace operations like opening, closing, and renaming files and directories.
7. It also determines the mapping of hdfs datanode blocks to individual datanodes.
8. The datanodes are responsible for serving read and write requests from the file system's directories from the clients.
9. The datanodes also can perform block creation, deletion, and replication upon instructions given by Namenode.

Starting from Hadoop V2.x, each & every datablocks in HDFS datanode should be size of 128 mb.

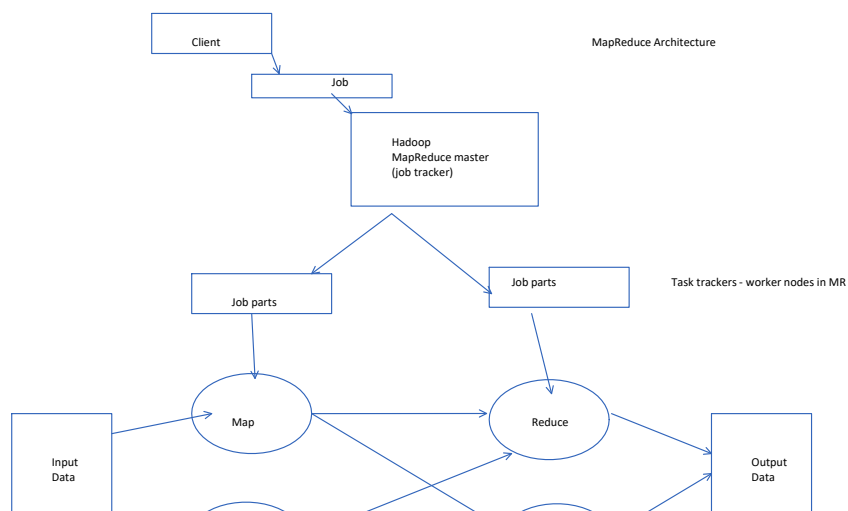
Input data file = 1GB = 1024 MB / 128 mb = 8 blocks of data would be stored in data node

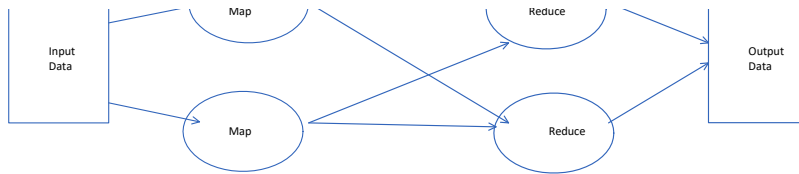
#### Features of MapReduce

1. MapReduce and HDFS are the two major components of Hadoop which makes it so powerful and efficient to use.
2. MapReduce is programming model used for efficient processing in parallel over large data sets in a distributed manner.
3. The input which is first stored within datanode blocks of HDFS, is first split and then combined to produce the final result.
4. The libraries for MapReduce is written in so many programming languages with various different optimizations.

#### Benefits of MapReduce

1. The purpose of MapReduce is to Map each of the jobs and then it will reduce the processing power.
2. The MapReduce task is mainly divided into two phases Map phase and Reduce phase.





1. Client - The MR client is one which brings the job for data computation for processing. There can be multiple clients available that can continuously send jobs for processing to the hadoop MR manager
2. Job - The MapReduce job is the actual work that the client performs which consists of so smaller tasks that client executes to process the data.
3. Hadoop MR master - it divides the particular job into subsequent job parts.
4. Job parts - The task of job parts which are obtained after dividing the main job. The result of all the job parts are combined to produce the final output
5. Input data - the data which is fed into MR for processing and transformation
6. Output data - refers to the final data i.e. processed and transformed data.

**Job Tracker** - the purpose of job tracker is to manage all the resources and all the jobs across the cluster and also to schedule each map on the task tracker running on the same data node. Since there can be hundreds of data nodes available in the cluster.

**Task Tracker** - The task tracker can be considered as the actual slaves which are working on the instruction given by the job tracker. This task tracker is deployed on each of the nodes available in the cluster which executes the Map and Reduce tasks as instructed by the job tracker.

**Job history server** - the job history server is a daemon process which saves and stores historical information about the task or application.

#### Big Data use cases

1. Product development - social media to develop social sentiment analysis
2. Predictive maintenance - Aerodynamics,
3. Customer experience - Telecom, Banking, Retail
4. Fraud Detection and Compliance - Banking, Financial
5. Drive innovation - BFSI, Healthcare and marketing domains

#### Apache Pig

1. Apache Pig can handle structured, semi-structured or unstructured data & stores the data into HDFS.
2. Every task or query running on pig is executed through MapReduce.

#### Benefits

- a) Ease of programming. We can write the queries in Apache Pig through scripting.
- b) Optimization of complex data processing (Load, Transform & DUMP)
- c) Flexible and with built-in operators (sort, filter, join, foreach)
- d) Supported for Avro based file formats (row wise)

#### Features of Apache Pig

- For performing several operations, apache pig provides rich sets of operators like the filters, join, sort etc.
- Joins operations are easier in Pig.
- Apache pig allows splits in the pipeline, data structure is multivalued, nested and richer.
- Pig can handle for analysis for both structured and unstructured data.
- 

#### Data types in Pig

1. Tuple - it's an ordered set of the fields
2. Bag - It's a collection of the tuples
3. Map - It's a collection of key/value pairs.
4. Atom - It's an atomic data value which's used to store as a string. It can be used as a number and string.

MapReduce	Apache Pig
It's a low level data processing	It's a high level data flow tool
Complex programs can execute through java/python	Apache pig we can execute simple pig/latin scripts which are simpler also consists of less no of lines.
Data operators are difficult compared to pig	Pig latin these built-in operators are available to work with
It does not allow nested data types	It provides the nested data types like tuple, bag and map

#### Apache Hive

It's data warehouse tool on Apache hadoop framework. Hive is used for data analysis in Big Data platform.

- a) Tools to enable easy access to data via SQL like query interface which is called HQL
- b) Enables us to perform data warehouse tasks such (extract, transform and load) ETL operations, analysis and reporting.
- c) Query processing , execution is done on hive through Mapreduce.
- d) In hive, we can query the data through SQL like query interface called as Hive QL.
- e) Hive supports various file formats like CSV/TSV, apache parquet (columnar file format which is performant), Apache ORC (row-columnar file).
- f) Hive queries can be executed based on query retrieval support when using Hive LLAP, Apache Yarn etc.
- g) Hive can operate with declarative SQL support.

Apache Pig	Apache Hive
Pig operates on the client side, since it's client side Scripting language	Hive operates on the server side of the cluster
Pig uses the pig-latin script for analysis	Hive uses the hive-ql language
Pig/latin is purely procedural language	Hive is declarative SQL language
Pig has support of AVRO file format (row-oriented)	Apache Hive is suitable for Parquet and ORC file format
Apache Pig is suitable for complex and nested data structure	Hive is suitable for match processing and OLAP system
Apache Pig does not support schema to store data	Apache Hive has the support for JDBC and ODBC
Pig does not have any metastore(metadata store)	Apache hive has the metastore support for DDL language (SQL, mySQL, postgresQL) etc.
Pig operates quickly and data loads quickly	Hive loads data slowly
.pig file is extension for pig-latin scripts	Any file formats(.hql) applicable for hive data query files.

#### Hadoop ecosystem components

- HDFS (data storage)
- Mapreduce (data computation)
- YARN (resource management / MR v2)
- Pig (data analysis)
- Hive (data warehouse)
- Kafka (real time data streaming)

#### Hands-on Lab

##### Create Azure Virtual Machine (on Linux)

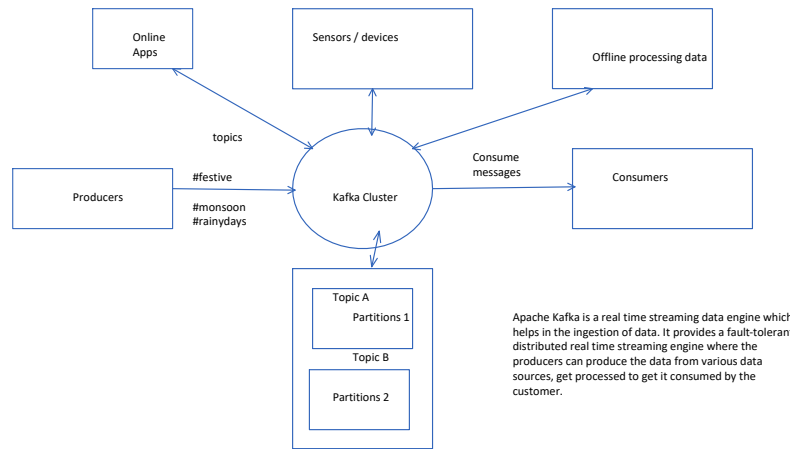
Pre-requisites

**Create Azure Virtual Machine (on Linux)**

## Pre-requisites

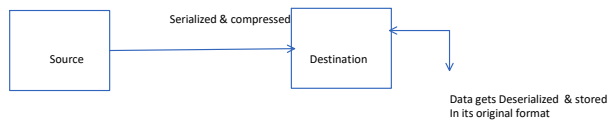
- Need to have the Azure Subscription (free trial)
- Resource provider, we've to register the resource "Microsoft.Compute"

- HDFS (data storage)
- Mapreduce (data computation)
- YARN (resource management / MR v2)
- Pig (data analysis)
- Hive (data warehouse)
- Kafka (real time data streaming)
- Spark (real time data transformation and processing)
- Sqoop (bidirectional transfer of data from hadoop / HDFS to any RDBMS (MSSQL mysql etc.)
- Flume (data ingestion tool)
- Oozie (workflow management tool)



## Serialization of Data

In HDFS for hadoop cluster



## HDFS federation &amp; High availability

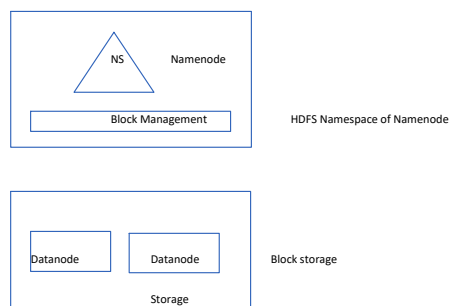
1. HDFS federation is the process through which the namenode failure can overcome in case of namenode is failed and the single point of failure happens.
2. HDFS has primarily two components

- Namespace - consists of files, directories and blocks
- Block storage service - which has two parts

## a) Block management

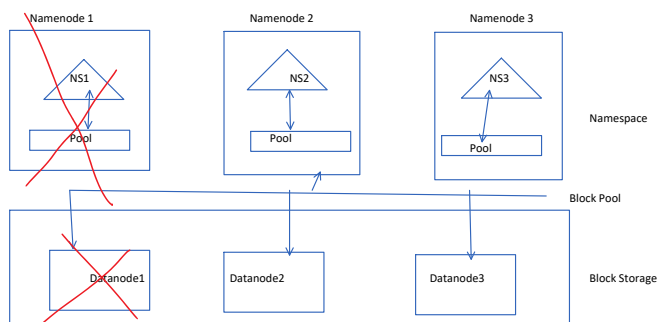
	Provides Datanode cluster membership by handling registration, also sends periodic heartbeats.
b)	Provides the block report and maintains the location of blocks.

Manage replica place, block replication for under replicated blocks, deletion of blocks which are replicated.



## Multiple Namespaces / Namenodes

As name service/namespace gets scaled horizontally, HDFS federation uses multiple namenode/ namespaces. The namespaces are federated, the namenodes are independent and do not require any co-ordination. The Datanodes are used as common storage for blocks by all namenodes. Each datanode gets registered with all of the namenodes in the cluster. Data nodes to send periodic heartbeats and send block reports. They can handle commands from the namenodes.





#### Block Pool -

It's a set of blocks which belongs to a single namespace. Datanodes store the blocks for all block pools in the cluster. Each block pool is managed independently. This allows a namespace to generate Block Ids for new blocks without the need for co-ordination with other namespaces.

A namenode failure does not prevent the datanode from serving other namenodes in the cluster.

A namenode and its block pool together are called namespace volume. A self contained unit of management. When a namenode/namespace is deleted, the corresponding block pool at the datanode is deleted. Each namespace volume is upgraded as unit during the cluster upgrade.

#### Benefits of Namenode HA & Federation

1. Namespace Scalability - Federation adds namespace horizontal scalability, large deployments or deployment using lot of small files benefit from namespace scaling by allowing more namenode to be added to the cluster
2. Performance - File system throughput is not limited by a single namenode. Adding more namenodes to the cluster scales the file system read/write throughput.
3. Isolation - A single namenode can offer no isolation in a multi user environment. It can overcome by using HDFS federation, where different applications can be isolated by a different namespace by using multiple namenode.

# Apache Hadoop Overview

12 December 2022 19:55

## HDFS Caching

Centralized Cache management in HDFS is an explicit caching mechanism which allows users to specify paths to be cached by HDFS. The namenode will communicate with datanode which has the desired blocks on disk, and instruct them to cache the blocks in off-heap cache.

### Centralized cache management in HDFS

1. Explicit caching of data, helps to evict the old data from memory. This is particularly important when the size of working set exceeds the size of main memory.
2. Because, datanode caches are managed by the namenode. Applications can query on the set of cached block of data when making task placement decisions. Co-locating a task with a cached block replica improves the read performance.
3. When a block is cached by the datanode, clients can use a new, more-efficient, zero-copy read API, since the checksum verification of cached data is done by the datanode.
4. Centralized cache can improve the overall cluster memory utilization, when relying on the OS buffer at each datanode, repeated blocks will result all n replicas of the block being pull down into the buffer cache. With the centralized caching, a user can explicitly specify the n replicas saving memory.

### Use cases

Centralized cache management is useful for files that being read/accessed repeatedly. For example, a small fact table in hive is often joins is good candidate for caching.

### Centralized caching,

For mixed workloads, with performance SLA. Caching is a working ser of high priority workload ensures that it does not content for disk I/O with a low priority workload.

## MapReduce Framework

1. A Mapreduce job usually splits the input dataset into the independent chunks of data which are being processed by the map tasks in a completely parallel manner.
2. This Framework sorts the outputs of the maps, which are then input to the reduce tasks. Typically both the input and output of the jobs are stored in a file system.
3. The framework takes care of scheduling tasks, monitoring them and re-executes the failed tasks.
4. Typically, the compute nodes and storage nodes are the same, MR framework & HDFS are running on the same set of nodes. This configuration allows the framework to effectively schedule tasks on the nodes where data is already being present resulting in very high aggregate bandwidth across the cluster.
5. The MR framework consists of single master ResourceManager, one worker NodeManager per cluster-node and MRAppMaster per app.
6. Minimally, apps can specify the input and output locations and supply map and reduce functions via implementation of appropriate interfaces and abstract classes. These and other job parameters are referred as the job configuration.
7. During mapreduce job execution, the hadoop client submits the job and shares the configuration with the ResourceManager which then assumes the responsibility of distributing the software/configuration to the workers.
8. Scheduling tasks and monitoring them provides the status and diagnostics information to the job-client.

## MapReduce Input and Output

The MR framework operates exclusively on <key, value> pairs, this framework to the job as a set of <key, value> and produces a set of <key, value> pairs as the output of the job.

The key and value classes have to be serializable by the framework and hence need to implement the Writable interface.

### Input and Output operation types of MR job

(input) <k1, v1> -> **map** -> <k2, v2>-> **combine** -> <k2, v2> -> **reduce** -> <k3, v3> (output)



# Apache Hadoop (Deep Dive)

12 December 2022 19:56

Partitioning of Hive - Hive Partition is a way to organize large tables into smaller logical tables based on values of columns. One logical table (partition) for each distinct value. In Hive, tables are created as a directory in HDFS. A table can have one more partitions which corresponds to a sub-directory for each partition inside the table directory.

Example: A student details table in Hive can be partitioned based on Departments.

The partitions can create the sub-directories which leads to the faster data processing.

Bucketing of Hive - It's technique to split the data into more manageable files.

## Features of Partitioning and Bucketing

Hive Partitioning	Hive Bucketing
It's the process of organizing tables into partitions for grouping same type of data together Based on a column or partition key.  Each table in the hive can have one or more partition keys to identify a particular partition. Using partition, we can make it faster to execute queries on slices of data.	In Hive tables or partition are subdivided into buckets based on hash function.  It's more efficient process because it gives extra structure to the data to be used for more efficient queries.
Pros - 1. It distributes the execution load horizontally 2. In partition faster execution of queries with the low volume of data takes place.	Pros - 1. It provides faster query response like partitioning 2. In bucketing, due to equal volumes of data in each partition, joins at Map side will be quicker.
Cons: 1. There is the possibility of too many small partitions creations - too many directories 2. Partition is effective for low volume of data. There're few queries which takes long time to execute. 3. There's no need for searching entire table column for a single record	Cons: We can define a number of buckets during the table creation.  But loading of an equal volume of data which has to be done manually by programmers.

## HDFS commands

1. Create a directory:

```
hdfs dfs -mkdir /user
```

2. Check the contents of the directory within HDFS

```
Hdfs dfs -ls /user
```

3. Copy the files from local drive to hdfs

```
Hdfs dfs -put /home/user/file.txt /user
```

```
Hdfs dfs -CopyFromLocal <local_file_path> <hdfs_file_path>
```

4. Create a new empty file

```
Hdfs dfs -touchz /user/username
```

5. Delete a file from the HDFS directory

```
Hdfs dfs -rm /demo/sample.txt <hdfs_directory_path>
```

6. Copy the file back to local drive from hdfs

Hdfs dfs -copyToLocal <local\_file\_path> <hdfs\_file\_path>

## Apache Sqoop

Sqoop is a tool which is designed to transfer the data between hadoop and relational databases. We can use sqoop to import data from a RDBMS such as mysql or Oracle, postgresQL, sqlite into the HDFS, transform the data in hadoop mapreduce, then export the data back to RDBMS.

Sqoop automates most of the process, relying on the database to describe the schema for the data to be imported. Sqoop uses MapReduce jobs to import and export the data, which provides the parallel operation as well as fault tolerance.

### Features of Sqoop :

1. With Sqoop, we can import data from a relational database system into HDFS. The input to the import process is a database table. Sqoop will read the table row-by-row basis into HDFS.
2. The output of this import process is a set of files containing a copy of the imported table. The import process is performed in parallel. The output will be in multiple files. These files may be in delimited text files (commas, tabs separated fields)
3. Sqoop uses jdbc driver in order to connect to the hive table or sql table in rdbms.
4. Using Conditional import, can specify which data to import like --table argument, to select the table to import. By default, all columns within a table can be selected for import.
5. Imported data is written to HDFS in its natural order, so that once the table containing the columns like A, B, C, it would result import of data like  
A1, B1, C1  
A2, B2, C2 ..

### Incremental Import in Apache Sqoop

Sqoop provides an incremental import mode which can be used to retrieve only the rows newer than the some previously-imported set of rows.

#### Arguments in incremental import of Sqoop

-- check-column -- specifies the column to be examined when determining which rows to import  
-- incremental (mode) -- specifies how sqoop determine which rows are new,  
--last-value (value) - specifies the max value of the check column from the previous import.

- a) Sqoop supports two types of incremental imports: append and lastmodified.
- b) We can use the -- incremental argument to specify the type of incremental import to perform.
- c) We should specify append mode when importing a table where new rows are continually being added with increasing row id values. The columns containing the row's id with --check-column.
- d) Sqoop actually can import the rows where the check column has a value greater than the one specified with --last-value.
- e) At the end of the incremental import, the value which should be specified as --last-value for subsequent import is printed on the screen. While running a subsequent import, we should specify --last-value to ensure only the new or updated data can be imported.

### File Formats in Sqoop

We can import data through sqoop in one of the two file formats - delimited text and SequenceFiles.

Delimited text is default import format. We can also specify it explicitly by using the --as-textfile argument. This argument will write string-based representations of each record to the output files, with delimiter characters between the individual columns and rows.

- The delimiters may be commas, tabs, or other characters.
- Delimited text is appropriate for most of non-binary data types. It also readily supports further manipulation by other tools (i.e. hive).
- Sequence files are binary format which store individual records in custom record-specific data types. These data types are manifested as java classes. Sqoop will automatically generate these data types. The format supports exact storage of all data in binary representations. It's appropriate for storing binary data, or data which will be manipulated by the custom MapReduce programs.

## Apache Flume

Apache Flume is a tool/service/data ingestion mechanism for collecting, aggregating and transporting large amounts of streaming data such as log files, events from various sources to a centralized data store.

Flume is highly reliable, distributed and configurable tool. It's designed to copy streaming data from various web servers to HDFS.



## Benefits of Flume

1. Using Apache Flume, we can store the data in any of the centralized stores (HDFS).
2. When the rate of incoming data exceeds the rate, at which the data can be written to the destination, Flume acts as a mediator between the data producers and centralized stores. Provides a steady flow of data between them.
3. The transactions in Flume are channel based where two transactions (one sender, one receiver) are maintained for each message. It also guarantees reliable message delivery.
4. Flume is reliable, fault tolerant, scalable, manageable and customizable.

## Flume Interceptors

1. Flume Interceptors are those who has the capability to modify/drop events in-flight.
2. There're different kinds of interceptors in flume -
  - a) host-flume interceptors - it inserts the hostname/IP address of the host where the agent is running on.
  - b) Regex-filtering interceptors
  - c) Remove-header interceptors
  - d) Static interceptors
  - e) Timestamp interceptors - it inserts the event headers , the time in which it processes the event.

# Azure Data Factory

12 December 2022 19:56

## ETL Definition

ETL is a data integration process which combines data from multiple data sources into a single, consistent data store which's loaded into a data warehouse / consumption layer or other BI systems.

ETL is the process for integrating and loading the data for computation and analysis, it's also the primary method to process data for data warehousing projects.

## Benefits

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

Azure Data Factory is a managed cloud service which's built for these complex extract -transform-load (ETL) , extract-load-transform(ELT) and data integration projects.

Data orchestration is the practice of acquiring, cleaning and matching, enriching and making data accessible across the technology systems. 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 a schedule.

## ETL vs ELT

ETL	ELT
ETL (extract, transform, load) can perform the end to end data transformation after loading the raw data into the transformation layer and captures the insights from the data.	ELT (extract, load, transform) 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 directly to the target data store to be transformed as required.
The ETL process, involves more structured datasets and data has to be relational 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 dataset and ideal for big data use cases.
In on-premise, ETL data pipeline is more common	On Cloud , Azure ELT pipeline is more popular

## ETL process

### Extract

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

- SQL / noSQL
- CRM / ERP
- Flat files
- Web pages

### Transform

In the staging area, the raw data undergoes the data processing, data is transformed and consolidated for its intended analytical use case.

- Filtering, cleansing, de-duplicating, validating, and authenticating the data
- Performing calculations, translations, summarizations of the 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.

### Load

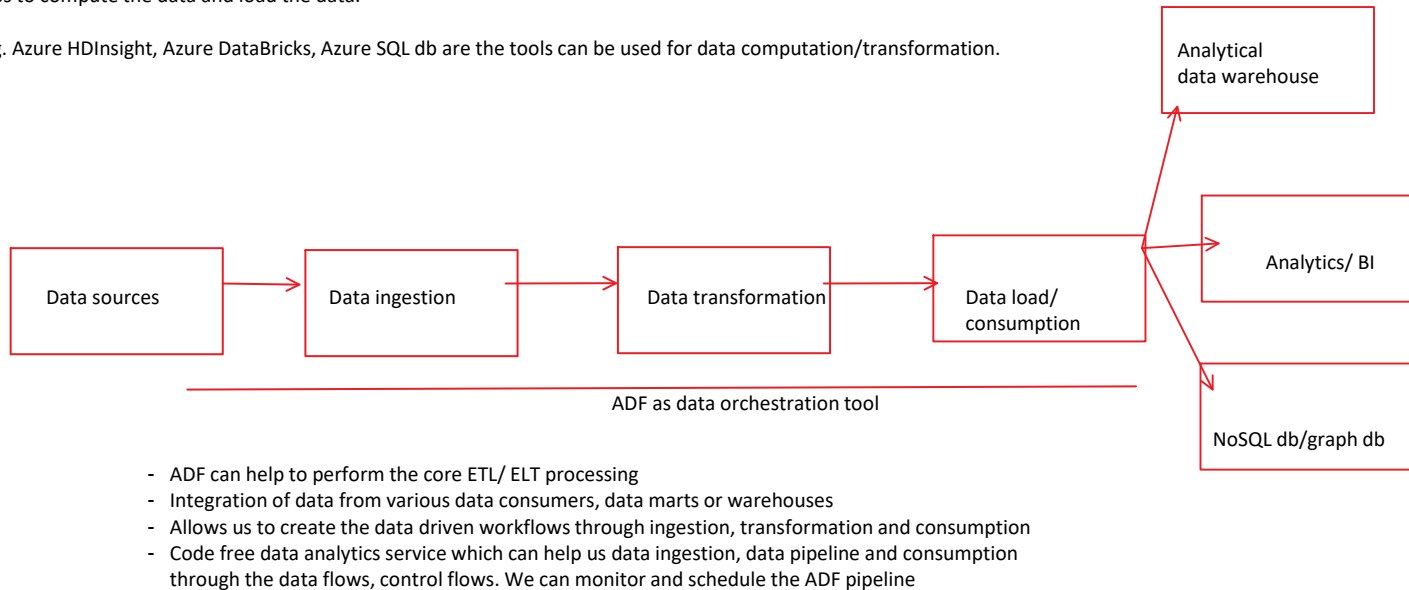
In the last step, the transformed data is moved from the staging area to a target data warehouse. Which involves loading of all the data, followed by the periodic loading of incremental data changes and full refreshes to erase and replace the data in the data warehouse.

## ETL Tool (Azure Data Factory)

1. Comprehensive automation support - leading ETL tools like 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 interface - the functionality can be used for specifying rules and data flows.
3. Support for complex data management - complex calculations, data integrations and string manipulations.
4. Security and compliance - the best ETL tool encrypt the data both in motion and at rest includes various regulations like GDPR, HIPAA etc.

Data factory is the cloud based ETL and data integration service allows to create data-driven workflows for orchestrating data movement and transforming the data at scale. Using data factory (ADF), we can orchestrate, create and schedule the data driven workflows (pipeline) which can ingest the data from disparate data sources and can build complex ETL jobs to compute the data and load the data.

e.g. Azure HDInsight, Azure DataBricks, Azure SQL db are the tools can be used for data computation/transformation.



## Azure Data Factory components

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 to the data factory.

e.g. db connection strings which define the connection information required for the service to connect to the external resource.

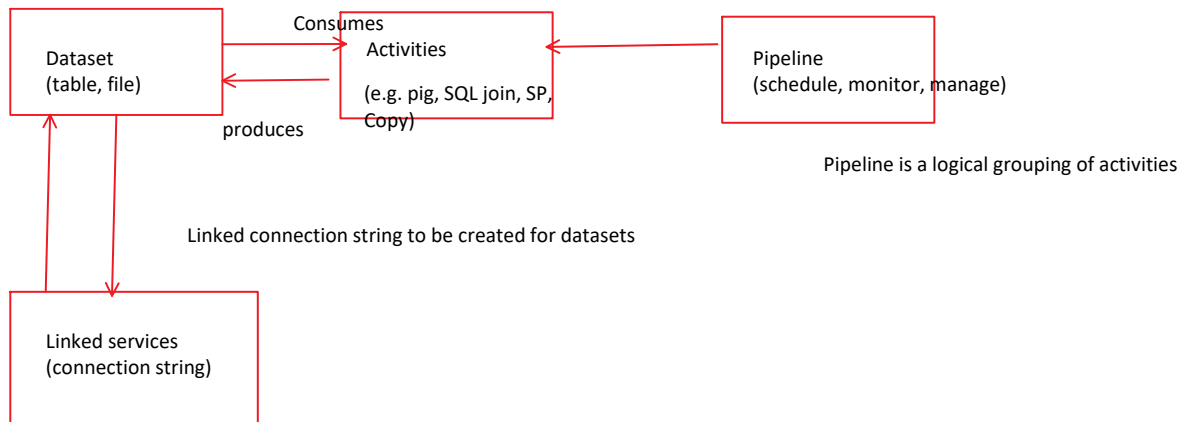


While copying/moving the data from Azure blob storage to Azure SQL db , the storage and the db linked service connection string need to be created.

2. Activity - It's the resource defines us which action / operation to be performed on the data.
  - Copy data,
  - Deduplicating the data
  - Formatting the data to remove all nulls, NaNs
  - Applying normalization, remove redundancies
  - Adding column headers and formatting
  -
3. Dataset - dataset represents the data source (sql table, flat files)
4. Pipeline - an integrated set of ADF activities

A data factory can contain one or more pipelines.

- a) A pipeline is the logical grouping of activities which together perform a task.  
For e.g. the 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.
- b) 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.
- c) The activities in a pipeline defines the actions to perform on the data.



### Hands-on Lab: 1

#### Pre-requisites : Azure Portal

#### Tasks

1. Create Azure Data Factory
2. Explore Azure Data Factory Studio
3. Explore ADF activities, data sources, pipelines
4. Monitoring and management capacity of ADF

### Hands-on Lab : 2

Transferring data from Azure SQL db to Azure Blob storage through ADF Copy Data Tool

#### Tasks:

1. Provision an Azure SQL db with sample AdventureWorks db
2. Provision an Azure Blob storage
3. Use Copy data tool to copy the data from Azure SQL db to Blob storage
4. Create the ADF pipeline
5. Monitor the pipeline

### Extract, Transform, and Load (ETL) process

Extract, Transform and load (ETL) is a data pipeline used to collect data from various sources. It then transforms the data according to business rules, and loads the data into the destination data store.

The transformation work in ETL takes place in a specialized engine, which often involves using staging tables to temporary hold data as being transformed and ultimately loaded to its destination.

### ETL process in ADF

Load





1. The data transformation which takes place usually involves various operations, such as filtering, sorting, aggregating, joining data, cleaning data, deduplicating data and validating data.
2. Often, these three ETL processes can execute in parallel to save time.  
e.g. while the data is being extracted, a transformation process could be working on data already received and prepare it for loading, and a loading process can begin working on the prepared data. Rather than just waiting for the entire extraction process to complete.

#### ELT in ADF pipeline (Extract, Load and Transform)

Extract, load and transform differs from ETL mainly in where the transformation takes place. In the ELT pipeline, the transformation occurs in the target data store instead of using a separate transformation engine.

In the ELT pipeline, the transformation takes place in the target data store. The processing capabilities of the target data store are used to transform the data. This simplifies the ELT architecture by removing the transformation engine from the pipeline.

1. Scaling of the target data store also scales the ELT pipeline performance. However, the ELT only works well then the target system is powerful to transform the data efficiently.



1. In Azure, the source flat files in scalable storage, such as HDFS, Azure blob storage or ADLS gen2 are all can be used as data storage layer for extraction.
2. For Spark, Hive, polybase, SQL query, Pig/latin scripts, MR jobs etc. can be used to query the source data.
3. In ELT pipeline, the key difference is that the data store used to perform the transformation is the same data store where the data is ultimately getting consumed. The data store reads directly from the scalable storage, instead of loading the data into its own proprietary storage. This approach skips the data copy step present in ETL which often just can be a time consuming operation for large data sets.
4. The data store only manages the schema of the data and applies the schema on read.
5. The final phase of ELT pipeline, is to transform the source data into a final format which is more efficient for the types of queries which need to be supported.

## 1. Ingest

- The data can ingest on multi-cloud and on-premise based on hybrid copy data.
- 100+ native connectors
- Serverless and auto scale
- Use wizard for quick copy jobs

## 2. Control Flow

- Design code free data pipelines.
- Generate pipelines via SDK
- Utilize workflow constructs - loops, branches, conditional execution, variables and parameters.

## 3. Data Flow

- Code free data transformations which execute in Spark
- Scaled out (add more instances) in Azure Integration Runtime
- Generate data flows via SDK
- Designers for data engineers and data analysts.

## 4. Schedule

- Build and maintain operational schedules for all of the data pipelines
- We can execute the ADF pipeline manually, programmatically, tumbling window and schedule basis

## 5. Monitor

- View active executions and pipeline history
- Detail activity and data flow execution status
- Establish alerts and notifications

### Schema Drift in ADF data flow

Schema drift is the case where the data sources often change metadata, fields, columns and types which can be added, removed, or changed on the fly.

Without handling for schema drift, the data flow becomes vulnerable to the upstream data source changes. Typical ETL patterns fail when the incoming columns and fields change because they tend to be tied to the source names.

- In order to protect against the schema drift ,
- Define the sources which has mutable field names, data types, values and sizes.
- Define the transformation parameters which can work with data patterns instead of hard-coded fields and values.
- Define expressions which can understand patterns to match incoming fields, instead of using named fields.
- Schema drift can be applied for both data source and sinks.

### Transform the drifted columns

When the data flow has drifted columns, we can access in the transformations with the methods like

- Use the 'byPosition' and 'byName' expressions to explicitly reference a column by name or position number
- Add a column pattern in a derived column or aggregate transformation to match any combination of name, stream, position , origin or type.
- Add rule-based mapping in a Select or Sink transformation to match drifted columns to columns aliases via the pattern.

### Control Flow Activity

Control Flow activities in ADF involves orchestration of pipeline activities including the chaining of activities in a sequence, branching, defining the parameters at the pipeline level.

The control flow activity also involves passing arguments while invoking the pipeline.

It also includes custom-state passing and looping containers.

The control flow activity defines also how control / sequence activities can pass through from one task to another.

Features:

1. **Control Flow activity is an orchestration of pipeline activities.**
2. **The orchestration includes the chaining activities in a sequence, branching and defining the parameters at the pipeline level.**
3. **It also helps to pass arguments while invoking the pipeline on demand or from a trigger.**

#### Examples:

##### 1. Lookup Activity in ADF

Lookup activity can retrieve a dataset from any of the data sources supported by data factory. We can use it to dynamically determine which objects to operate on in a subsequent activity, instead of hard coding the object name, some objects examples like files and tables.

Lookup activity reads and returns the content of a configuration file or table.

- It also returns the result of executing a query or stored procedure.
- The output can be a singleton value or an array of attributes, which can be consumed in a subsequent copy, transformation or control flow activities like ForEach activity.

#### Features of Lookup activity

1. The Lookup activity can return upto 5000 rows, if the result set contains more records, then the first 5000 rows will be returned.
2. The lookup activity output supports up to 4 MB in size, activity will fail if the size exceeds the limit.
3. The longest duration for Lookup activity before timeout is 24 hours.

#### Type properties in Lookup Activities

a) Dataset - provides the dataset reference for the lookup. Get details from the dataset properties section in each of the corresponding connector,

- b) Source - Contains the dataset-specific properties. The same as the Copy activity source.
- c) firstRowOnly - indicates whether to return only the first row or all rows .

2. **Foreach Activity** - It's a control flow activity type in ADF. The foreach activity can define the repeating control flow in ADF pipeline. This activity can be used to iterate over a collection and executes specified pipeline activities in a loop.

The loop implementation for the activity is similar to Foreach looping structure in programming.

#### Type properties

- Name
- Type
- isSequential (specifies whether the loop should be executed sequentially or in parallel. Max of 50 loop iterations can be executed at once in parallel.) If we have a ForEach activity iterating over a copy activity with 10 different source and sink dataset, with **isSequential** set to false, all copies are executed at once.

If 'isSequential' is set to false, ensure it that there's a correct configuration is being executed to run multiple executables. Otherwise, this property should be used with caution to avoid incurring write conflicts.

- batchCount
- Items - expression which returns a JSON array to be iterated over
- Activities - the activities are to executed.

#### Hands-on Lab 3

1. Create the ADF
2. Create the ADLS gen2 storage
3. Create Azure SQL db (with AdventureWorks db)
4. Providing the Linked Service connection through managed identity from ADLS to ADF
5. Copy the data through Copy data activity in ADF pipeline
6. Test, Validate and Monitor the pipeline

#### Hands-on Lab 4:

##### Mapping Data Flow

#### Hands-on Lab 5:

Move/Copy data through ADF from on-premise SQL Server database to Azure Blob storage

Pre-requisites - SQL authentication is required to be enabled for on-premise SQL server database

#### Hands-on lab 6:

Building a reusable pipeline in ADF

#### Hands-on Lab 7:

Monitoring, validating the ADF pipeline

3. Switch Activity - This Switch activity provides the same functionality which is a switch statement in programming. It evaluates a set of activities corresponding to a case which matches the condition evaluation.

4. Execute Pipeline Activity - This activity allows the data factory to invoke another pipeline.

property

Name	Name of the execute pipeline activity
Type	Must be set to : Execute Pipeline
pipeline	Pipeline reference to the dependent pipeline which the pipeline invokes. A pipeline reference object has two properties,  a) referenceName - the referenceName property specifies the name of the reference pipeline b) Type - the type property should be set to pipeline reference.
parameters	Parameters to be passed to that invoked pipeline
waitOnCompletion	Defines whether activity execution waits for the dependent pipeline execution to finish. Default is true.

5. The Custom Activity - There are two types of activities in ADF, to use in Data Factory.

- Data movement activities to move data between the supported data source and sink data source
- Data transformation activities to transform the data using compute services (such as Azure HDInsight)
- To move the data to/from the data store which the service does not support, or to transform/process data in a way not supported by the service, we can create a custom activity with the data movement or transformation logic and use the activity in the pipeline.

#### Overview of Azure Data Factory Control Flow Activities

Control Flow Activity Name	Purpose / Definition
Append Variable	Append variable activity could be used to add a value to an existing array variable defined in ADF pipeline
Set variable	Set variable activity can be used to set the value of an existing variable of type string, bool, or array defined in a ADF pipeline
Execute Pipeline	The execute pipeline activity allows ADF pipeline to invoke another pipeline
If condition	If condition activity allows directing pipeline execution, based on evaluation of certain expressions.  The If condition activity provides the same functionality that an if statement provides in programming.  It executes a set of activities when the condition evaluates to true and another set of activities when the condition evaluates to false.
Get Metadata	Get Metadata activity can be used to retrieve the metadata of any data in ADF
The Foreach activity	The activity defines the repeating control flow in the ADF pipeline. This activity is used to iterate over a collection and executes specified activities in a loop.  The loop implementation of this activity is similar to Foreach looping structure in programming.
Lookup	Lookup activity can retrieve a dataset from any of the ADF supported data sources

Filter	Filter activity can be used in a pipeline to apply a filter expression to an input array
Until	Executes the set of activities in a loop until the condition associated with the activity set to true.
Wait	Wait activity allows pausing pipeline execution for the specified time period
Web Activity	Web activity can be used to call custom REST endpoint from an ADF pipeline
Azure Function	Allows to run the Azure Function in the ADF pipeline
Validation Activity	We can use the Validation in a pipeline, to ensure the pipeline only continues execution once it has validated the attached dataset reference exists, that it meets the specified criteria, or timeout has reached.
Webhook activity	It can control the execution of pipelines through the custom code. With the webhook activity, code can call an endpoint and pass it a callback URL. The pipeline run waits for the callback invocation before it proceeds to the next activity.

Webhooks are automated messages sent from apps when something happens. Webhooks have a message, or payload and sent to the unique URL,

- We've to get the webhook URL from the application to send the data to.
- Use the URL in the webhook section of the application we want to retrieve the data from
- Choose the type of events we want the application to notify about

Hands-on Lab 4:

Pre-requisites

1. **Create the ADF resource**
2. **Create the Azure SQL db (with sample AdventureWorks db) (input data source)**
3. **Create the ADLS gen2 storage (data sink)**

#### 4. **Create the linked services**

- 4.1. Linked Service for Azure SQL db
- 4.2. Linked Service for ADLS gen2

#### 5. **Create the datasets**

- 5.1. ADLS dataset
- 5.4. Azure SQL for product table
- 5.5 Azure SQL for ProductCategory table

#### 6. **Three pipelines**

- 6.1. copy\_product
- 6.2 copy\_productCategory

1. Start the Lab 4 for copy data and use the data flows.

### **Mapping Data Flow Activities**

1. Mapping Data flow activities are visually designed data transformations in Azure Data Factory. Data Flows allow data engineers to develop data transformation logic without writing code.
2. The resulting data flows are executed as activities within ADF pipelines which use scaled-out Apache Spark cluster, Data flow activities can be operationalized using existing Azure Data Factory scheduling, control flow and monitoring capabilities.
3. Mapping data flows provide an entirely visual experience with no coding required. The data flows run on ADF-managed clusters for scaled out data processing. Azure Data factory handles all the code translation, path optimization and execution of the data flow jobs.

Data flow data types

1. Array
2. Binary
3. Boolean
4. Complex
5. Decimal (includes precision)
6. Date

7. Float
8. Integer
9. Long
10. Map
11. Short
12. String
13. Timestamp

Mapping data flow allows to interactively see the results of each transformation step applied while building and debugging the data flows, the debug session can be used both in when building the data flow logic and running pipeline debug runs with data flow activities.

#### Schema Drift

Schema drift allows us to define data sources which should have mutable field names, data types, values and sizes

- Define transformation parameters which can work with data patterns instead of hard-coded fields and values
- Define expressions which understand patterns to match incoming fields, instead of using naming fields.

Azure Data Factory natively supports flexible schema which change from execution to execution which we can use to build generic data transformation logic without the need to recompile the data flows.

- There's an architectural decision in the data flow to accept the schema drift throughout the data flow.
- With schema drift, can protect against the schema changes from the sources.
- Through ADF, it treats the schema drift flows as late-binding flows, so while building the transformations, the drifted column names aren't being available in the schema views throughout the flow.
- Columns which are coming to the data flow from the source definition are defined as 'drifted' when they're not present in the source projection.
- We can view the source projection from the projection tab in the source transformation.
- when we select a dataset from the source, the service will automatically take the schema from the dataset and create the projection from the dataset schema definition.
- When the schema drift is enabled, all incoming fields are read from the source during execution and passed through the entire flow to the Sink.
- By default, all newly detected columns, known as drifted columns which arrive as a string data type, of we wish the data flow to automatically infer data types of drifted columns, check the drifted column types in the source settings.

#### Transforming drifted columns

When the data flow has the drifted columns, we can access it in the transformations with the method,

- Use the byPosition and byName expressions to explicitly reference a column by name or position number.
- Add a column pattern in a derived column or aggregate transformation to match on any combination of name, stream, position, origin or type.
- Add rule based mapping in a select or sink transformation to match drifted columns to columns aliases via a pattern.

#### Integration Runtime in ADF:

Integration runtime provides the compute capacity the ADF resources. Cpu, memory and disk utility for the infrastructure compute unit within the ADF integration runtime.

1. Azure IR (AutoResolveIntegrationRuntime) - all of the data copy, data movement, core Azure data transformation can be executed.
2. Self-hosted IR (Self-hosted Integration Runtime) - whenever, on-premise data sources or data sinks are required to be orchestrated in the ADF pipeline.

e.g. data copy from on-premise SQL db to Azure SQL db

Data copy from Azure Blob storage to on-premise SQL db

3. Azure - SSIS - used while transformation, copying data from on-premise SSIS db to Azure SQL db.

## **Hands-on Lab 5:**

### **Create reusable pipeline in ADF**

#### **Pre-requisites**

1. Create the ADF resource
2. Create Azure SQL db (with AdventureWorks sample db)
3. Create Azure Data Lake Storage (ADLS Gen2)
4. Create Linked Service for Azure SQL db and ADLS gen2

#### **Tasks:**

4. Create data flows using dynamic contents with reusable datasets
5. Explore the options to build the ADF reusable pipeline with dynamic contents to retrieve the data from SQL dataset (10 different tables) and sink to ADLS storage.
6. Used the generic pipeline to extract ten different Azure SQL tables with a single pipeline containing only three activities (Lookup, Foreach -> Execute Pipeline activity)

#### **Azure**

# Azure Data Lake Gen2

12 December 2022 19:56



# Azure SQL Database

12 December 2022 19:57

# Azure Blob Storage

12 December 2022 19:57

# Azure Analysis Services

12 December 2022 19:57

# Azure Synapse Analytics

12 December 2022 19:57

# Case Study

12 December 2022 19:57

# Apache Spark

12 December 2022 19:58

# Python Programming

12 December 2022 19:58

# Azure Databricks

12 December 2022 19:58



# Overview of AWS

12 December 2022

19:58

# Overview of GCP

12 December 2022

19:58

# Case Study

12 December 2022

19:58

# L1 Preparation

12 December 2022 19:59