



Azure Data Fundamental DP-900

HASSAN RAFIQUE MSc Data Science

- Describe core data concepts
-
- Identify considerations for relational data on Azure
-
- Describe considerations for working with non-relational data on Azure
-
- Describe an analytics workload on Azure
-
- Azure Synapse Analytics
-
- Azure cosmos DB
-
- Azure Data Factory
-
- Stream and Batch processing
-
- Data warehouse

Azure DP-900

1 - Explore Core Data Concept =

Key-value, Document DB, Graph DB

Transaction DB

Ex: Banking data.

It must deal to ACID.

✓ Atomicity:

- Execute in single unit
Fail or Pass.

✓ Consistency:

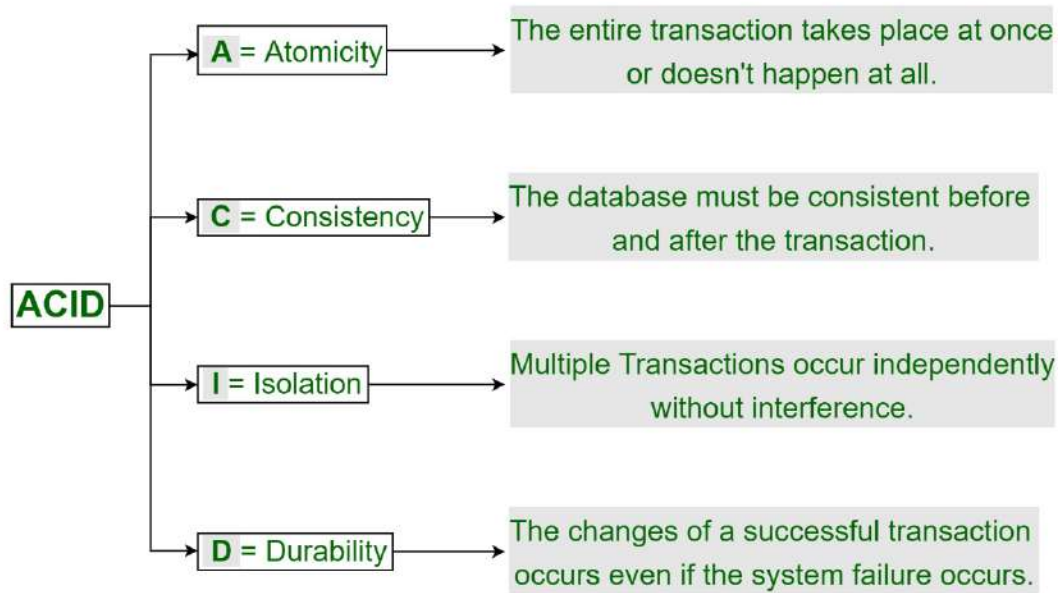
Take data from valid
state to another state.

state-1 → state-2

✓ Isolation: multiple transaction occur independent

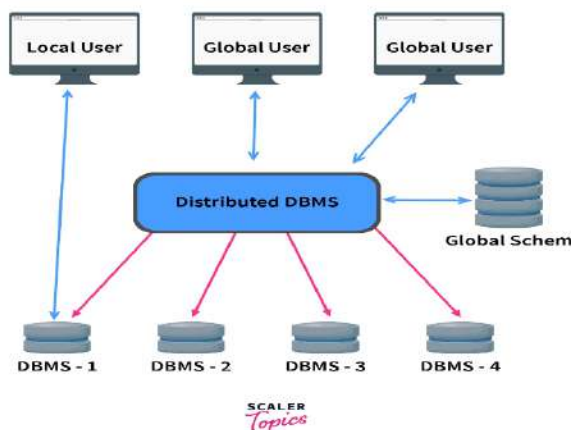
✓ **Durability:-** Committed transactions remain committed even system failure occurs.

ACID Properties in DBMS



Distributed DB.

Data store in different location. on same or different network.



Batch and Stream Process:-

★ data in the form of group is Batch.

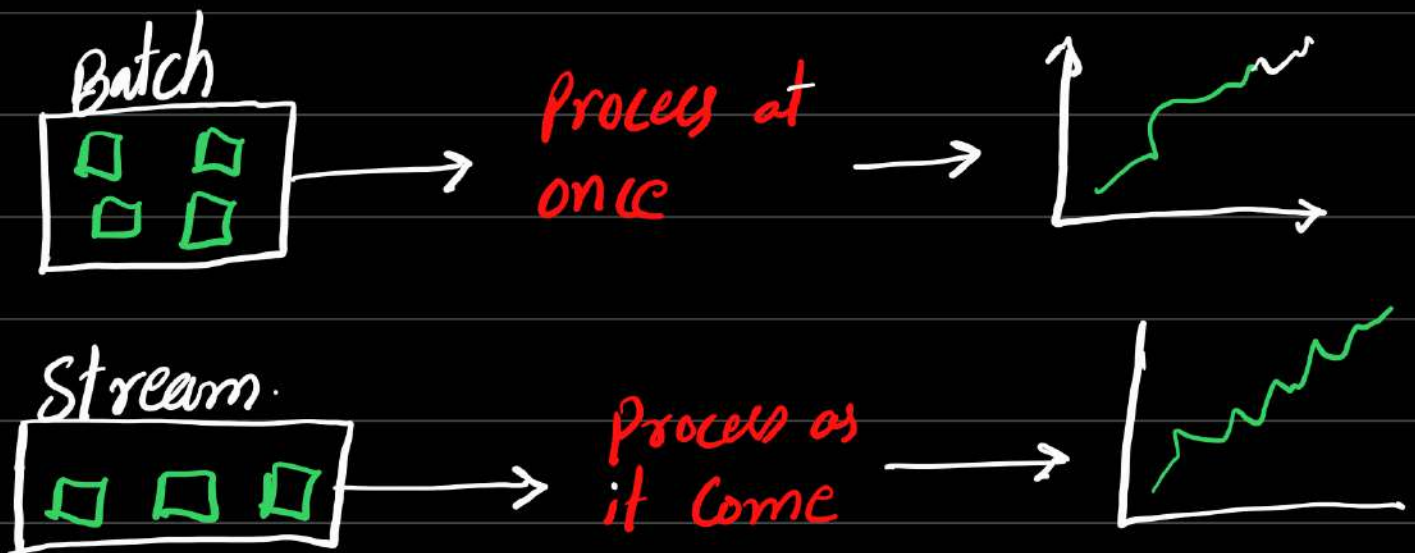
- apply OLAP.

- Data warehouse solution, monthly analysis

★ Stream OLTP. transactional data. New incoming data.

Ex: Stock Price

online gaming company
Real-estate website.

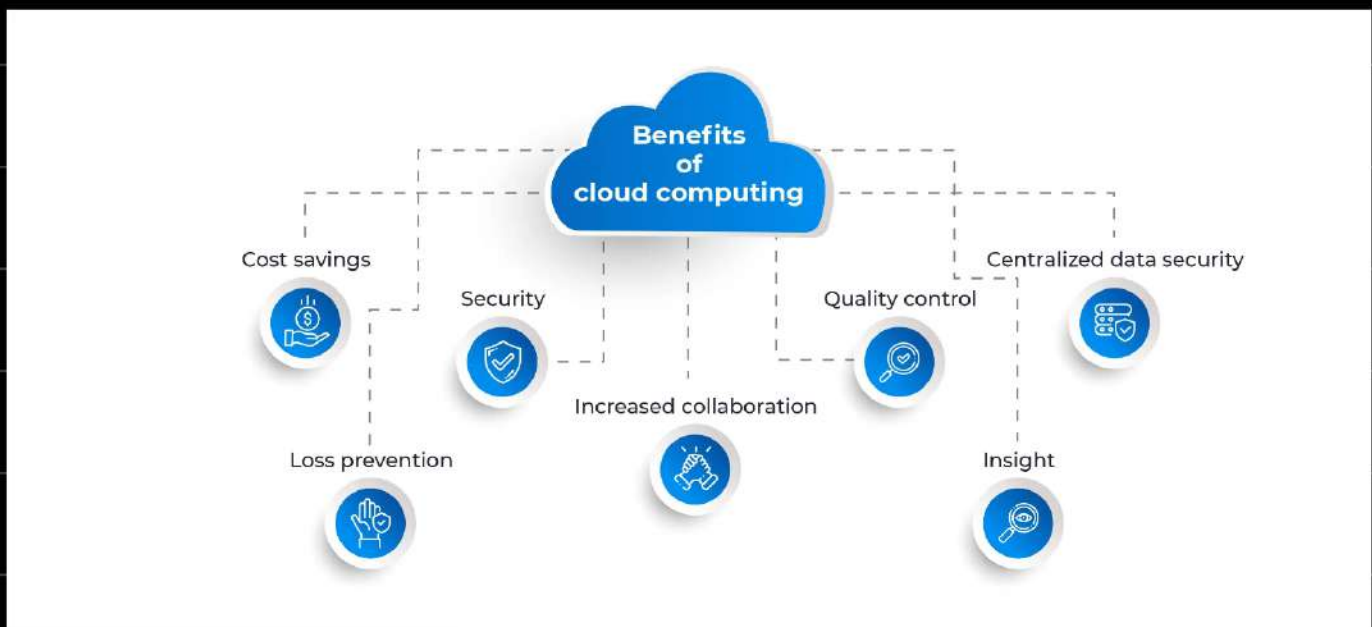


1 **Role and Responsibilities in Data:**

- Data Administrator
- Handle data
- Modify DB structure

- Control and monitoring
- Backup database
- Data replication.
- Install and upgrade DB server.

• EX:- Azure Data Studio
SQL server management studio.



Relational and Non relational data.

Entity, field, Table-Name
Relation DB structure is fixed format.

★ Primary key : Identity



* Foreign Key:-

Primary key in an other table called foreign key.

* Cluster Index / Non Cluster Index.

- Insert data in table with order Formate
- Fast search
- Primary key create cluster automatically.

ID	Name
1	A
2	B
3	C

* View:-

Immediate access to relative data.
A virtual Table.

```
"Create VIEW as "myview"  
Select "Name" From mytable;
```


Non-Relational

Azure Cosmos DB.

- semi-structure
- unstructured data } data structure.

JSON - semi-structure, XML,

• Avro - Apache -

Row based format. Header store JSON information store in Binary.

• Apache ORC

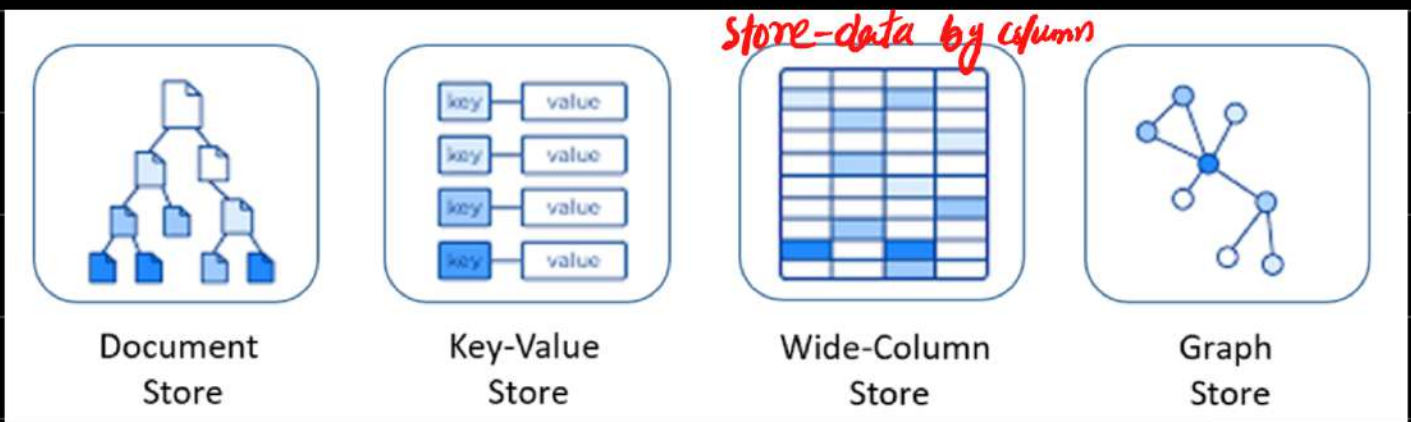
• Parquet :-

Column data base.

TYPES

Azure Cosmos DB.

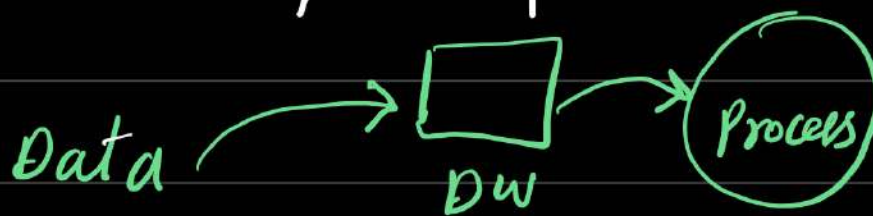
- key-value — Azure Table Store.
- Column — Column family — Apache Cassandra
- graph — neo4j,
- Document — Cosmos DB.



Ingestion → Process → Explore

ETL, ELT.

ELT use for complex data model.



Azure data Factory.

Cloud base - data transform tools.



2- **SQL:**

DBMS:-

on premiss — Cloud.

SQL server.

Stored procedure

linked server (SQL server we join to one server to another server)

PaaS - Azure data Services are available to create.

< SQL, Maria DB, Postgress > Relational-DB
my SQL server are popular tools.

On-premises Data $\xrightarrow{\text{migrate}}$ Azure Cloud Data

- Azure resources managed by server.

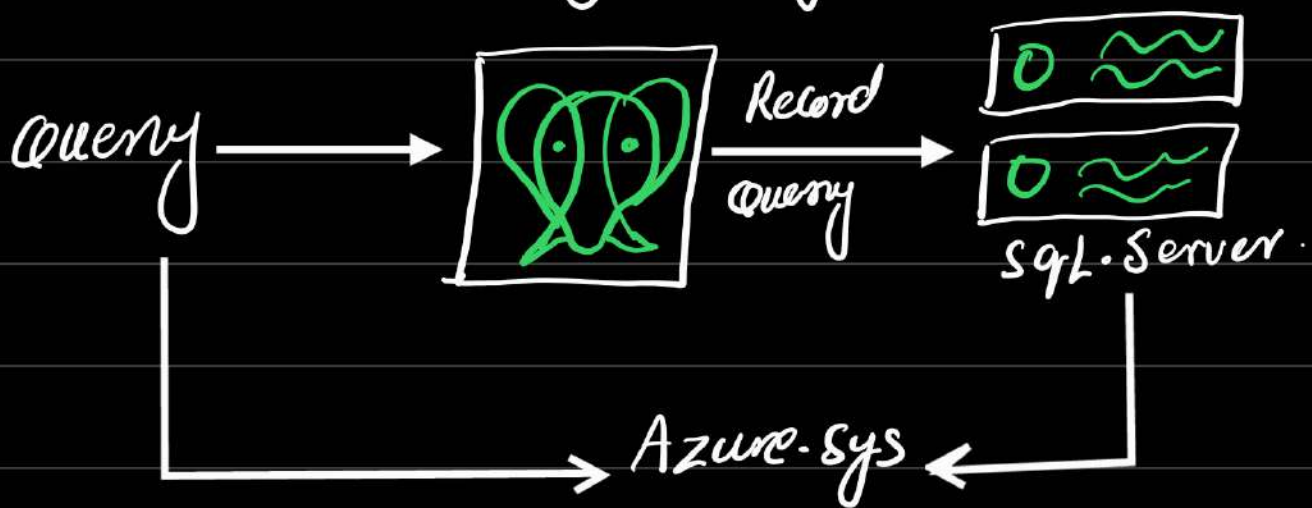
{ Availability
Scalability
Security, [AD, key vault]
threat detection
Automatic update

SQL server login $\xrightarrow{\hspace{1cm}}$ Azure AD

SQL server managed is good for connecting more server.

✓ **Postgres** (Pg-admin tool)

- Hybrid relational-object database
- Relational and non-relational store
- Store geometric data [2d-data]
- Pgsqll own Query Language.



- Support ultra high performance.
- Create ^{read} Replicas in Azure to 5 number **Virtual Network**.

Azure VM \longleftrightarrow $\langle \dots \rangle$ \longleftrightarrow Azure Service
AVN

RBAC:- Role base access control.

allow resources who are use Resource.

owner, Reader, Contributor.

Azure DB communicate over Port 1433.



3 Azure Cosmos DB

Azure database services for non relational.

- (1) Az Cosmos DB.
- (2) Table Storage
- (3) Blob Storage
- (4) File Storage [Continuous Availability]
[transfer legacy system]

✓ Table Storage.

- un-structure data store
- Easy and Scalable.
- Partition key, Rowkey.

✓ Blob Storage.

image, video, Storage

- ① - Block blob
- ② - Page - 8Tb - VM

③ - append -

Blob → Container → Blobs

✓ File Share:-

- Data replicated
- Scalability
- move data for legacy system to Cloud.
- automatically handle delta.
- Back-up File

✓ Cosmos DB -

- Multimodel NoSQL system. | store document data.
- API Cosmos DB support.
- SQL API
- Mongo db API
- Table API
- Cassandra API
- Gremlin API (graph API)

Document in Cosmos db

Partition-key — Container.

- automatically allocating space.
- Height availability
- Replicas

Non-Relational.

Provision:-

creation of providing and supply something like : Azure-portal

Cosmos DB → Container → data

Data Lake Storage:

- use to store large amount of data.

Security Component.
Firewall

Active directory (AD)

- user authentication.
- multi-user authentication.
- etc = mobile
- Add user.
- Access Control.

COSMOS-DB

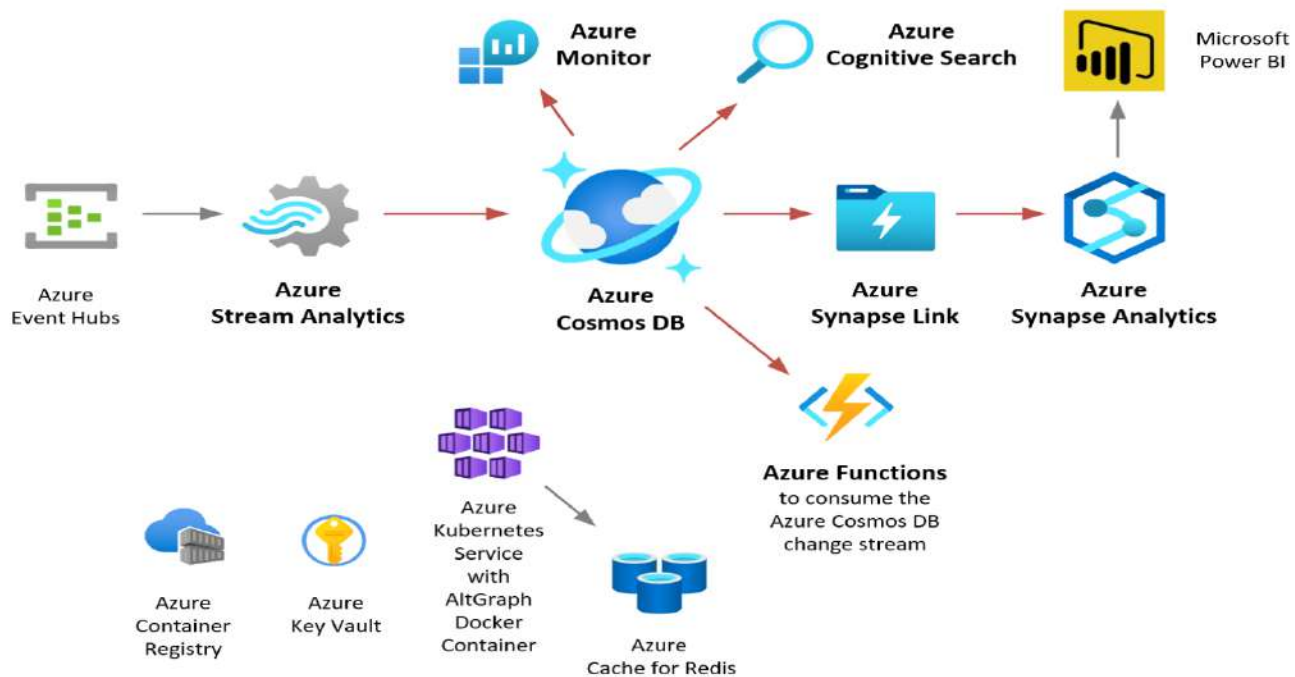
- Replication
 - in region - 4 time
 - outside region

(★)

- It manage data as set of document.
- select portion to get all related containers.

Portion key	container
<div style="border: 1px solid black; padding: 5px; display: inline-block;">1</div>	" _____ "
<div style="border: 1px solid black; padding: 5px; display: inline-block;">1</div>	" _____ "

(★) Native API of Cosmos DB is SQL API.



Data Warehouse:-

- Data ingestion
- Store
- PowerBI
- Data Factory. (handle ETL service)

Synapse:-

Data warehouse + Bigdata Analytics.

4 Azure service to build datawarhous
Solution.

Data Factory, Data bricks, Synapse
Power BI.

Stream and batch data
Data warehouse Store Structure data.



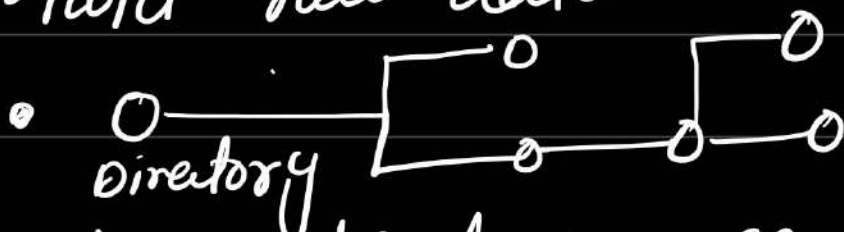
① Data Factory:

- ETL / ECT Process.
- Ingest data
- Extract Interesting Data.
- Data Integration Service.

② Data Lake

Azure data lake Gen(2)

- hold raw data.



Store data in
Folder in Folder.

— hierarchical Name Space.

Note - Blob only mimic data

- HDF

Data lake → Factory → Synapse → Report

③ Data bricks.

is an apache spark environment running on Azure, use for machine learning.

Data processing.

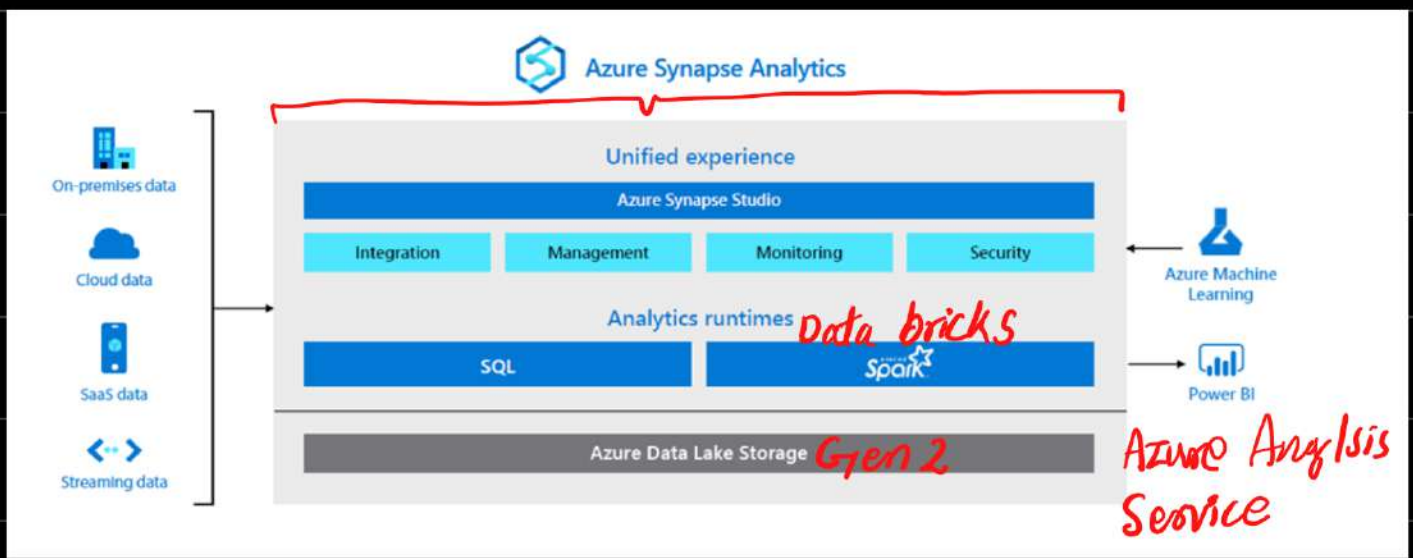
④ Synapse..

- ingest data from different source
- Data Aggregation



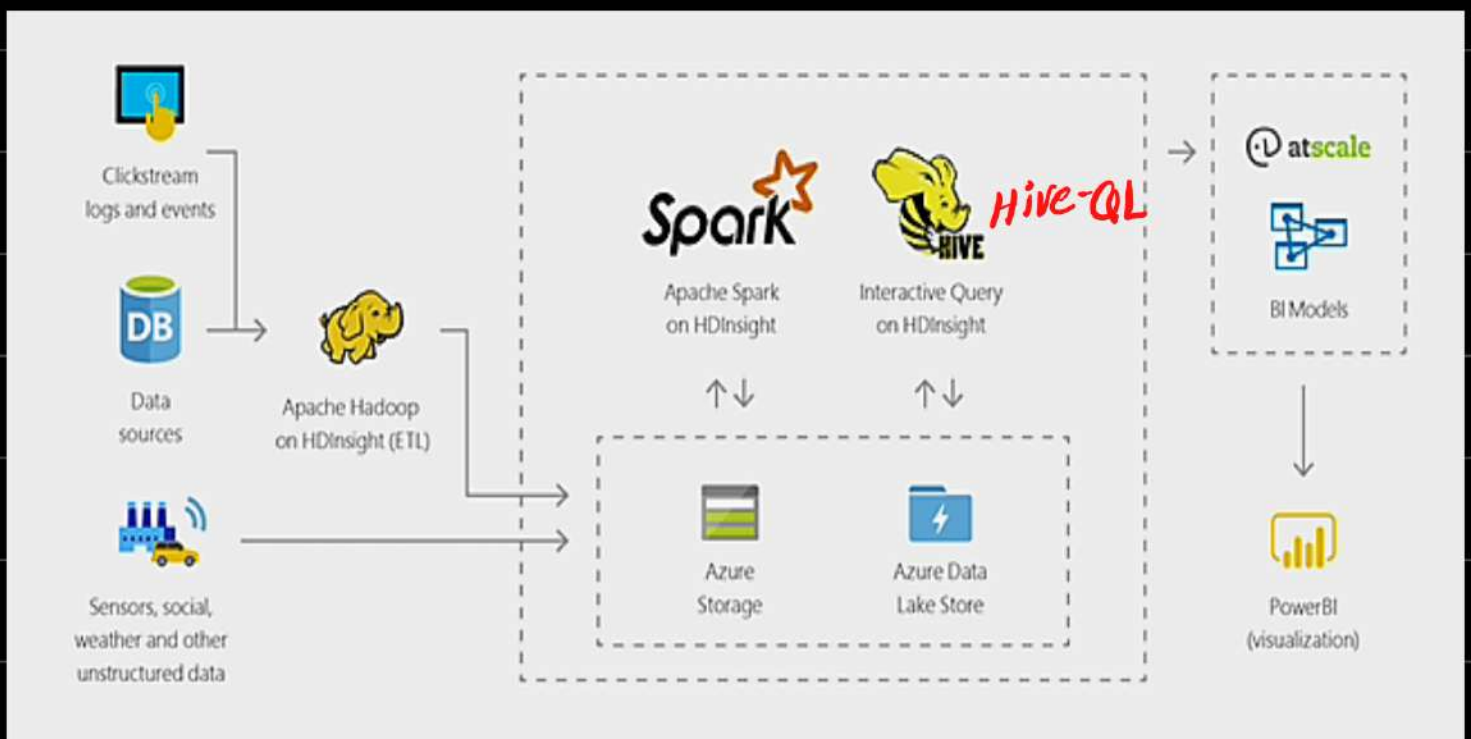
(★) Spark is in memory computation tool.

- ETL operation



Azure HD insight

- Big data processing
- implement cluster model.
- Support Streaming technology such as (Kafka).



Azure Data bricks:-

- parallel process.
- Streaming data
- use Apache Spark to distribute data.

Data Storage and processing:

- ☆ Azure Synapse Analytics.
 - ☆ Data Factory
 - ☆ Data bricks
 - ☆ Data Lake
- } Data process

☆ PolyBase:-

- Polybase is technique to make external data → look like sql table.
- Run Queries against these table directly.

[Pool use]

- Complex data
- data injest from many source

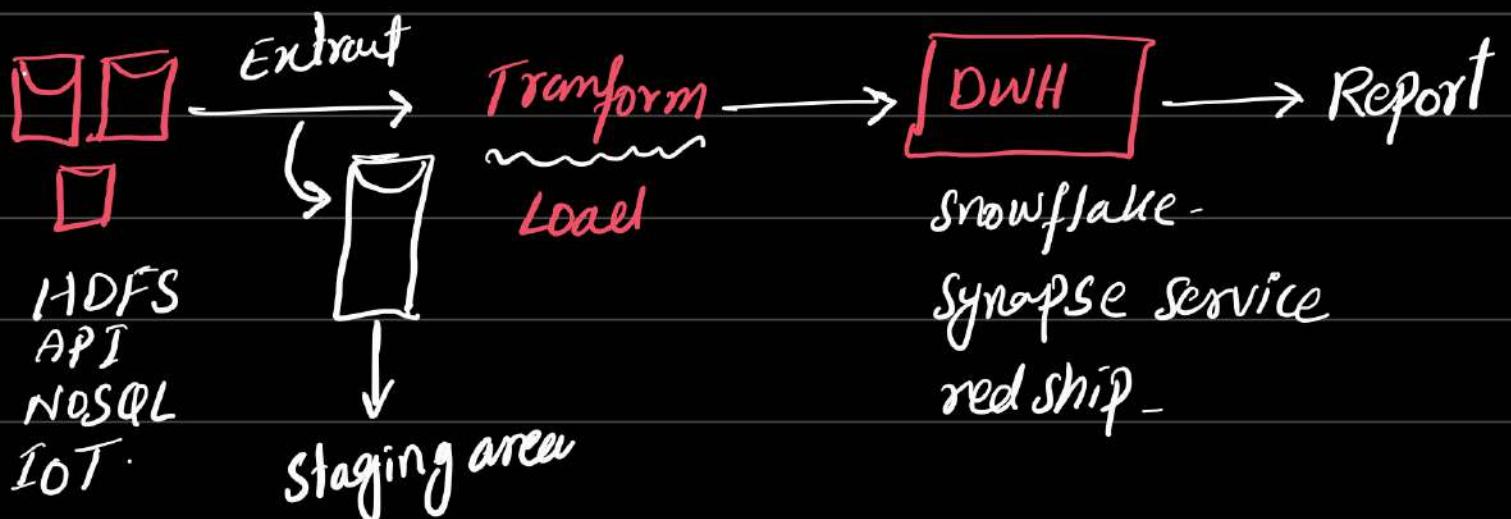
- **Synapse Spark** to write code base on Apache Spark.
To train ML model.
- In memory computing more faster.

Component of Synapse

Pipeline



Ineuron



[ELT]



{ HDFS , Amazon S3 }
{ Azure Blob }

Data Lake:-



Categorical data in DWH for specific purpose called data-mart.

OLAP / OLTP

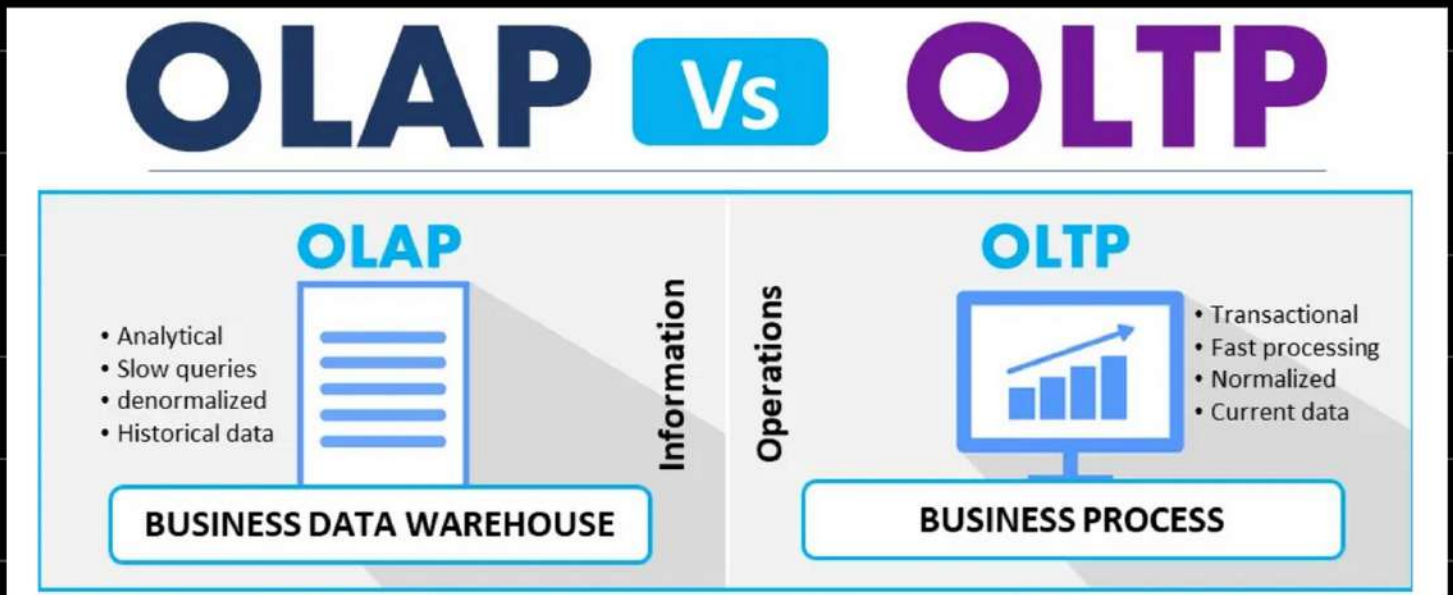
OLTP:-

It support ACID property and

- main goal to handle realtime data.
- RDMS
 - msql, oracle.

OLAP:-

- to handle large data.
- Analytics on History data
- Data warehouse
Hive, Azure synapse.

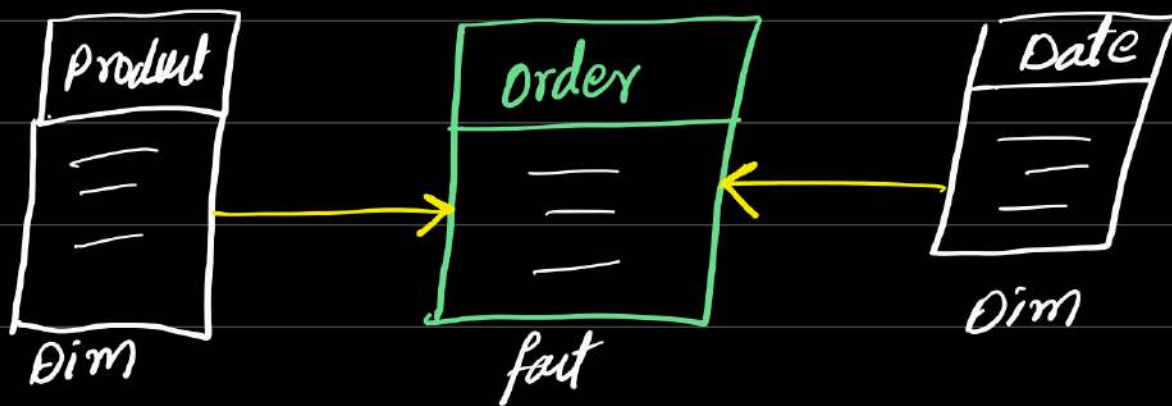


Important Terminology in warehouse.

- Dimension table
- Fact table

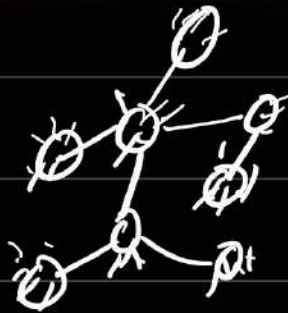
• Star & Snowflake.

Dimension table connected with fact table using foreign key.



Dimension is de-normalize table.

Star Schema vs Snowflake.



- Star Schema denormalize
- Snowflake normalized & less space because remove redundancy.
- Query Fast in Star Schema because less number of joins.