

## Traditional Data Warehouse

- ✓ Uses structured relational tables.
- ✓ Data is queried using SQL
- ✓ Enforcement of rigid schema-on-write (Data must conform to predefined structure)
- ✓ Optimized for structured data & complex queries
- ✓ Generally more expensive for large-scale storage.

Combines Benefits of both DW & DL.

## Data Lakehouse (modern Approach)

Microsoft Fabric

files are stored in Parquet Formats in Table as

### Benefits

- ✓ Allows row file storage & structured querying
- ✓ More flexibility
- ✓ more structured than pure data Lake.

Storage Layer: OneLake (ADLS Gen2)

Metadata Layer: Uses Delta Lake format

- ✓ schema enforcement
- ACID transaction
- Data Versioning (Time Travel)
- schema evolution
- SQL Query capabilities

## Data Lake Evolution

- ✓ Emerged as a response to "big data" challenge

3 V's

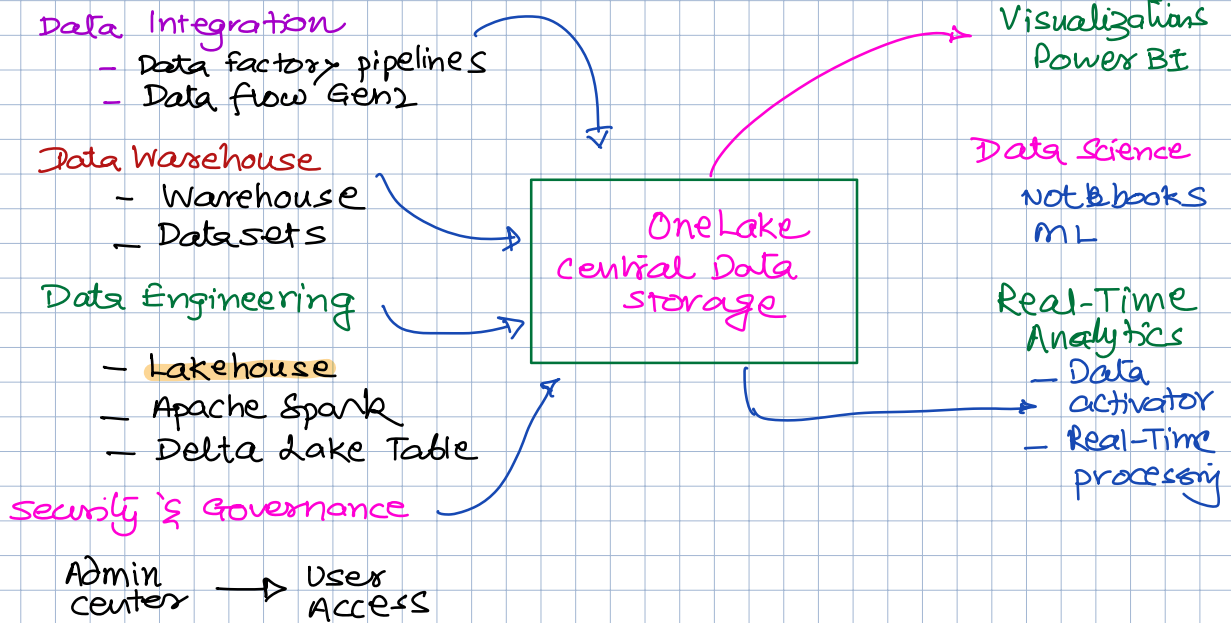
Volume:

Variety

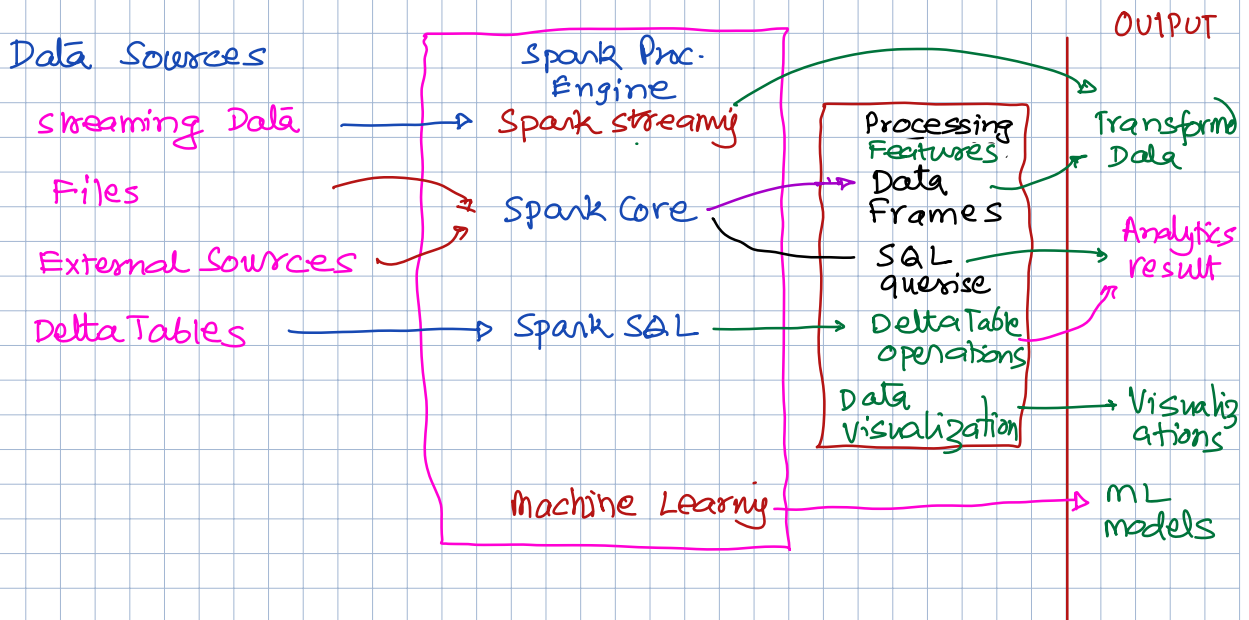
Velocity

- ✓ stores data in raw format.
- ✓ schema-on-read (No fixed schema)
- ✓ More cost effective for large scale storage
- ✓ Flexible but rather hard to query.

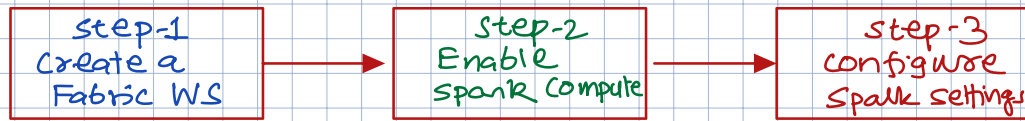
## Microsoft Fabric



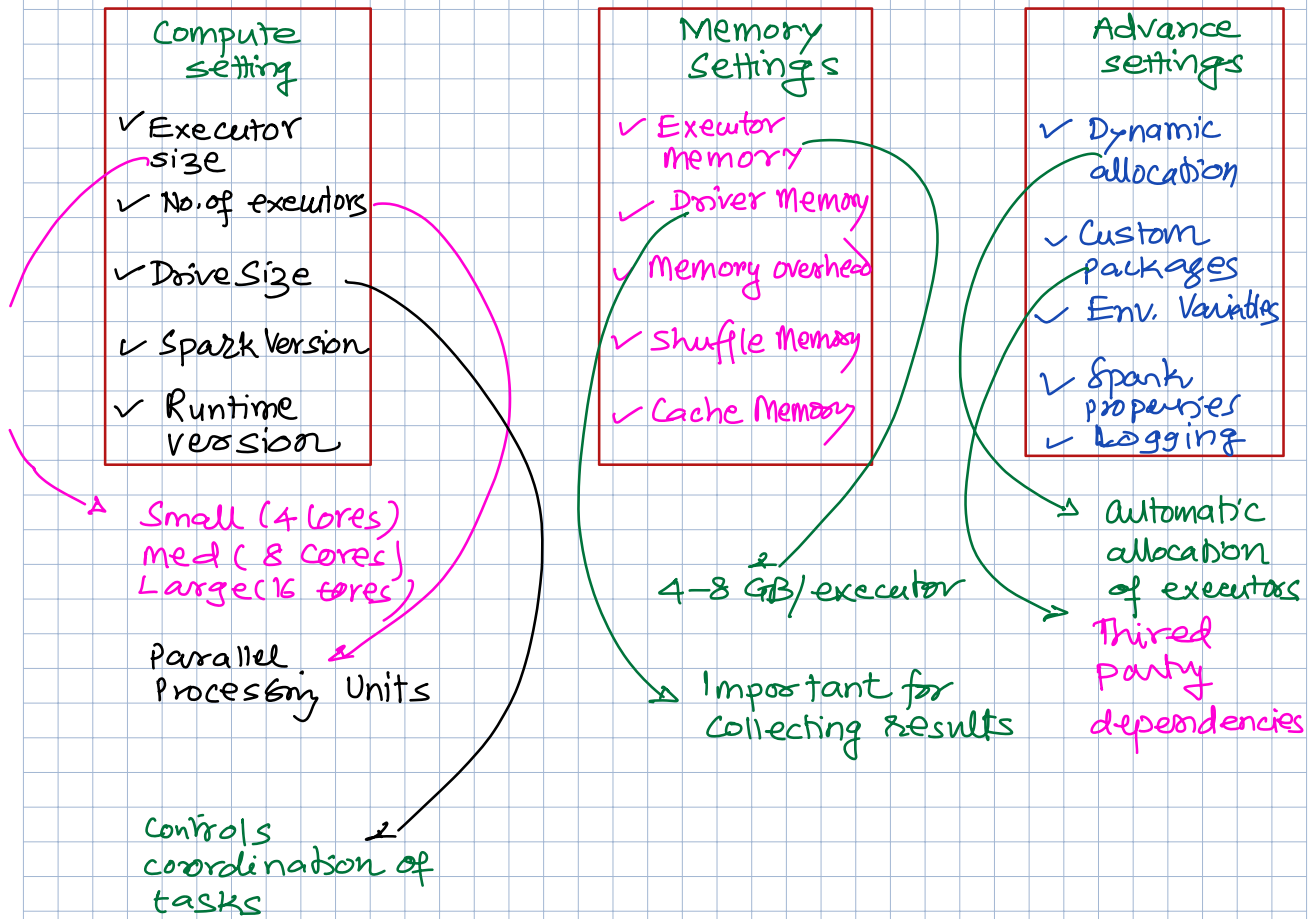
## Apache Spark on Fabric



## Spark Configuration



## Key Config. Settings



## Spark Notebooks Vs Jobs.

Exploratory work

**Notebook**

✓ Data Exploration

✓

✓ Model Development

✓ Documentation

Production-ready scheduled execution.

**Jobs**

Production ETL

✓ scheduled data processing

✓ Large-scale transformations

✓ Automated workflows.

Batch Processing

model Deployment

Integration  
works  
with many  
tools

Performance  
upto 10x faster  
than traditional

Key  
Benefits  
Using  
Spark

Flexibility  
multiple Data  
formats & sources

Scalability  
handle  
large data  
volumes

## Spark Data Frame

ID	Name	Age	City
1	Vijay	34	Blore
2	Neha	28	Chennai
3	Amit	36	Hyderabad

- ✓ Distributed collection of data organized into names & columns
- ✓ Immutable & lazily evaluated
- ✓ Supports - structured as well semi-structured data.

### Original Data Frame

ID  
Name  
Age  
City

### Filter

df.filter(age > 25)  
df.where(city == 'chennai')

### Select

df.select('name', 'age')  
df.selectExpr('age + 1')

### Group By

df.groupBy('city').  
agg(avg('age'))

### JOIN

df1.join(df2, df1.id == df2.id)

Output

DF