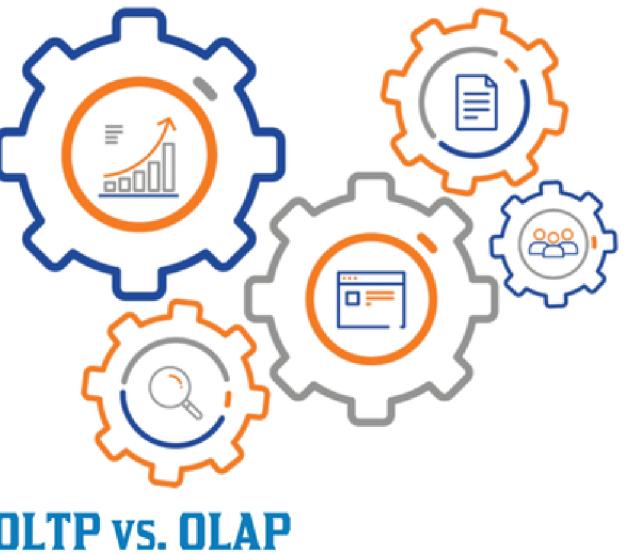




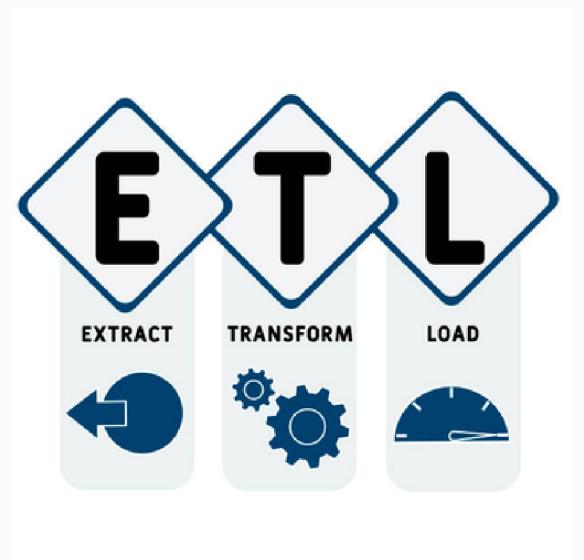
databricks



# DATA STORAGES, AND BASIC ETL PROCESS

Mentor/Intern: Nguyen Ngoc Thien / Tran Quoc Hai

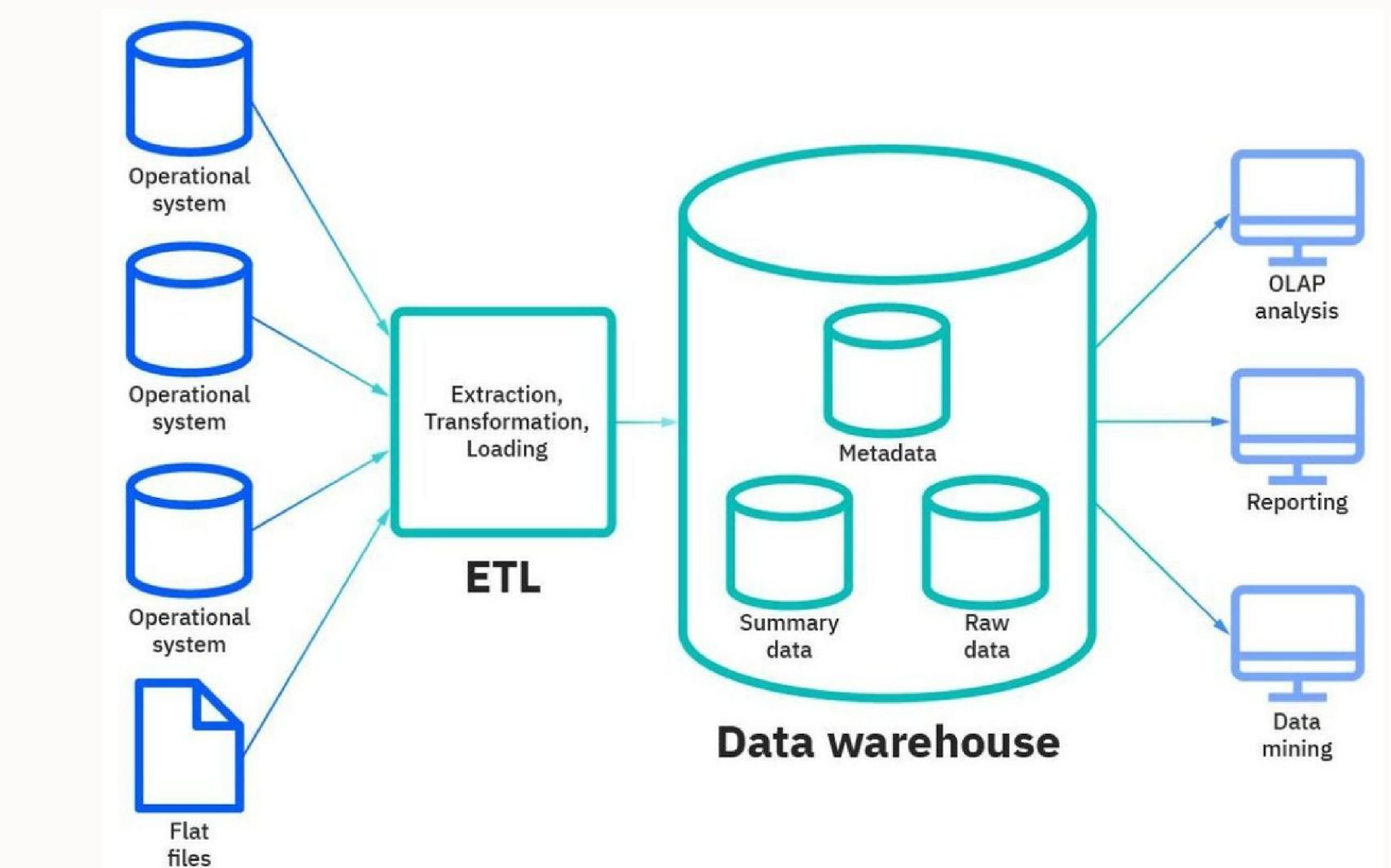
Le Trung Viet / Tran Thi Anh Thu  
Phan Thien Huu



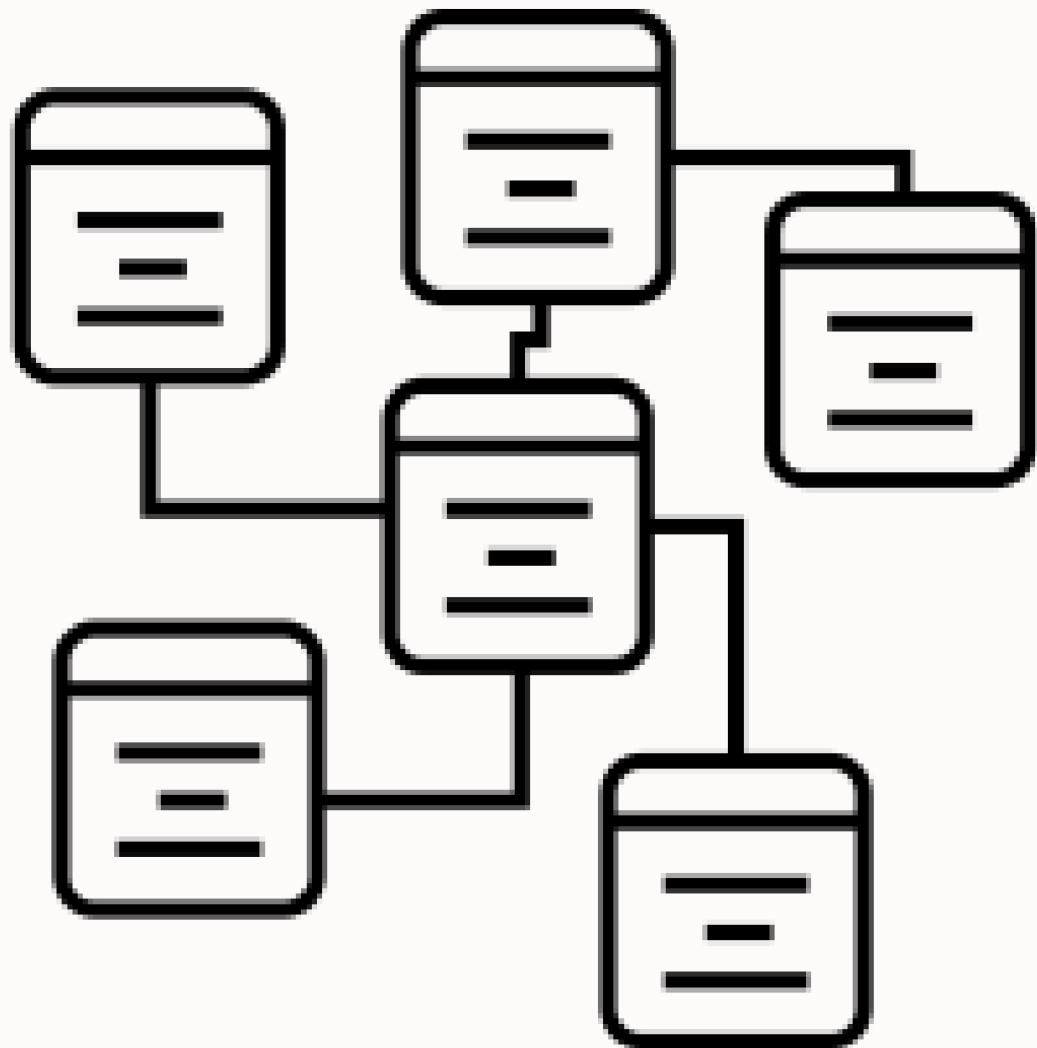
# TABLE OF CONTENTS



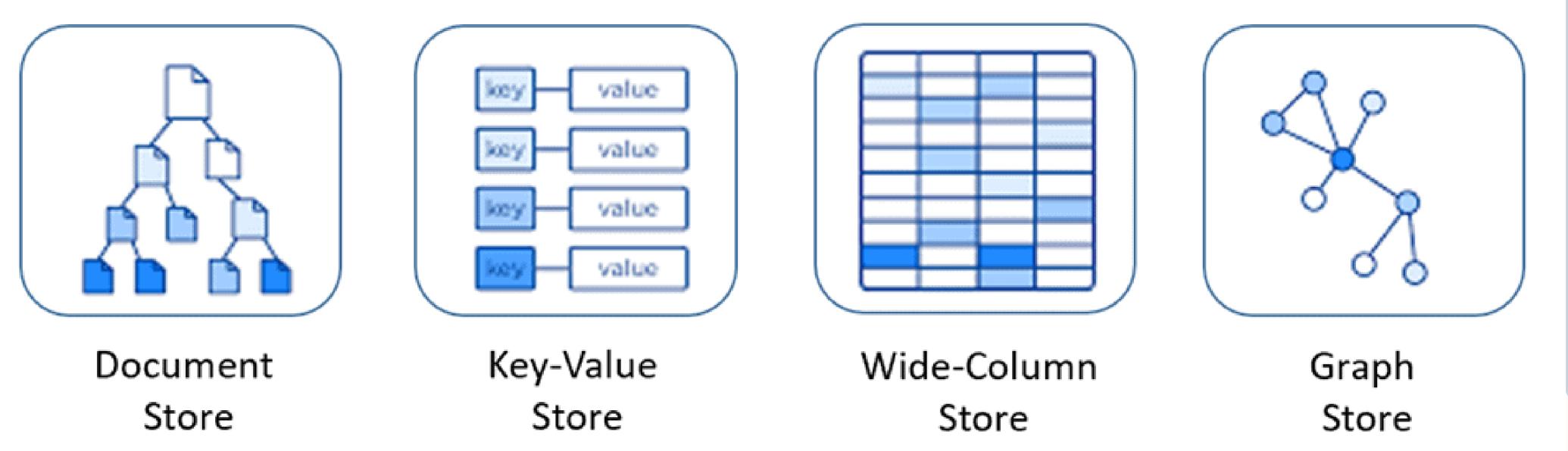
- 1 Database & Data Warehouse & Data Lake
- 2 OLTP & OLAP
- 3 ETL & ELT
- 4 DEMO



# DATABASE



Relational Database

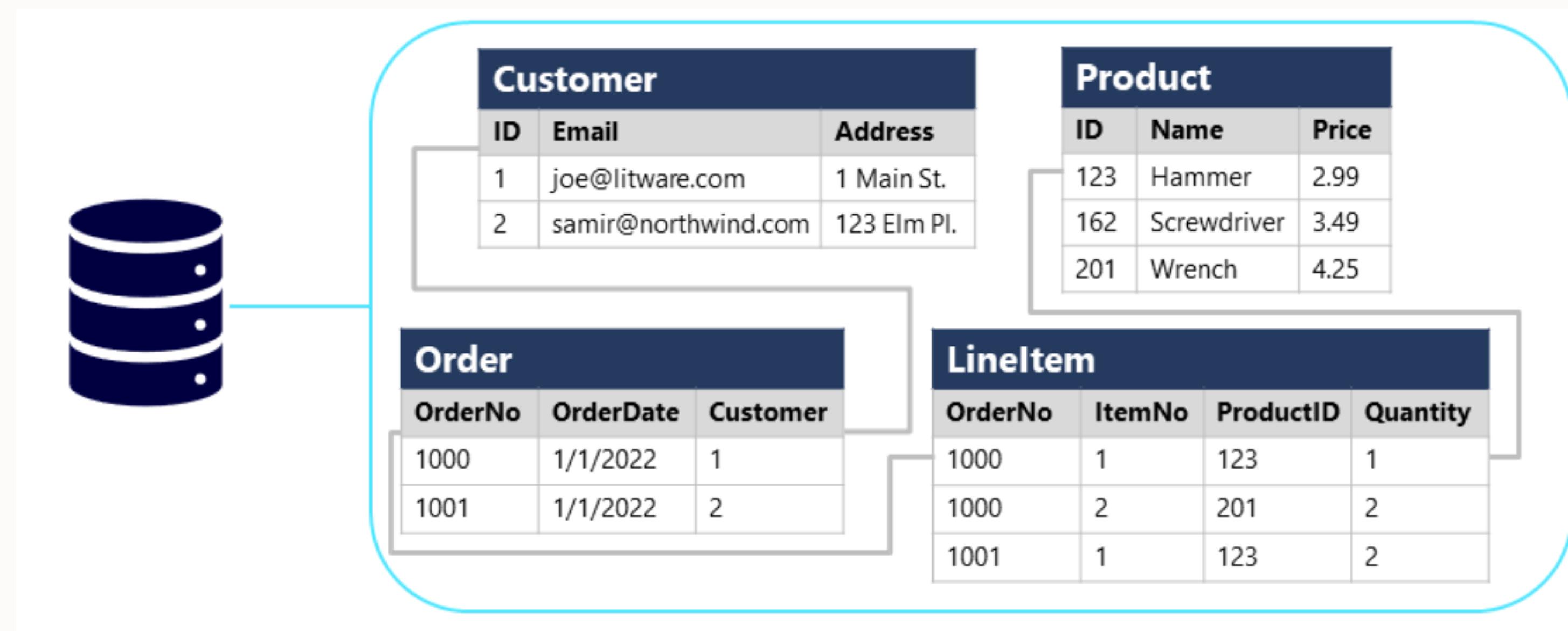


Non-Relational Database

# DATABASE



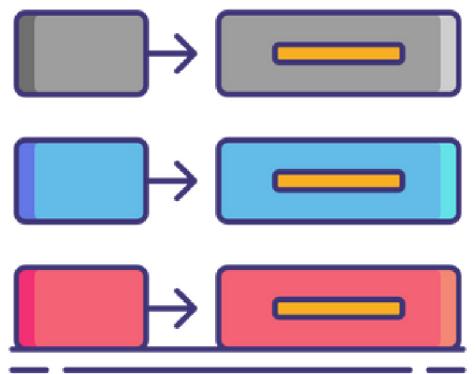
## Relational Database



# DATABASE

## Non-Relational Database

### 1. Key - Value



Products	
Key	Value
123	"Hammer (\$2.99)"
162	"Screwdriver (\$3.49)"
201	"Wrench (\$4.25)"

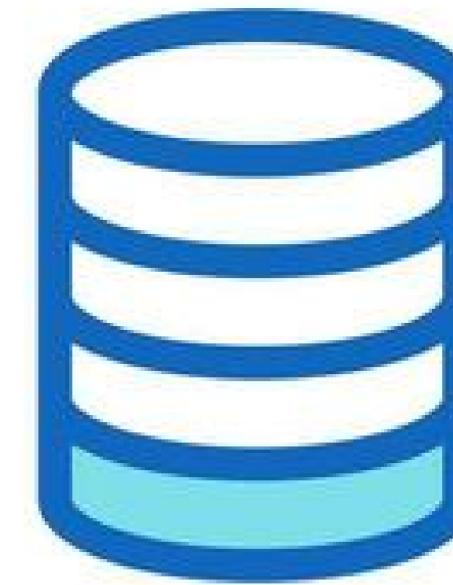
# DATABASE

## Non-Relational Database

### 2. Document

Customers	
Key	Document
1	{ "name": "Joe Jones", "email": "joe@litware.com" }
2	{ "name": "Samir Nadoy", "email": "Samir@northwind.com" }

# DATABASE



## Non-Relational Database

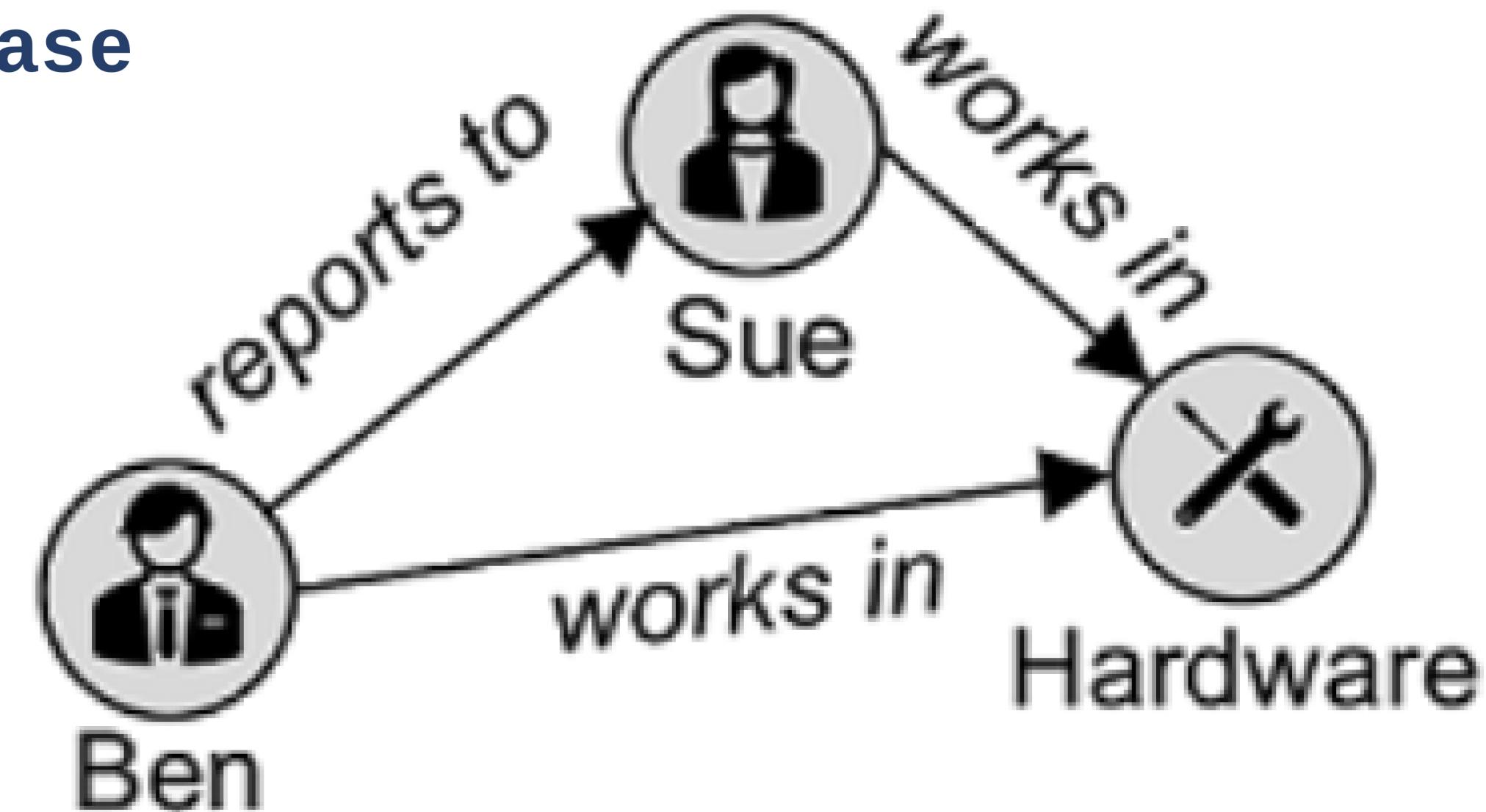
### 3. Column family

Orders				
Key	Customer		Product	
	Name	Address	Name	Price
1000	Joe Jones	1 Main St.	Hammer	2.99
1001	Samir Nadoy	123 Elm Pl.		4.25

# DATABASE

## Non-Relational Database

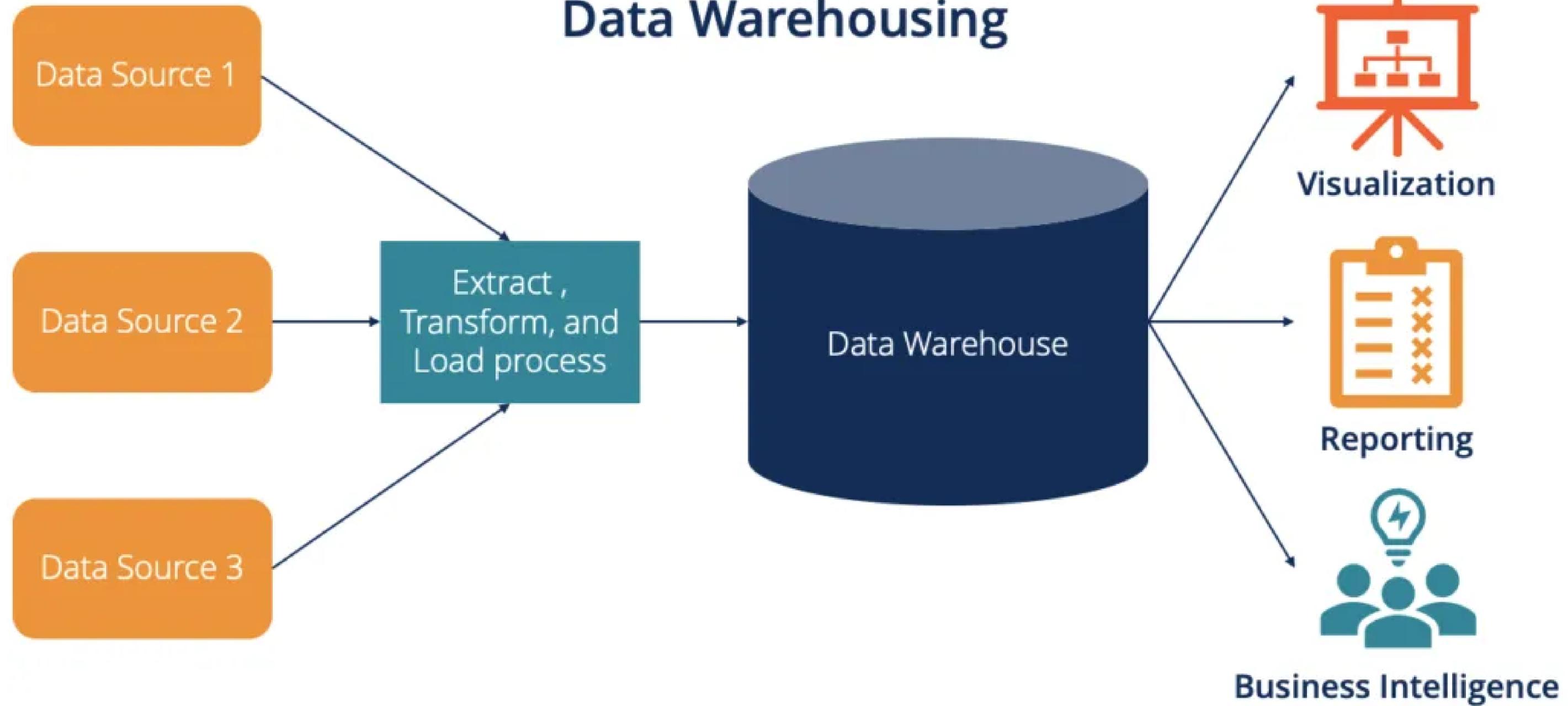
### 4. Graph



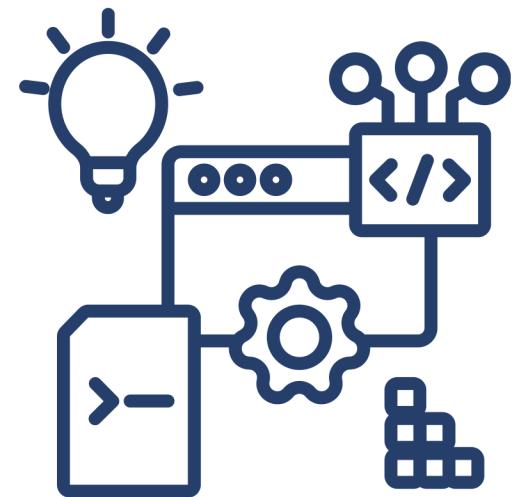
# DATA WAREHOUSE



## Data Warehousing



# DATA WAREHOUSE



**1**

Integrated

**2**

Subject-oriented

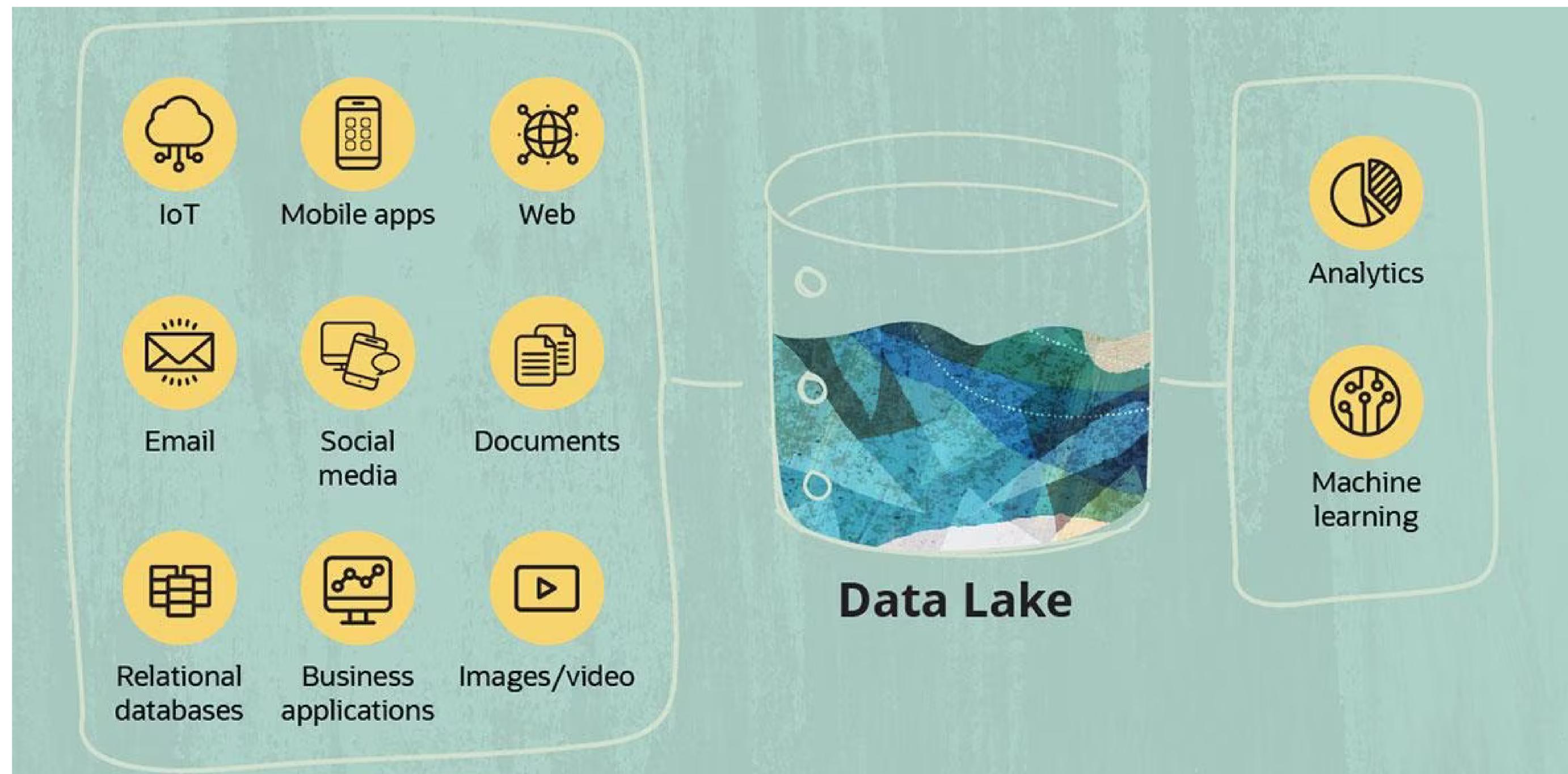
**3**

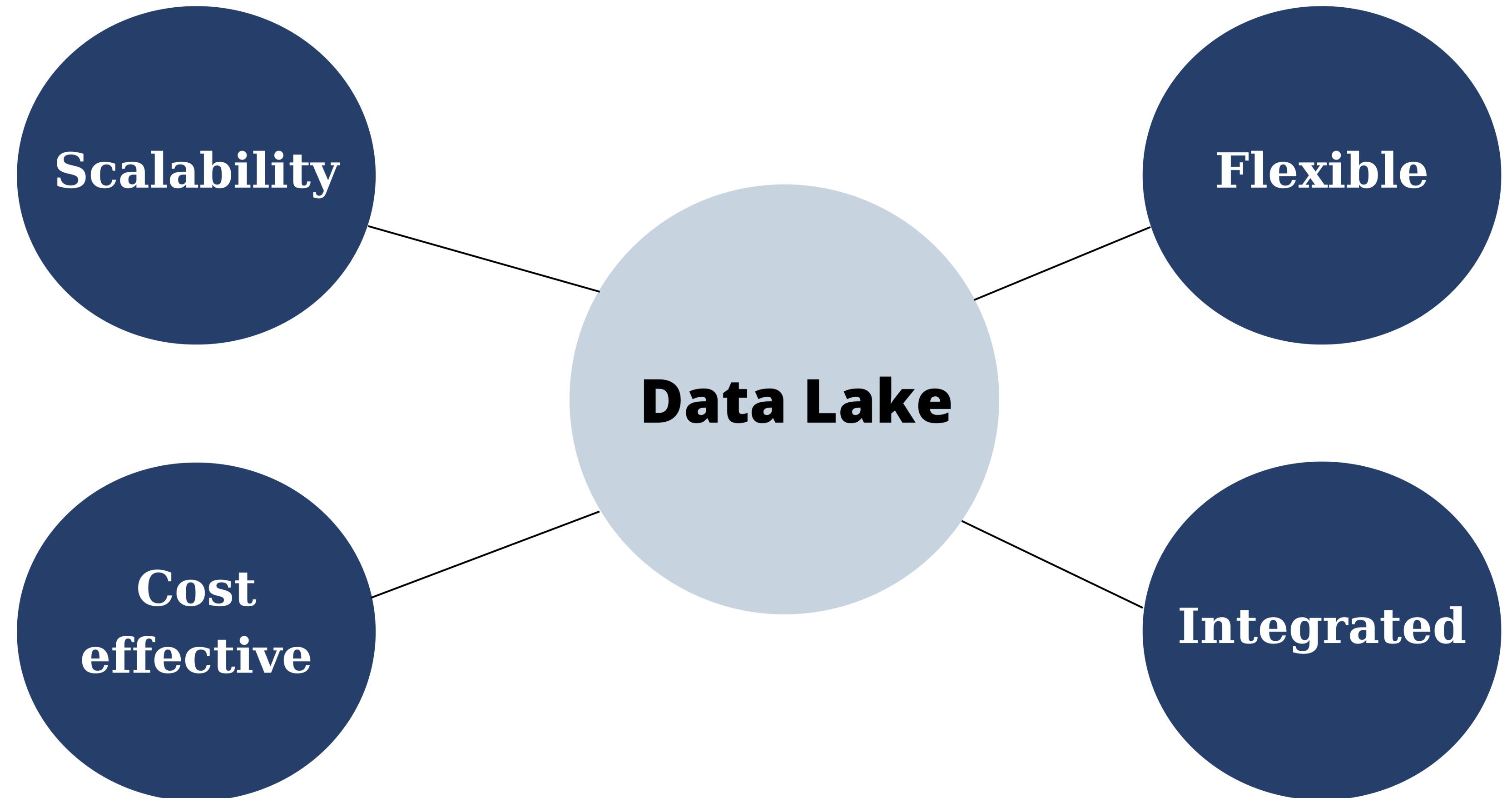
Time variant

**4**

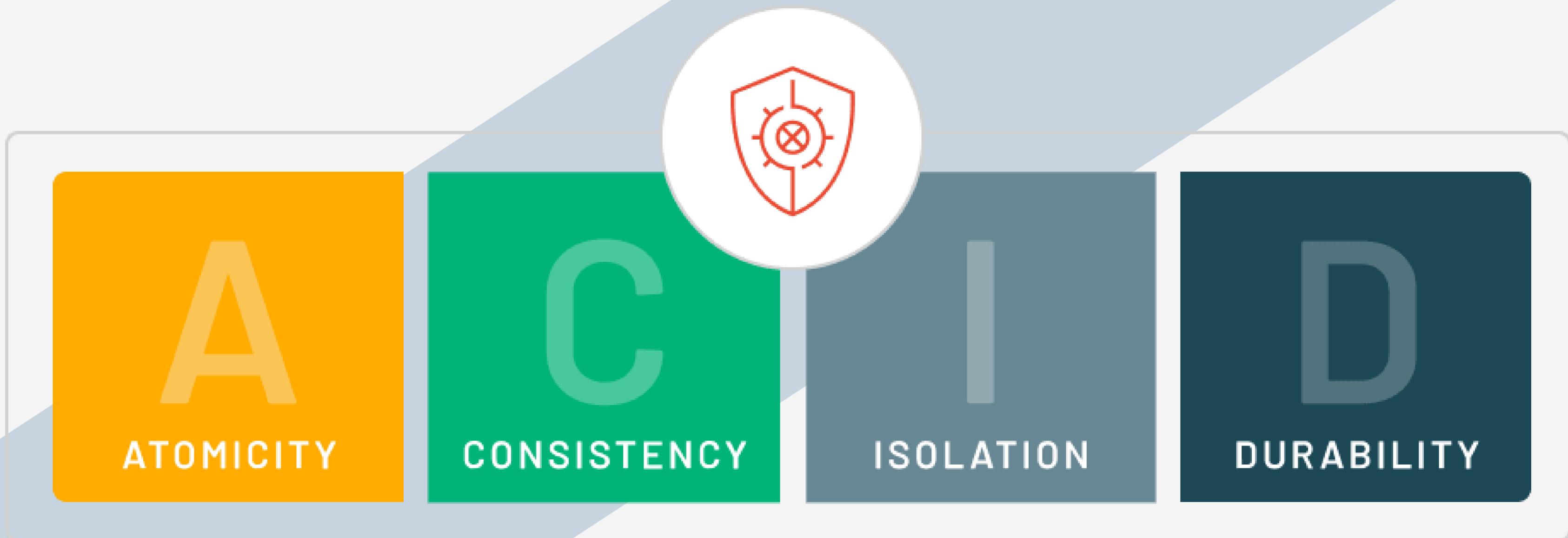
Non-volatile

# DATA LAKE

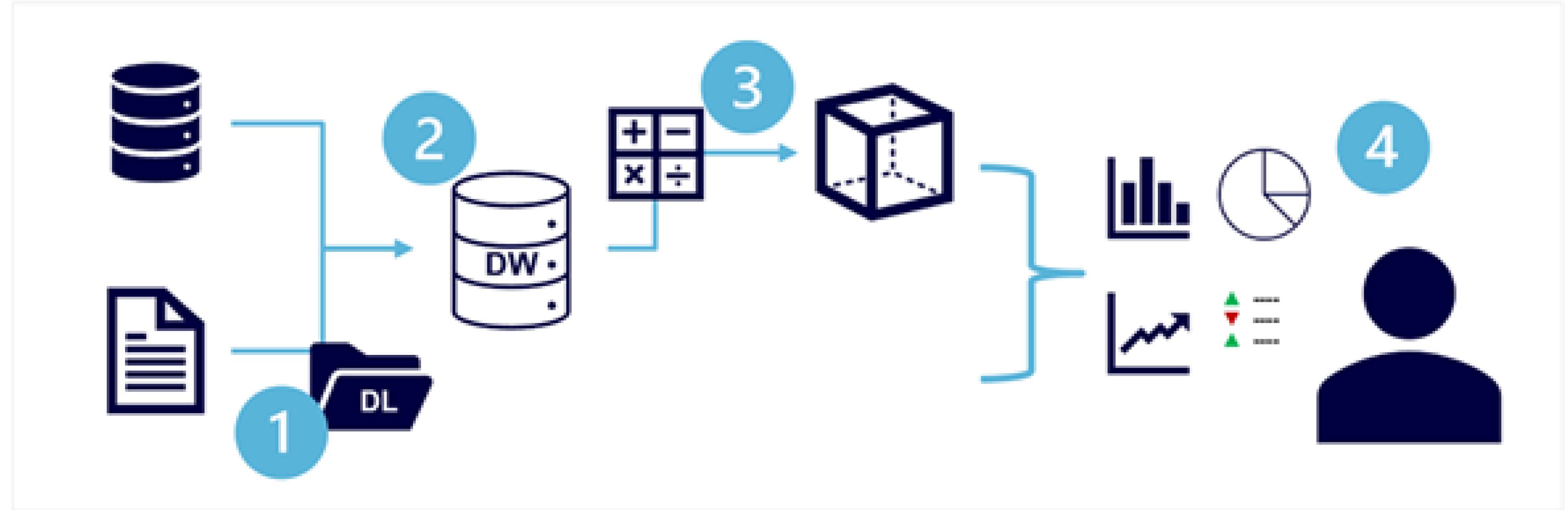




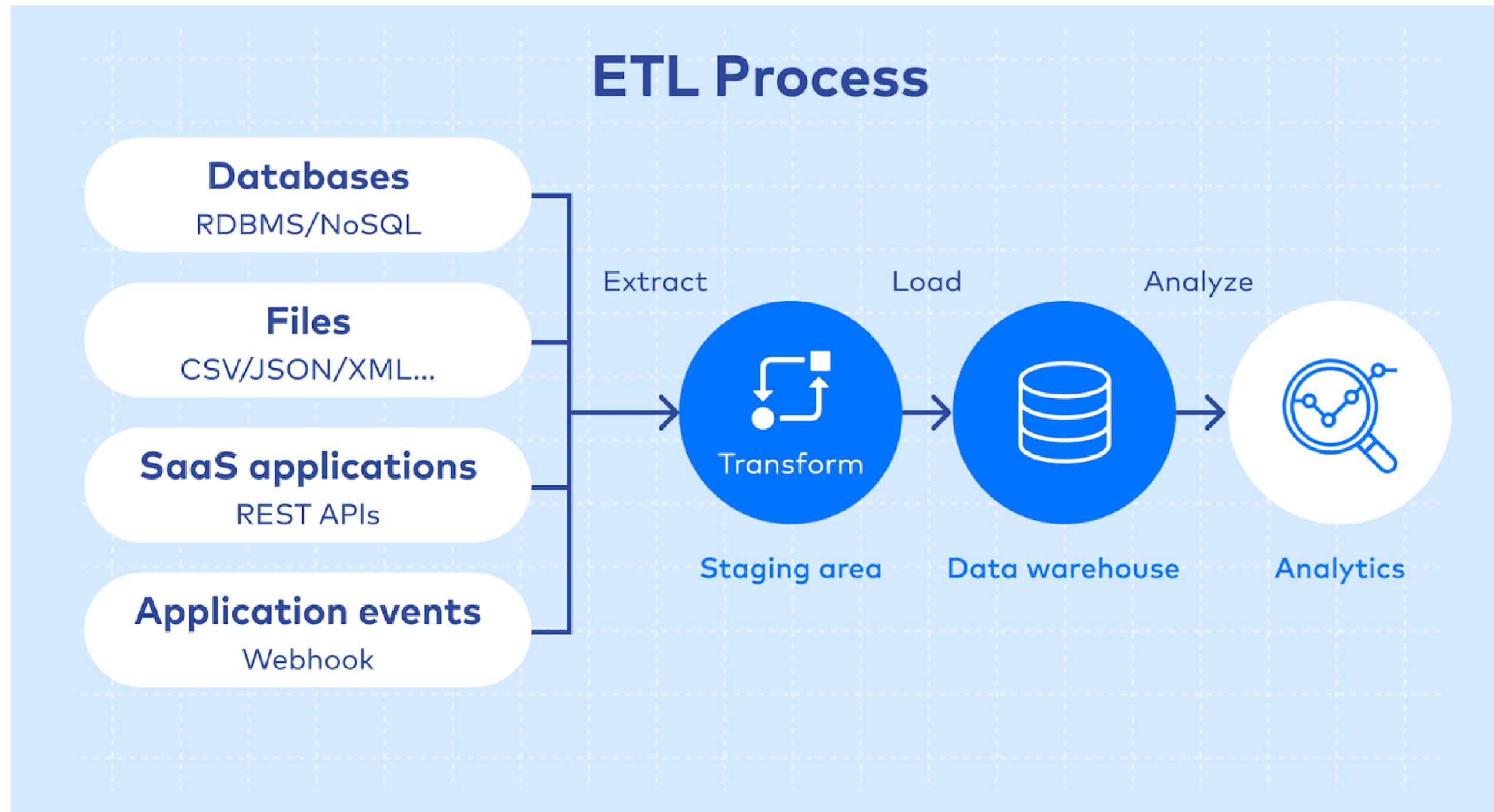
# OLTP



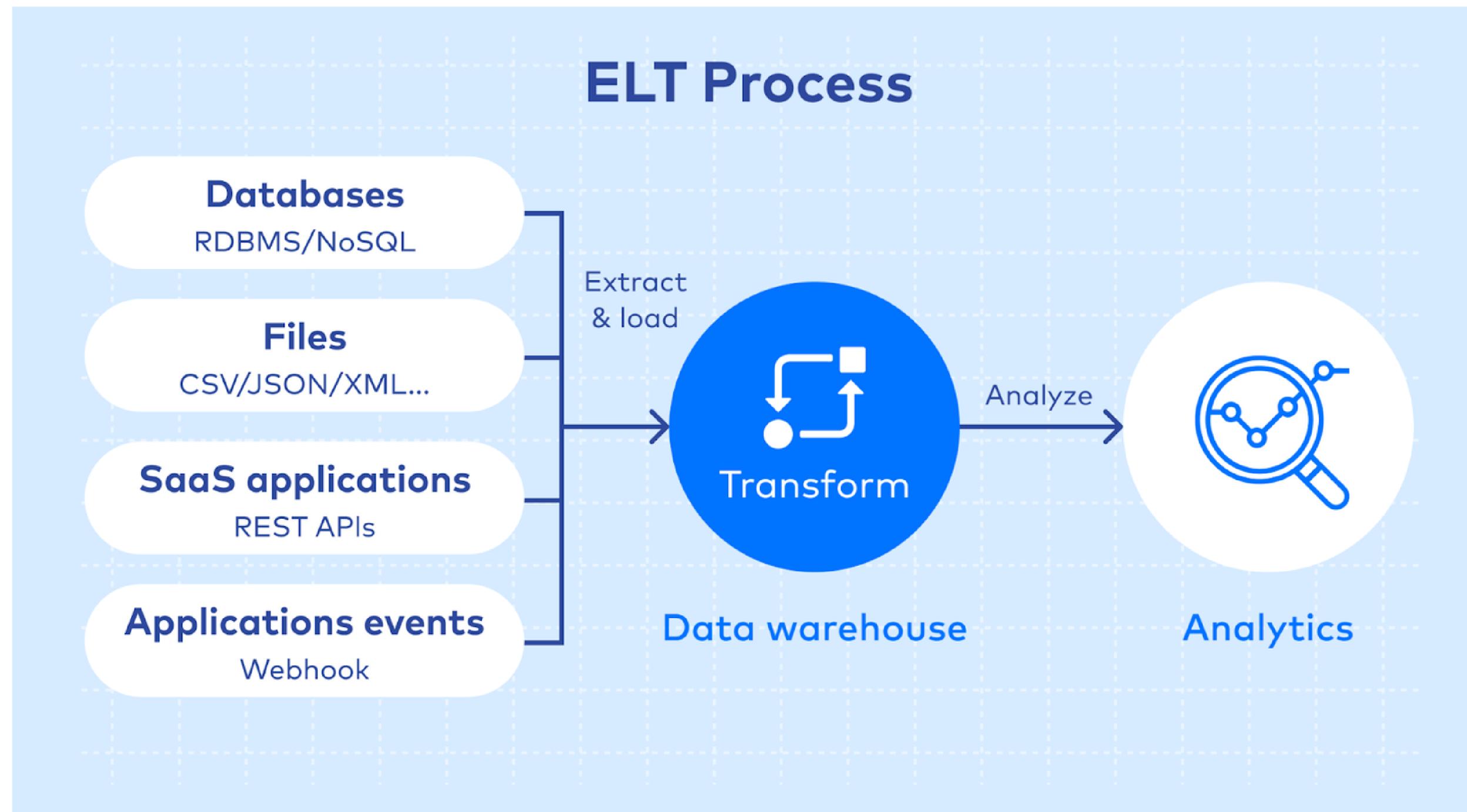
# OLAP



# ETL PROCESS



# ELT PROCESS



# Brief Comparison

---

	ETL	ELT
<b>Privacy</b>	Yes	No
<b>Load speed</b>	Slow	Fast
<b>Powerful target system</b>	No	Yes
<b>Capital for hosting required</b>	Yes	No
<b>Data Lake support</b>	No	Yes

# DATA

## Tokyo 2020 Olympic Summer Games



Athletes



Coaches



Medals

# Athletes

name	short_name	gender	birth_date	birth_place	birth_country	country	country_code	discipline	discipline_code	residence	residence_code	height_m	url
AALERUD	AALERUD	Female	12/4/1994	VESTBY	Norway	Norway	NOR	Cycling Road	CRD				.../..../en/i
ABAD Nes	ABAD N	Male	3/29/1993	ALCOI	Spain	Spain	ESP	Artistic Gymnastics	GAR	MADRID	Spain	1.65/5'4"	.../..../en/i
ABAGNALI	ABAGNALI	Male	1/11/1995	GRAGNAN	Italy	Italy	ITA	Rowing	ROW	SABAUDIA	Italy	1.98/6'5"	.../..../en/i
ABALDE A	ABALDE A	Male	12/15/1995	FERROL	Spain	Spain	ESP	Basketball	BKB			2.00/6'6"	.../..../en/i

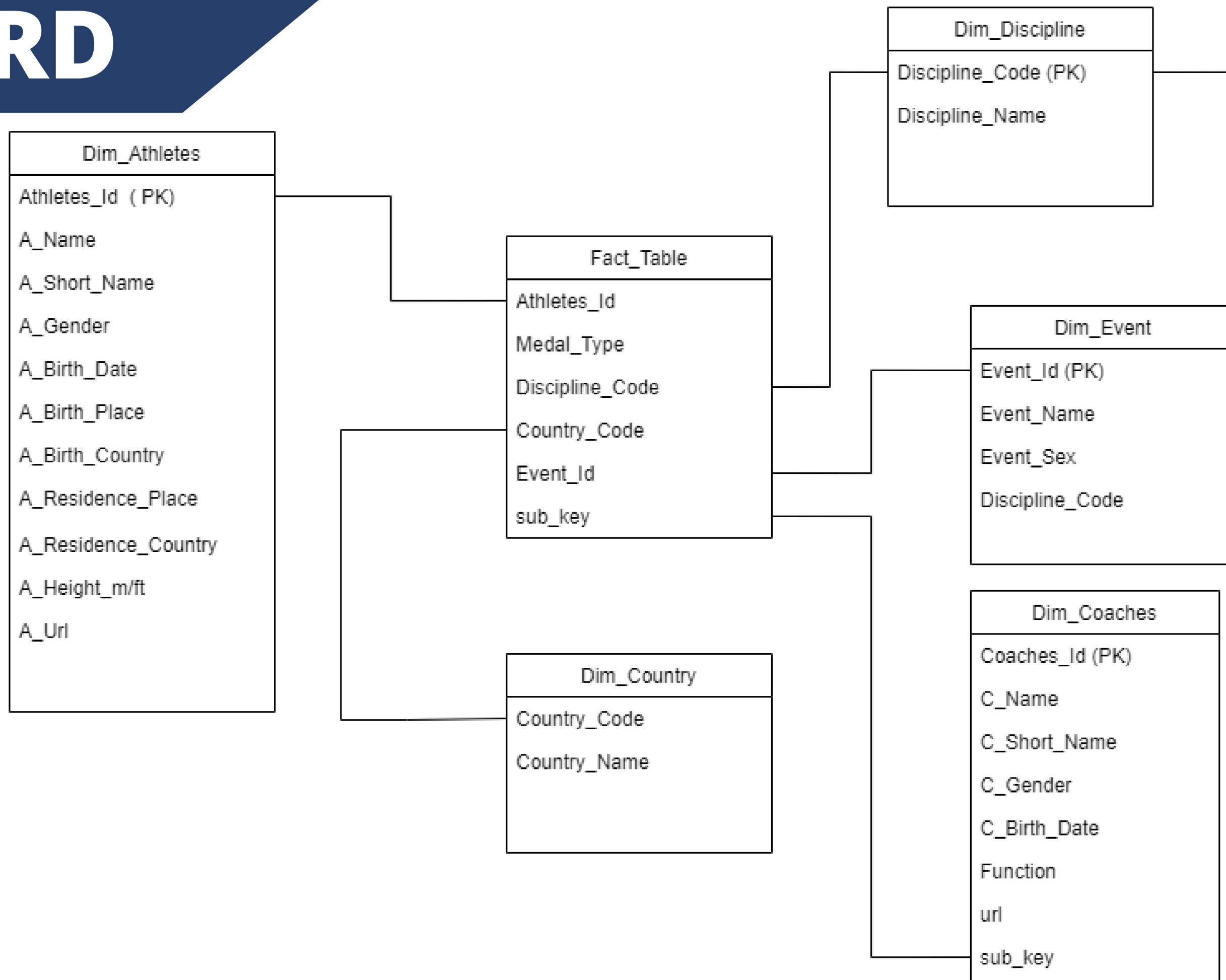
# Medals

	medal_type	medal_code	medal_date	athlete_short_name	athlete_name	athlete_sex	athlete_living	country_code	discipline_code	event	country	discipline
0	Gold Medal	1	00:00:00	KIM JD	KIM Je Deok		.../..../en/i	KOR	ARC	Mixed Team	Republic of Korea	Archery
1	Gold Medal	1	00:00:00	ANS	AN San	X	.../..../en/i	KOR	ARC	Mixed Team	Republic of Korea	Archery
2	Silver Medal	2	00:00:00	SCHLOESSER G	SCHLOESSER Gx		.../..../en/i	NED	ARC	Mixed Team	Netherlands	Archery

# Coaches

name	short_name	gender	birth_date	country_code	discipline	function	event	url
ABDELMA	ABDELMAGID W	Male	8/2/1982	EGY	Football	Head Coach		.../..../en/i
ABE Junya	ABE J	Male	7/25/1990	JPN	Volleyball	Head Coach		.../..../en/i
ABE Katsu	ABE K	Male	9/23/1979	JPN	Basketball	Coach		.../..../en/i

# ERD



```

1 import pyspark
2 from pyspark.sql import SparkSession
3 from pyspark.sql.functions import *
4 from pyspark.sql.types import *
5 spark = SparkSession.builder.getOrCreate()
6
7 df = spark.read.csv('/FileStore/tables/medals2.csv', header=True, inferSchema=True)
8 display(df)

```

▶ df: pyspark.sql.dataframe.DataFrame = [0 rows, 10 columns, 100% complete, 100% memory] [more fields]

# DEMO

Table ▾ +

_c0	medal_type	medal_code	medal_date	athlete_short_name	athlete_name
1	Gold Medal	1	2021-07-24T00:00:00.000+0000	JUD	KIM Je Deok
2	Gold Medal	1	2021-07-24T00:00:00.000+0000	SAN	AN San
3	Silver Medal	2	2021-07-24T00:00:00.000+0000	SCHLOESSER G	SCHLOESSER Gabriela
4	Silver Medal	2	2021-07-24T00:00:00.000+0000	WIJLER S	WIJLER Steve
5	Bronze Medal	3	2021-07-24T00:00:00.000+0000	ALVAREZ L	ALVAREZ Luis
6	Bronze Medal	3	2021-07-24T00:00:00.000+0000	VALENCIA A	VALENCIA Alejandra
5	Gold Medal	1	2021-07-24T00:00:00.000+0000	CARAPAZ R	CARAPAZ Richard

# THANK YOU

