

Why ?

- ↳ operational data keeping (OLTP) online Transaction Processing
 - ↳ Analytical decision making (OLAP) online Analytical Processing
-

Requirements

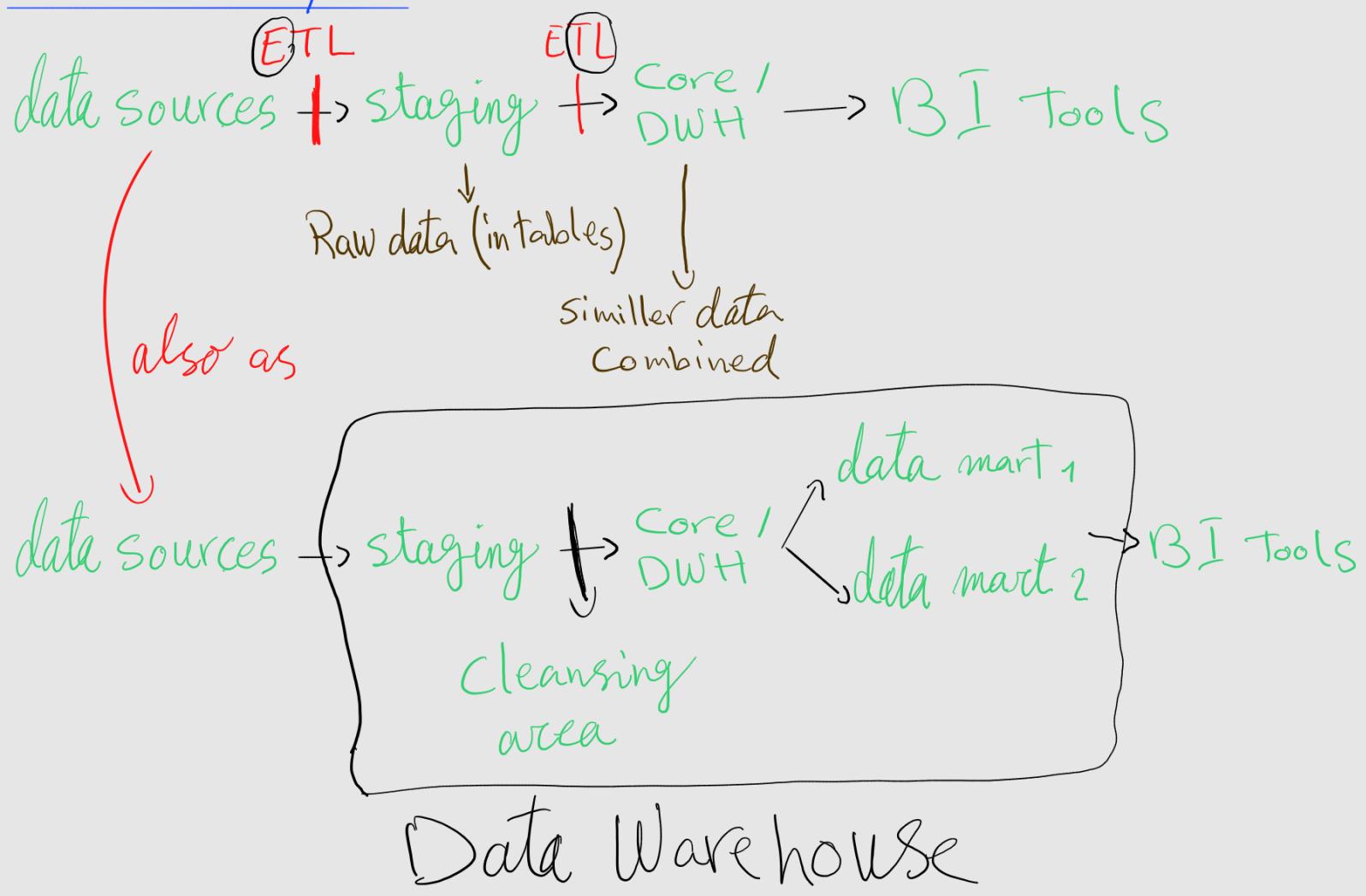
- OLTP:
 - one record at a time
 - Data input
 - no long history
 - OLAP:
 - Thousands of records at a time
 - Fast query Performance
 - Historical context
 - Usability
-

Why data lake?

- ↳ stores raw data without process
 - ↳ " Big data"
 - ↳ unstructured
 - ↳ not ready to be used
-

DWh architecture ↴

DWH layers



Staging area

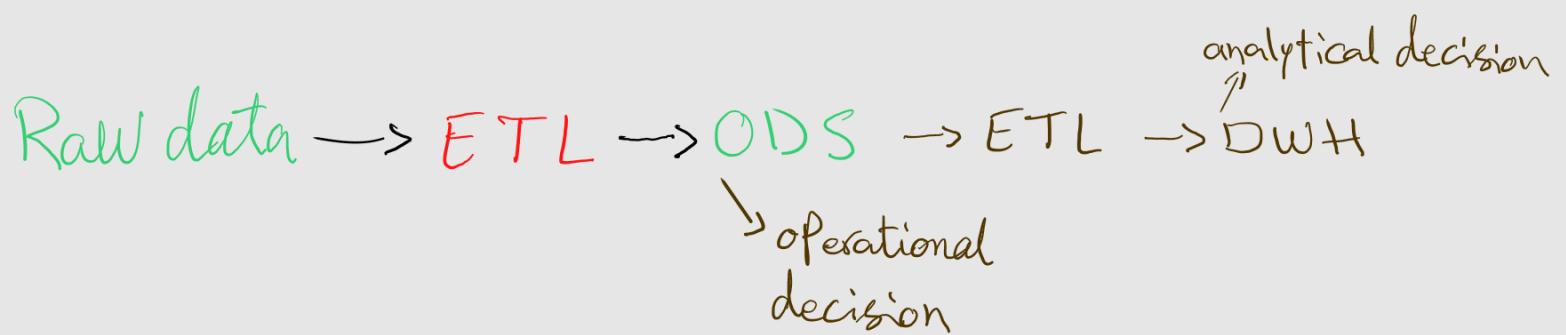
- ↳ quickly extract
- ↳ move the data into relational database
- ↳ start transformations from there
- ↳ landing zone extracted data
- ↳ data in tables and on a separate database
- ↳ temporary or persistent

In-memory databases → used usually with data marts



Operational data storage (ODS) (uncommon)

↳ Used for operational decision making



Dimensional modeling ↴

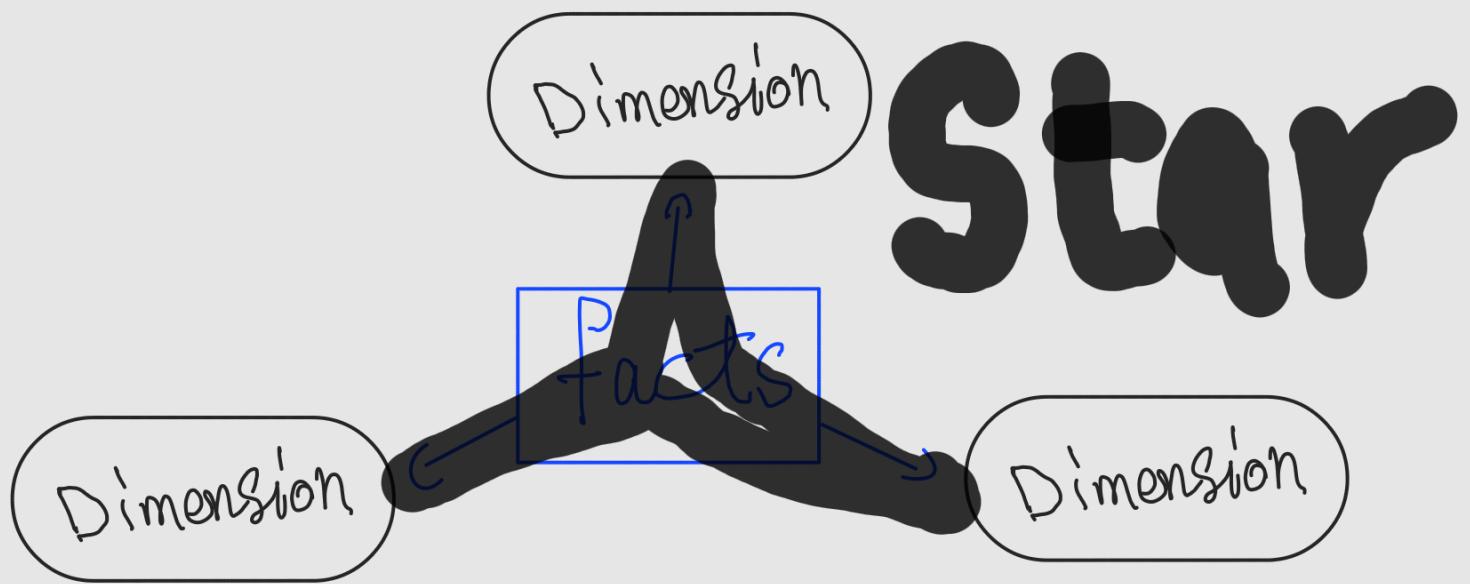
↳ method of organizing data (in data warehouse)

↳ facts

↳ Measurements like
Profit

↳ Dimensions

↳ context like category



↳ Unique technique of structuring data

↳ Commonly used in DWh

↳ Optimized for faster data retrieval

↳ Oriented around Performance & usability

↳ Designed for Reporting / OLAP

- Facts

- ↳ Foundation of DWH
- ↳ Key measurements →
 - ↳ aggregated and analyzed
 - ↳ aggregatable (numerical values)

لعن هو الاحاجي او المعاشر
لـ (الـ Dim_Tables) عـ (الـ Ids) او
(Tables) الـ (cues) تـ (المـ اـ (بيانـ اـ (دـ لـ (ـ

- fact table : PR, FR & facts
- Grain : the most atomic level facts are defined

- Dimensions

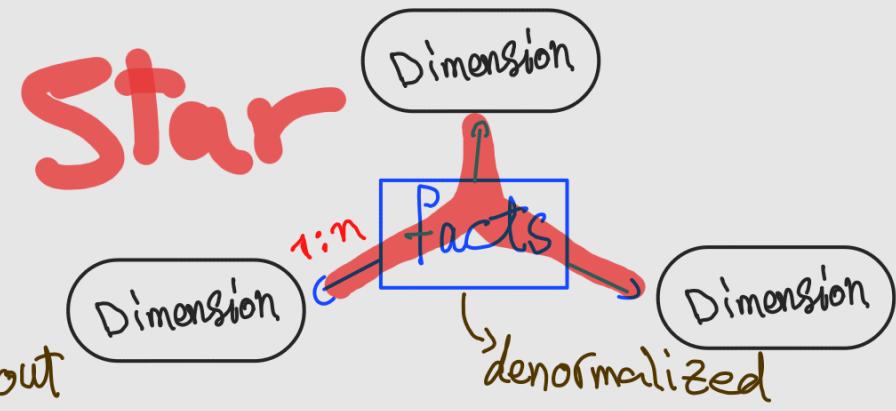
- ↳ categorizes facts
- ↳ supportive & descriptive
- ↳ filtering, Grouping & labeling
- ↳ non-aggregatable
- ↳ (more) static

- Dimension table : PR , Dimension , (FR)
- People , products , places , time

Star Schema

↳ Denormalized

- There is data redundancy
- optimized to get data out
- query performance (read)
- user experience



↳ Normalized

- Technique to avoid redundancy
- minimizes storage
- performance (write / update)
- many tables → many joins necessary

Snowflake Schema

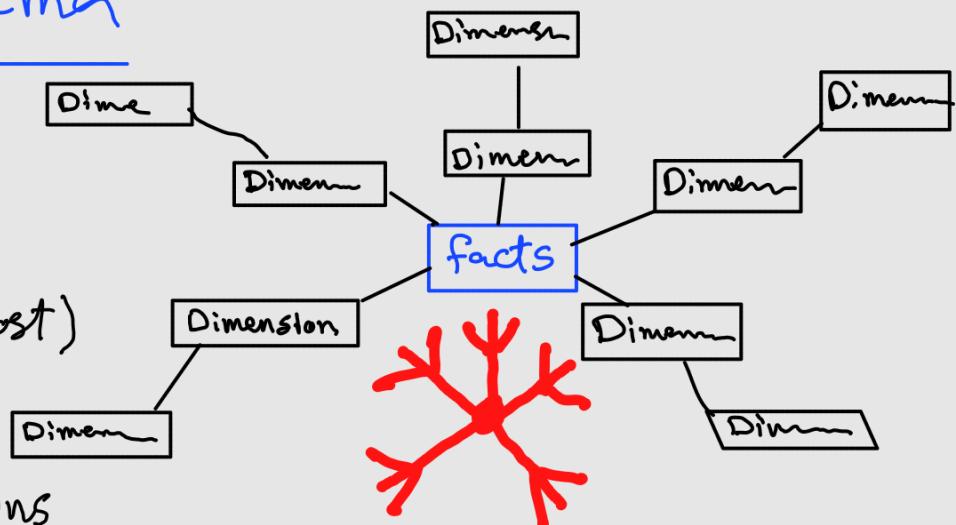
↳ more normalized

- Advantage

↳ less space (storage cost)

↳ less redundant data

↳ solves write slow downs



- Disadvantage

↳ more complex

↳ more joins (more complex SQL queries)

↳ less performance Data marts

Fact table types

1- Transactional fact table

↳ 1 row = measurements of 1 event / transaction

↳ 1 transaction defines the lowest grain

(FK) Keys will be unique and single value per row *

↳ Tend to have a lot dimensions associated

↳ can be enormous in size

2- Periodic snapshot Table

↳ 1 row = summarizes measure of many events / transaction

↳ summarized of standard period (1 day, 1 week, etc)

↳ lowest period defines the grain

(FK) Keys will be unique and single value per row *

↳ Tend to be not as enormous in size

↳ Tend to have a lot of facts and fewer dimensions associated

↳ no events = null or 0

3 - Accumulating snapshot fact table

- ↳ 1 row = summarizes measure of many events / transactions
- ↳ Summarized over lifespan of one process (order fulfillment)
- ↳ definite beginning & definite ending (& steps between)
Child steps of Order in Activity Log Table *
- ↳ Multiple Date / Time foreign keys (for each process step)
- ↳ Date / Time keys associated with role-play dimension

4 - Factless fact table

Activity log analysis will belong to Table *

employee registration will be

• Steps to create a fact Table

1- Identify business process for analysis

↳ Sales, order processing

2- Identify the grain → what one record in a fact table represents

↳ Transaction, order, order lines, Daily,
Daily + Location

3- Identify dimensions that are relevant

↳ What, when, where, How and why

↳ Time, location, Products, customers

4- Identify facts for measurements

↳ defined by the grain & not by specific use-case

○ Natural Keys

↳ Come out of the source system

○ Surrogate Key

↳ Artificial keys

↳ more beneficial

• integer number

• -pk or -fk suffix

• created by the database / ETL Tool

● Dimension Tables

- Hierarchies in dimensions

↳ preferred to be denormalized / flattened

1- Date dimension

* Replace nulls with descriptive values

2 - Conformed dimensions

↳ Shared by multiple fact tables / stars

↳ used to compare facts across different fact table

Shared attributes *one fact table only belongs to one Dimension Table*

3-Degenerate dimension

↳ dimension key without an associated dimension

order Id like Customer Table Columns

4-Junk dimensions (Transactional indicator dimension)

Yes / No or other types of attributes Table *
in table like

5-Role-play dimensions

↳ Dimension that is referenced multiple times by a fact

like Order Date has one attribute

-Slowly changing dimensions

Type 0 : Retain Original

- ↳ There won't be any changes
- ↳ Date Table (except for holidays etc.)
- ↳ original
- ↳ Very simple and easy to maintain

Type 1 : Overwrite

- ↳ Old attributes are just overwritten
- ↳ Very simple
- ↳ No fact table needs to be modified

Problem:

- ↳ History is lost
- ↳ insignificant changes
- ↳ might affect /break existing queries

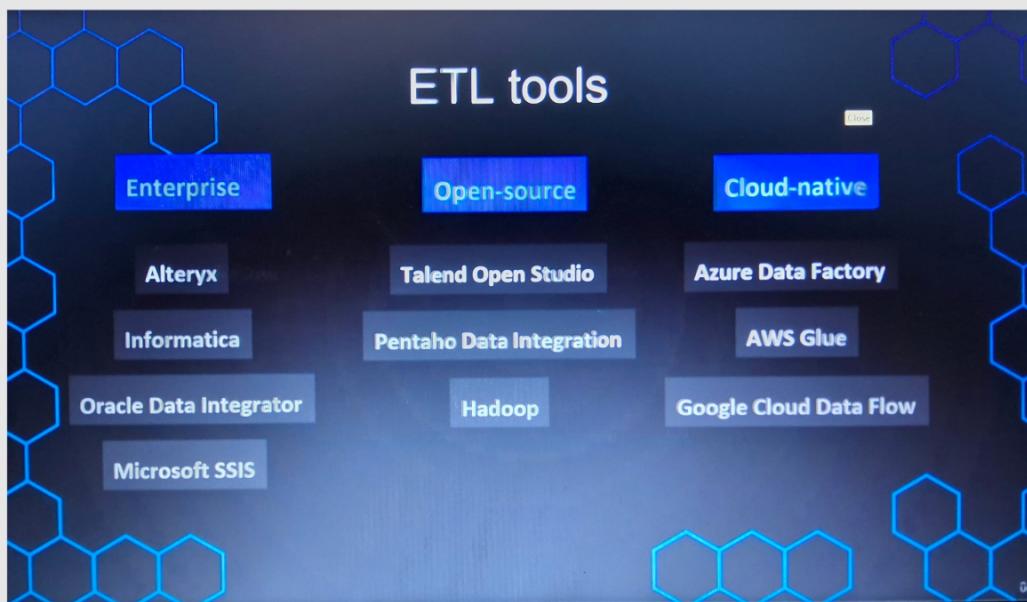
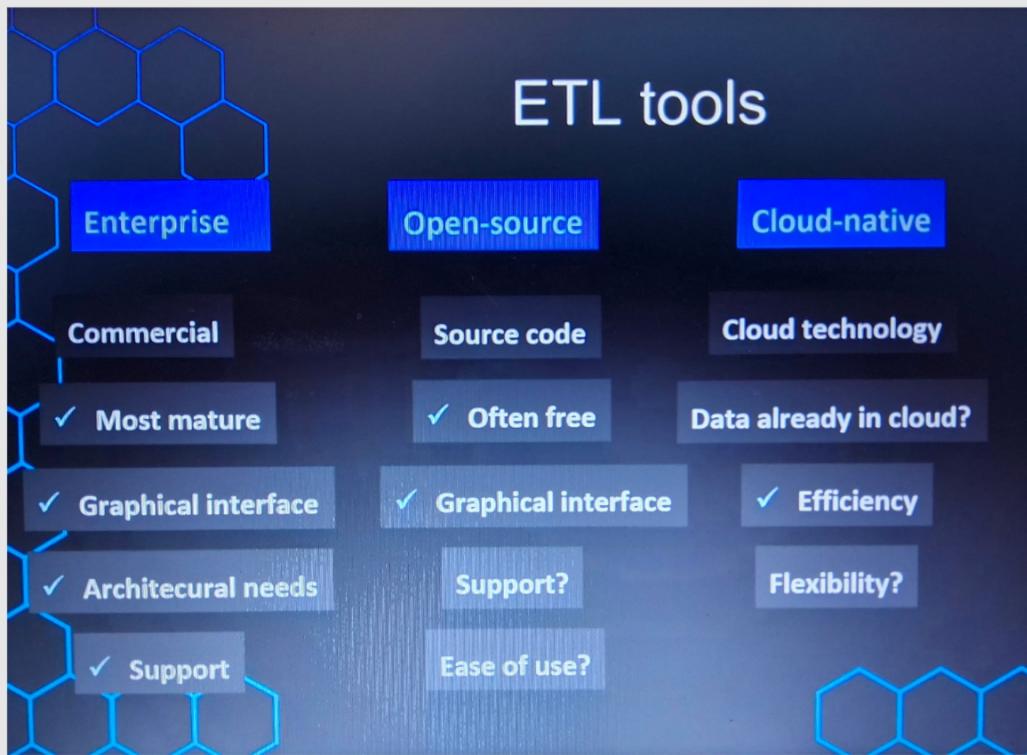
Type 2 : New row

- ↳ Changes reflected with history

Type 3 : Additional attributes

- ↳ instead of adding a row, we add a column
- ↳ Typically used for significant changes at a time

ETL



- ETL v.s. ELT

- Query optimization
and database internals