

chp # 2.

- * Major steps in ETL process.
 - The process of extracting data from source systems & bringing it into the data warehouse is commonly called ETL.
 - Extraction
 - transformation
 - Loading
 - It is a process in which ETL tool extracts the data from various data source systems, transforms it in the staging area & finally, loads it into the Data W. system
 - ETL process requires active inputs from various stakeholders, including
 - developers
 - analysts
 - testers
 - top executives
 - Business DW technique needs to change with business changes
 - ETL is a recurring methods (daily, weekly, monthly) of a data Warehouse system & needs to be agile, automated, & well documented

ETL tools: sybase

Oracle Warehouse Builder

Informatica ETL

Mark Logic

Steps:

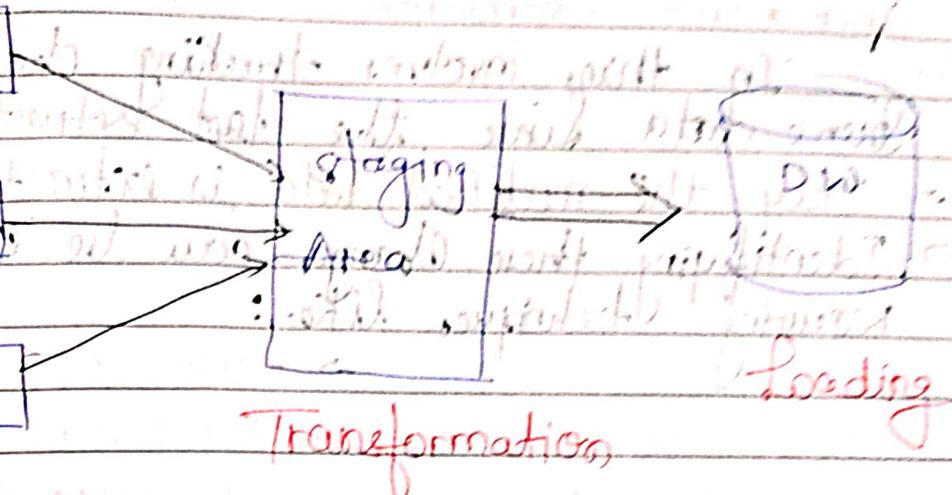
Source

RDBMS

SQL server

Flat files

Extraction



Extraction:

- First step of ETL process is extraction. In this step, data from various source systems is extracted which can be in various formats like relational DB NO, SQL, XML & flat files into the staging Area.
- It is imp to extract the data from various systems & store it into the staging area first & not directly go into the DW as it may be incorporated during logical extraction.

Data Extraction Techniques

(Logical) Physical Extraction

① full extraction: ~~logical~~ physical

This data extraction method involves taking all data from a source system without needing additional logical information. The data is provided as-is & there's no tracking of changes in the



Eg:

source system.

Exporting an entire table as a flat file.

Incremental Extraction:

- In this, involves tracking changes in the source data since the last extraction.
- Only the modified data is extracted & loaded
- Identifying these changes can be complex, requiring techniques like:
 - timestamp comparison ; or
 - maintaining 'a' change table in source system

Physical extraction

i) **Online Extraction**: running SQL for gate drift information. In this process, extraction process directly connects to the source system & extracts the source data with help of DBMS, ORACLE, DB2, etc.

ii) **Offline Extraction**: All the data is stored in form of files. The data is not extracted directly from the source system, but is staged explicitly outside the original source system.

Common structure of in offline extraction

① flat file : Generic format

② dump file : Database specific file

Transformation:

Page No.

Date

The second step of the ETL process is transformation. A set of rules, or functions, are applied on extracted data, to convert it into single, standard format. Involves following tasks:

① Filtering:

Involves selecting & loading only specific attributes or data elements into the data warehouse. → focus on relevant info for analysis.

② Cleaning:

Addresses data quality issues by handling NULL values & inconsistencies.

Eg: mapping variations like "USA", "United States", & "America" into a standardized format.

③ Joining:

Combines data from multiple attributes or tables into a single, unified dataset.

④ Splitting:

Divides a single attribute into multiple attributes based on defined rules.

useful when piece of info needs to be broken down into its constituents parts.

⑤ Sorting:

Arranges data tuples based on specified attributes, typically a key attribute.

Data Transformation Techniques

(1) Data Smoothing

- This method is used for removing the noise from a dataset. Noise is referred to as the distorted & meaningless data within a dataset.
- Method: Utilizes algo's to highlight significant features & detect small changes revealing special patterns.

(2) Data Aggregation

- Collects data from various sources in a unified format.
- Data is analyzed, summarized & presented in reports, facilitating the gathering of extensive info about specific data clusters.

(3) Discretization

- Converts continuous data into intervals by substituting continuous attribute values with interval labels for easier study & analysis.

(4) Generalization

- Transforms low-level data attributes into high level attributes using concept hierarchies.

Eg: Age data (20, 30)

Transformed into categorical value (young, old)

(5) Attribute Construction

- Creates new attributes from an existing set.

Eg: Employee dataset attributes (Name, ID, address) used to construct a dataset containing employees

who joined in 2019 enhancing mining efficiency.

Normalization

- Data pre-processing
- Transforms data to fall within a specified range, addressing challenges in data modelling & mining posed by attributes in different scales.

Loading is the third step of ETL process

→ The third & final step of the ETL process is loading. In this step, the transformed data is finally loaded into the DW.

→ Sometimes the data is updated by loading into the data DW very frequently, sometimes it is after a longer but regular interval.

→ Rate & period of loading depends on Requirements & varies from system to system.

→ Loading can be carried in 2 ways:

① Refresh:

When data is completely rewritten. This means that old file is replaced. Refresh is usually used in combination with the data extraction to populate data warehouse initially.

② Update:

Only those changes applied to source information are added to DW. It is carried out without deleting or modifying pre-existing data.

~~OLAP VS OLTP~~

OLAP VS OLTP:

(we can divide IT systems into transactional (OLTP) & analytical (OLAP). In general, we can assume that:

OLTP → provides source data to DW

OLAP → helps to analyze it.

OLTP (Online Transaction Processing):

It is characterized by large No. of short

online transactions (INSERT, UPDATE, DELETE)

→ Put on very fast query processing, maintaining data integrity in multi-access environment. A. effectiveness measured by No. of transactions per second.

→ mostly there is detailed & current data

→ Schema used — store transactional DB in the Entity model (BNF)

OLAP (Online Analytical processing):

→ relatively low volume of transactions

→ queries often very complex & involves various aggregations involving multiple facts.

→ OLAP systems, response time is an effectiveness measure

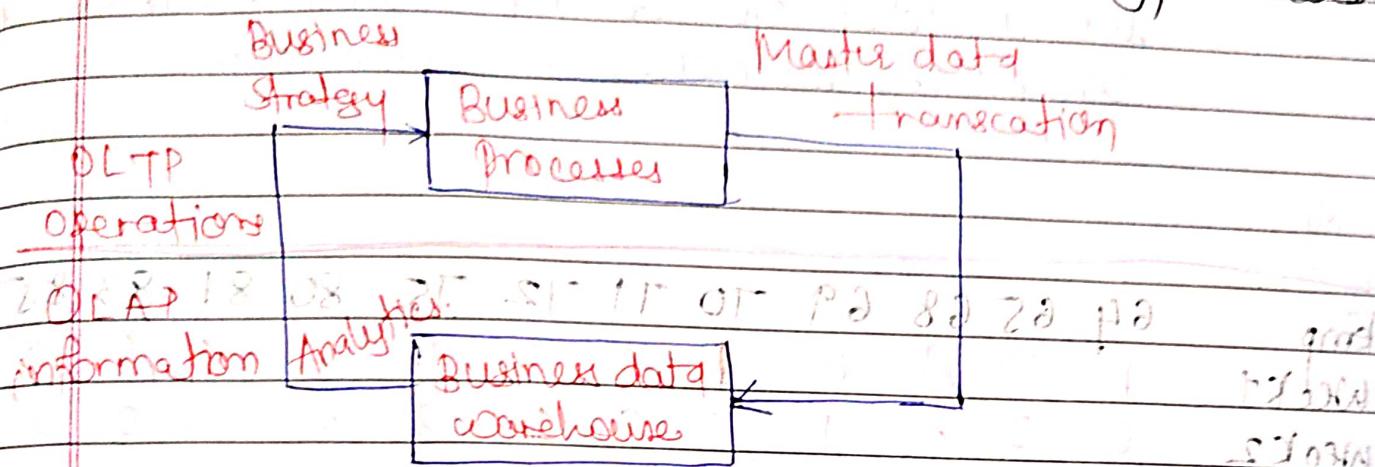
→ OLAP DB's is used by Data Mining

→ In OLAP, DB is aggregated, historical data stored in multidimensional schemas (usually star) or star with

→ divided into 4 for one or more cubes.

Cubes

- designed in such a way that creating & viewing reports become easy.
- OLAP cube is the data structure optimized for very quick data analysis.
- OLAP cube is also called hyperscubes.



OLAP operations :-

(P) -> (1) Pass (2) -> (3) Filter (4) Task

- In multidimensional model, data are organized into multiple dimensions, & each dimension contains multiple level of abstraction defined by Concept hierarchies.
- This organization provides users with the flexibility to view data from different perspectives.
- OLAP provides a user-friendly environment for interactive data analysis.
- A number of OLAP data cube operations exist to materialize different views of data, allowing interactive querying & analysis of the data.

1) Roll-up operation is done at higher level of abstraction.
 2) The roll-up operation is (i) concept hierarchy
 Platforms aggregation, and a concept hierarchy either
 by climbing up or via a concept hierarchy
 for a dimension or by climbing down a
 concept hierarchy, ie dimension

Eg:

Temp
Week 1
Week 2

64	65	68	69	-70	71	72	75	80	81	83	85
1	0	1	0	1	0	0	0	0	0	1	0
0	0	0	1	1	0	0	1	0	1	0	0

Consider (30-20) up levels hot (80-85), mild (-70-75) cool (64-69)

→ We have to group up the column and add up the values according to the concept hierarchies. This is called roll-up.

Temp
Week 1
Week 2

cool	mild	hot
1	1	1
2	1	1

DRILL DOWN

17/11/23

- The drill down operation (also called roll-down) is the reverse operation of roll-up.
- Drill down is like zooming in on the data cube.
- Navigates from less details to more detailed data.
- Drill down can be performed by either
 - Stepping down a concept hierarchy for a dimension or adding additional dimensions.
- Bcz as a drill-down adds more details to the given data, it can be performed by adding a new dimension to a cube.

Eg:

Temp	/0/0			weather
	cold	mild	Hot	
Day 1	0	0	0	sunny, part
Day 2	0	0	0	part
Day 3	0	0	0	cloudy
Day 4	0	1	0	cloudy
Day 5	1	0	0	cloudy
Day 6	0	0	0	rainy
Day 7	1	0	0	part
Day 8	0	0	0	rainy
Day 9	1	0	0	rainy
Day 10	0	1	0	rainy
Day 11	0	1	0	rainy
Day 12	0	1	0	rainy
Day 13	0	0	0	rainy
Day 14	0	0	0	rainy

SLICE

Page No. _____
Date _____

- A slice is a subset of the cube corresponding to a single value for one or more members of the dimensions.
- For example, a slice operation is executed when the user wants a selection on one dimension of a three-dimensional cube resulting in a 2-D slice.
- So, the slice operations perform a selection on one dimension of the given cube, thus resulting in a sub-cube.
- If we select,

temperature = cool

Temperature	cool	hot	warm	normal
Day 1	0			
Day 2	0			
Day 3	0			
Day 4	0			
Day 5	1			
Day 6	1			
Day 7	1			
Day 8	1			
Day 9	1			
Day 10	0			
Day 11	0			
Day 12	0			
Day 13	0			
Day 14	0			

DICE

- (-) The DICE operation describes a sub-cube by operating a selection on two or more dimension.

28.3.1932. 9/10

9/10 part 1)

For Eg. if we want to select a subcube

Selection (Time = day 3) OR (temp = cool)

OR (temp = cool) AND (Time = day 4)

(Time = day 4) AND (temp = cool)

(Temperature is cool) OR (temp = hot)

to the original cube we get the foll. sub-cube (still two-dimensional).

temp	Cool	hot
Day 3	0	1
Day 4	0	0

log result

- PIVOT (Rotation) was covered in 9/10 part 1
 PIVOT operation is also called as transpose operation
 → Pivot is a visualization operation, which rotates the data axes in view to provide an alternative presentation of the data.
 → Contains swapping the rows & columns or moving one (or more) row(s) dimension into the column dimension of matrix.

POLAP (Relational)

OLAP SERVERS

MOLAP (Multi-dimensional)

① ROLAP

- Rational On-line-Analytical Processing (ROLAP) work mainly for data that resides in a relational DB, where the base data & dimension tables are stored as relational tables.
- ROLAP servers are placed between the relational back-end server & client front-end tool.
 - use RDBMS to store & manage warehouse data.
 - use OLAP middleware to support missing pieces.

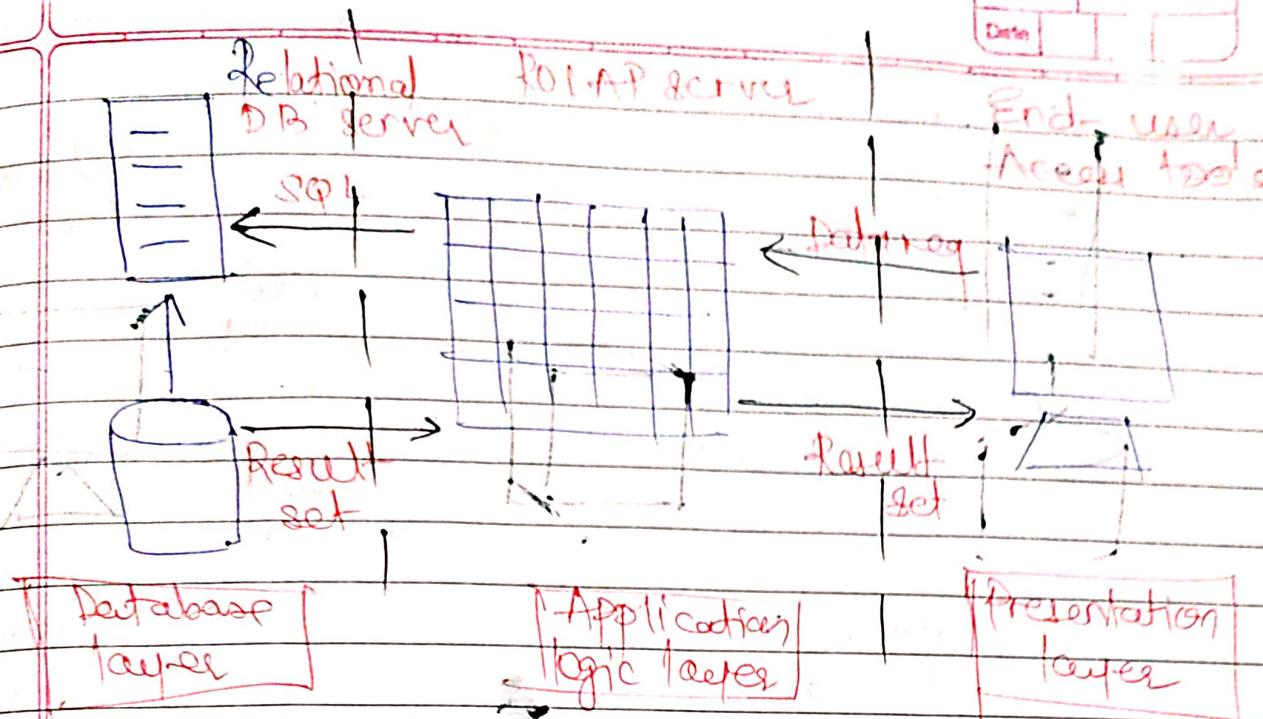
Eg: DSS Server of Microstrategy.

Advantages

- ROLAP can handle large amounts of data.
- Can be used with data warehouse & OLTP systems.

Disadvantages

- Limited by SQL functionalities
- Hard to maintain aggregate tables.



MOLAP

- MOLAP supports multidimensional views of data through array-based multidimensional storage engine
- With multidimensional data store, the storage utilization may be low if the data set is sparse
- Eg! Oracle Essbase.

Advantage:

- optimal for slice & dice oper
- performs better than ROLAP when data is dense
- Can perform complex calculations.

Disadvantage

- Difficult to change dimension without re-aggregation
- MOLAP can handle limited amount of data