

## What is Datawarehouse

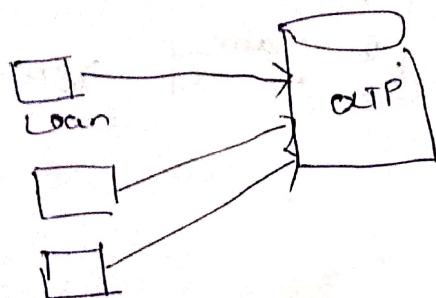
DB : All companies maintaining databases to store only transactional data.

Bank : Online, store, organization.

⇒ They are storing day-to-day transactions.  
Ex: in bank depositing, money transferring

⇒ This database is known online transaction processing (OLTP) because they are storing day-to-day transactions.

⇒ The source for OLTP data are the applications.



Applications are generating data and data are stored into OLTP DB.

⇒ Now days every organization maintains two databases OLTP and OLAP.

⇒ Online Analytical processing (OLAP) is also known as Data warehouse (DW) OR DSS (Decision support system).

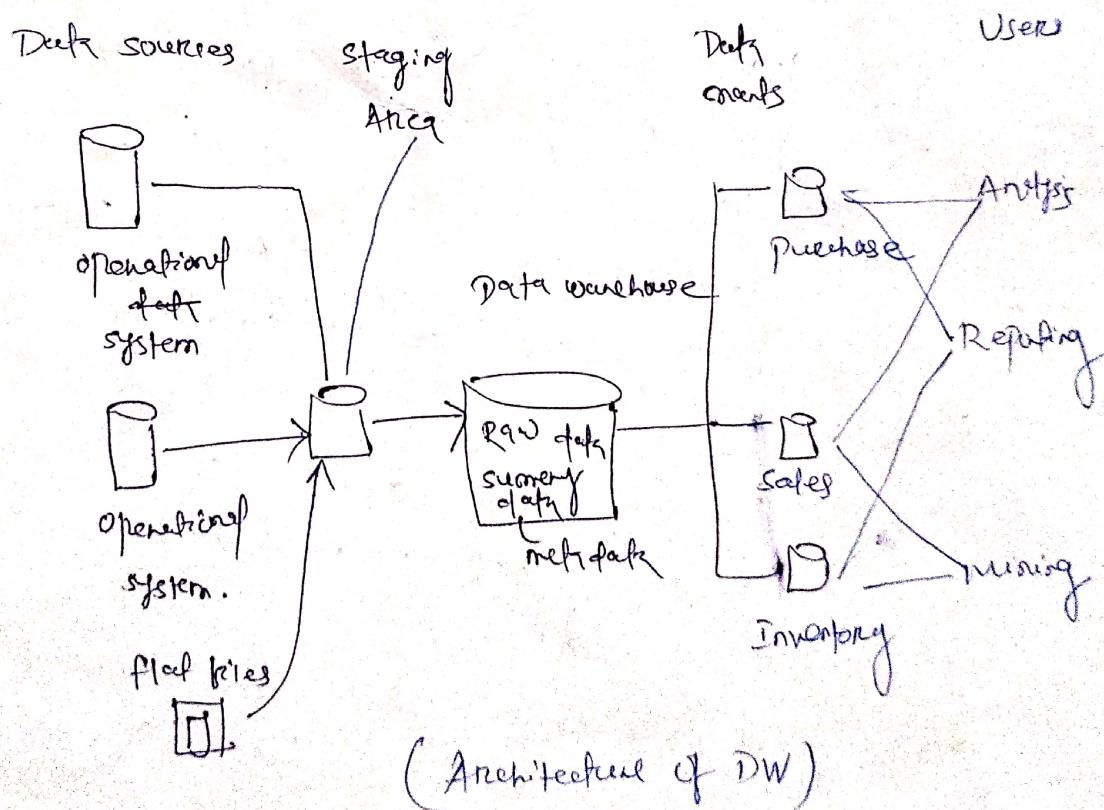
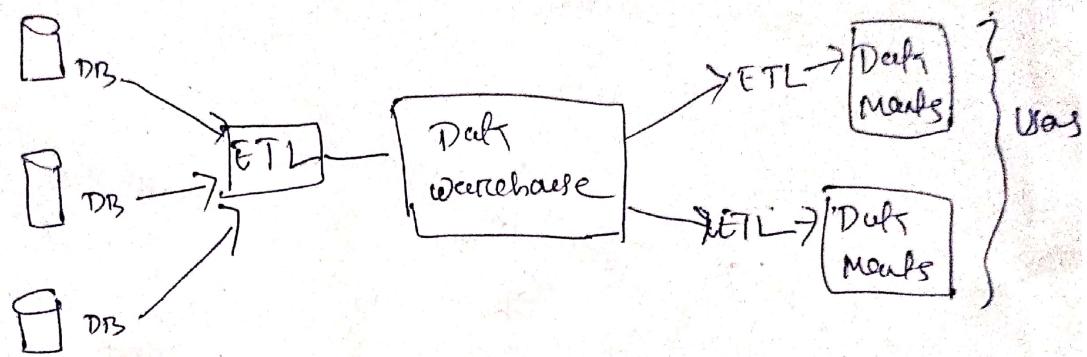
⇒ OLAP is used for analysis of data. Because we need historical data. That historical data are placed in OLAP.

Ex: If ten years data are available then we can find who is the best player.

⇒ Source of OLAP are OLTP.

The process of moving the data from ETP to OLAP is done by ETL (Extraction, Transformation, Loading) process.

- Data warehouse is the process of collecting and managing data from various sources to provide meaningful business insights.
- DW is the central repositories of integrated from one or more sources.



Staging Area → intermediate Area. used for data processing during ETL process.

Data marts → is a simple form of a DW that focus on a single subject. Ex Sales. These are built and controlled by a single department within an organization.

### ETL Process

- Extraction: - Data from various various source systems are extracted which can be in different formats.
- Like Relational DB, MySQL, XML and flat files.
- collected and stored into staging area.
- first collect data from different sources. then store into staging (intermediate) area.
- not directly to data warehouse because data may be corrupted.
- Transformation: - In the process set of functions or rules are applied on the extracted data.
- that function convert all data into a single standard format.
- sub processes: filtering, joining, splitting / sorting
- Filtering: - reading some of attributes into the data warehouse.
- Cleaning: - filling the NULL values with some default values, delete incorrect records.
- splitting: - splitting a single attribute into multi-field attributes.

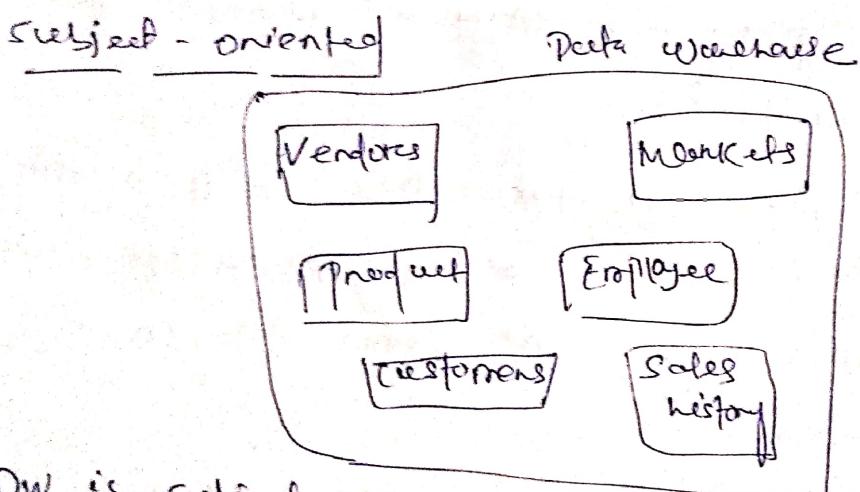
→ Sorting: - sorting type on the basis of core attribute (generally key attribute)

→ Loading: - the transformed data is finally stored / loaded into Data Warehouse.

→ ETL tools examples:- sybase, oracle, vertica  
- ss builder, marmalade.

### Characteristics of DW

- Subject-oriented
- Integrated
- Time-variant
- non-volatile



→ DW is subject oriented because it provides information in subject wise.

→ Means data are stored in Data warehouse in subject wise so that data can analyze easily.

## Integrated

→ Data are collected from different sources and then combined and placed in DW. (Operational DB)

## Time-Variant

→ Historical data is kept in DW.

→ for analyzing data DW must contain historical data.

→ if one can refine data from 3 months to 6 months, 12 months older data from a DW.

## Non-volatile

→ once data is in data warehouse it will not change.

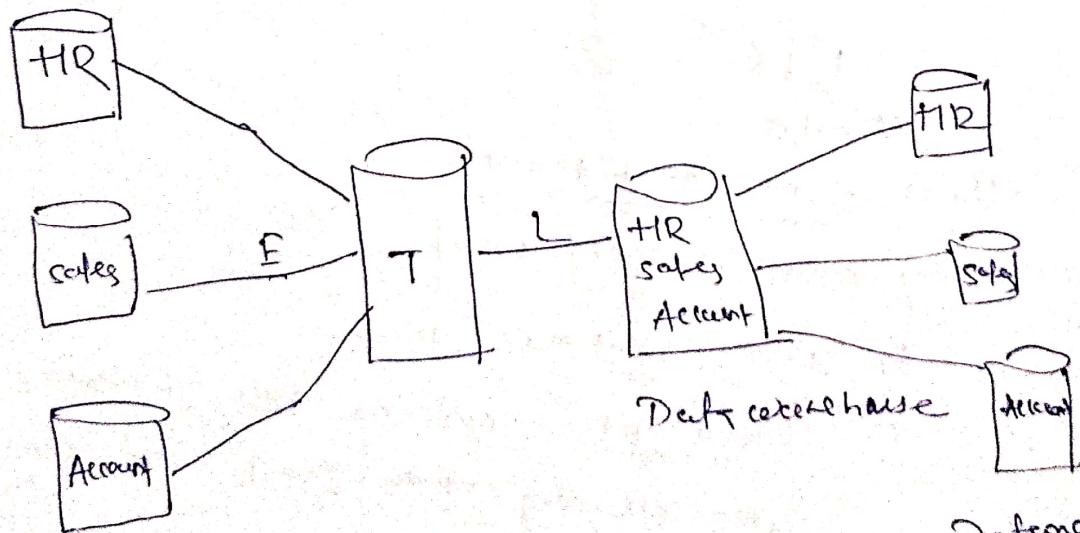
→ so historical data in a DW should never be altered.

→ The new data can be added into the DW.

→ non-volatile means we can't update the historical or any data which is present in DW.

## Data marts

→ subset of Data warehouse is known as data marts.



- It focuses on single subject.
- It is subject-oriented and it is designed to meet the needs of a specific group of users.
- Data marts are fast and easy to use, as they make use of small amount of data.

### Components of Data warehouse

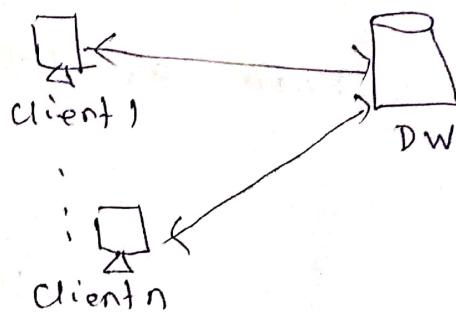
#### DW Architecture

##### Single-tier Architecture

- objective to minimize the amount of data.
- goal is to remove redundancy.
- this architecture is not used in practice.

##### Two-tier Architecture (Client/Server)

- it separates physically available resources and data warehouse.
- it can't support more number of users.

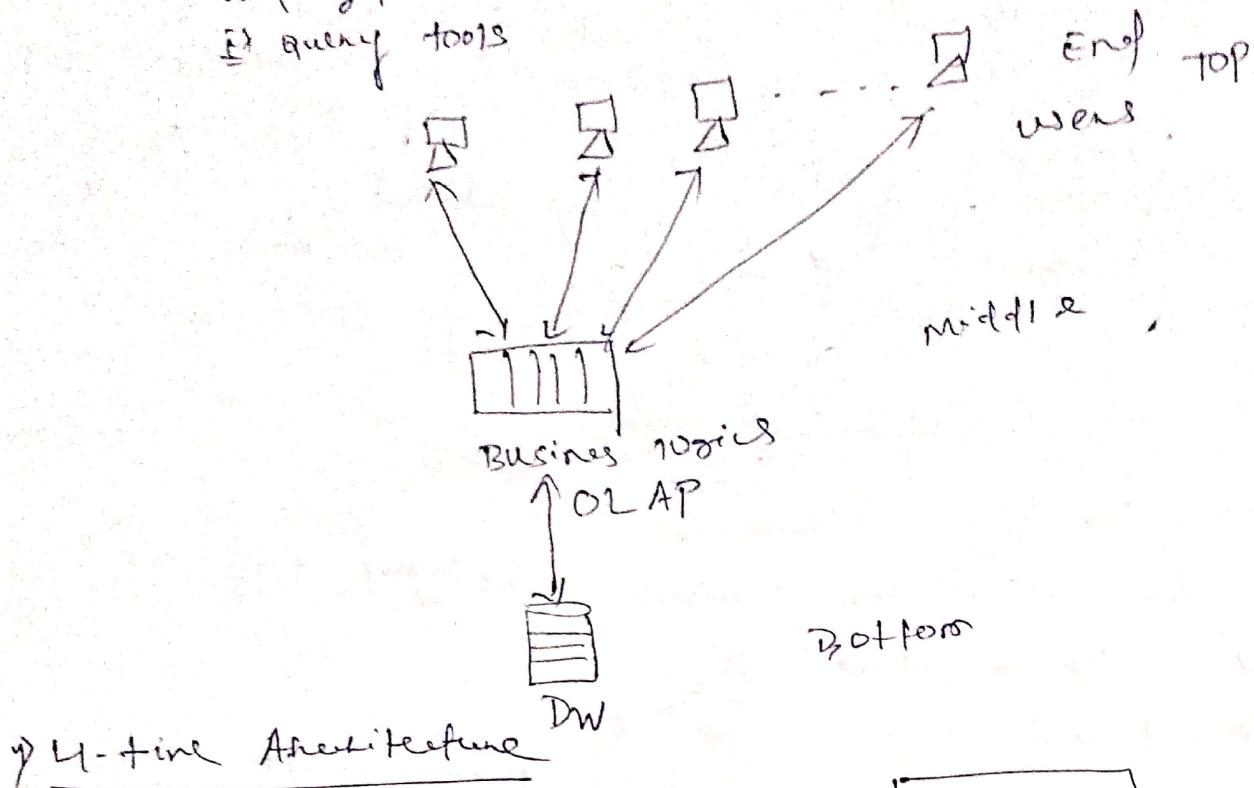


##### Three-tier Architecture

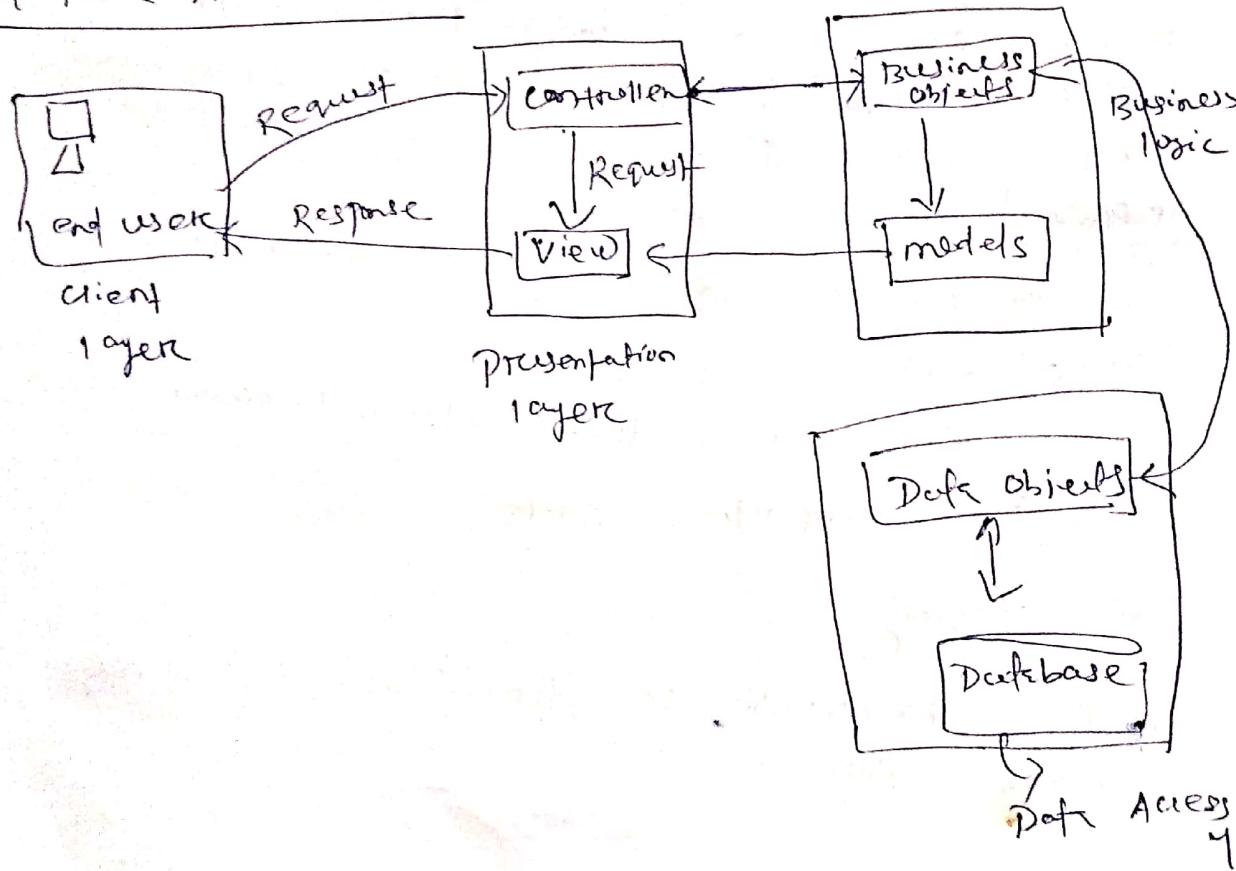
- most widely used.
- contains three tiers.
- Bottom tier contains DBs (RDBMS). After cleaning, transformed and loaded into bottom tier.

→ middle tier contains OLAP server. It gives an abstract view of database.

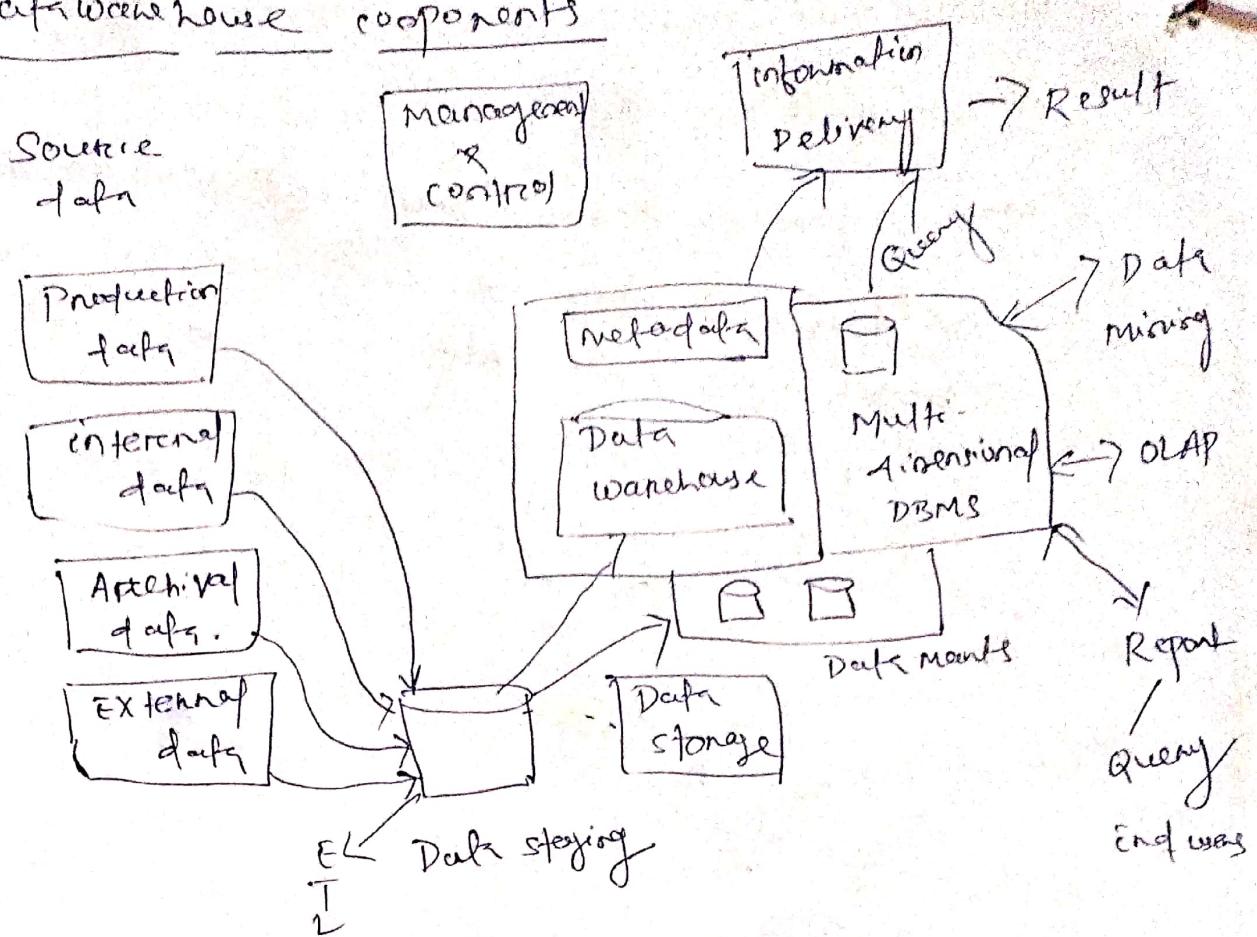
→ top layer/tier is a front-end client layer. Here tools and APIs are available that connect and get data from data warehouse. It query tools



#### 4-tier Architecture



## Data warehouse components



## Production Data

data comes from various operational task cores from the enterprise.

## Internal Data

→ organization / client keeps private spreadsheets, reports, customer profile. These are the internal data.

## Archival Data

It contains old and historical data.

## External data

→ data comes from external source.  
Ex: competitors data

## Data staging

→ Apply extraction, transformation and loading (ETL) process.

After ETL process store the data into some intermediate storage known as data staging.

### Metadata

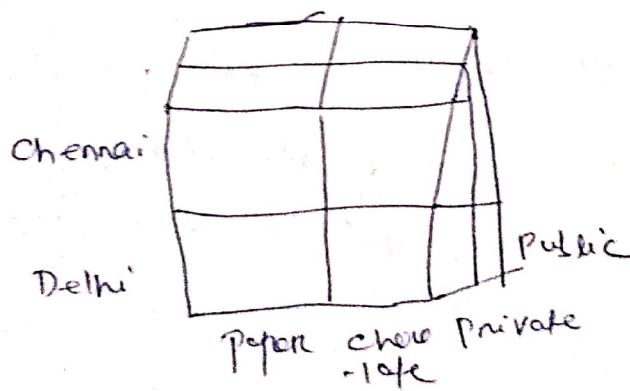
- It is a facts dictionary on the data catalog.
- It keeps information like indexes, record and address and internal details.

### Data Marts

- Data marts are sub parts of the data warehouse.
- It is designed for a particular line of business such as sales, marketing.
- Data mart is an access layer which is used to get data out to the users.
- Access data from data marts take less time and takes less cost to build because it is derived from data warehouse.

### Multidimensional Database

- It is used mostly for OLAP (online Analytical Processing)
- It allow users to quickly get answers to their requests by generating and analysing the data.



## Information delivery

- It is used to enable the process of subscribing lot of data warehousing.
- It will distribute warehouse stored data and other informations to offline data warehouse or end user's product such as spreadsheets and user databases.

## Management & control

- It will coordinate the service and functions with in the DW.
- The control component controls the data transformation and load management into DW.
- Manages the data delivery to clients.
- Security management.

## DW Access tools

- Different tools are available at end user site by which they can perform different operations.
- Query & Reporting:- generate business reports for analysis.
- Application developer:- application also developed.
- Data mining:- used to find patterns and correlations in large amounts of data.
- OLAP tools:- helps to build multi-dimensional DW  
Analyze the enterprise data from multiple perspective

## Metadata in Data warehouse

- Data dictionary contains metadata.
- Metadata means facts about facts.
- It contains
  - information about the logical data structure
  - information about files and addresses
  - information about indexes.
- It stores all information one stored in one place.  
Means collect the previous records.
- Metadata component serves as a dictionary of the contents which are present in data warehouse.

### Types of metadata

- operational metadata
- extraction and transformation metadata
- end user metadata.

### Operational Metadata

- Data comes from several operational systems of enterprise.
- These sources contain different data structures.
- The elements selected for data warehouse from different sources have different fields and different length, different data types.
- we are splitting the records, combining the records from different sources files and deal with multiple coding schemes and field lengths.

When the data is delivered to the end user  
that data must be in original form which  
was in the source of data sets.

### Extraction & transformation metadata

ETL contains

- 1) the extraction of data from the source system
  - the extraction frequencies
  - extraction methods
  - Business rules for the data extraction.

- 2) Also contains information about all data transformation that take place at the data staging area.

### End-user metadata

→ It enables the end users to find data from the data warehouse.

→ It helps end users to understand the various types of information resources are available in data warehouse.

→ These resources can take many forms like data elements, queries, reports and published documents.

## Data cube / OLAP cube

- When data is grouped on combining in multidimensional measures need is called as data cubes.
- Data cube is a structure that enable OLAP operation to achieve the multidimensional functionality.
- It is the extension of 2-D.
- whenever there are lots of complex facts to be aggregated and there is a need to extract the relevant / important data.

⇒ Data cube is created from a subset of attributes in the database.

⇒ Specific attributes are selected to be measured attribute.

Ex xyz create a sales fact warehouse to keep records of sales for the dimensions time, item, branch & location.

Each dimension may have table known as dimension table. Ex:- dimension table for item containing item-name, brand and type.

Ex 2-D view of data

Electronic item sales. data for the computer in the city store.

location = "Hawaii"		item (type)	
time (quarter)	home entertainment	computer	phone
Q1	605	825	14
Q2	680	952	31
Q3	812	1023	50
Q4	715	250	28

in dollars

Row  
item

column,  
quarter

Q-D

### 3-D cubes

→ To view data according to time, item, location  
for cities Chicago, New York, Toronto, ~~Mumbai~~.

location = "chicago"			location = "new york"			location = "toronto"			location = "mumbai"		
item			item			item			item		
time	home	comp. phone	time	home	comp. phone	time	home	comp. phone	time	home	comp. phone
ent.	ent.	ent.	ent.	ent.	ent.	ent.	ent.	ent.	ent.	ent.	ent.
Q1	605	892	89	Q1	1087	968	38	Q1	818	746	43
Q2	680	890	54	Q2	1130	1024	41	Q2	894	769	52
Q3	812	924	59	Q3	1084	1048	45	Q3	940	995	58
Q4	715	992	23	Q4	1142	1091	54	Q4	978	864	59

Here we are viewing data in three dimensions like time, items and locations.

Location (cities)

	Chicago	New York	Boston
Time (Quarter)	Q1	Q2	Q3
Item (Types)	home	computer	phone
entertainment - rent	605	825	19
Q1	1081	968	38
Q2	818	746	43
Q3	892	1023	30
Q4	927	1038	38

## OLAP operations

4 types of operations are there

- Roll-up

- Drill-down

- Slice and dice

- Pivot (Rotate)

### Pivot

→ Here, rotate the data axes to provide a substitution of data.

→ Here, rotation is done.

Ex

New jersey

Los Angeles

pearth

sydney

locations  
(cities)

PC			
Book			
Shoe			
Clothes	605	825	14
			400

PC Book Shoe Clothes

(item types)



pivot



PC				605
Book				825
Shoe				14
Clothes				400

New Jersey Los Angeles Perth Sydney

Jersey Angeles

Location (cities)

Roll up

- It is also known as aggregation.
- Roll up operation can be performed in 2 ways.
  - Reducing dimension
  - climbing up concept hierarchy
- concept hierarchy is a system of grouping things based on their own level.

Locations (cities)	USA			
Australia				
Time (quarter)	Q1	1000		
	Q2			
	Q3			
	Q4			
	PC	Book	Shoe	Clothes
		Item types		

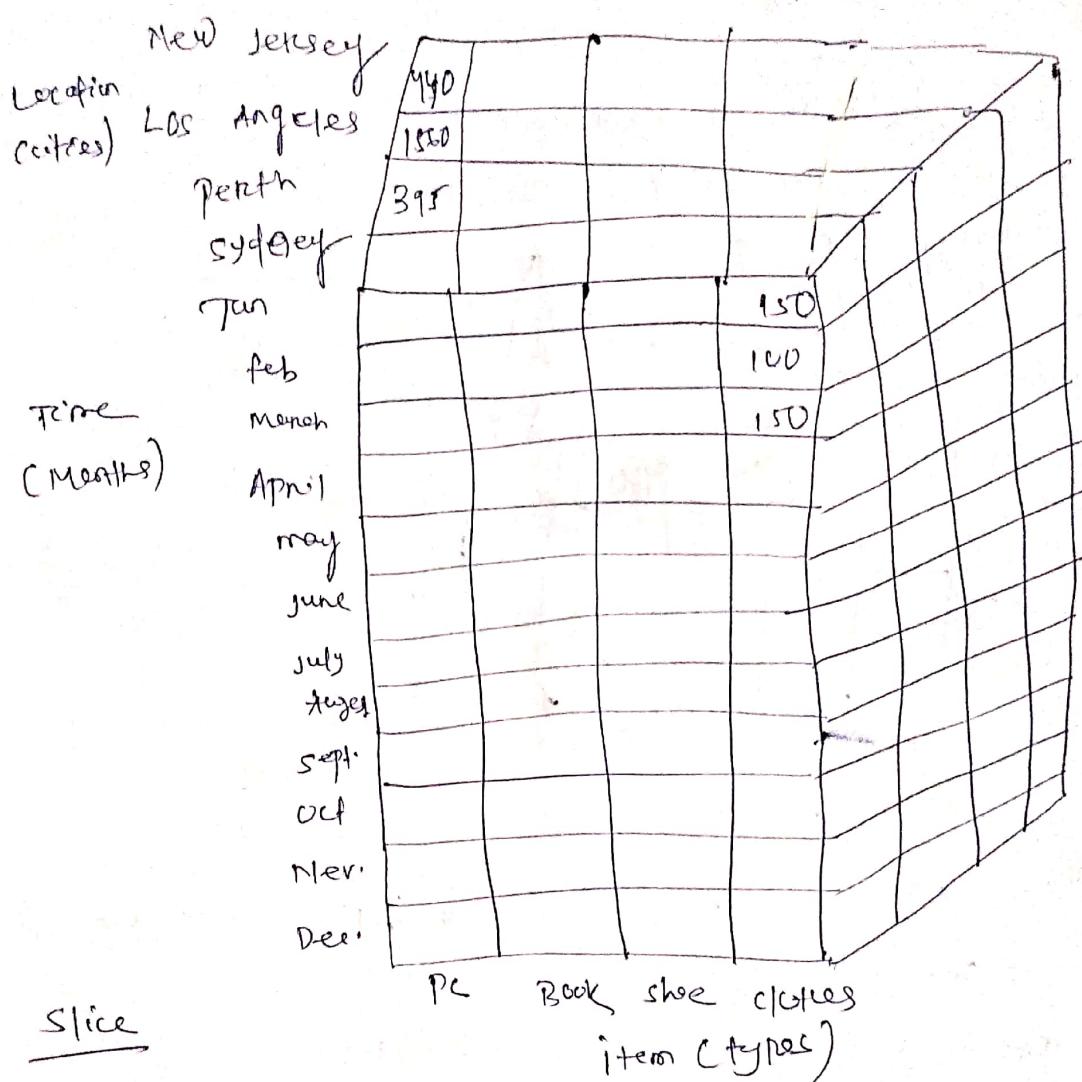
Roll-up on location  
(from cities to countries)

Locations (cities)	USA				
New Jersey					
Los Angeles					
Perth					
Sydney					
Time (quarter)	Q1	605	825	14	400
	Q2				
	Q3				
	Q4				
	PC	Book	Shoe	Clothes	
		Items types			

- Here, cities New Jersey, Los Angeles are aggregated into one country USA.
- Perth, Sydney into Australia.
- Data is location hierarchy moves up from city to country.

### Drill-down

- Here, data is fragmented into smaller parts.
- It is the reverse of the roll up process.
- It can be done by
  - moving down the hierarchy
  - increasing a dimension.



→ Here, one dimension is selected and a new subcube is created.

→ Here, time dimension is sliced with G1 castle kitten.

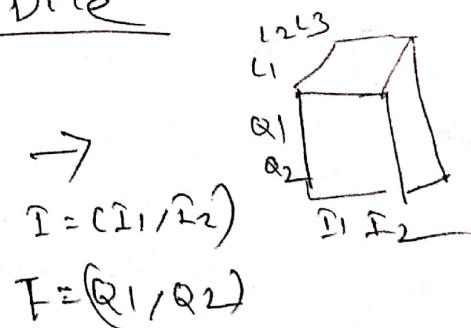
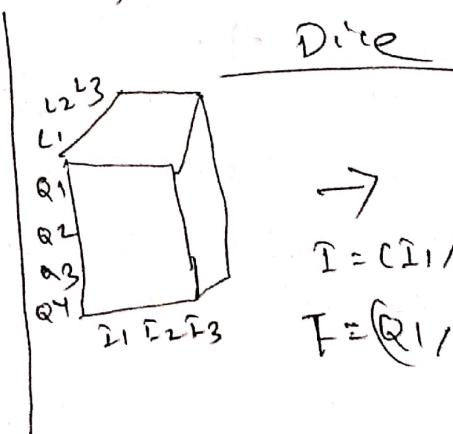
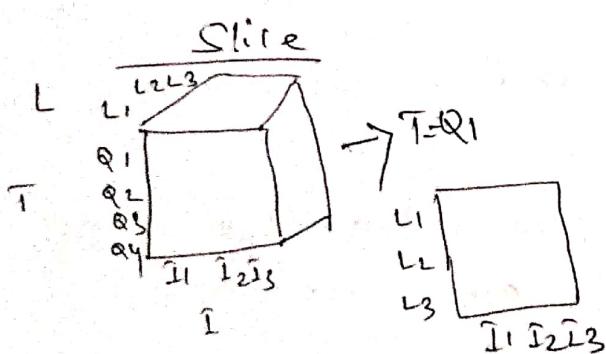
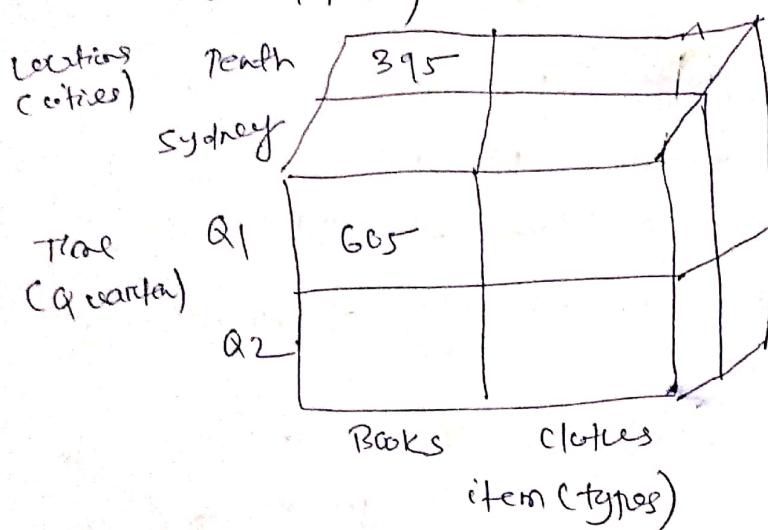
New Jersey			
Los Angeles			
cities	Tenth		
Sydney	605	825	141

PC Book Shoe clothes

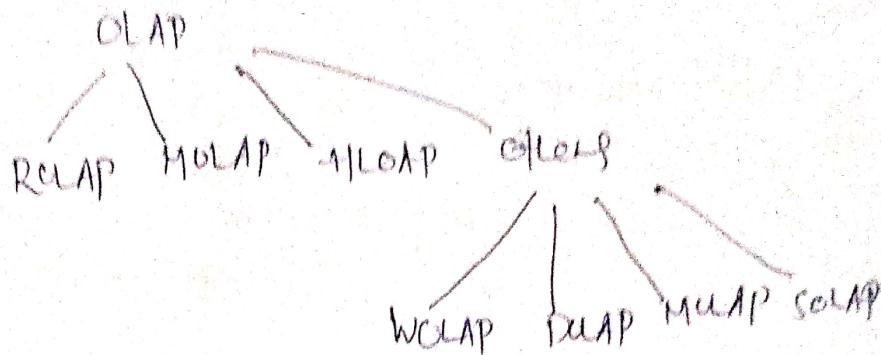
Item types.

Dice:- this operation is similar to a slice.  
only difference is dice selected 2 or more dimensions  
that gives a new sub-cube.

Ex Dice like (location = 'peach' or 'sydney')  
and (time = Q1 or Q2) and (item = 'books'  
or 'clothes')



## OLAP TYPES



→ ROLAP:-

- × extension of RDBMS.
- × Highly scalable.
- × Multi dimensional data mapping to perform the standard relational operation.

→ MOLAP (Multidimensional OLAP):

- × It implements operation in multidimensional data.
- × Information retrieval is fast.
- × It doesn't contain detailed description.

→ Hybrid OLAP:-

- × Aggregated totals are stored in a multidimensional DB.
- × Details are stored in RDB.
- × Offers efficiency of ROLAP and MOLAP.

→ Desktop OLAP:- User downloads a part of the data from DB locally and analyze it.

→ Web OLAP:- The OLAP which is accessible via web browser.

→ Mobile OLAP:- User will access and analyze OLAP data using mobile devices.

→ spatial OLAP:- It is created to facilitate analysis.

- consist of both spatial and non-spatial data in geographic information system.

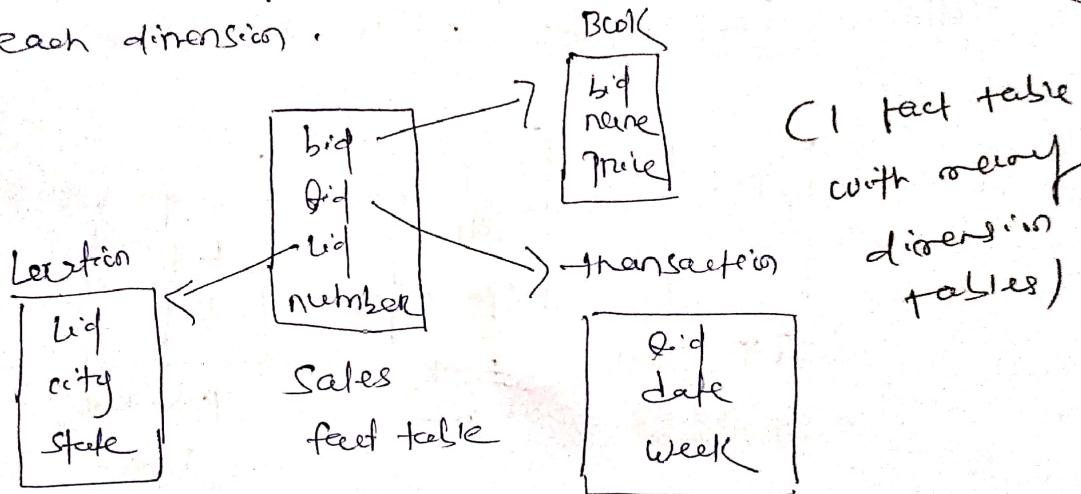
• Spatial data:- large amount of nistion of data

### Data warehouse schema

→ schema is a collection of database objects including tables, views, indexes etc.

#### 1) star schema

It consists of fact table with a single table for each dimension.

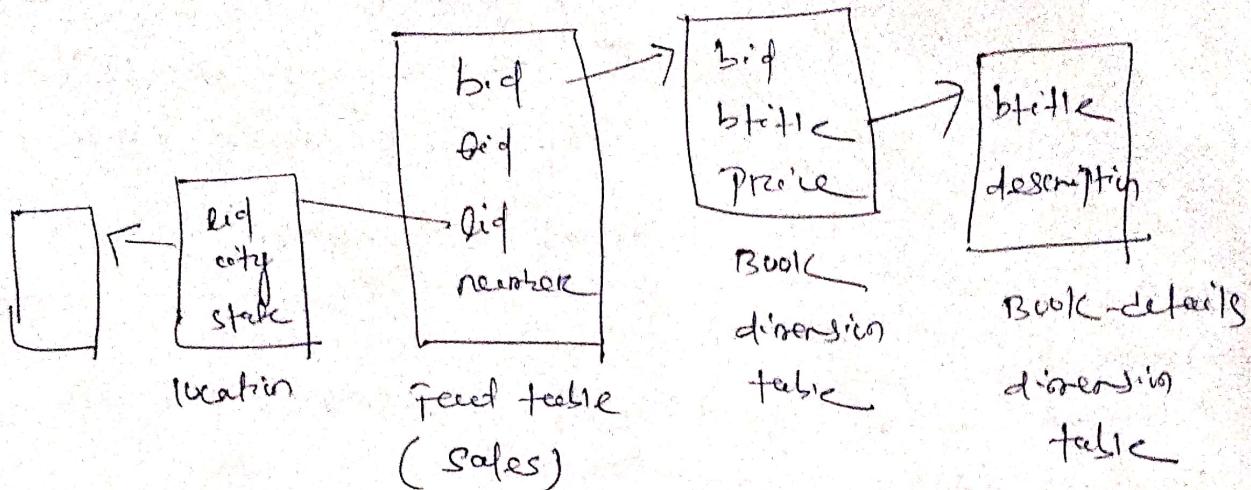


#### 2) snowflake schema

→ It is a variation of star schema which have multi-level dimension tables.

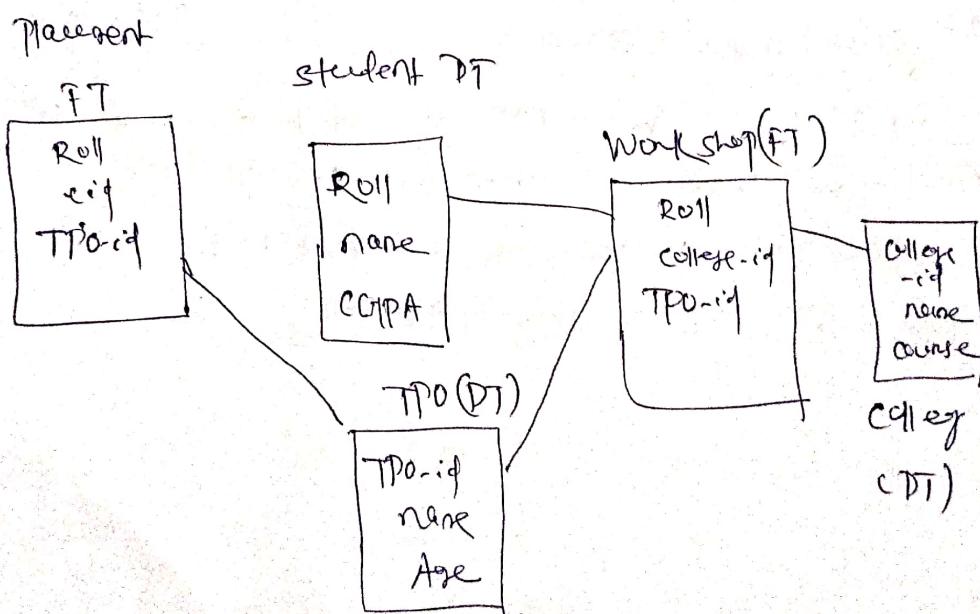
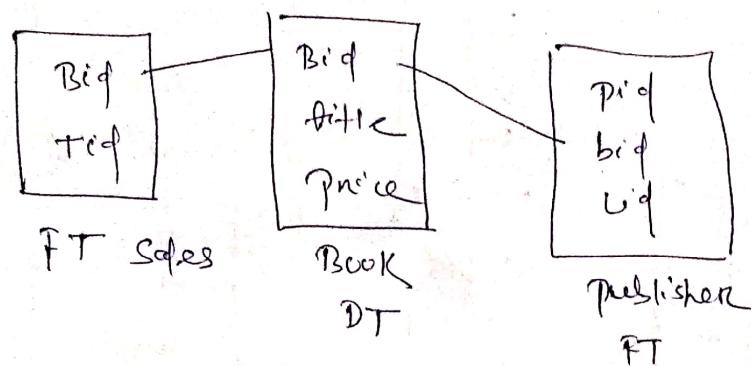
→ Dimension tables are normalized.

→ It is like ER-Diagram.



### fact constellation schema

- Also known as galaxy schema.
- Multiple fact tables share the dimension tables.



## OLAP vs OLTP

→ OLAP (Online Analytical Processing) is a category of software tools that analyze data stored in a data warehouse.

Ex A company might compare their mobile phone sales in September with sales in October.

→ OLTP (Online Transaction Processing) supports transaction-oriented applications. OLTP administers day-to-day transactions.

Ex ATM center  
Withdraw of amount / deposit  
Daily transactions of an organization.

→ It controls

-Zation.

### OLTP

→ It manages DB modification. → OLAP is an analysis and data refining process.

→ It consists of large number of short online transactions. → It is characterized by large volume of data.

→ It is an online DB modifying system.

→ It is an online query management system.

→ It uses DBMS.

→ It uses Data warehouses.

→ Databases are normalized.

→ Tables in OLAP DB are not normalized.

### OLTP

- Different transactions are the courses of data.
- It helps to control and run fundamental business task.
- Simple queries are used.
- Processing time is less.
- OLTP must maintain data integrity constraint.
- It focus on insert, delete, insert information from the database.

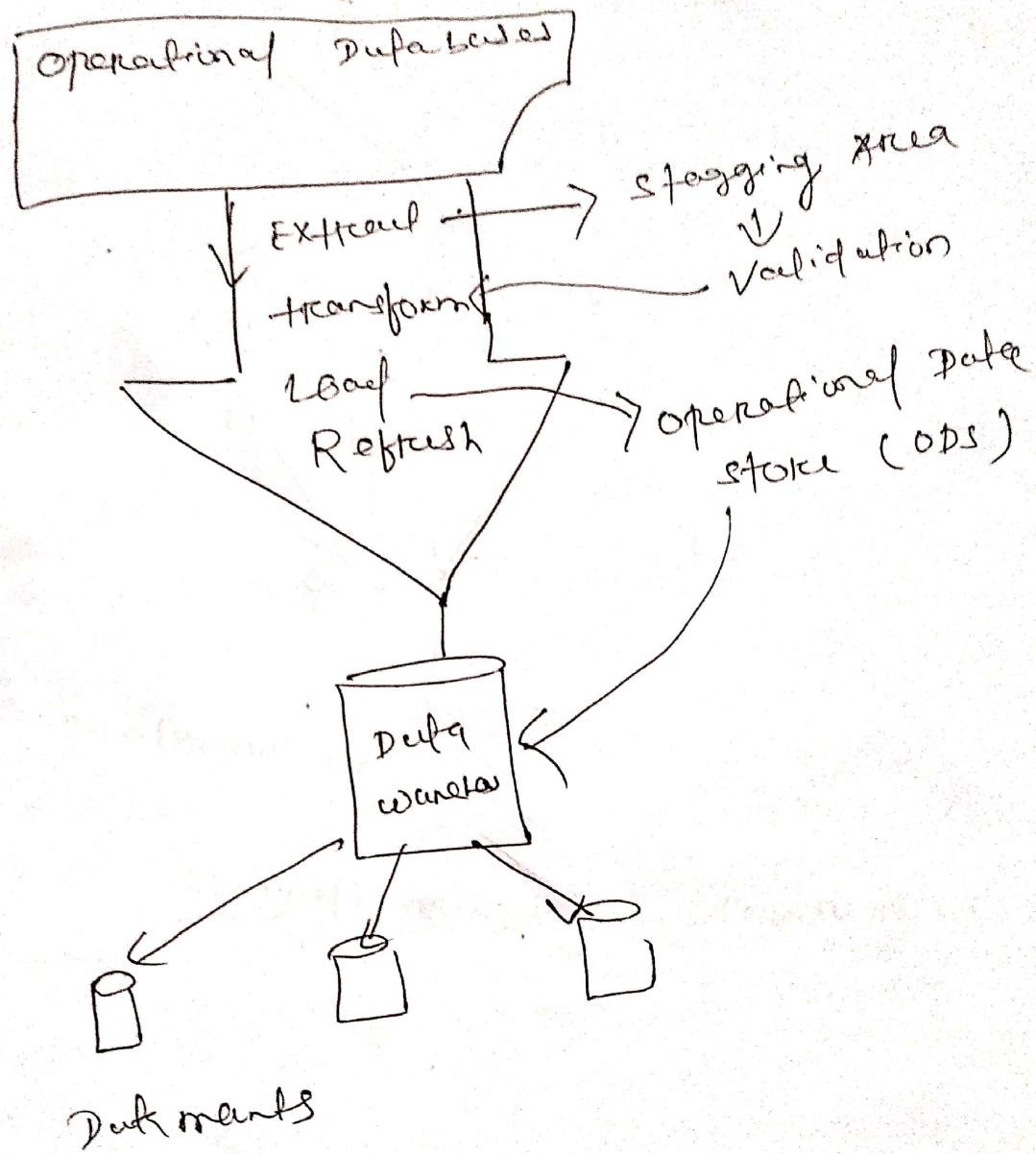
### OLAP

- Different OLTP database uses are the sources of data for OLAP.

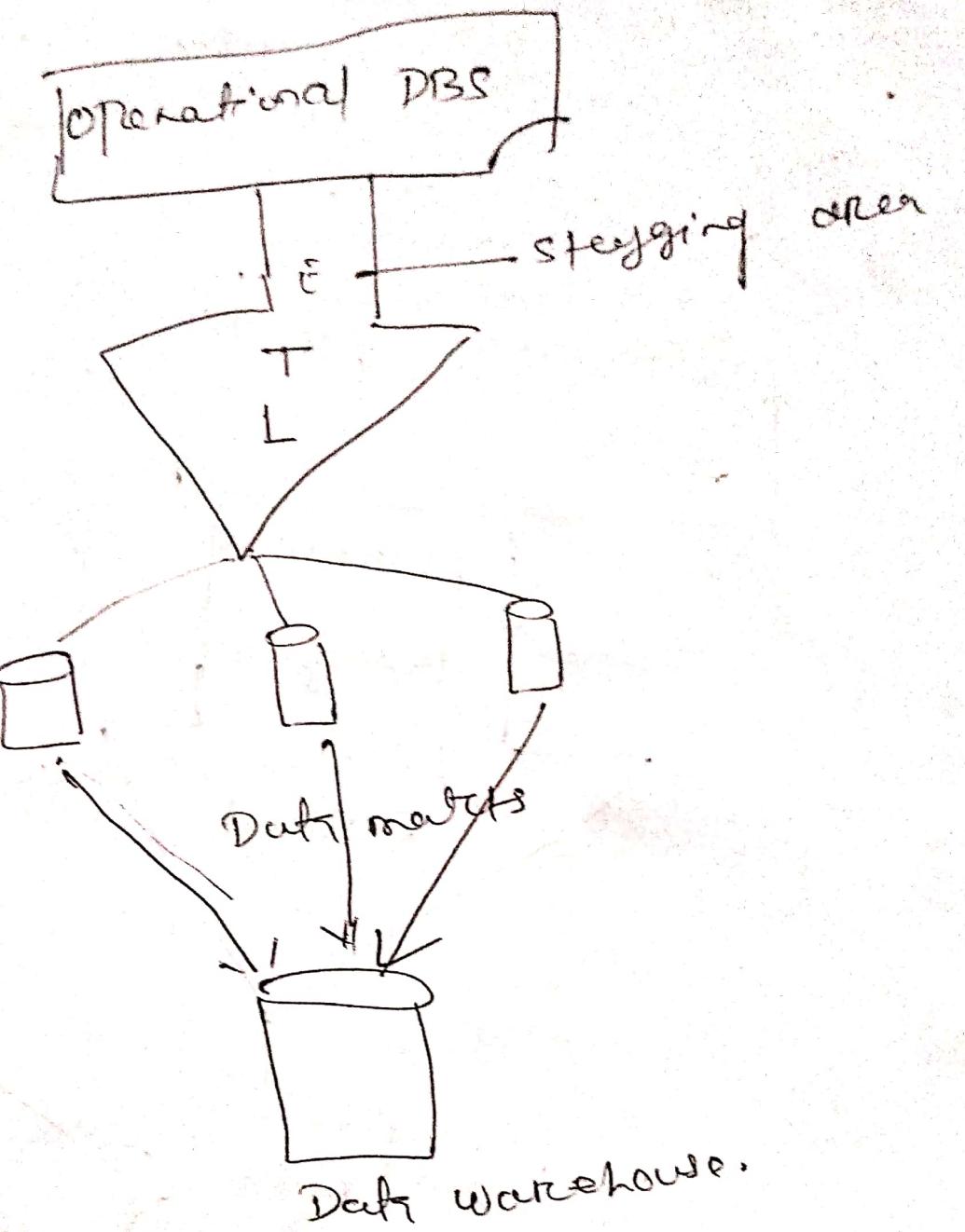
- It helps with planning, problem-solving and decision support.
- complex queries, involving aggregations.
- Processing time is more.
- OLAP DB doesn't get frequently modified. Hence, data integrity is not affected.
- Extract data for analyzing that helps in decision making.

## Approaches for Building DW

### 1) Top-down approach



### 2) Bottom-up Approach



## TOP-down Approach

- It is considered as data driven approach.
- first data is collected & integrated.
- After integration data marts are created.

### Steps

- Raw data are collected from external sources.
- store the data into staging area.
- Apply ETL operation.
- Then load into Data warehouse.
- Then we can create data marts which stores the information of subject wise.

### Advantages

- Developing new data mart from Datawarehouse is very easy.
- Data marts provides consistent dimensional view of data marts.

### Disadvantages

- we can't change the departmental data when required.
- The cost of implementation is high.

## Bottom-up Approach

- first data is extracted from external sources.
- Then store into staging area and loaded into Data marts.

- Then first data marts are created.
- Now, data marts are integrated into Data Warehouse.
- Then apply data mining techniques.

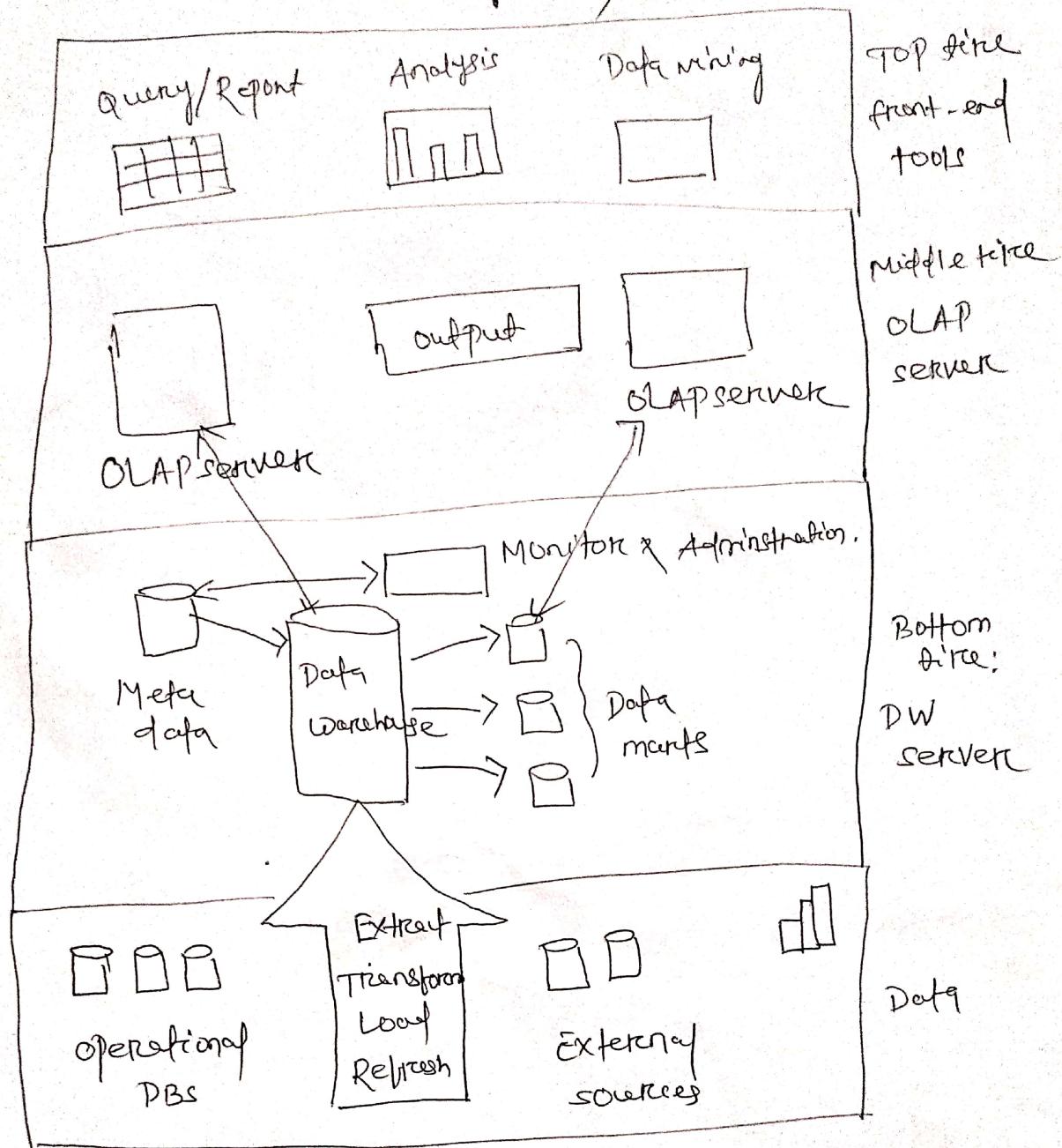
### Advantages

- Reports are generated quickly from the data marts.
- Data warehouse can be extended by adding more number of data marts.
- Cost & time taken in designing this model is low.

### Disadvantages

- The data marts are not consistent, so the dimensional view of data are inconsistent.

# Data warehouse Architecture (3-tier Architecture of DW)



Bottom tier :-

- The database of datware house servers are the bottom tier.
- It is usually a relational Database systems.
- Here, Data is cleaned, transformed & loaded into this layer.

## Middle layer:-

- This layer contains OLAP servers which is implemented using ROLAP or MOLAP.
- This layer presents an abstract view to the user.
- This is an interface between end-user and the data warehouse.

## TOP-layer:-

- This layer is for end users.
- It contains tools, APIs that can help to refine data from DW.
- The tools available at this layer are
  - Query & Reporting tools.
  - Analysis tools
  - Data mining tools
- Basically it is a front-end client layer.