

DWDM

Unit-1

Data :- Data is a collection of a distinct small unit of information. It can be used in a variety of forms like text, numbers, media, bytes, etc. It can be stored in pieces of paper (or) electronic memory etc.

DB :- A database is an organised collection of data, so that it can be easily accessed and managed. You can organize data into tables, rows, columns and index it to make it easier to find relevant information.

DWB :- A Data warehouse is a Subject oriented, integrated, time variant and non-volatile collection of data from heterogeneous sources to provide meaningful business insights.

A DW is a separate from DBMS, it stores a huge amount of data, which is typically collected from multiple heterogeneous sources like files, DBMS, etc.

need for DW :- An ordinary Database can store MBs to GBs of data and that too for a specific purpose. For storing data of TB size, the storage shifted to Data warehouse.

Eg :- Social media websites :- like Facebook, Twitter etc. These are based on analyzing large data sets.

These sites gather data related to members, location etc.

* Banking - Most of the banks these days use warehouses to see the spending patterns of account/cardholders.

Government → Tax payments.

Difference b/w DW and operational DB :-

	DW	ODB
Basic	A dw is a repository for structured, filtered data that has already been processed for a specific purpose	These are those database where data changes frequently.
DS	DW has de-normalized schema	It has normalized schema
	It is fast for analysis query	It slow for analysis queries.
	It is used for OLAP.	It is used for OLTP

Subject oriented → customers, emp, products

Characteristics
of
DW

Integrated → Combined → It allows for easy access and analysis of the data.

Time-Variant → data is stored with a time dimension, It allows easy access to data

Non-volatile → never updated (&) deleted only added to specific periods, such as last quarter (&) last year.

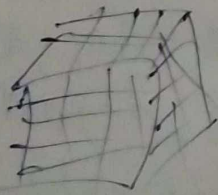
Benefits & DW:-

- 1) Better business analytics:-
2. Faster Queries - It is designed to handle large queries that's why it runs queries faster than the DB.
3. Improved data Quality:-
4. Historical Insight:-

Top layer:-

Front end tools:- Users & DW & use these tools
 meta data \rightarrow Data about the data.
 Data mart \rightarrow Sub part of the DW/
 Sequence of the DW.

OLAP \rightarrow Online Analytical process
 It is used for Analyze the multidimensional data.



Top layer:- It is communicate
 b/w the users and DW.

Single Tier architecture:- It is not periodically used in practice. Its purpose is to minimize the amount of data stored. Only layer physically available is the source layer. In this method, dws are virtual. This means the dw is implemented as a multidimensional view of operational data created by specific middleware (or) an intermediate processing layer.

Two-tier architecture:-

↓
The requirement for separation plays an essential role in defining the two-tier architecture for a dw.

→ It is typically two-tier architecture to highlight a separation b/w physically available sources and dw.

4 ^{How} stages

① Source layer:- A dw system uses a heterogeneous source of data that data is stored initially to corporate relational databases.

② Data Staging:- The data stored to the source should be extracted, cleaned to remove inconsistencies and fill gaps and integrated to merge heterogeneous sources into one standard schema.

③ DW layer:- Information is saved to one logically centralized individual repository: a dw. The dw can be directly accessed, but it can also be used as a source for creating a m, which partially replicates dw content and designed for a specific enterprise department. Meta data repositories store information on sources, access procedures, data staging, users, data mart schema and soon.

Analysis:- In this layer, integrated data is efficiently and flexibly accessed to issue reports, dynamically analyze information and simulate hypothetical business scenarios.

Feature aggregated information navigators, complex query optimizers and customer-friendly GUIs.

Reconciled layer is that creates a standard reference data model for a whole enterprise & it is also directly used to better some operational tasks such as producing daily reports.

Components of the dw: - A database serves as the foundation of your dw. Traditionally, these have been standard relational databases running on premise (or) in the cloud. But ~~be~~

because of Big data, need for true, real-time performance, drastic reduction in the cost of RAM, in-memory db are rapidly gaining access tool. It allows users to interact with the data ^{in PoPs} _(table)

in your dw. Eg: - query & reporting ~~time~~ tools, Appl'n development tools, dm tools, and OLAP tools

ETL process: - ~~Extract~~ Extract: - The 1st stage in the ETL process is to extract data from various sources such as transactional systems, spreadsheets and flat files.

This step involves reading data from the source systems and storing it in a staging area.

② Transform: - The extracted data is transformed into a format that is suitable for loading into the dw. This may involve cleaning and validating the data, converting datatypes, combining data from multiple sources and creating new data fields.

③ Load: - After the data is transformed, it is loaded into the dw. This step involves creating the physical data structures and loading the data into the warehouse.

Adv:-

Improved the quality.
Better data integration
Increased data security
Improved scalability
Increased automation

Disadvantages

- ① High cost
- ② Complexity
- ③ Limited flexibility
- ④ Limited scalability
- ⑤ Data privacy concerns

Disadvantage:-

Data modeling:-

Process of designing schema of detailed and summarized information & dw.

Designing Schema

need:- In order to improve the efficiency & a sys.

Support complex queries.

3 levels:-

1. Conceptual:- It is mainly explain the semantics & the data ^{meaning}
2. Logical:- It defines ^{all} the information ^{data} in structural format.
It contains such as DS, procedures, rules, relⁿ so many things.
3. physical:-

It is describes how the data is presented in dw?
like tables, columns, foreign keys.

DW Models:- ① Enterprise warehouse → ^{go to the} implementing mainframes, super servers.

It contains all information/data about the subject related to entire organisation (detailed and summarized)

② Data mart:- It contains data specific to group & users (not entire organisation) only summarized data

↳ Independent ② dependent

③ Virtual warehouse

It contains data copied from multiple sources during a production.
→ data can be easily search.

TKR COLLEGE OF ENGINEERING & TECHNOLOGY

② It is based on multidimensional data model and allows the user to query on multi-dim.

OLAP ^{cube} ~~operation~~: It is the technology that allows users to analyze information from multiple systems at the same time. These cubes are known as Hypercube.
OLAP is a data structure for quick data analysis.

ETL \Rightarrow data is loaded onto OLAP cube.

after that different operations performed on the OLAP cube: - analysis of the data.

It can be performed even using the spreadsheets also.

Generally, operations can be performed by the spreadsheets in two-dimensional data but DWH is multi-dimensional data.

OLAP operations: - 4 operations: -

- ① Roll-up (drill-up)
- ② Drill-down
- ③ Slicing & Dicing
- ④ Pivot

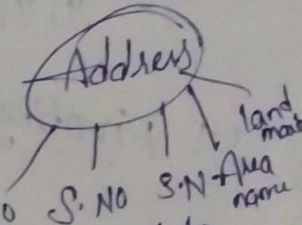
- It uses the concept of hierarchy

\rightarrow high data abstraction

③ When Roll-up is performed, more dimensions from the data cube removed.

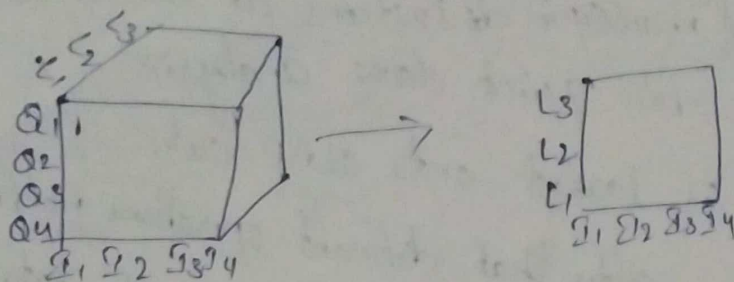


② Drill-down: - Reverse of roll-up operations. data is drilled down. low data abstraction.

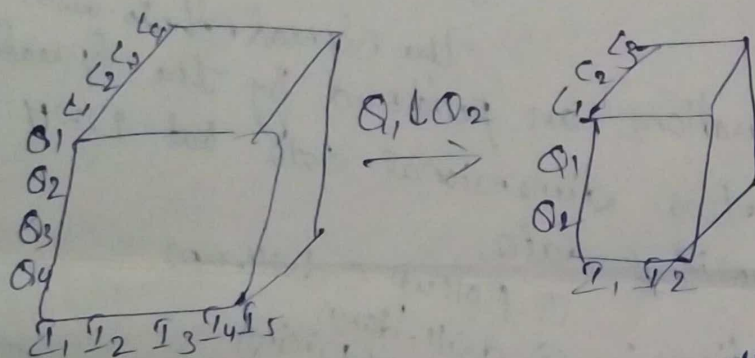


3. Slicing and dicing -

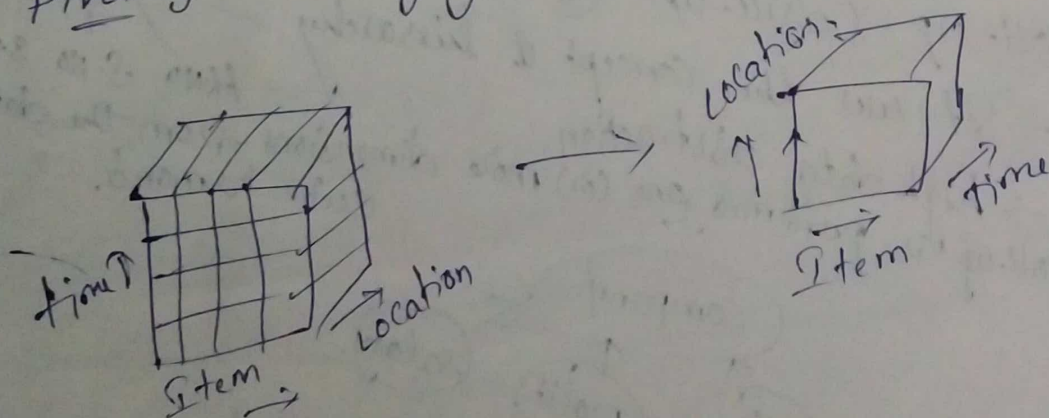
- It is the selection technique.
Data cube is sliced and information is divided and new cube is formed.



Dicing - more than one dimension one is considered.



4) Pivoting - changing the axes of the data (Rotation).



OLAP based on the multidimensional data model.

OLAP implement the multidimensional analysis of business information.

OLAP support the capability to complex estimation treat analysis and sophisticated data modeling.

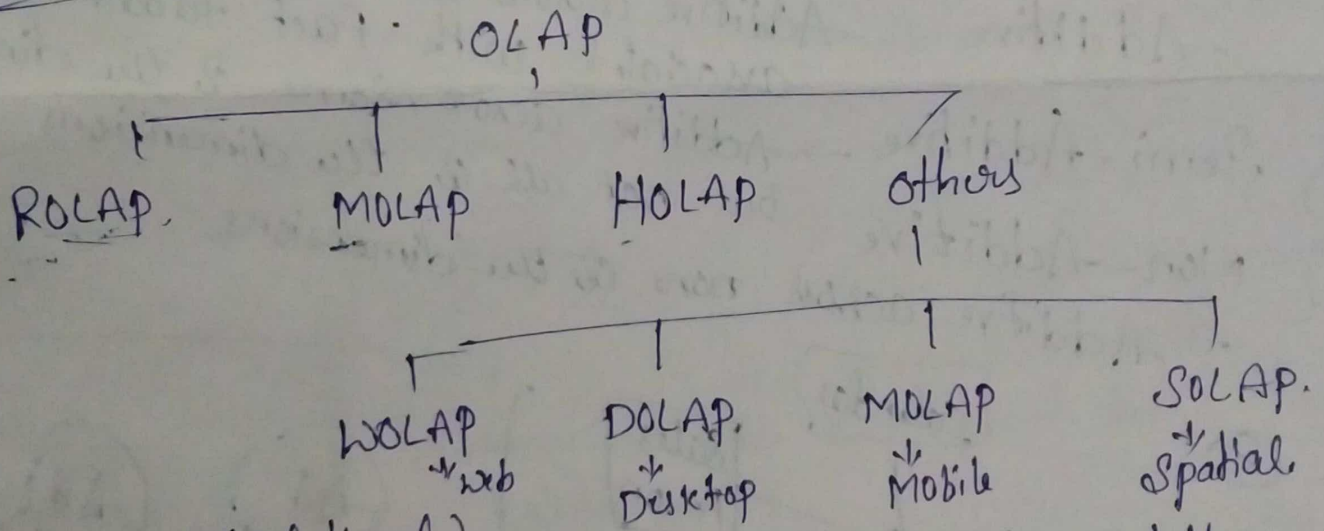
OLAP Applications are used in:

- ① Finance and Accounting
- ② Sales and marketing
- ③ production

OLAP characteristics:

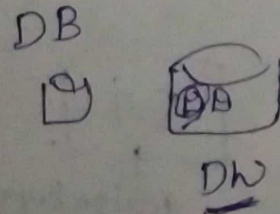
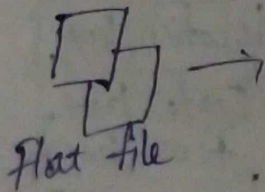
1. useful in analyzing the business.
2. It focus on information out.
3. contains historical data.
4. No. of users is in hundreds.
5. no. of record accined is in millions.
6. Database size is from 100 GB to 1TB.
7. Highly flexible.

OLAP Architecture:



ROLAP: (Relational)

- extension of RDBMS. → will be represented in tables.
- capable of large amount of data analysis.
- can store and analyse ~~large~~ changing amount of data also



2. MOLAP - (Multidimensional)

- Limited amount of data.
- Arrays are used.
- Fast information retrieval.
- Complex calculations.

It is not extend of RDBMS

↓
1GB
2GB

3. HOLAP (Hybrid) :-

ROLAP + MOLAP.

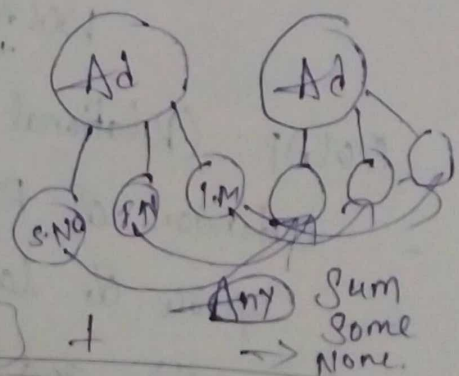
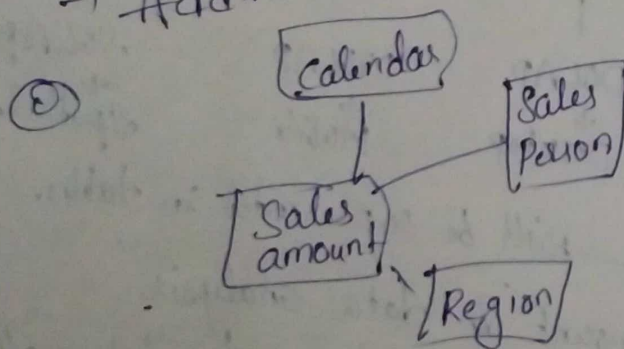
Here, Detailed data and fast information retrieval.
It is large amount of data.

Types of measures in fact table :-

Additive :- Additive across all the dimensions associated with fact table.

Semi-Additive → Additive across some of the dimensions but not all of the dimensions.

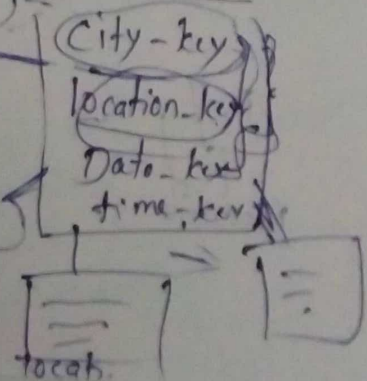
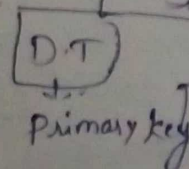
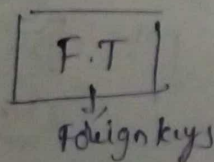
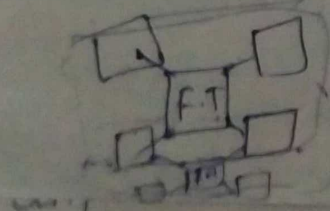
Non-Additive → Additive across none of the dimensions.



Star Schema →

Snowflake →

Fast Constellation :-



will be represent in structural format

→ Collection of database objects.

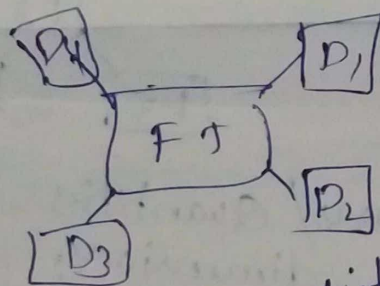
Schema :- It is a logical description of the entire database. It includes the name and description of records of all record types.

① Star Schema.

② Snowflake

③ fact constellation.

④ 2 Important components - FT & DT
only 1 fact table & many dimensional tables.
will contains foreign keys → primary keys
→ every dimensional table.



④ foreign keys.

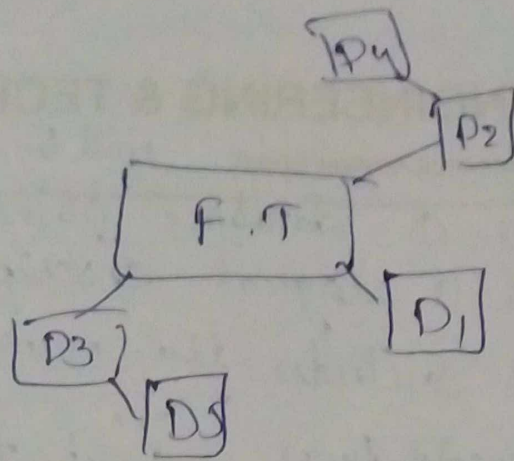
Each foreign key associated with dimensional table.

Each dimensional table ~~having~~ will have the primary key.

Here, fact table can be normalized → removing the redundancy in the data.
Dimensional table can't be.

② It is a type of Star schema where dimensional tables are also normalized.

How → By existing dimensional table you will be deriving the new D.T.
→ It will manages size of the F.T.



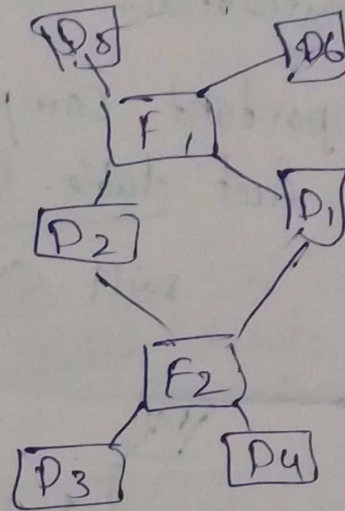
D2, D1, D3 is the existing ones.
 D5 is ~~is~~ derived out of D3.
 D4 is derived out of D2.

3) Fact Constellation Schema:-

It is Group of Star Schemas.

- Complicated design.
- Hard to understand and implement.

Many fact table & DT.



Fact:- numerical measures/quantities by which can analyze relationships b/w dimensions.

Dimensions:- collection of logically related attributes used to modelling the data.

Fact table:- relations b/w multi dimensional data.

Dimension table:- Tables related to each dimension and helps in describing dimension further.

Types of