

ASSIGNMENT - 01

Q1. Explain features of Datawarehouse in detail.

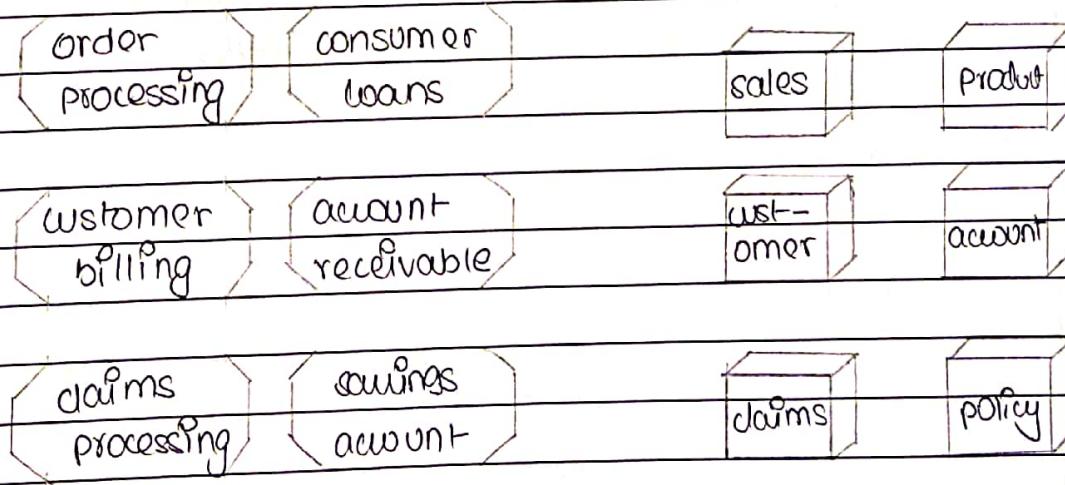
- Data warehouse is a subject oriented, integrated, non-volatile and time variant collection of data in support of management decisions.

Q2. Subject oriented -

- In database, data is stored in the form of application whereas in datawarehouse, data is stored in the form of subjects. As every thing is available in the form of subjects, any complex query can be answered with less processing time and it will provide 360° view of the organisation.

Operational Applications

DW subjects



- As shown in the figure, in the operational systems shown data for each application is organized separately by application.

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

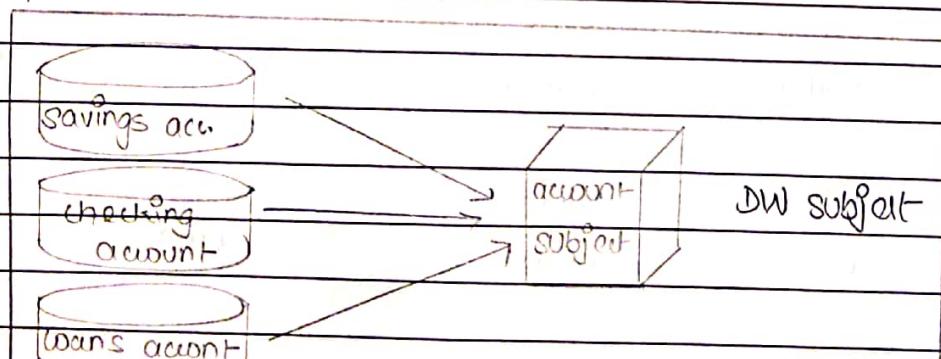
JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

- whereas in data warehouse, data sets contain data needed for functioning of the particular applications as subjects.

2o Data integration?

- For proper decision making, all the relevant data from various applications are to be pulled together.
- The data in the DW comes from various & operational systems. Source data are in different database files and data segments.
- Before data from various sources can be usefully stored in the DW, the inconsistencies have to be removed. and data elements have to be standardized. Before moving data into the DW, a process of transformation, consolidation and integration of data source is done.
- Standardization includes - naming conventions, codes, data attributes, measurements.

example,



data from
applications.

3. Time variant data :

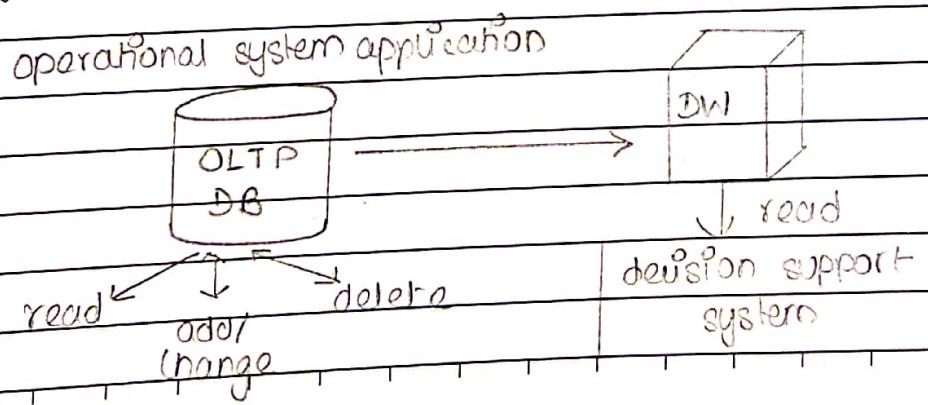
- for an operational system, stored data contains the current values.
- A DW because of its very nature of its purpose has to contain historical data, not just current values.

Data is stored as snapshots over past and current periods. Every data set in DW contains the time elements.

4. - For example, in DW containing units of sales, the quantity stored in each file record relates to the specific time element. Depending upon the level of detail, the sales quantity in a record may relate to a specific date, week, month, or year.

4. Non-volatile data :

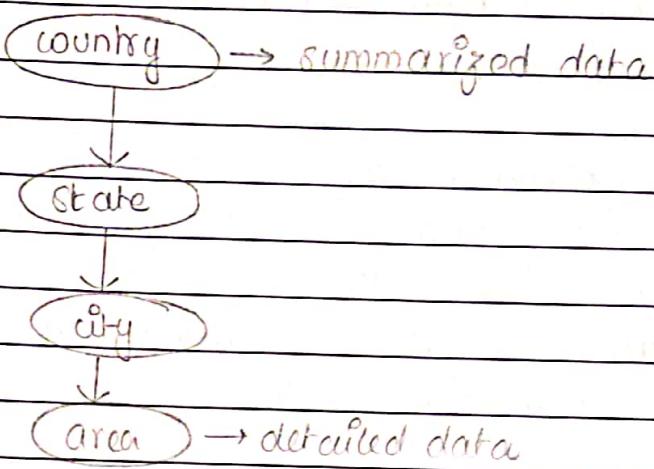
- The data in the DW is not intended to run the day to day business.
- Data from operational systems are moved into the DW at specific intervals depending upon the requirements.
- In DW only a single operation i.e. read operation is performed and therefore data items are preserved for future analysis.



5. Data granularity:

- Data granularity in DB refers to the level of detail. The lower the level of detail, finer the data. granularity.

- example,



- The granularity levels has to be decided on the basis of data types and the expected system performance for queries.

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI
MCT
 MANJARA CHARITABLE TRUST
 JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

Q20

Differentiate between :

- 10 Top down and bottom up approach.

Top down

- DW is created first and then individual data marts are created.

- Answers Answering the question is easy as compared to bottom up approach.

- As DW is created first, technical/ professional skills are needed.

- It is time consuming as iterations are needed to build DW.

- Risk failure is higher

Bottom up approach

- Individual data marts are created first and then the whole view of DW is created.

- Answering the question takes more time as compared to top down.

- Users with less cross-functional skills can build data marts.

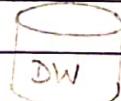
- It is less time consuming as data marts are created first.

- Risk failure is lower.

operational

SOURCES

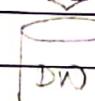
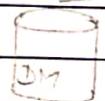
↓ ETL



operational

Data source

↓ ETL



2. ER modelling and dimensional modelling.

ER modelling

Dimensional modelling

- It is used for designing tables from database.
 - It is used to specify relationship between entities.
 - The main aim of ERM is to remove redundancies by normalising tables.
 - It is well suited for answering queries at transactional level.
- It is used in order to create DW.
 - It is used to focus on business analysis.
 - Most of the tables of dimensional modelling are denormalised.
 - It is used to answer queries at strategic level.

3. OLTP and OLAP:

OLTP

(online transactional processing)

- It is used for performing day to day operations.
- It is application oriented.
- Detailed data is stored.

OLAP

(online analytical processing).

- It is used for analysis and decision purpose.
- It is subject oriented.
- Summarized data is stored.

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

MANJARA CHARITABLE TRUST
JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

Page No. 4

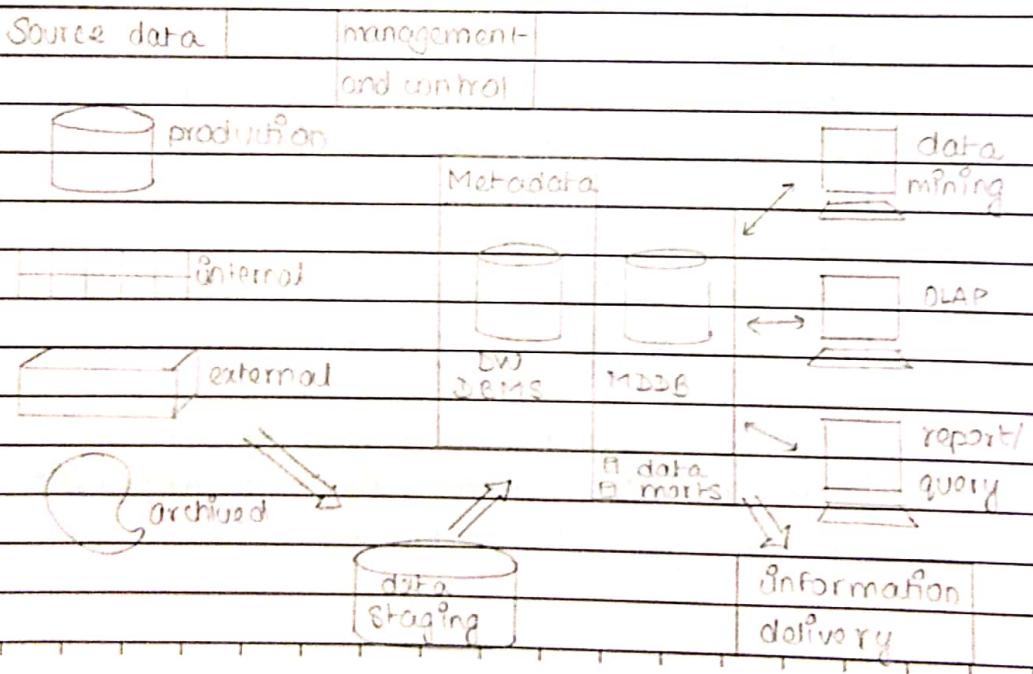
OLTP

- Read, write and update operations can be performed.
- It is volatile.
- Technical users are involved in this.
- Data here is not time variant.
- Tables are normalized.
- Size of OLTP is below 100MB - Size of OLAP ranges from 100 GB to few TBs.
- Queries are simple.

OLAP

- Only read operation can be performed.
- It is non-volatile.
- Executives and knowledge workers are involved.
- Here, data is time variant.
- Tables are in denormalized form.
- Queries involved in this are comparatively complex.

Q3: Describe architecture of Data warehouse with diagram.



- Source data: Input to the DW is taken from various source applications and categorised into 4 parts which are as follows:

1. Production data - In this, data is collected from various operational systems of organisation.
e.g. - employee table may be stored in account department, HR department and payroll department. For data to be stored in DW, it must be collected from all these systems. This is production data. Problem arises when related data is not in standard form. So, standardisation is done before production data is stored in DW.
 2. Internal data - It is available within the organisation. This is information which must be kept confidential.
e.g. personal details of employees.
 3. Archived data - Operational system is used to run daily transaction. So, old data from it is periodically taken out and stored in archived files. Thus, archived files are created for storing historical data.
 4. External data - It is the data which is not available within the organisation and needs to be gathered from outside sources for analysis.
- Data staging component - Data staging area is the place where all the extracted data is temporarily stored and prepared for loading in DW.

Three major functions are performed:

1. Extraction - Source applications are extracted from and placed in data staging area.
2. Transformation - data cleansing and standardisation is done here.
3. Loading - Initial data is loaded in a data warehouse in a single run and then updated periodically.

- Data storage component: It is a separate repository. In this large amount of historical data is stored for analysis. It is designed specifically for analysis purpose and not for quick retrieval of information.

- Data marts: Data mart is a decision support system that stores number of subject areas based on the needs of the users in that department.

Thus, data mart can be thought of as a subset of the enterprise wide DW.

Eg: Finance has its own department, marketing has its own dept. and so on. Thus, each dept. will have its own Data mart.

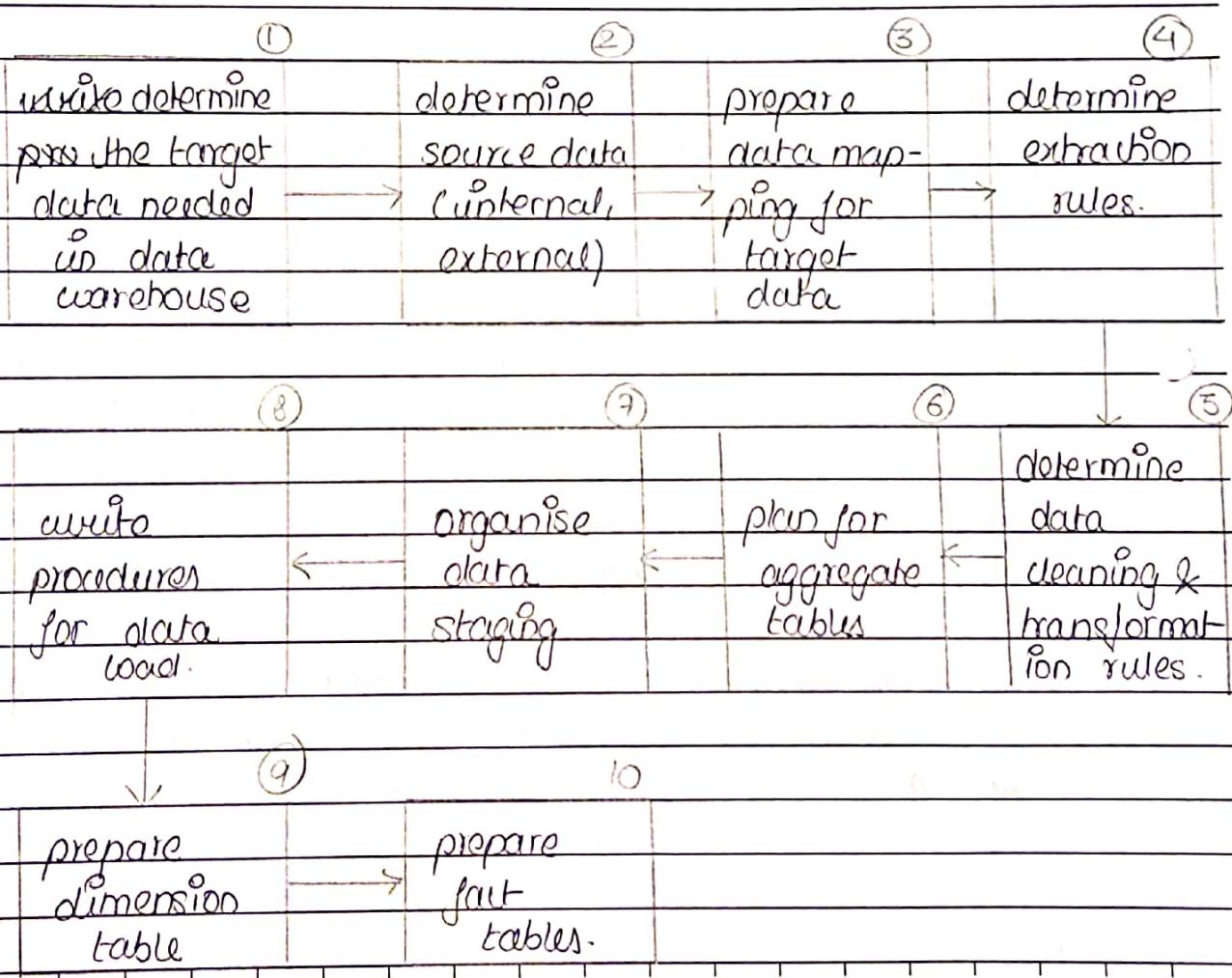
There are 2 approaches -

- a) top down approach
- b) bottom up approach.

- Metadata component: It is data about data in DW. It is similar to the data dictionary or data catalog in DBMS.

- Management and control component: It sits on the top of all components and manages all the activities of DIO. It helps data staging area to perform ETL activities.
- Information delivery component: It is used to provide info. to the wide variety of users with the help of IDC, info. of DW can be used for various purposes like solving complex queries to do prediction and pattern analysis.

Q4. Explain ETL process.



MGT
MANJARA CHARITABLE TRUST
RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI
JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

- The ETL (Extraction Transformation and Loading) process plays an important role in DW architecture. Without ETL process there would be no strategic information present in the DW for future analysis.
- ETL process starts with determining the target data (what output we want for future analysis).
- Once, the target data is set, source identification takes place, in which internal & external sources are identified from where input must be taken.
- The input is then mapped with the output.
- Once, mapping is done, extraction methods are applied.
- After extraction, data is cleansed and transformed into standardized format.
- In next step we decide whether to keep single fact table or plan for aggregate fact tables.
- After planning dimensional modelling, data staging area is organized.
- Different procedures are written to load data in DW.
- All the subjects are defined in the terms of dimensional tables and fact tables.

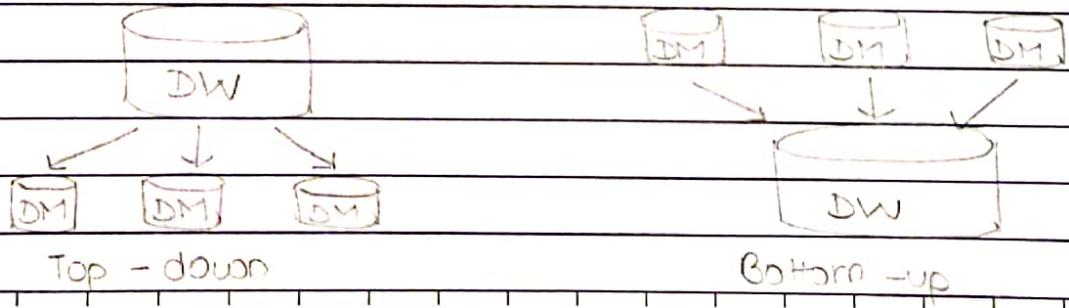
RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

Q5. Write short note on the following-

a) Data mart :

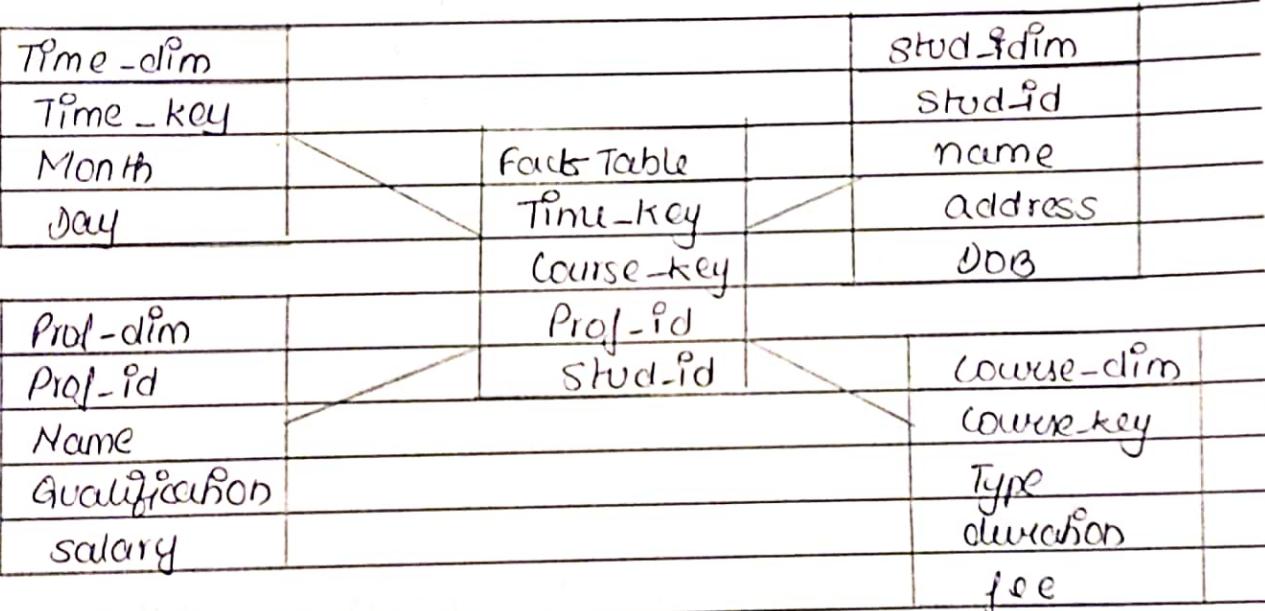
- Data mart is a decision support system that stores number of subject areas based on the needs of the users in their department.
 - Data mart can be thought of as a subject of the enterprise wide data warehouse.
 - e.g. the finance department has their own data mart, likewise every department will have its own data mart.
 - There are two approaches for building data mart -
 - top down approach
 - bottom up approach
- a) Top down approach - DW project team looks at the larger picture of the organization and builds huge DW first that will feed the individual data marts.
- b) Bottom up approach - DW project team provides the requirements of individual department and builds data mart first that will feed data to the data warehouse. In this data marts are created first.



b) Fact less fact table -

- In factless fact tables, no facts are present in the fact table. Only the primary keys of all the dimension tables are included.
- The job of factless fact table is to record event occurrences or describe conditions.

- example,



c) Star schema and keys involved in star schema:

- Star schema is an arrangement where fact table is placed in the middle and all the dimension tables are directly connected to the fact table.
- In this schema all the dimension tables will have equal chance in query analysis as no intermediate tables are created.

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

example, OTs → ^{prod}Prod-dim, Time-dim, cust-dim, store-dim
 FT → Order-fact

Prod-dim			Time-dim
Prod_id			Time_id
Name	Order-fact		Year
Quantity	Prod_id		Month
Brand	Store_id		Day
	cust_id		
Store-dim		Time_id	cust-dim
store_id		Total-sales	cust_id
Name		Profit-dollars	Name
address			Address

* Keys involved in star schema -

1. Primary key:

In dimension table, keys are used to identify each record separately. In FTs, we have concatenated keys which acts as primary keys for FTs.

2. Foreign key:

Primary key of OTs act as normal columns in FT transformed as foreign key for FTs.

3. Surrogate key:

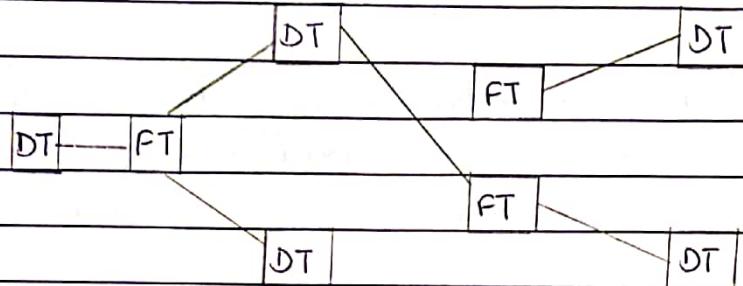
These keys are generated automatically by the software. Primary keys of database and data warehouse cannot be same because,

primary key of DB keeps on changing rapidly.

If conflicts may be present in the DB.

d) Fact constellation schema -

- In fact constellation schema, more than one fact tables are present in the organization.
- No two fact tables are directly connected to each other although they are connected using dimension tables.



- There are four ways by which fact constellation schema can be created.

10. Aggregate FT and derived DT:

In aggregate FT, pre calculated summaries of FTs are stored.

The dimensions which are shared between base FT and aggregate FT are termed as derived DTs.

2. Core and custom table:

Core FT is basically used to store values which are similar in nature whereas custom table is used to store the values which are dissimilar in nature.

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI
JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

5. Snapshot and transactional FT:

In snapshot FT, information based on particular point of time is stored whereas in transactional FT, full view of the organization is stored.

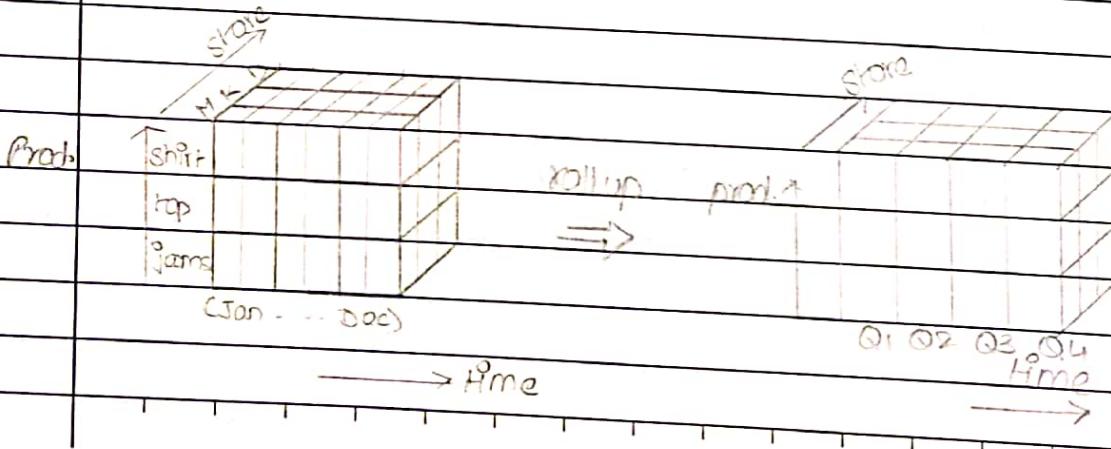
4. Supporting enterprise value chain:

As product goes from raw parts to finished product, it has to go through various stages. In each stage an individual stage schema is maintained.

Q6. Explain OLAP operations in detail.

- There are five types of OLAP operations -
 - 1. roll up
 - 2. drill down
 - 3. slice
 - 4. dice
 - 5. pivot / rotation

10. Roll up: In roll up operation (also known as drill up), we move upwards in the hierarchy from less detailed to more summarised data.



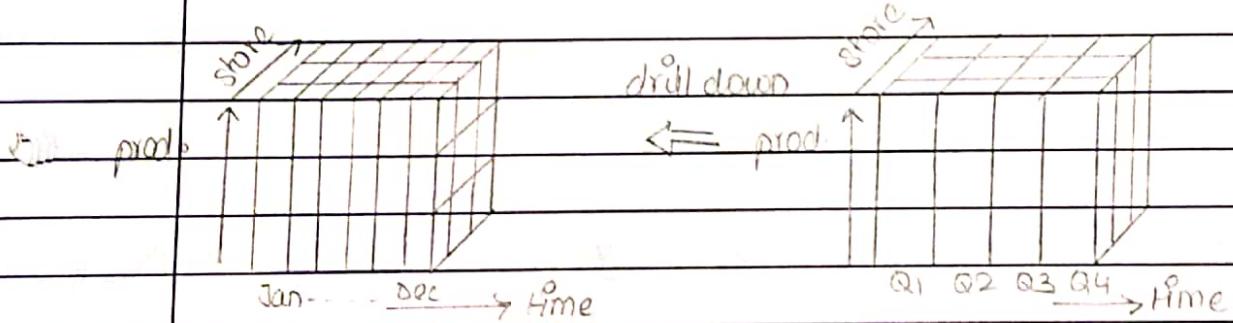
Eg store Mumbai.

Shirt	900				Shirt			
Top	450				Top			
Pants	200				Pants	1300		
	Jan - Dec					Q1	Q2	Q3 Q4

- In this example, time dimension is represented using roll up operation, where time dimension is represented in quarter using the hierarchy as,

Day < month < quarter < half year < year.

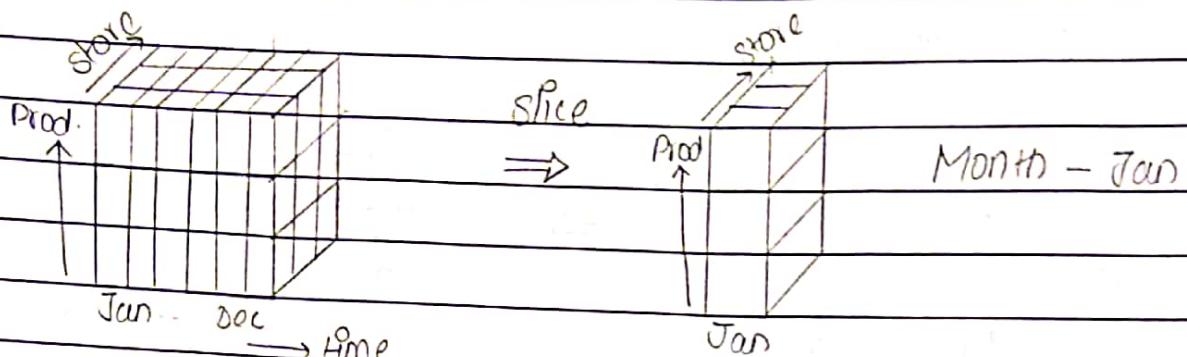
20. Drill down: In drill down operation (also known as roll down), we move downwards in the hierarchy that is from summarised data to detailed data.



30. Slice operation:

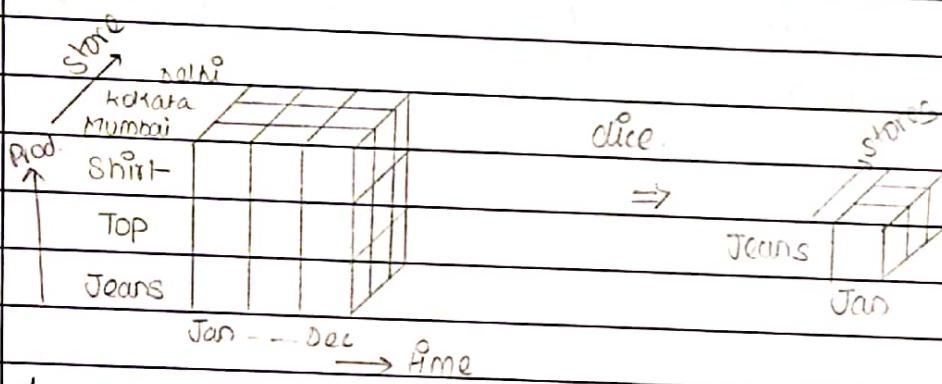
If we apply selection operation on only one dimension, it will result in a sub cube. This operation is termed as slice operation.

MGT
MANJARA CHARITABLE TRUST
RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI
JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.



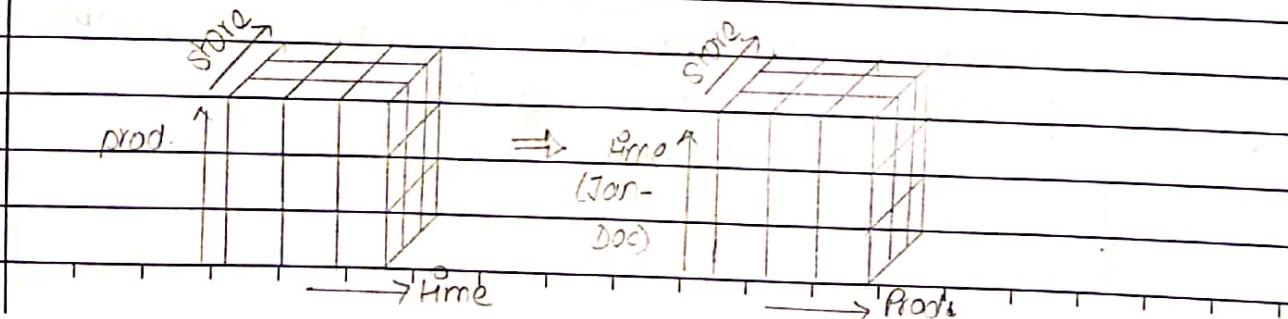
In above example, products from all the stores for the month of January are selected.

4. Dice operation: If we apply selection operation on more than one dimension, it results in a subcube and this operation is termed as dice operation.



In the above example, selection of Jeans for the month of January is shown.

5. Pivot operation: It is obtained by interchanging the dimensions.



RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

MANJARA CHARITABLE TRUST
JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

Q7.

Explain metadata, types of metadata and example describing metadata in detail.

- Metadata is known as data dictionary or data catalog as it describes all the entities of a data warehouse briefly.
It is also called as data about data.
- In most of the organisations, data warehouse is created externally and therefore the users of DW are not aware about the contents of it.
And therefore, along with DW, data dictionary is given which not only describes the entities but also describes the syntaxes and semantics of DW.
- Types of metadata -
 - 1. Operational metadata - It is created for source application and details of input of the DW.
 - 2. Extraction and transformation metadata -
It gives details about the extraction method, frequency and transformation rules.
 - 3. End user metadata -
It is created for the stakeholders using DW for analysis.

4. Business metadata -

It is created for the managers and the analysts.

5. Technical metadata:

It gives the technical details required by the developers.

example,

Sample metadata for customer.

Entity name : Customer

Definition : A client, a person or an organisation that purchases goods or services from company.

Remark : It includes current and past customers

Source systems : Finished good orders, online sales, maintenance contracts.

Create date : 25th April 2015

Last update date : 22nd September 2016

Update cycle : yearly

Data quality reviewed : 20th September 2016

Responsible user : John M.

Q8.

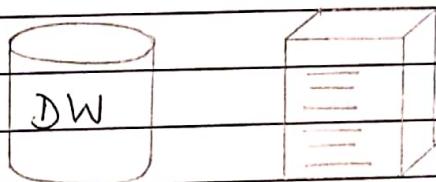
Explain different OLAP models.

There are four OLAP models -

1. MOLAP
2. ROLAP
3. HOLAP
4. DOLAP

10 Multidimensional OLAP (MOLAP):

- In MOLAP, pre calculated cubes are created by the MOLAP engine depending upon the data present in DW. MOLAP engine creates as many cubes as possible and therefore retrieval of information is rapid.
- As many outcomes are created, complexity of MOLAP is more.



data
layers

application
layer

presentation
layer

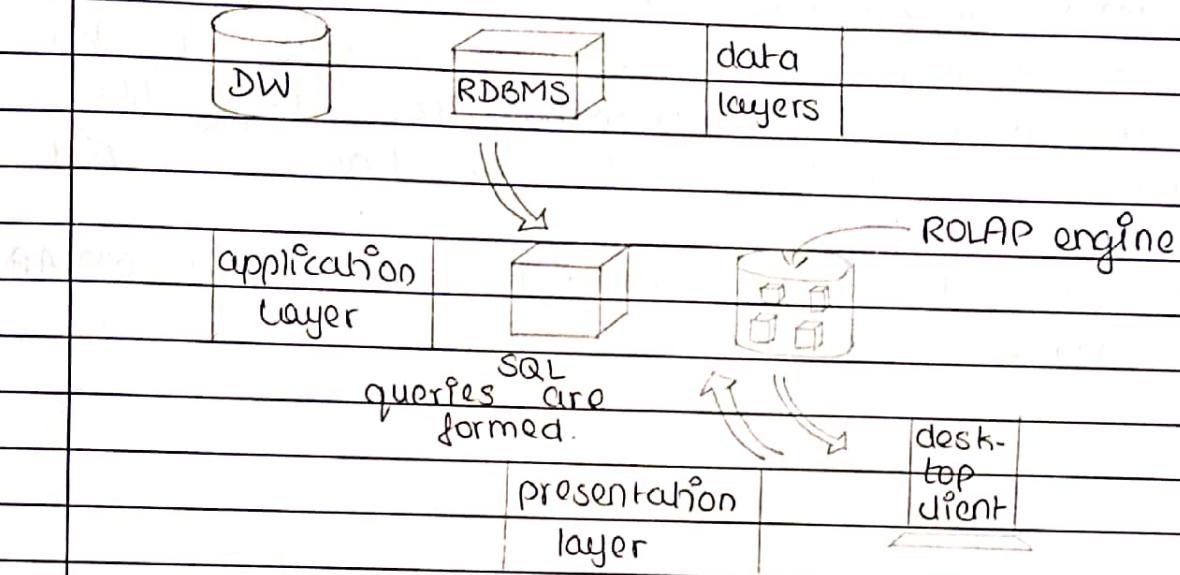
MOLAP engine

desk-
top
client

(Arrows indicate flow from DW to cube, cube to application layer, application layer to presentation layer, and presentation layer to desktop client.)

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI
JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.2o Relational OLAP (ROLAP):

- In ROLAP, pre calculated cubes are not created and cubes are created depending upon the requirements of the desktop client.
- SQL queries are formed in order to get information from DW. ROLAP engine will create cubes based on the requirements, therefore, information retrieval is slower.



3o Hybrid OLAP (HOLAP): It takes the advantage of both the OLAP models i.e. MOLAP & ROLAP and removes the disadvantages.

4o Desktop OLAP (DOLAP):

It programme performs the task of sending the cubes of requested data to the client.