

# ETL

Puja Patel  
10

Kalyanika  
10

Srijanika  
10

CLASSmate  
10-8-22

It is used for decision system.

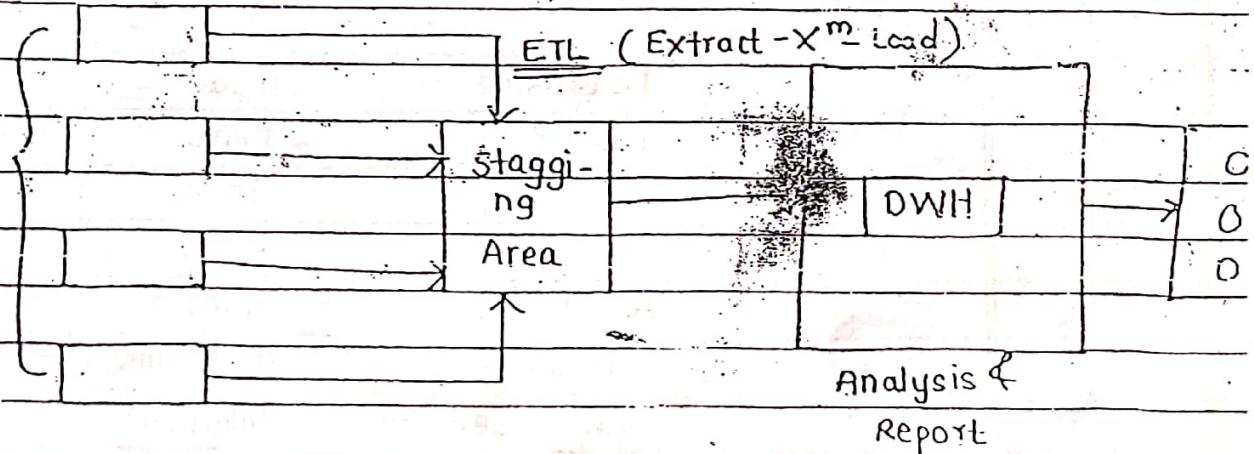
\* DWH - For Analysis & reporting purpose.

Defn → DWH is a DB designed for query & analysis rather than transaction processing.

1. It contains historical data derived from transactional data (operational/current)

2. DWH separates analysis workload from transactional workload and enables to consolidate from various DB (sources)

Defn → DWH is a subj. oriented, integrated, nonvolatile, time varying coll'n of data in support of management decn making process.



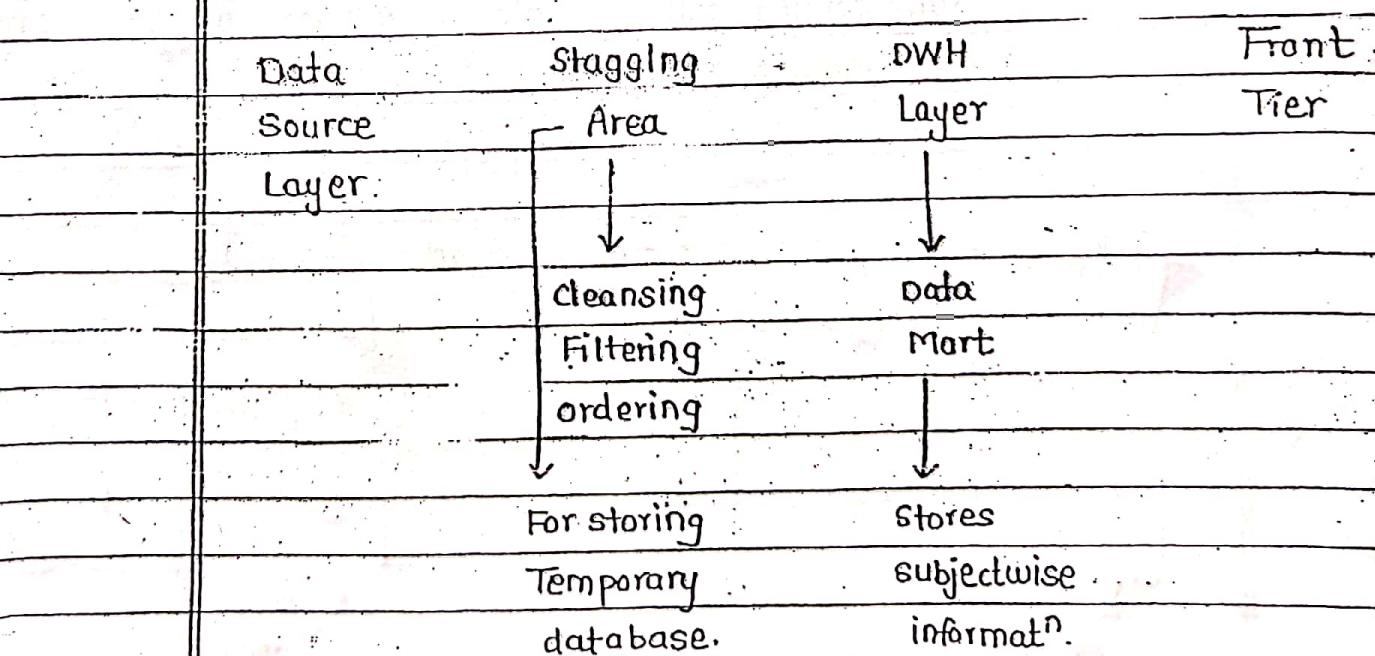
\* Data Ware House  
Business Ware House  
Decn support syst  
Business Intelligence soln

Various names.

IMP

## \* \* ARCHITECTURE OF DWH →

Internal Data Sources	CRM Syst.	S	Meta Data Repository		Reporting Tool	
	ERP Syst.	T				
		G	Data	Data		
External Data Sources	SAP Syst.	G	Mart	Mart	Data mining Tool	
		I				
	Flat files	N	Data	Data		
		G	Mart	Mart		
		DB				



### 1. 1<sup>st</sup> layer — Data source Layer.

- which refers to various data sources in multiple format like relational dB, flatfiles & others.
- contains Transactional / current / operational data.

### 2. 2<sup>nd</sup> layer — Staging Area

- Intermediate area stage betn source to target where required business rules [Transform] are applied.
- This layer takes care of data processing methods like
  - 1. Data cleansing
  - 2. Filtering
  - 3. Merging
  - 4. Splitting, etc.To avoid duplicate data.
- After this all data put it into DWH i.e. DWH Layer/ data store/ Data meta Repository.

### 3. 3<sup>rd</sup> layer — DWH Layer

- In this layer :- cleaned  
Integrated  
Transformed  
ordered data is present in multidimensional environment

#### 4. 4th Layer -- Front Tier

- In front tier data in DWH layer is used for reporting analysis with the help of reporting & data mining tool.

#### \* Staging Area →

- Also called as landing zone.
- Intermediate storage area used for data processing during ETL process.

#### \* Need of Staging Area →

- In the absence of staging area data load will have to go from OLTP to OLAP syst. directly.
- This can hamper the performance of OLTP syst.
- This is the primary reason of existence of Staging area.

#### \* Data Mart →

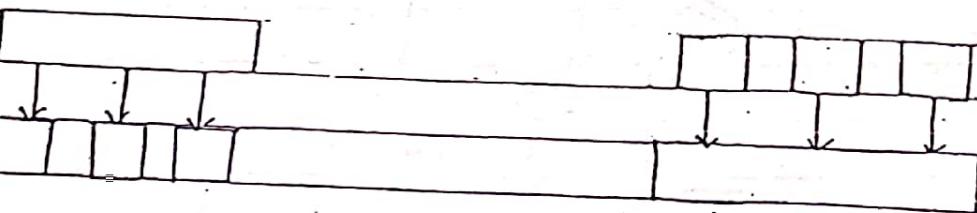
- is a subset of DWH which is limited to specific fun area (subj. area) or group of users.

It is a condensed version of DWH focussed on single fun area at a time.

2 Approaches → Top Down  
Bottom Up

Top Down

Bottom Up



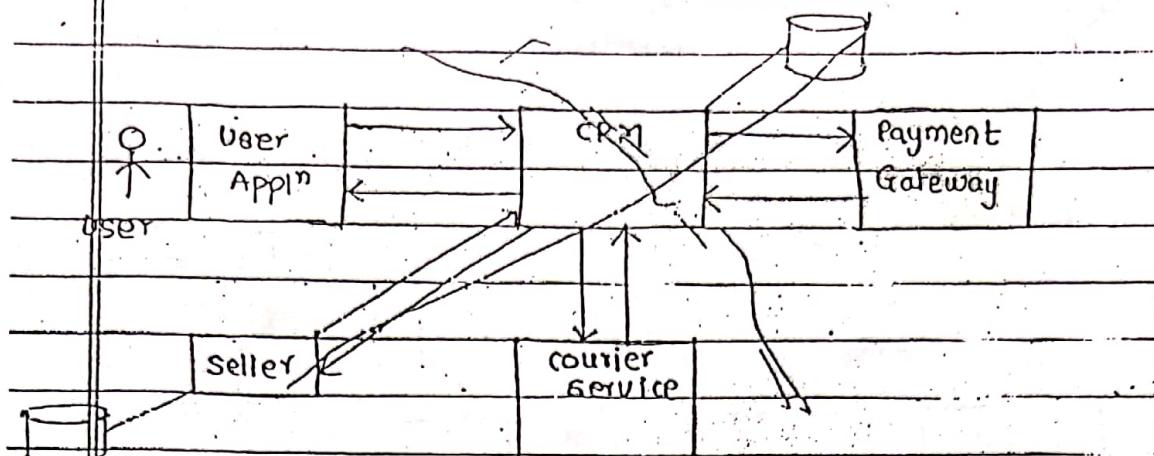
- Data Mart can be build later to form DWH.
- High cost
- Less Time consuming.
- More Planning & designing.
- Data Mart can be build before or in line with DWH.
- Low cost
- less time consuming.
- Less Planning & designing.

Data Ware Housing

Data Mart

- |                                  |                                       |
|----------------------------------|---------------------------------------|
| 1. Union of all data.            | 1. Data of single subj. area.         |
| 2. Implement^n - Time consuming. | 2. Implement^n - less time consuming. |
| 3. Enterprise View               | 3. Departmental view                  |
- ↓  
No. of user more
- ↓  
No. of users less

## \* OLTP VS OLAP —



### I. OLTP :- online Transact" Processing.

- OLTP DB → detailed and current data is stored , of transactional syst. Mostly data is in 3NF form.
- It is mainly used for data processing not for analysis.

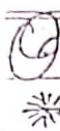
### II. OLAP :- online Analytical Processing.

- It allows user to analyze info. from multiple DB syst. at once.
- Mainly used for data analysis not for processing.
- It provides single platform for all platforms of business needs which include
  - ↓ planning, designing, analyzing info & reporting.

## OLTP

## OLAP

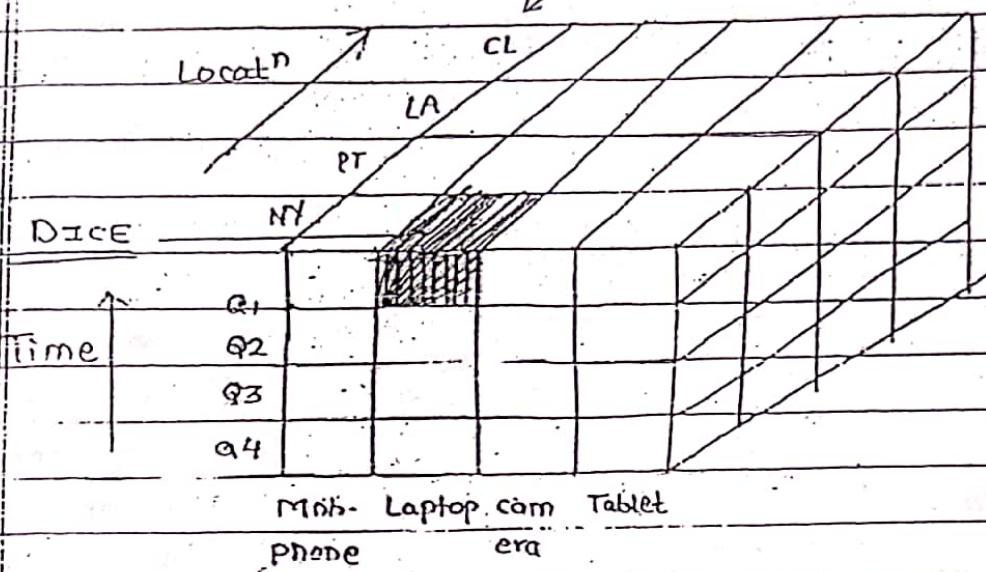
- |  |   |
|--|---|
| 1. Used for Transaction processing.  | 1. Used for Query & Analysis.   |
| 2. Data <ul style="list-style-type: none"> <li>1) current data/operational.</li> <li>2) Data is in normalised form.</li> <li>3) Stores all data.</li> <li>4) Volatile data.</li> </ul> | 2. Data <ul style="list-style-type: none"> <li>1) Historical data / statistics.</li> <li>2) Data is in denormalised form.</li> <li>3) stores only relevant data.</li> <li>4) Nonvolatile data.</li> </ul> |
| 3. To run fundamental business task.   | 3. To help with planning prob. solving, & decn support.   |
| 4. No. of users are more.  | 4. comparatively less users.  |
| 5. Appn driven.  | 5. Analysis driven.   |
| 6. Used to store the data into db.   | 6. used for reading data from DWH.  |



## OLAP Operations :-

Date  
Page

citynames.



### operations :-

Or navigate data from less detail to high detail

1. Drill down. - Add dimension from data cube.
2. Roll up. - Remove dimension from it.
3. Dice
4. Slice
5. Pivot

1. Drill Down :- physically data is present in quarter wise user want it month wise

2. Roll up :- physically data present city wise user want region wise.

User wants

3. Dice → Particular data e.g. user wants data (laptop) in 1st quarter.

4. Slice → Detailing data in 2D. Q1, Q2, Q3, Q4

3D → 2D

	Q1	Q2	Q3	Q4
	NY	PT	LA	CL

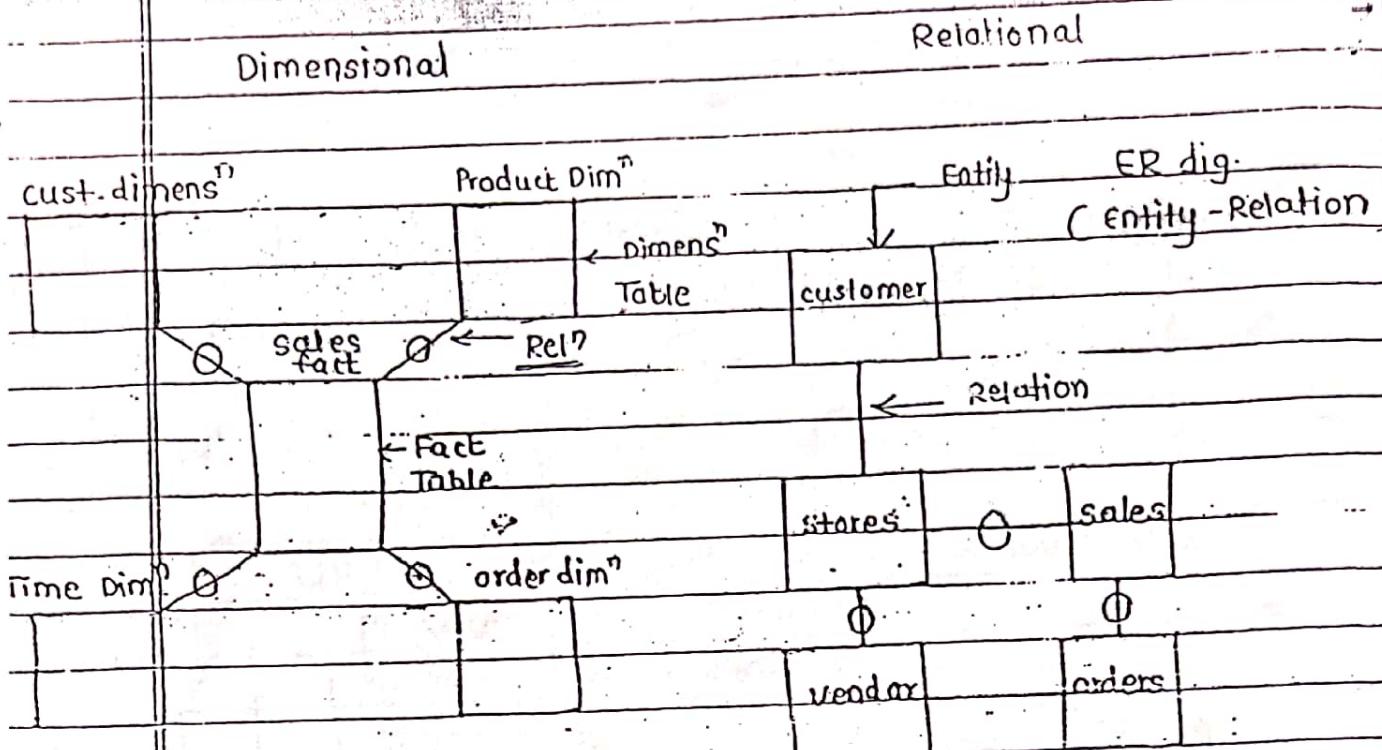
5. Pivot → Rows becomes columns & columns city becomes rows.

NY	Q1	Q2	Q3	Q4
PT	-	-	-	-
LA	-	-	-	-
CL	-	-	-	-

①      ②      ③      ④  
⑤      ⑥      ⑦      ⑧  
⑨      ⑩      ⑪      ⑫  
⑬      ⑭      ⑮      ⑯  
⑰      ⑱      ⑲      ⑳  
⑳ PT LA CL      Quarter

\* Dimensional Model & Relational Model →

- |             |            |
|-------------|------------|
| Dimensional | Relational |
|             |            |
- 1. Design of data for business processing. (Analysis)
  - 2. Captures the fact along with their dimensions.
  - 1. Design of data for data processing.
  - 2. Removes the data redundancy & ensures the data consistency & integrity.



\* Design : models →

1. conceptual
2. logical
3. physical

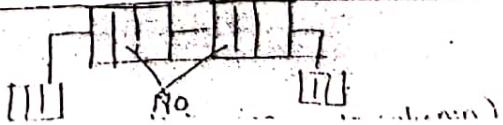
1. conceptual Model —

→ Important entities & their relation



→ No attributes, No datatype (defined) in entity  
defined

→ High level design of dB



8 Retesting with diff. data  
Regression --

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

2. Logical Model  $\rightarrow$  It defines the data as much as possible to show & how they can be physically implemented in dB.
- $\downarrow$
- $\rightarrow$  Imp. entities & rel'n & their possible Attributes
- $\downarrow$
- 3) Primary key of each entity is specified.
- 4) Foreign key may be specified.

customer		sales:	
PK	custid	PK	salesid
	custid		sales-name
	cust-name		

3. Physical Model  $\rightarrow$  It defines how the models are physically exists in syst.

- $\Rightarrow$  Displays all the required constraints
- $\Rightarrow$  data types - also shown.

like int

varchar2(20)

\*

Diff. betn

— conceptual model & logical model & physical Model

→ 1) conceptual Model → High level represent<sup>n</sup> of entities which is used as a part of dB design.

2) logical Model → It is a next level represent<sup>n</sup> of conceptual model in which entities & attributes are specified & relat<sup>n</sup> betn entities.

3) physical Model → It is a granular level of represent<sup>n</sup> of dB design of all entities, column, constraints present in detail.

\*

Facts

→ Prod.dim

cust-dim

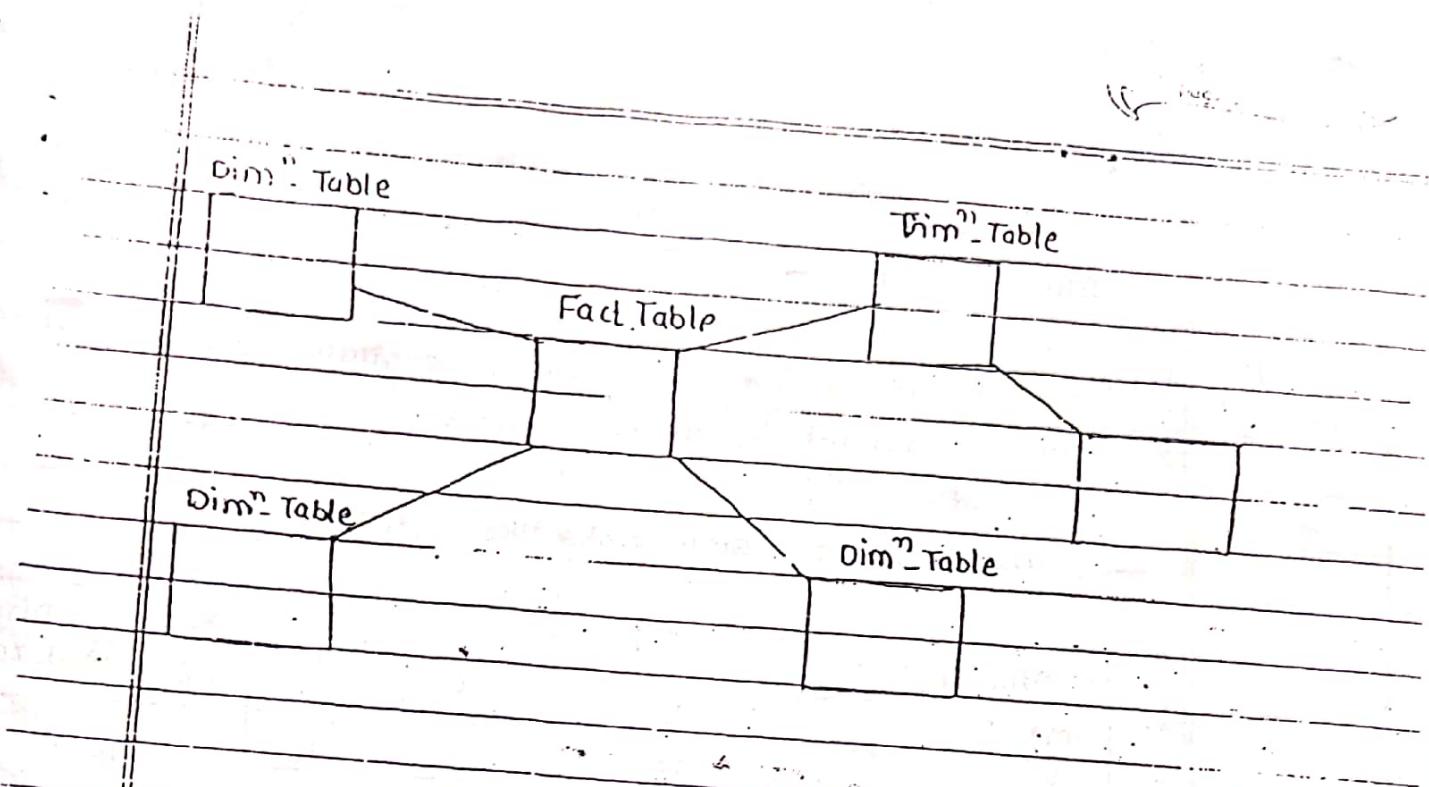
sale-fact

PK salesID

Time.dim

order.dim

PK



Advantages of Star Schema →

1. Easy to understand.
2. Provides Better Performance.
3. can easily handle future changes.

3. Galaxy Schema →

— Also known as Galaxy or fact constellation schema.

— Multiple fact table share the same dim<sup>n</sup> table viewed as a min of star schema is called as a Galaxy schema or Fact constellation.

★ 1.

Star Schema →

design

- It is a db engine which contains centrally located fact table surrounded by dimension table.
- The database design looks like a STAR.

Time-Dim<sup>n</sup>

PK Time-ID

Day

Month,  
Year

Quarter

Sales fact Table

PK

Sales-ID

Co!

FK

Time-ID

FK

Prod-ID

FK

Branch-ID

FK

Locn-ID

Unit sold

Avg sales

Product - Dim

PK Prod-ID

P-Name

P-Type

P-Bran

FK Supplier-

Branch-Dim<sup>n</sup>

PK Branch-ID

Branch-Name

Branch-Addr

Branch-Type

Location-Dim

PK Locn-ID

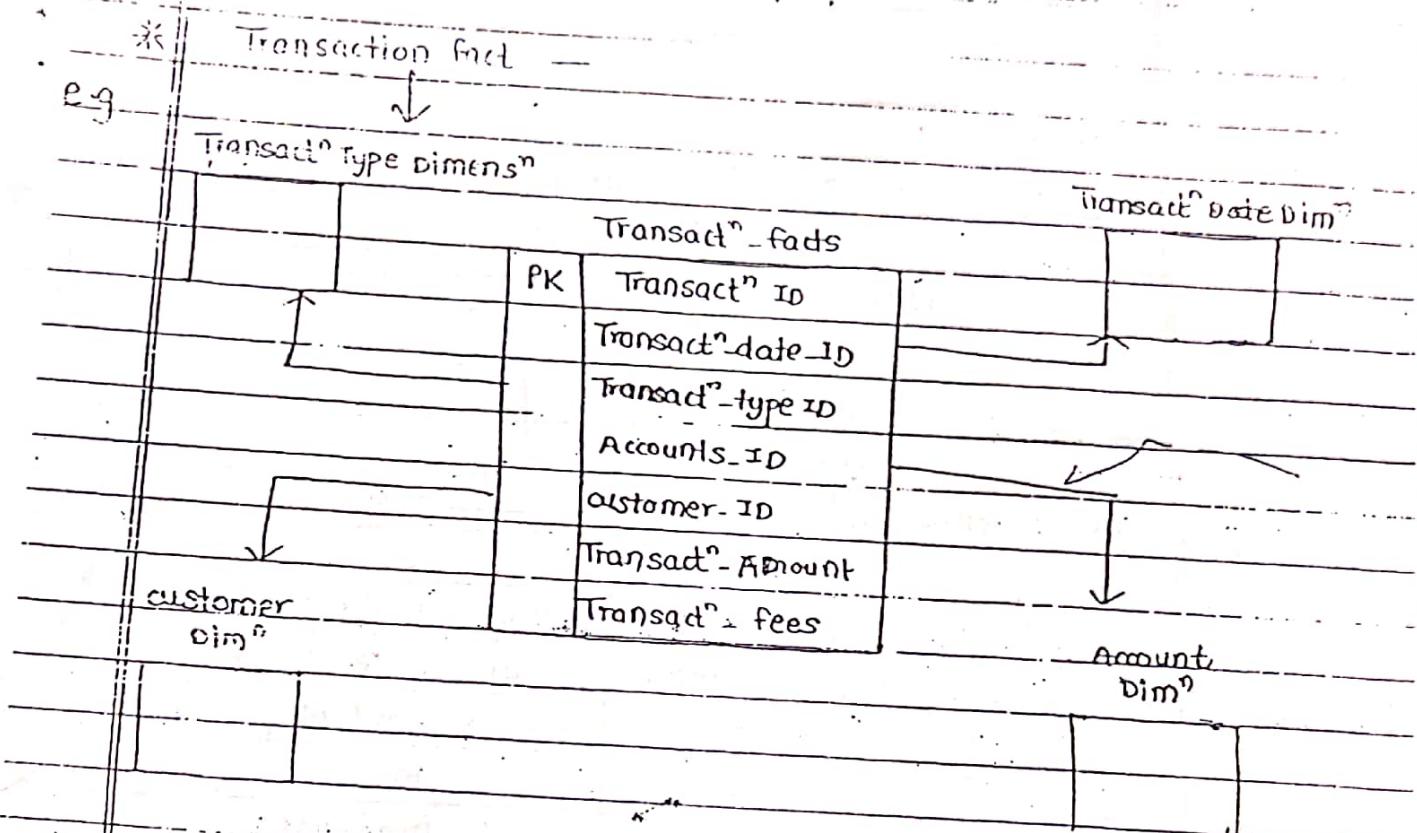
FK Store-ID

city-ID

2. Snowflake Schema →

- It is an extension of Star schema.

- Dimension Table in Snowflake schema are normalized & the process of normalising dim<sup>n</sup> table is called Snowflaking.



\* How Data is stored in the fact Table?

Trans-ID	Trans. DateID	Trans-Type	Acc-ID	Cust-ID	Trans-Amt	Trans-fees
1	25	2	47	51	5000	0.00

\* SCHEMA :-

Skeleton structure that represents logical view of entire database.

Schema defines how the data is organised in database.

\* e.g. fact & dimens' Table?

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

Facts :-

- It is a counted or measured event.

\* Dimension → contains referential info. about fact.

	No. of Students	57	
Dimension	↑	↑	Facts

\* Fact Table → 1) central table in dimens' model surrounded by dimens' tables,

contents,

facts,

Measures

2) contains foreign keys of dimens' tables.

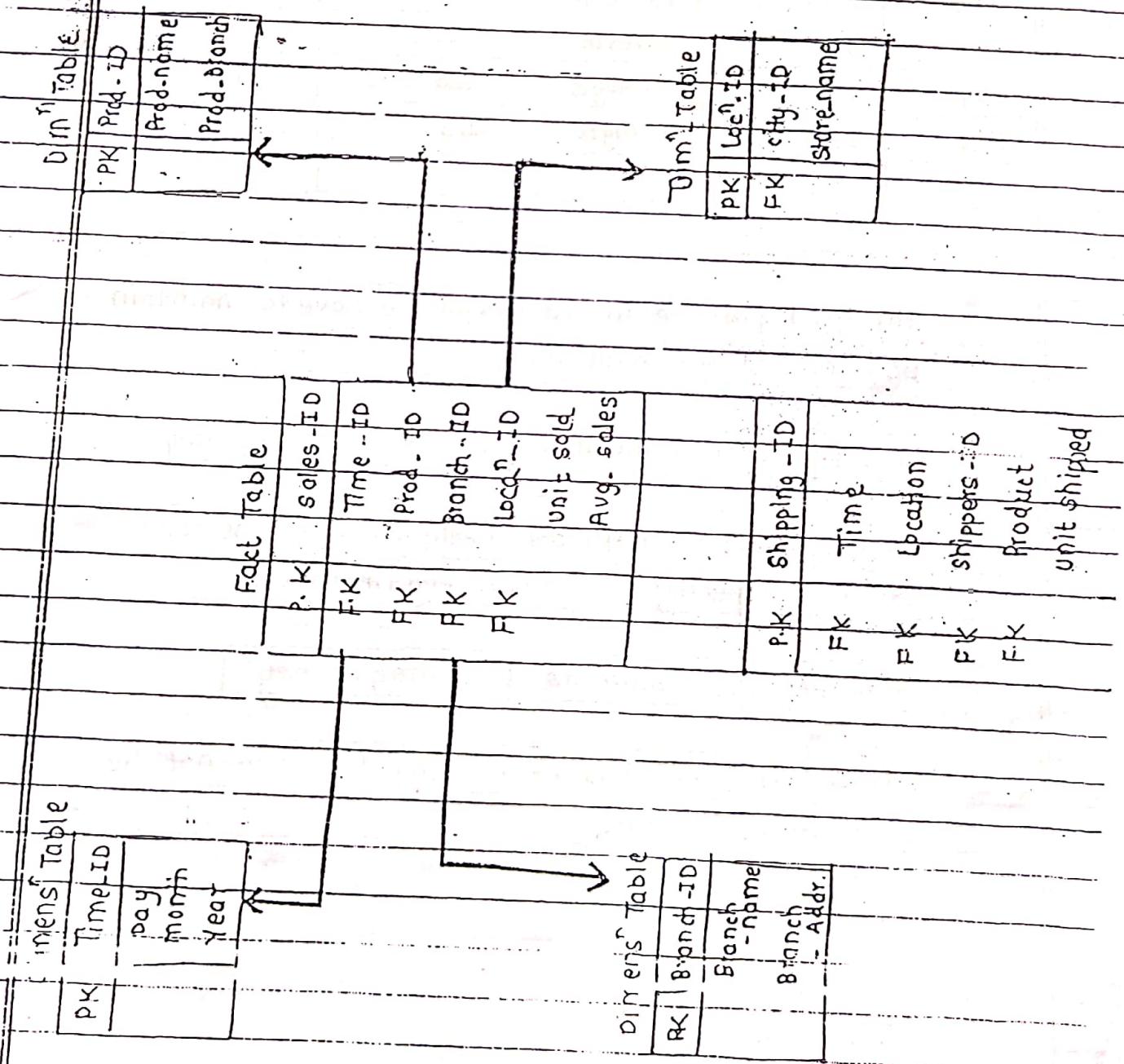
\* Dimension Table → 1) Dimension Table contains

dimens keys

f Attributes,

values

# \* GALAXY SCHEMA \*



Q \*

## Surrogate Key :-

ID(S-K)	Start Date	End Date	Material Rate	Material Name
101	1 Jan 18	12 Feb 18	50	Steel
102	15 Feb 18	28 Feb 18	175	Cement
103	1 Jan 18	28 Feb 18	125	Bugash
104	1 Jan 18	30 Aug 18	100	coalash
105	1 March 18	30 Aug 18	175	Cement

— For good practice in dB design we have to maintain primary key for each table.

1. Use part of a data in a table as a primary key
2. Use new field with artificially or autogenerated value to specify a primary key in a table.

This key is known as 'Surrogate Key'.

— This key itself has no meaning & it may not be visible to end user.

Syntax :- 1) col-name identity (start-val, increment val)

primary key (col-name)

e.g. id identity (100, 1)

SQL  
server

Primary key (id)

2) create sequence sq-name

min value = 100

increment by 1

max value = 1000;

Oracle

id sq-name (100, 1)

Primary key (id);

Key

In fact  
Table

Date 20-8-17  
Page .....

\* Composite key :-

Primary key Primary key

order-ID	Date	Product-ID	Product-Name	Price \$	Quantity
E101	22/11/17	P125	Modem	50	1
E101	22/11/17	P16	Router	100	5
E102	23/11/17	P125	Modem	100	2
E102	23/11/17	P51	Printer	25	2

— composite key refers to case where more than one column is used to specify primary key in a table.

Syntax :- create table table-name

    column1 integer

    column2 varchar2 (20)

    column3 varchar2 (20)

    Primary key (column1, column2)

);

JMP  
Page \*

### Types of Fact :-

1. Additive
2. Semi-Additive
3. Non Additive

( Additive )

1. Additive →

- Additive facts can summed up across all dim<sup>n</sup> in a table.

	Date-ID	Sales Amount	
	Store-ID	↓	
Additive →	Product-ID	Day 1 :-	100 \$
	Sales-Amount	Day 2 :-	+ 250 \$
Non →	Profit-margin	Day 3 :-	+ 300 \$
Additive			650 \$

2 Non-Additive →

- cannot summed up with any dim<sup>n</sup> in a table.

Profit Margin

↓

Day 1 :- 5%

Day 2 :- + 8%

13% X (not possible)

### 3. Semi Additive →

- can summed up with some dim<sup>n</sup> in a fact table but not with others.

Date ID	
Amount ID	
current_Balance	]

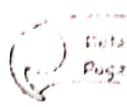
- current balance can summed up with account ID dim<sup>n</sup> but not with Date-ID dim<sup>n</sup>.

## ① \* Types of Fact Tables :-

1. Transactional Fact Table
2. Periodic Snapshot Fact Table.
3. Accumulated Fact Table.
4. Factless Fact Table.

### 1. Transactional Fact Table :-

- Fact Table that represents an event that occurred at instantaneous pt. of time.  
mostly Additive facts



## 2. Periodic Fact Table →

- Fact Table that describes set of things in a particular instance of time.
- Time period is predictable or regular.

Dimens <sup>n</sup> Table				Fact Table		Dimens <sup>n</sup> Table	
Date-ID	Day	Months	Year	BatchID	→	Moool	JAVA
1	14	4	2018	BatchID	→	Moool	.Net
2	23	6	2018	No. of Students		Noool	Manual

Fact Table	Batch-ID	Date-ID	No. of students
	moool	2	55

## 3. Accumulated Fact Table →

- used to show activity of process that has well defined beginning & end.

Step 1 - | order Date | NULL | NULL |

Step 2 - | order Date | shipping Date | NULL |

Step 3 - | order Date | shipping Date | delivery Date |

### \* DIFFERENCE Betw →

	Transactional	Periodic	Accumulated
--	---------------	----------	-------------

e	Unpredictable	Regular or predictable	undetermined time span
---	---------------	------------------------	------------------------

E	Insert	Insert with Insert or Update
---	--------	---------------------------------

in	one row per Transaction event	one row per period	one row per life.
----	-------------------------------	--------------------	-------------------

Date 21-8-18  
Page No. 2

#### \* 4. Factless Fact Table →

- It is a fact table which does not contain — Numeric Fact OR Measures.  
& contains only dimension keys.
- It captures an event that happened only at info. level not at calculation level.
- It captures many to many relationships bet<sup>n</sup> dim<sup>n</sup> table but contains no numeric facts.

#### \* Types of Dimensions →

1. Slowly changing dim<sup>n</sup>
2. confirmed dim<sup>n</sup>.
3. Degenerated dim<sup>n</sup>.
4. Junk Dim<sup>n</sup>.
5. Role Playing dim<sup>n</sup>.
6. Rapidly changing dim<sup>n</sup>.

#### ★ 1. Slowly changing dim<sup>n</sup> →

- Dim<sup>n</sup> attribute that changes slowly over a period of time rather than changing regularly.

\* Types of Slowly changing Dim<sup>n</sup> (SCD) →

1. SCD Type 1 :-

cust-ID	cust_Name	Year	Location
101	xyz	2005	London

PIP →

cust-ID	cust_Name	Year	Loc <sup>n</sup>
101	xyz	2008	Paris

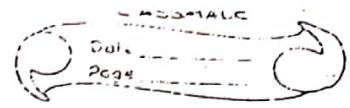
Type I :- replaces old entry with new one.

:- only new data is present & old data is completely lost.

2. SCD Type 2 :-

cust-ID	cust_Name	Year	Location
101	xyz	2005	London
101	XYZ	2018	Paris

long  
key  
will  
Sync get  
key



Type II :- Stores new data as well as old data.

:- New record & old record present in a same table.

3. SCD Type 3 :-

e.g.

cust-ID	aust-Name	old Yr	old Loc <sup>n</sup>	old NewYr	New Loc <sup>n</sup>
101	xyz	2005	London	2008	Paris

Type III :- creating new fields in a table & maintain old as well as new data in a same record.

2. confirmed Dimension →

calender Year  
(week - mon-sat)

Budget Year  
(week - sun-sat)

Time

Time

HR

Finance

Emp

Product

- We can call Time dim<sup>n</sup> as a confirmed dim<sup>n</sup> when it contains same descript<sup>n</sup>, same contents & one is a subset of another.

ADVANTAGE :- 1) Less Maintenance cost.

2) Easy development.

3) Efficient ETL Work.

4) It can be reused whenever needed.

#### 3. Degenerated Dim<sup>n</sup> →

- It is dim<sup>n</sup> and attribute stores as a part of fact table not in a separate dim<sup>n</sup> table.

#### 4. Junk dimension →

- It is a single table with comb<sup>n</sup> of different & unrelated attributes to avoid having large no. of foreign keys in a fact table.

#### 5. Role Playing Dimension →

- It is where the same dim<sup>n</sup> key along with its associated attributes, can be joined to more than one foreign key in the fact table.

\* Oracle → Information Schema

\* Examples → Additive  
Non Additive  
Semi Additive

- CLASSmate

Date 23-8-18

Page

### 6. Rapidly changing Dimension →

- Dim<sup>n</sup> attribute that changes frequently is called as a rapidly changing.
- If you need to track the changes using standard slowly changing dim<sup>n</sup>, this tech. can result in a huge inflation of size of the dim<sup>n</sup>.
- one sol<sup>n</sup> is to move attributes to its own dim<sup>n</sup> with a separate foreign key in a fact table.  
This new dim<sup>n</sup> is called Rapidly changing Dim<sup>n</sup>.

### \* DB Testing

L. Primary Goal

Data validation &

Data Integration

### ETL Testing

Data Extract<sup>n</sup>, transform  
& loading for reporting  
& analysis.

Applicable syst.  
where business flow  
occur.

syst. containing  
historical data not  
business flow occur.

Need  
Ensuring data  
integrity from multiple  
appln.

Ensured the req.  
data is moved from  
source to target.

Modelling  
ning

E-R  
diagram

multi-dimensional.

Normalised data -  
more no. of tables

Denormalised data -  
few no. of tables.

### \*

What is diff. betn DWH & Data Mining.

DWH is a huge concept as compared to data mining.

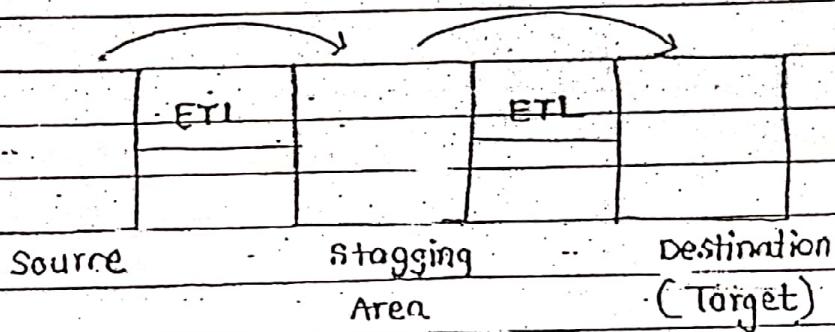
DWH is a db & data mining is a process.

Data mining involves extracting info. from DWH  
& interpret it for future predict<sup>n</sup> (Analysis, planning,  
designing, reporting)

development

## \* E - testing :-

### \* what is ETL ?



ETL defines mechanism of dataflow from source syst. to target syst.

It is a process of extracting data from source syst, transformed w.r.t set of rules, after getting req-data data is loaded to target.

### \* ETL Testing →

defn → ETL Testing is done to ensure that the data that has been loaded from source to target (DWH) after applying set of rules (Business logic) is correct or not.

\* (ETL - product reconciliation)

↳ Table balancing

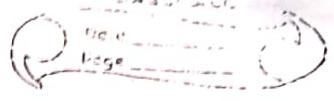
 Data Validation

Table structure valid"

Data completeness

Data Duplication

Finding Invalid data

ETL Testing also involves verifica<sup>n</sup> of data in various diff. stages that has been used bet<sup>n</sup> source to target.

### \* Challenging during ETL Testing →

Issues

1. Data loss during ETL process.
2. Incomplete, Incorrect, Invalid, Duplicate data is present at target.
3. DWH contains historical data, data volume is too large.
4. Tough to generate test cases and finding scenarios because the data vol. is large.

### \* Types of ETL Testing →

- |                                 |   |
|---------------------------------|---|
| 1. constraint                   | 5. Incremental & historical process Testing |
| 2. source to target             | 6. data completeness                        |
| 3. Source to target data valid" | 7. Data Transform <sup>n</sup> Testing.     |
| 4. Data integrated Testing      |   |

## 1. \* Constraint Testing →

— During this testing test-ER identifies whether the data is mapped from source to target or not.

constraint testing.

Not NULL

Default

unique

check

Primary key

Foreign key

} Constraints

1. NOT NULL — ensure that the column cannot have null value.

2. Default — ensure the default value for a column whether the none is specified.

3. Unique — ensure that all values in column are different.

4. Check — make sure that all values in column satisfies certain criteria.

5. Primary — used to identify uniquely row in the table.

6. Foreign — used to ensure referential integrity of data.

~~que~~ select count \* from lamp-table

Date 27-8-18  
Page

## 2. \* Source to target count Testing —

- During this test ER ensure that count of source & target data is expected or not.

## 4. \* Data Integration syst. is same as SIT.

## 5. \* Incremental & Historical Process Testing. →

- In this type of testing we are going to verify new data (changed data) or historical data in target syst.

- There are 2 types of loading in ETL process.

1. Incremental Loading → only loads the data that changed in source syst.
2. Full Loading:



When

It truncates all existing tables and reloads all the data in source syst.

## 6. \* Data completeness Testing —

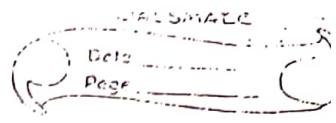
- Data completeness testing ensures all the desired data is completely loaded to the target.

2. validate parent to child relationship of new field.

5) Data quality -

1) data format check — 2) ( precision value check )

e.g. 19.2356



### Duplicate check Testing —

- In this testing we are going to identify duplicate data in target syst.
- Duplicate data may arise
  - because
  - 1. Primary key not defined.
  - 2. Wrong development.
  - 3. Environmental issues
- when there is huge amount of duplicate data in target syst. that may results incorrect analysis & reporting.

### \* Types of ETL Bugs. :-

1. Table str. Issue
2. Issue with data & source syst.
3. Data count not matching betw? source & target.
4. Duplicate data loaded issue.
5. X? rules issue.
6. Data format issue
7. Index not create after jobrun.

### 8. Performance Issue.

#### \* S/w Testing

1. S/w testing  
carried out prior to  
deployment of s/w.

2. Source code  
specific.

3. Focussed on  
Used cases which  
contains various  
Test cases.

#### DWH Testing

1. DWH Testing  
carried out post  
deployment.

2. Content  
specific.

3. Focussed on  
querying the test  
data loaded by  
ETL Testprocess.

#### \* Mapping Document —

Date 8-9-18  
Page

### \* Mapping Document :-

source	Target	X <sup>n</sup> logic
S- Name Tbl-name col-name data type	T-name Tbl-name colname data type	
SystDB material Material varchar & M- type (50) name	material material varchar & type (50) material which is applicable for weight	
dBname user user Date DM user user user -details -DOB time Name details age varchar (545 -DOB 365		

Que. What is Mapping Document ?

Mapping document describes relationship betn extreme starting pt. & extreme end pt. of ETL process which is used in DWH.

Que. Why tester needs mapping doc. ?

→ 1. Because it contains details about each & every table which is a part of ETL process from source to target.

2. While working on target tables we have to refer tables in source syst. & mapping doc. contains

Date  
Date

complete info about tables in source syst,  
target syst along with its X^n logic.

\*\*\*

Sr.No. — 1

Test case ID — TC-001

Priority — High

Reference — SRS doc

Title/ summary — structure valid" Table-name

Pre-cond — source exists

Test data —

Action/ Description — To validate table constraint against  
com-mapping doc.also validate col-name

Expected Result —

Actual Result

Defect ID

Select month between (01/01/2018, 01/05/2018),