

- * Database lang \rightarrow SQL (technology)
structural query lang.

operation perform - (CRUD) \rightarrow Delete
 Create \downarrow Read \downarrow Update

SQL is technology

f Database testing, ETL- BI + testing \rightarrow Implementation.

* Database testing

- \rightarrow In Database testing, we are going to check impact of Backend due to frontend operations of users, and operations are - CRUD.
purpose \rightarrow validation of data

+ Domain \rightarrow

1. Transactional app

→ Run fundamental process

e.g Amazon - steps

a) Amazon open

b) login

c) search product

d) Add to cart

e) Address

f) payment

g) Product delivery

Gpay

- money transfer
- Recharge, bill

} middle man.

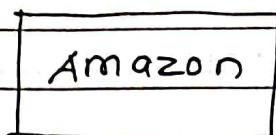
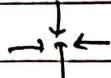
What's app - data transfer

2) Analytical applications

→ Analysis

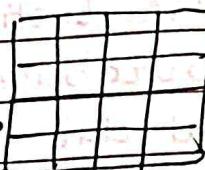
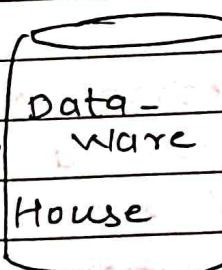
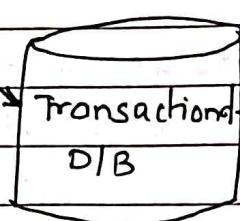
understand business, improve -

23/02



Account create
address

create
update
delete.



Data-
ware
House

Report

• Drawback of transactional system.

1. Degrade performance
2. time consuming because of multiple location
3. Not fit for analysis purpose

24/02

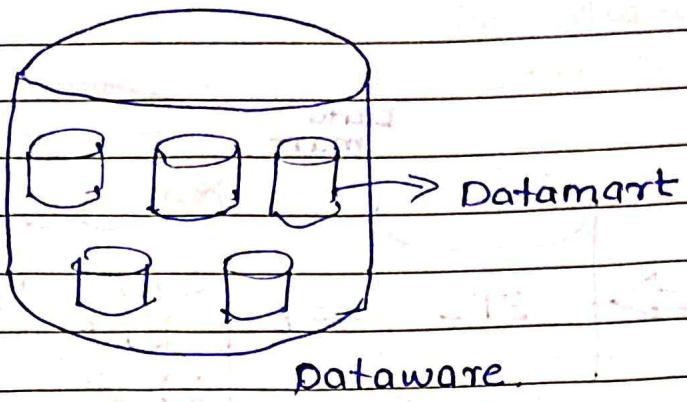
*

Datawarehouse

- It is a database that is designed for querying and analysis rather than transactional processing
- It separates analysis workload from T.P
- This help in
 - i. Maintaining historical records
 - ii. Analyzing the data to gain a better understand of the business and to improve the business.

IMP → Datawarehouse is a Subject-oriented, integrated, time-varying, non-volatile collection of data in support of the management's decision making process

- subject oriented
 - subject wise data store



- Integrated

- save in one format, from diff sources

Transactional

pdf, word, Excel, ppt, etc
ppt, EDI, text

diff format

Batch



one format

e.g. Grader

- Time-varying

- Historical data is maintained over time

e.g. Employee info., year wise

101 Pooja 2010 Trainee

101 Pooja 2010 Trainee A.P.

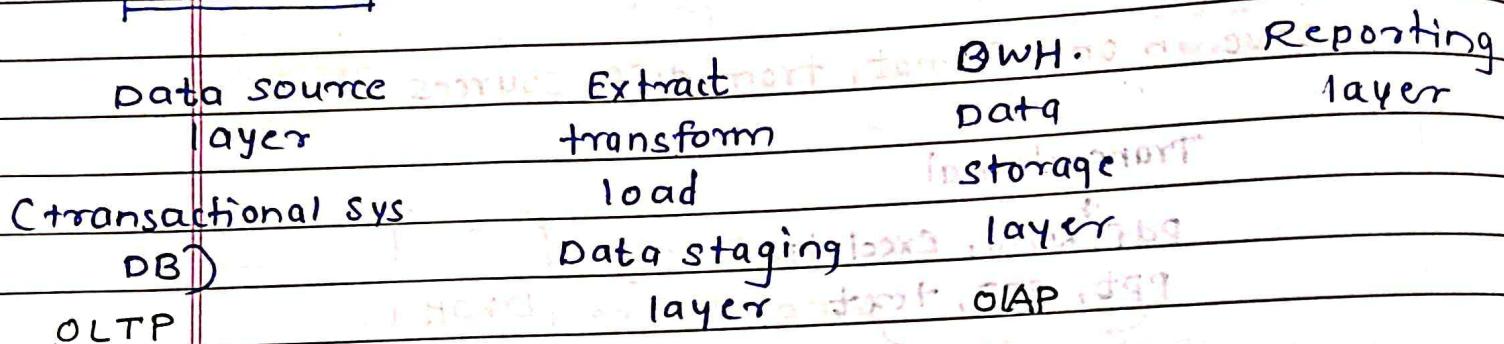
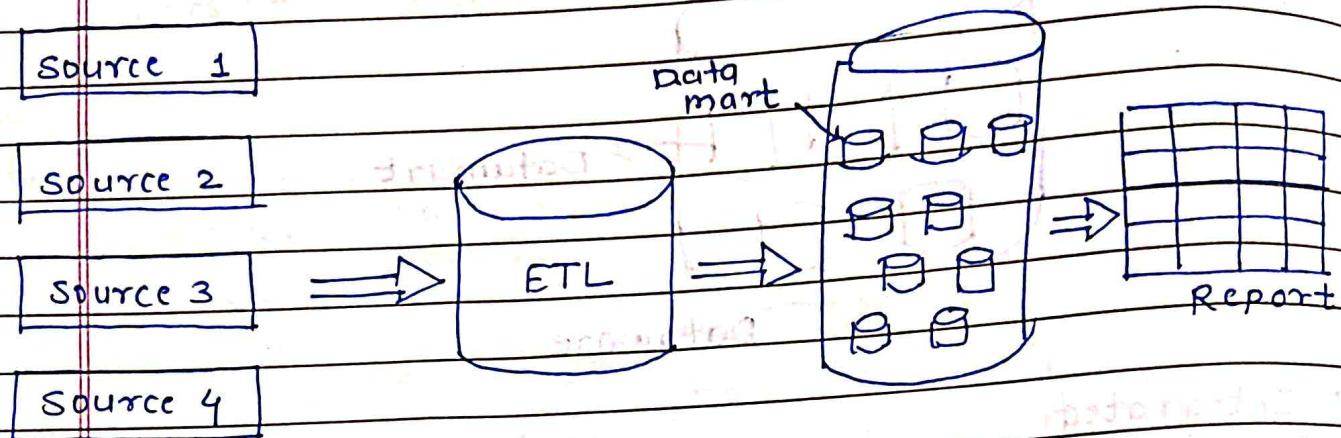
101 Pooja 2011 Tester

101 Pooja 2012 Manager

- Non-volatile → once dataware is placed in dwh, it can't be changed. (No-Read, No-Update, No-delete)

25/02

* Architecture of DWH



1. Data source layer

- It is nothing but transactional system Database
- stores various data, business data like sale, customer, product, purchase, delivery (in different format)

2. Data staging layer (Landing zone)

- ETL

Extract → Reading data from source database

Transform → converting extracted data in Req^ form

Load → Writing the data into target database

- It is process which defines how data is loaded from Source system to target system. (Respective data mart)
- Transactional data → Analytical Data
- minimise data loss, temp data processing during ETL

3. Data storage layer

- Data warehouse

- the place where the successfully cleaned, integrated, transformed & ordered data is stored in a multi-dimensional environment

- Now, data is available for querying & analysing.

Data mart - subset of DWH

- subjectwise data stored, single functional area

Meta data - it is data about data.

Meta data repository is used to store metadata of data which is actually present in DWH.

4. Reporting layer

- Data, in data storage layer is used to create various type of management reports from where user can take business decisions for planning, designing, forecasting etc

* OLTP - online Transactional processing system - source layer

OLAP - online Analytical processing system - storage layer

• use
OLTP
OLAP
used for transaction processing
Query processing

• Data holds current data

current / History

stores all data

only relevant data

small database

large database

volatile data

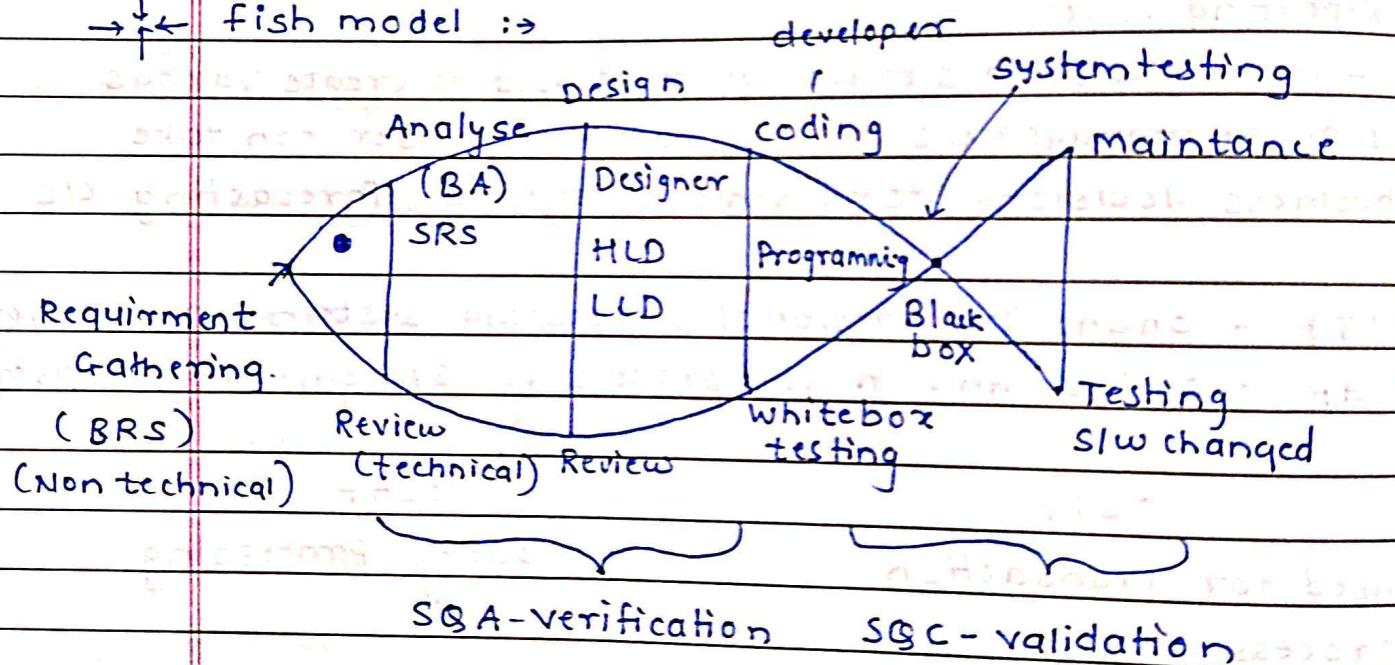
Non volatile

(CRUD)

(only Read)

• source	original source of Data.	Data comes from Various OLTP sys.
• purpose	To control & run fundamental business task	To help with planning, problem solving, decision support
• queries	std & simple SQL queries	complex SQL queries.
• DB design	Highly normalized with many tables (3NF)	De-normalized with fewer tables
• users	Many users	few users

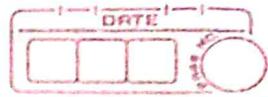
→ ↕ fish model :>



27/02 Normalization

carid	Model	color
101	Bmw 1	Red
102	Bmw 2	Blue
103	Bmw 3	Green
104	Bmw 4	Blue
105	5	Red
106	6	Red

→ Model 1

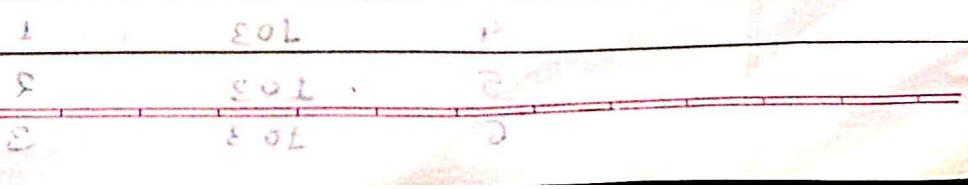


model 2 consists of a Primary key.

Card-id	model	color-id	other	color-id	color-name
101	1	1	1	1	Red
102	2	2	2	2	Blue
103	3	3	3	3	Green
104	4	4	4	4	Yellow
105	5	1	1	1	Red
106	6	2	2	2	Blue
			foreign key	101	Red
				102	Blue
				103	Green
				104	Yellow
				105	Red
				106	Blue
size	large		model	2	2
data			size	small	small
facing	simple		data	addr	addr
DB			facing	addr	addr
design	simple		DB		
performance	low		design	complex	complex
			performance		
(Denormalization)				Better	Bad
				red	blue
				data	data
				freq.	freq.
				add	update
				/ delete	

- Normalization is the process of efficiently organizing the data in the DB.
- used to minimize redundancy.
- used to eliminate the undesirable characteristics like insertion, updation, deletion anomalies.
- Divides larger table into smaller table & links them using relationship.

Normal form: 1NF, 2NF, 3NF, 4NF, BCNF, 5NF, 6NF, 7NF.



full \rightarrow first name last name
 partial \rightarrow name qualification



e.g.	ID	Name	qualification
	101	Sumedha	BE
	102	Vipin	BE, ME
	103	Viha	BE, M.E., PhD.

28/02

LNF	ID	Name	Qu.
	101	Sumedha	BE
	102	Vipin	BE
	102	Vipin	ME
	103	Viha	BE
	103	Viha	ME
	103	Viha	PhD

- Attributes of table can't hold multiple values

RNF

- should be in LNF

- Should not have partial dependency

Dependency - full, Partial, transitive

ID	Name	Qualification
101	Sumedha	BE
102	Vipin	ME
103	Viha	PhD

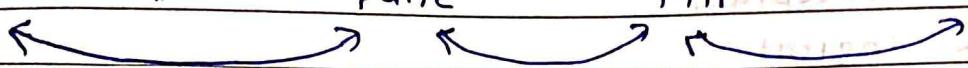
Name partially dependent on qualification

id-quid	ID	qu-id
1	101	1
2	102	1
3	102	2
4	103	1
5	103	2
6	103	3



3NF → should be in 2NF at first, then, the subset, place, etc.
Doesn't have transitive dependency.

ID	Name	city	state	country
101	S	Newyork	US	US
102	P	london	UK	UK
103	X	pune	MH	IN.



ID	Name	city-id
101	S	1
102	P	2
103	X	3

city-id	city	state	country
1	Newyork	US	US
2	london	UK	UK
3	pune	MH	IN
4	mum	MH	IN
5	Mg	MH	IN



2NF

state id	state	country
1	US	US
2	UK	UK
3	MH	IN

country → country-id is good, country-id and country word

1	US
2	UK
3	IN

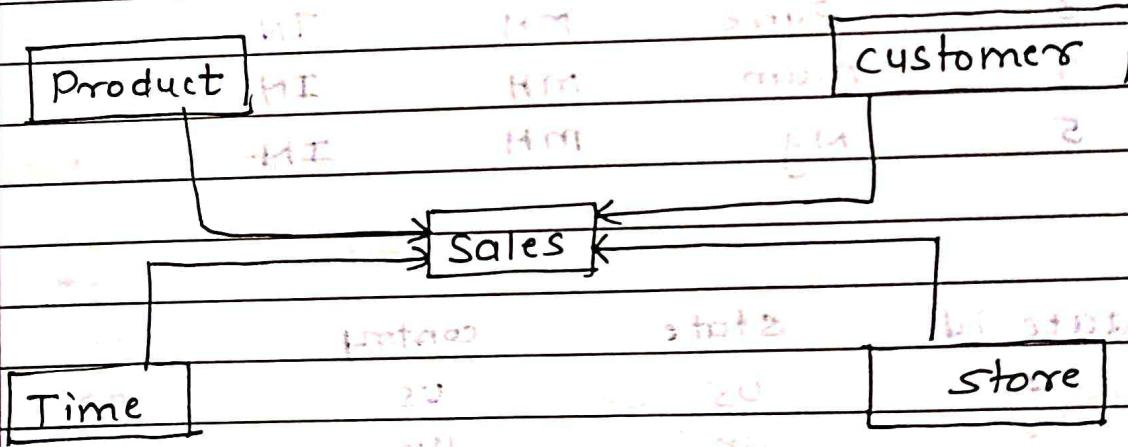
03/03.

Data Models.

- Data models tells how the logical structure of a DB is designed.
- it is defined as how data is connected to each other & how it will be proceed and stored inside the system
- Types of data model
 - 1. conceptual
 - 2. logical
 - 3. physical

1. conceptual design / model

- it is high level design of DB
- features :-
 - Display imp entities & relationship among them
 - No attributes is specifies
 - No primary key is specifies
 - Level 1



2. logical data model.

- It is define the data as much as possible to show they can be physically implemented in DB
- Detail info
- table name, col name, datatype, Pk.fk

* logical data model

→ design by Architect.

colname datatype null constraints sales manmt
↓ ↓ ↓ ↓ ↓

1. include all entities & relationship among them

2. All attributes for each entity are specified.

3. primary key, foreign key specified. also - constraints.

→ Physical model

- Actual implementation of logical model into database
is called physical model.

Diff betⁿ Relational Datamodel & Dimensional Datamodel

- Suitable for OLTP application.
- many number of tables with chain relation among them
- CRUD activity
- Normalised.
- Relational DB is suggested for DWH applications.
- fixed star fact table, dimension table
- only Read.
- De-Normalised.

04/03



Fact Dimension

Batch number of student

T ₉	31
T ₁₀	29
T ₁₁	30
T ₁₂	40
T ₁₃	41

Total

171

order NO.	Date	Product	city	qty	Total price
981	16/02	Laptop	Pune	9	11400 \$
982	17/02	PD	Mum	2	1300 \$

- 11 - 12700 \$

fact table

Primary key present in dimension table for their foreign key present in fact table

fact table has primary key to measure value

It is counted or measured event wise

Dimension table present in database of DB

It contains preferential info about fact table

fact table

consist of measurements of facts of a business process

it is central table in dimension model surrounded by dimension table

- A fact table typically has two types of col.
 - i. Those that contain facts
 - ii. Those that are a foreign key to dimension table

→ Dimension table
used to described dimensions.

06/08

- Types of Dimension model.
 - i. Star (schema) → logical arrangement)
 - ii Snowflake
 - iii. Galaxy
1. Star
- it is simplest form
 - In this, central table is called fact table & radically connect other table are called dimension table
 - Entity relationship diagram is look like star
 - De-normalized form
 - It is good for datamart with simple relationship.

2. snowflake

- The process of normalizing dimension table is called snowflaking
- Dimension table is in normalized form
- It is extension of Star schema
- ER diagram is look like snowflake

3. Galaxy | fact constellation

- contain two or more fact table that share same dimension.
- looks like collection of star

• Diff betn Star & Snowflake

star

snowflake

1. contain less no. of rows
2. Good for datamart
3. Simple DB design
4. less easy to maintain
5. contain only single dimension
6. fact & Dimension both are in de-normalized form
1. contain more no. of rows
2. Good for DWH
3. complex
4. Easier to maintain
5. more than one
6. Dimension - Normalized fact

* Types of fact :-

- i. Additive - Additive facts are facts can be summed up through all of the dimensions in the fact table
- ii. Non-add - that cannot be summed up for any of the dimensions present in the fact table
- iii. semi-add - that can be summed up for some of the dimensions in the fact table, but not with all.

07/03

* Types of Dimensions :-

(i) slowly changing dimension

- Dimensions that changes slowly over a period of time, rather than changing on regular schedule
- is a dimension that stores & manages both current & historical data over time of in a DWH
- It is considered of implemented as one of the most critical ETL task in tracking the history of dimension record

→ types SCD

SCD - 0 - Passive Method

SCD - 1 - overwriting the old value

SCD - 2 - Creating New additional record

SCD - 3 - Adding New col

SCD - 4 - using Historical table

SCD - 5 - combine 1,2,3

- * SCD 0
 - NO Special Action is performed upon dimensional changes.
 - Dimension data that is remain same as it was first time inserted

Before

101 Rohit paris 2011

After

101 Rohit New York 2008

*

SCD-1

- old value is simply replaced

- only new value maintained

- No history

- easy to maintain, e.g - correction of spellings, etc

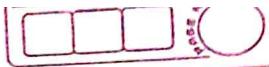
Special characters.

Before

101 Rohit paris 2011

After

101 Rohit paris 2011



* SCD-2

- old & New value is present in same table
- New row is inserted for new value.
- All history maintained

Before	101	Rohit	Newyork	2008
After	101	Rohit	Newyork	2008

Before	101	Rohit	Newyork	2008
After	101	Rohit	Paris	2011

* SCD-3

- old and new value is kept in same table, same row
- New value located into 'new' col and old one into 'previous' col.
- History is limited
- least commonly needed technique

Before	101	Rohit	Newyork	2008
After	101	Rohit	Paris	2011

* SCD-4

- Separate table are there for old value & new value
- Separate history table is used to track all History
- Main table keeps only new data.

Before	101	Rohit	Newyork	2008
After	current - 101	Rohit	Paris	2011
History	- 101	Rohit	Newyork	2008

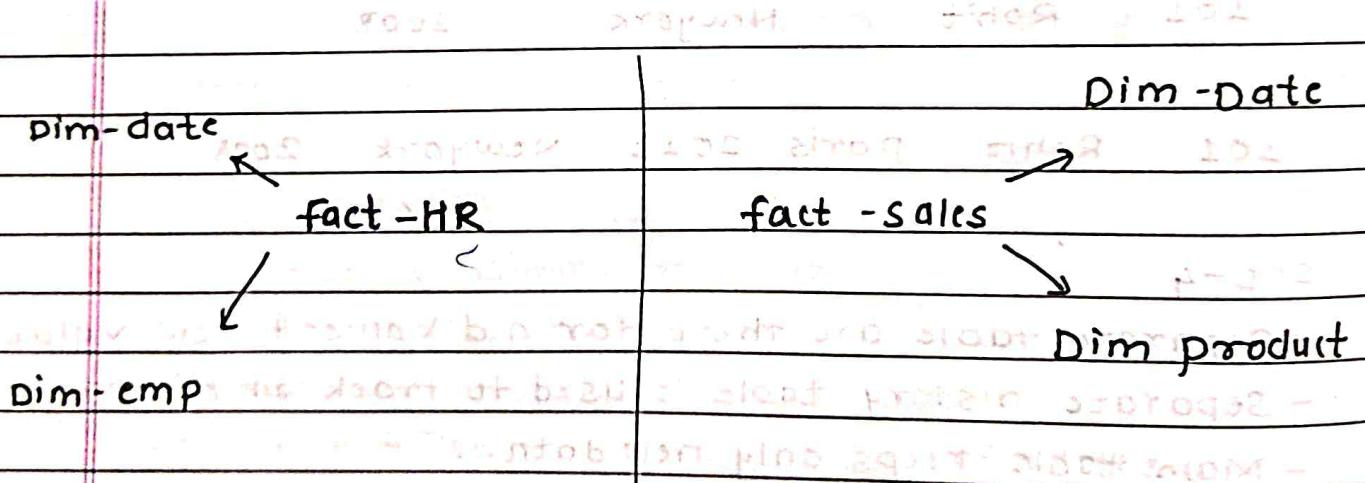
* SCD - 6

- combine approaches of 1, 2, 3, 4, 5 & 6
- ↳ current_address, current_year - for keeping current value
- :- previous_address, previous-year - History value
- :- current_flag - for keeping information about the most recent record.

ID	Name	add	year	prev.add	year	flag.
101	Rohit	Mum	2020	Newyork	2008	N
101	Rohit	Mum	2020	Paris	2011	N
101	Rohit	Mum	2020	Pune	2015	Y

(ii). conformed dimension.

- conformed dimension are those dimensions which have been designed such way that the dimensions can be used across many fact-tables (Data mart) in different subject areas of DWH.



- calendar year

Jan.

Monday - Friday

English

- year - budget year

April

Sun - Sat

french

- instead of creating two dimensional table, create only one

(iii) Degenerated Dimension

- which are directly present in fact table not in separate dimension tables

(iv) Junk Dimension.

- when group of independent dimensions stored in one dimension table that dimension table is called as Junk dimension.

obj/03

Surrogate Key

- primary key made up of real data in table, that called as Natural key.

- sometimes in database we cannot make primary key from real data

- we add one artificial col in table which is unique & not null

- This primary key which is generated from artificial col is called Surrogate key.

Mapping Documents

- It define relationship b/w source data field to their related target data field which are involved in ETL process.

- Map b/w source data & target data in ETL process.

- ETL tester needs mapping document because, while testing data on target table we have to refer data in source tables & mapping document having detail info about each and every table that is part of ETL process from source to target along with transformation logic.

ETL testing.

ETL - Extract - Transform - load

Reading data
from dB

↓
converting the extracted
data to required form

Writing the data
target

- It is a process which defines how data is loaded from the source to the target system (DWH).
- All the data migration and data movement is done by ETL process.
- We are going to validate data that has been loaded from source system to the data warehouse, is uniform in term of quality, quantity and format.

purpose → To identify and mitigate general data errors during ETL process.

→ we need to ensure data accuracy at target system because if the data is not accurate, then the business decision will be wrong.

ETL Testing vs Database Testing.

ETL Testing

Database.

1. We are validating data which is loaded from Transactional System to analytical system.

1. We are validating data in dB which is created due to front end operations of application user

- i.e. - We are going to check impact of backend due to front end operations of user.

2. ETL testing is normally performed on data in a DWH system.
3. DB testing is commonly performed on transactional system.
4. We have to validate data in both ETL testing and dB testing but approach and steps taken for completion are different.

* Types of ETL Testing

1. Metadata Testing - F - Verification phase.

- We have to validate physical model against its logical model
- involves verification of:
 - Table name
 - column name
 - col datatype
 - col datalength
 - constraints.

2. Data completeness testing - F - (Validation phase)

source - count - level 1 - After build.

target - count

also count(null) - check.

- We are going to ensure that all the expected data is loaded in target (DWH) from the source system.
- it involves - checking of validating record counts, null count between source & target.
- also - aggregate function, filter, against incremental load of Historical data.

- Initial load - When first time inserting data into target
- Incremental - insert only newly added or modified data from source.
- full load - Truncate all data at target and then insert all source data at target.
(History not maintained)

(iii) Data transformation testing. -f - val.

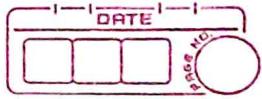
Source EM101 Patil Sumedha S

target Em101 patilsumedhaS

- format different - we can't compare. transform source format like target or vice versa.
- We are going to assure source data which is converted as per business rule/reqt is loaded correctly in target
- Data transformation can be
 - concatenation
 - Joins, splits
 - conditional
 - Aggregation
 - Data conversion.

(iv) Data quality testing -f - val.

- we are going to ensure accuracy of data in target system.
- involves checking validation of
 - Duplicate data
 - Rejected Record
 - Data validation rule.
 - Data integrity.



- Duplicate data - Duplicate value
- Rejected - present at source, but because of transformation not present at target.
- Validation rule - red + green lines highlighted TEST ↓
 - constraints e.g. from date < TO date
- Data Integrity - check pk, fk

↓ ↲ Test scenarios - QA team starts STS + transformation

- statement describing the functionality of application to be tested
- Test condition / Test possibility
- what we have to test
- High level test case
- one test scenarios cover one or more test cases.

e.g. login page, - Homepage, id, password page.

↓ ↲ Test cases.

- set of positive and negative executable steps of all test scenario which has a set of objective, actions, expected result and actual result
- How to be tested
- To validate test scenario by executing set of steps
- Good test cases
 - should be short
 - easily understandable
 - strong objective
 - must have expected result

- Steps.

1. understanding the req^x
2. Test scenario
3. Test cases
4. Test execution will start after getting build (code) from dev. team
 - i. Run ETL job/ ETL code
 - ii. validate data at target system

Development → ETL code → job → Tester.

ETL code - SSIS ETL packages run/ some of the ETL codes are in scripts and we have to execute that for loading data from source system to target system.

Example → Test scenario of case

e.g. login page.

. Test scenario

1. Verify user login

. Test cases.

id	type	objective	Action	Expected result	Actual result	status
1	TC	Verify user login	Enter id and password	User should be able to login and go to home page	Not showing to homepage	fail.

→ ↵ Tester roles & responsibility

1. understand logical flow of the appn[↑]

(Data → source → target)

2. understand and review database design

(logical model vs physical model)



3. understand and review source to target data mapping document. (what transformation logic used)
4. Analysis on business rules / data transformation/validation rules provided by client.
5. create test plan (created by team leader)
6. Identify test scenarios from mapping document for ETL testing.
7. After that we have to developed test cases from scenario
8. Create test data and write sql queries for all test cases
9. Review test cases.
10. After getting build from development team we are going to start executing test cases.
11. Run ETL job to load source data in target system.
12. Document test result and log the code defect for failed test cases.
13. Retesting for failed test cases.
14. Regression testing for code changes.
(Defect fixed, Enhancement, New feature)
15. Test closure
16. Sign off from testing team to project manager and Development team.

* challenges

1. Data loss during ETL process.
 $S \rightarrow T$ (except to dev. team)
(invalid data, time matter, wrong implementation of code)
2. Large volume of data.
(time max req for large data)
Soln → divide data.
3. Invalid and duplicate data at target system.
extra time, efforts req.
waste of time → stop testing.

4. Many NO. of source system.
- Data at multiple places - difficult.
5. source to target mapping information may not be provided to ETL tester.
- old document.
6. Does not have permission to execute ETL code/Job.
7. unstable testing environment.
(Remote testing)