

Data Warehouse

Why do we need Data warehouses?

Operational purposes

→ receive orders

→ react to complaints

→ fill up stock (demand - supply)

basically to keep the company running

→ refer to OLTP (online transactional processing)

requirements

- ① one record at a time
- ② no longer history (only current data is maintained)
- ③ data input

(Older course)

Data warehouse concepts

② Analytical purpose.

→ # of sales this month vs last month

→ what is best category?

→ what can be improved?

to evaluate performance & make better decision

→ refer to OLAP (online analytical processing).

① Analyze 1000 to millions of records all at the same time.

② fast query performance

③ to analyze data over time, history data is maintained.

Data warehouse (DW) is there to address analytical data needs.

It is used for reporting and data need.

- ✓ user friendly (name readability)
- ✓ fast query performance (to pull & process lot of data)
- ✓ enabling data analysis.

Operational data systems

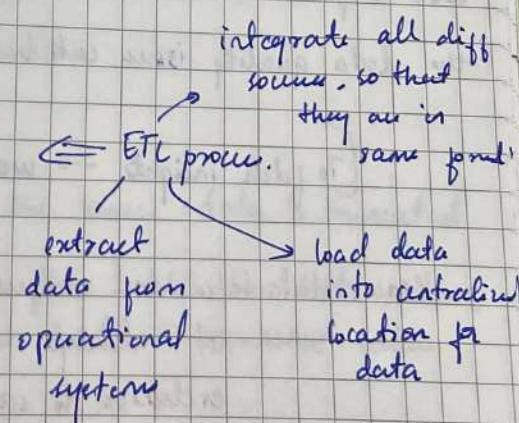
Sales data

HR data

CRM data

Other data

Data warehouse



Goals of Datawarehouse:

- ① Data is available in centralized system in a consistent manner
- ② Data must be fast accessible (query performance)
- ③ User-friendly (easy to understand in term of modeling the data)
- ④ must load data consistently & repeatedly (ETL)
- ⑤ reporting & build data visualizers on top of it.

3 Layer

Business Intelligence: Transform raw data to fetch meaningful insights & make better decisions.

Data lake.

- ① own centralized location for data storage as well.
- ② store raw data
- ③ Big data technology
(CSV, JSON, images, video files - store unstructured data)
- ④ usage is not defined yet
- ⑤ user are data scientist
- ⑥ needs cleaning

disadv.: data quality is not ensured, do not

know the purpose of data but by using cloud technologies which are used to process this data & are scalable the data quality issue can be reduced.

[To fetch insight → use datawarehouse → on some parts of datalake]

Hence, data lake and datawarehouse are not always exclusive to each other.

[we use ETL process to get data out of data lake & load into datawarehouse which is user friendly to fetch insights]

Data warehouse

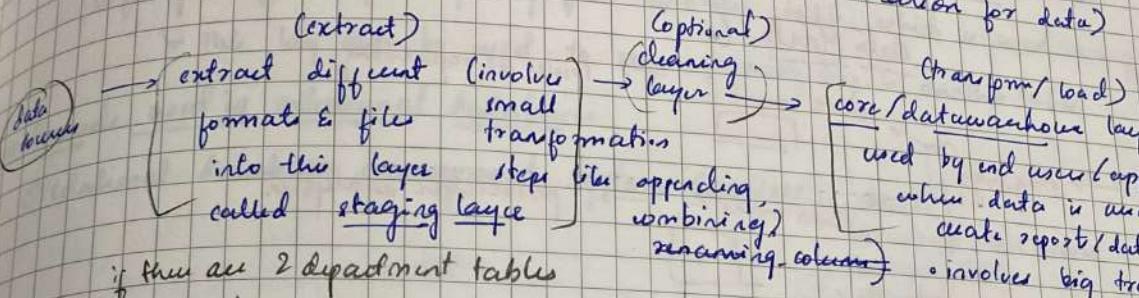
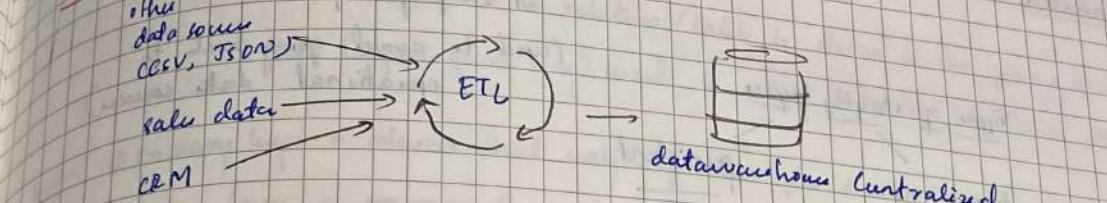
- ① centralized location for data storage.
- ② store processed data (user friendly based on we can clean data)
- ③ database (tables, schemas)
(clean tables, structured data)
- ④ specific & ready to be used
- ⑤ Business user & IT
- ⑥ highly user friendly

Data source

eg: me dep in

Module 3 Data warehouse Data warehouse Architecture

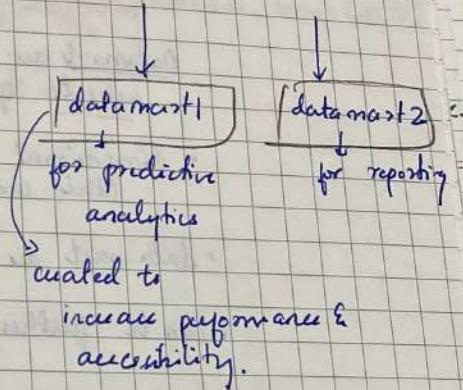
3 layers



if there are 2 department tables

	dept1	dept2
	id name desc	id name desc
merge 2 department tables		
into 1 table called dept.	id name desc	id name desc

e.g. merge 2 department tables into 1 table called "dept".



Staging Area

Purpose: extract data from data sources and place in staging area.

Why do we need this layer?

① Spend less time on source systems to avoid huge transactions

② Quickly extract (quick read access \rightarrow pull data \rightarrow place it in tables)

③ move data into relational database in the form of tables

④ Start transformations from here

- Staging layer stores temporary data, i.e. after every ETL the layer's data is truncated.
- To check what is new in data sources, build delta logic in place which has delta column, for ex: if delta column's max value was 5, then load all rows w/ values greater than 5.

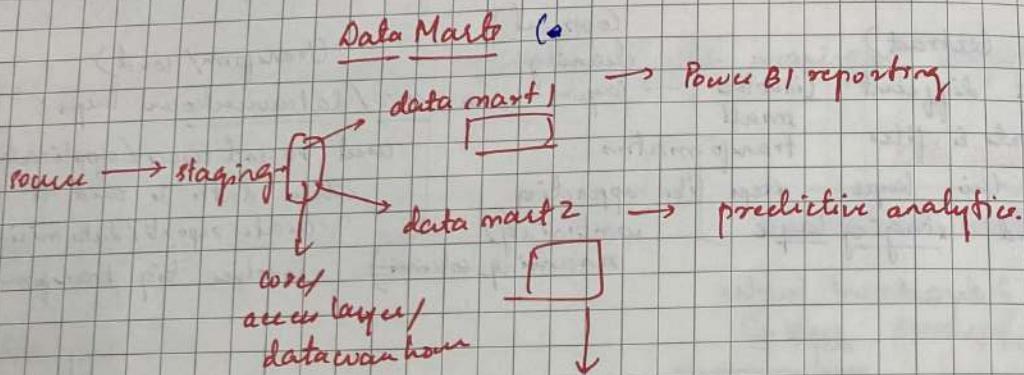
usually date is used as delta column.

• Append data using transformation logic in DW.

Alternative: Persistent staging layer: used as a rollback option when we are never truncating data in staging layer.

(used to avoid going back to operational data source)

Type of staging layers
temporary persistent



Datamart are:
• subsets of datawarehouse

• dimensional model (facts & dimensions)
(this can be done in DW layer as well)

• datamart are built for specific use case

• can be further aggregated based on requirement

Benefits of datamart: ① usability + acceptance (data acceptance for use cases)

ability to work w/ only specific data

② used to improve performance : data is modeled in a dimensional way and we can use specific tools like in-memory DB / datacubes.

When to use datamarts / use cases: ① data marts can be created for diff tools

like DM for reporting w/ diff data
DM for predictive analytics

② DM's for diff department

③ diff DM's for diff regions.

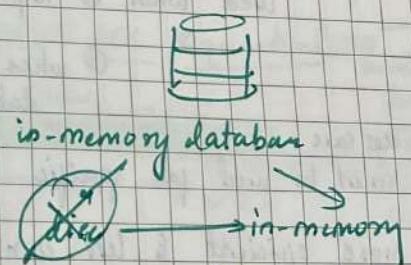
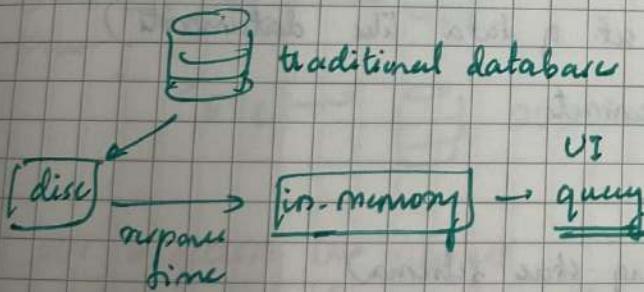
④ diff DM's for diff use cases.

- Relational Databases: A database where data is stored in the form of tables.
 (Tables are also called relations).
- data in a relation / table is stored in the form of rows/columns.
 - SQL (Sequential query lang) used to query data in relational DBs.
- Keys:
- ① Primary key: a column that uniquely identifies every row in a table.
 - ② Foreign key: a primary key from another table, which is used in a given table
 ↳ This key can be used to refer primary key in another table
- * Products used in relational databases:

- ① Relational database management system (RDBMS): Oracle, Microsoft SQL Server, PostgreSQL, MySQL, Amazon relational database service (RDS), Azure SQL DBs.

In-Memory databases: (to use only relevant data)

- ① highly optimized for query performance
- ② good for analytics/high query volume
- ③ usually used for data mart
- ④ relational & non-relational



the response time coming from disk is eliminated

- ⑤ columnar storage (data is scanned through columns)
- ⑥ parallel query plan:

larger query can be broken down into multiple parts and then executed.

- Ex:
- ① SAP HANA
 - ② Oracle in-memory
 - ③ MS SQL server in-memory tables
 - ④ Amazon Memory DB

disadvantage: ⑤ durability: loses all information when device loses power or is reset.

resolution: ⑥ snapshots / images of current state of data can be maintained if we need to roll back.

⑦ cost: the amount of data is ↑ than to store them.

⑧ traditional DBs are also reducing usage of disc.

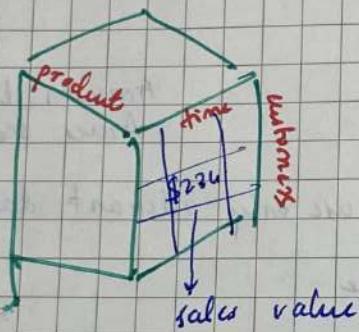
Cubes

OLAP cubes :

- In traditional DW, data is organized based on relational RDBMS (ROLAP)
- Data is organized non-relational in cube (MOLAP)
cube = multi-dimensional dataset.
- Arrays are used instead of tables
- Used for fast query performance.
- works well with many BI solutions.

sales data into multiple dimensions

* pre-calculated
aggregated values.



- MDX language :
Multi-dimensional expression from Microsoft can be used to query cubes.

Used when ① high performance is needed

② when hierarchy plays a role, & when one needs to slice & dice the data.

use case :

- * must be used for specific use case (specific set of data like datamarts)
- * more efficient & less complex w/ separate datamart
- * good for interactive queries with hierarchies.
- * can use data
* can be used from relational DB (built using star schema)

Operational Data Storage (ODS)

- similar to data warehouse
- data is fetched from multiple data sources & using ETL, data is stored in operational data storage (ODS).
- ODS is used for operational decision making

use cases / used when:

- ① you don't need long history of data
- ② can be very current or real-time.

e.g.: In a finance company,

customers can invest in ETF, stocks, crypto, etc.

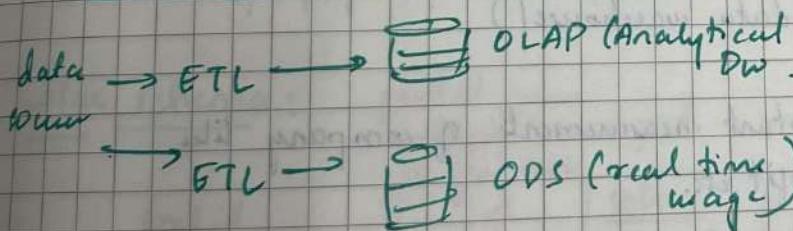
- we need to combine overall money invested from different investments & show their current balance

↓ company
this helps customers in deciding whether the company should give credit to customers or not.

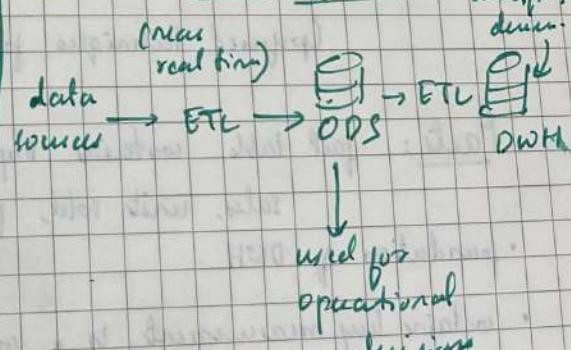
- ① Data is not appended / history is kept. But it is removed after usage

(updated)

ODS parallel

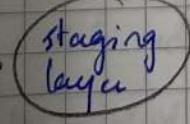


ODS sequential



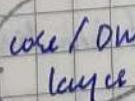
ODS is getting less relevant because of advancement of how to new tech for real time fast data.

Summary



(landing zone)

- minimal transformation
- stage data in tables



data marts

- access layer
- made for specific use case
- optimized for performance

- Business logic / single point of truth
- can be some times access layer for end user

Dimensional Modeling [unique way to structure the data]

[method of organizing data (in a data warehouse)]

facts

- measurement like profit.

dimensions

- context like category (period).

e.g. profit by year → gives meaningful insight for data.

- fact is usually in middle connected by multiple dimensions

why dimensional Modeling?

goal: fast data retrieval

• oriented around performance & usability.

→ wide tables w/ many columns might have duplicate data.
This is not feasible for query performance.

Hence we create separate customer, product table &
retain only their IDs as foreign key in
fact table

↓
this results in ↑ performance & ↑ usability.

(preferred technique for data warehouse!)

Facts: fact table contains important measurements of company like sales, unit sold, profit etc.

- foundation of DWH
- contains key measurement in a company.
- aggregated & analyzed.

facts are:

- ① aggregatable (numeric values, they are additive / can be added up)
- ② measurable
- ③ they are event- or transactional data.
- ④ date / time in a fact table.

Fact table contains:

- ① PK that uniquely identifies each record in a table
- ② FK's which are used to refer diff dim tables
- ③ fact → sales / profit / budget / KPIs.

- fact tables are defined by grain.

e.g:

id	date_id	region_id	profit
1	2022/1/12	1	\$23
2	2022/4/9	2	\$100

most atomic level are defined.

for specific region @ specific date, profit is defined.

∴ date is the most atomic value in this fact table.

- different types of fact.

Dimensions: In star schema, dimensions are clustered around fact.

- their function is to categorize the facts.

- character of dimensions: supportive & descriptive.

- used to filter, group & label our data.

→ dimensions are: ① not aggregatable (e.g. date cannot be added $\frac{2019+2020}{2}$)

no meaning

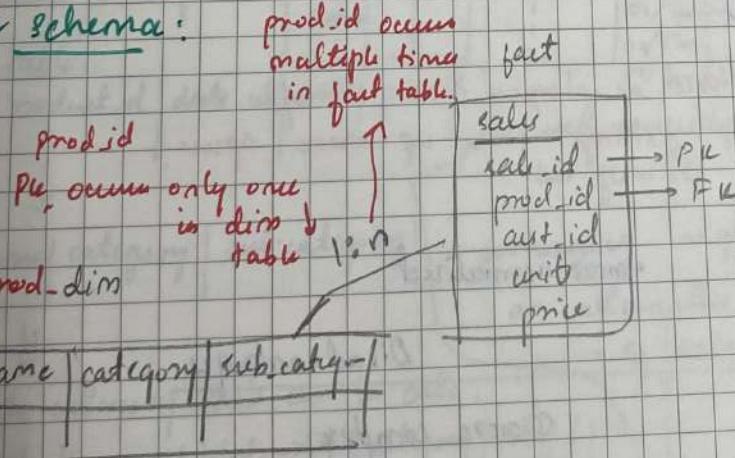
② descriptive in nature

③ more static (product name/categories) are not changing more than facts.

Dim table: contains → PK, FK & dimensions.

e.g.: people, product, place, time.

Star schema:



- there is data redundancy in dimension tables (data is denormalized & is not normalized)

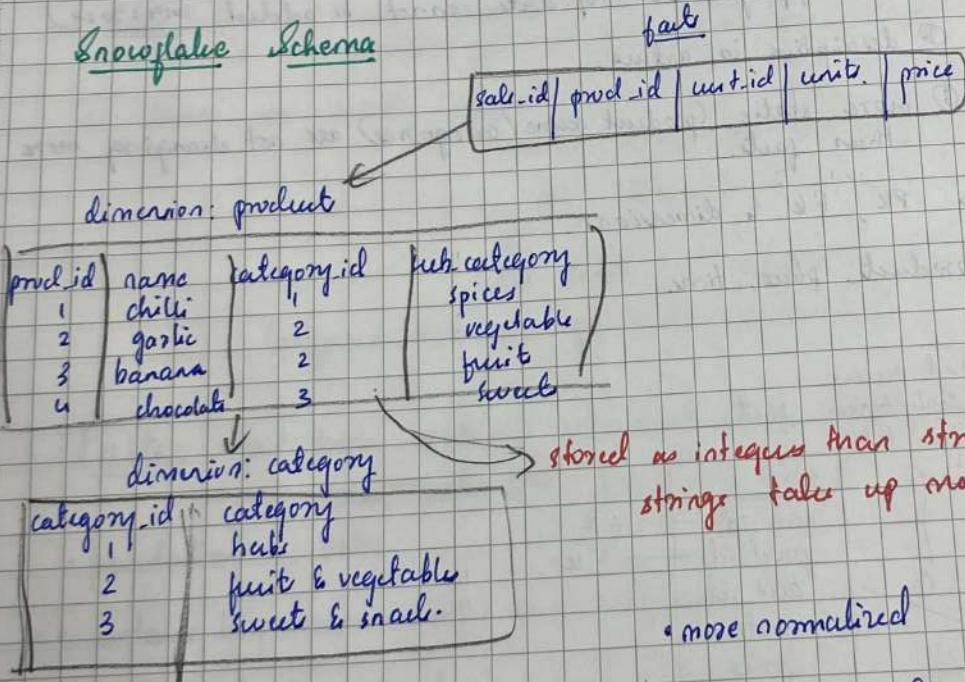
prod_id	name	category	sub-cat
1	chilli	herbs	spices
2	garlic	fruit & veg	vegetable
3	banana	fruit & veg	fruit
4	chocolate	sweet & snacks	sweet
5	chips	sweet & snacks	snacks

- optimized to get data out
- query performance (read)
- less expensive.

- Normalized :
- ① technique to avoid redundancy
 - ② minimizes storage.
 - ③ write/read operations are easier
update
 - ④ many tables
 - ⑤ many joins & complex queries

- Star schema is
- ① most common schema in data mart (usable & high performance)
 - ② simplest form (vs. snowflake schema)
 - ③ works best for specific needs (simple set of queries vs. complex queries)
 - ④ usability + performance for specific (read) use-case.

Snowflake Schema

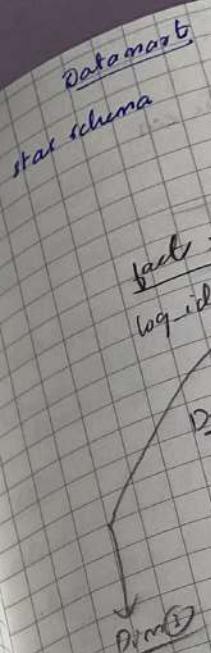


Advantage

- ① less data redundancy
- ② less disk space → less storage cost
- ③ easier to maintain/update data,
less risk of corrupted data. due
to no redundant data.
- ④ solves write/update slow downs

Disadvantages

- ① more complex
- ② more joins (more complex SQL queries)
- ③ less performance datamart/cubes



Data Mart

star schema

fact:

log-id	update-date ^{1D}	hrs-logged	project-ID	emp-ID	div-ID
--------	---------------------------	------------	------------	--------	--------

Dim ①: project ID

project-name

project-priority

core (DW)

star schema (default)

may be snowflake schema (if freq update is high)

thus storage is concurrent

Dim ②:

Employee-ID	emp-name	div-holiday	div
-------------	----------	-------------	-----

Dim ③

Dim ④: date-ID

update-date	year	month	quarter	day
-------------	------	-------	---------	-----

Facts

Additivity in facts

Additive

- can be added across all dimensions

- most flexible & useful

Eg: sales-table

sal-ID	product-ID	date-ID	unit	price
1	2	3	4	5

prod-ID | name | category | sub-category |



- across each category units can be calc.

date	total units	name	total units	category	total units
1	2	3	4	5	6

sub-category | total-unit |

Semi-additive

- can be added across a few dimensions.

Eq:	bal-ID	post-folio-ID	date	balance
bal-ID	1	1	01/Jan	50
	2	1	02/Jan	100
		2	03/Jan	100

- used carefully & less flexible

- averaging might be an alternative.

Eq: balance

Non-additive

- cannot be added across any dimension.

Eg: price of a unit in sales table.

cannot get total price of category; it is unit price.

i.e. we need to multiply units * price = revenue.

Eg: price percentage ratio (what is the inventory in warehouse).

- they have limited analytical value.

↓
store underlying values like numerator & denominator which can be calc from fact table.

Nulls in fact table

- do not use null values as is, if present in foreign keys as it can be an error to join main table w/ referenced dim table

↓

Hence, replace null w/ any default value like '999' / -1

Ex:	<u>portfolio_id</u>	<u>portfolio_id</u>
	1	1
	2	2
	null	999
	1	1
	null	999

- replace null w/ 0 if calculating avg, sum, min

Year-to-date facts

- often requested by business users
- MTD, QTD, fiscal YTD etc.

↓
these calculations are not in the defined grain[↑] fact table.

fact[↑] table is @ a defined grain @ daily level.

solt: store underlying values in defined grains

↓
like date value / month/year value.

↓
calculate to-date variations in BI tool

<u>Types of fact</u>	
<u>transactional</u>	<u>Periodic</u>
1 row = measurement of event / transaction • taken place at a specific time. • 1 transaction defines the lowest grain.	1 row = summarize measure of many events / transactions. • summarized over std. period (e.g.: 1 day, 1 week etc.) • lowest period defines the grain.
<u>sales</u> <u>dim's</u> <u>sale_id</u> <u>prod_id</u> <u>date_id</u> <u>units</u>	<u>weekly sales table</u> <u>weekly_id</u> <u>revenue</u> <u>sales</u> <u>cost</u> measure
• typically additive • have lot of dims (FK's) associated • enormous in size most common & very flexible.	<ul style="list-style-type: none"> not large in size aggregated table (not very fine grained) growth is not rapid but controlled. typically additive data. has lot of facts & very few dims. if, no event = null or 0

Types of fact table

Type	<u>transactional</u>	<u>Periodic</u>	<u>Accumulating snapshot</u>
① grain	1 row = 1 transaction	1 row = 1 defined period (plus other dim)	1 row = lifetime of proc/ud
② Date dimensions	1 transaction date	snapshot date (end of period)	Multiple snapshot date
③ No. of dimensions	high	lower	Very high
④ facts	measures of transactions	accumulative measure of transactions in periods	measure of process is life
⑤ size	largest (most detailed grain)	middle (little less detailed)	lowest (highest aggregate)
⑥ Performance	can be improved w/ aggregation	Better (less detailed)	good performance.

Factless Fact Table

A fact table can comprise of fact (measures) → numeric dimensions (PK, FK)

Defn: Sometimes only dimensional aspects of an event are recorded. → Factless fact table

Eg: reg_id | entry_date | dep_id | region_id | manager_id | Pos_id |

↓
has only events & no metrics.

① How many employees have been registered last month?

② How many employees have been registered in certain region?

Steps to create a fact table:

① Identify business process for analysis.

Eg: sales, order processing

② Define the grain

Eg: transaction, order, order lines, daily, daily + location
↓
rather we move fine grain (higher atomic level)

↓
as it leaves option
to analyze from all aspects of data
where as pre-aggregated data can be
used only upto limited use-cases.

③ Identify dimensions that are relevant
(what, when, where, how and why) → used for filtering
and grouping.

Eg: time, location, product, customer.

④ Identify facts for measurements.

Surrogate keys: artificial keys (integer no. assigned to every single row).

→ Natural keys: come out of source system / source generated.

Eg: sal_id → PX 0123 → alpha numeric values that take heavy space

→ Surrogate keys: are generated by database / ETL tool

Benefits → ① much smaller in size

② improve performance (less storage / better joins)

③ handle dummy values (nulls / missing values)
by using default values like 999 or -1.

④ easier to administrate / update

⑤ integrate w/ multiple source systems.

⑥ sometimes there are even no natural keys available.

⑦ always
⑧ both for fact
⑨ optionally keep

Case
Corporate IT

E-commerce

⑩ websites

⑪ each website
on multiple pl

Step 1: sales

Step 2:

Step 3:

Step 4:

plus fact table

- ① Always use surrogate keys as PK or FK.
- ② Both for fact & dimension (except date dimension)
- ③ optionally keep natural keys

Y guideline

Case Study: E-commerce

corporate IT

E-commerce company

① 3 websites

- ② each website is operated independently by multiple department

Step 1: Identify the Business Process.

sales transaction - which product sold
what is sale profit
sales of each website
performance over time
sales over time

Step 2: Define the grain

→ level of detail

most analytical value comes w/ atomic grain.
what 1 row of fact table should represent.

Step 3: Define / Identify dimensions

e.g.: male, date/time, file
customer
product

Step 4: Identify the facts must comply w/ grain.

to check • if facts are additive (so that can be agg across all dim)

* can use absolute discount price

* do not store discount percentage

* calc profit

- Tables used to slice and dice data.

- Always has a primary key (PK)

ProductID	Name	category
P001	abc	A
P002	def	B
P003	ghi	C

- use surrogate key (Eq: prod-PK \rightarrow 1, 2, 3)

↓
used to reduce space
usually integers

Data Dimensions

contain date related features: Year, Month (name & number), day, quarter, week, weekday (name & number).

- meaningful surrogate key: YYYYMMDD

Eg: 2022-04-02 \Rightarrow 20220402

- Extra row for no date (null) (source) \Rightarrow 1900-01-01 (dim)

- Time is usually a separate dimension.

- can be populated in advance (Eq: for next 5 to 10 years)

Date features:

① Numbers & text (Eq: January 1)

② Long & Abbreviated (Jan, January - Mon, Monday)

③ Combinations of attributes (Q1, 2022-Q1)

④ Fiscal dates (fiscal year etc)

⑤ Flags (weekend, company holidays etc).

Nulls in Dimensions

- Null must be avoided in PK's \rightarrow as they break referential integrity

- They don't appear in T-SQL.

Eg: Promotion table

Promo_Pk	Promo_name
Social media	1
Website	2
No promo	-1

date table \rightarrow	
Promo_Pk	date
20220101	1/1/2022
19000101	1/1/1900

- Replace null w/ descriptive values.

but this
bc as
we're

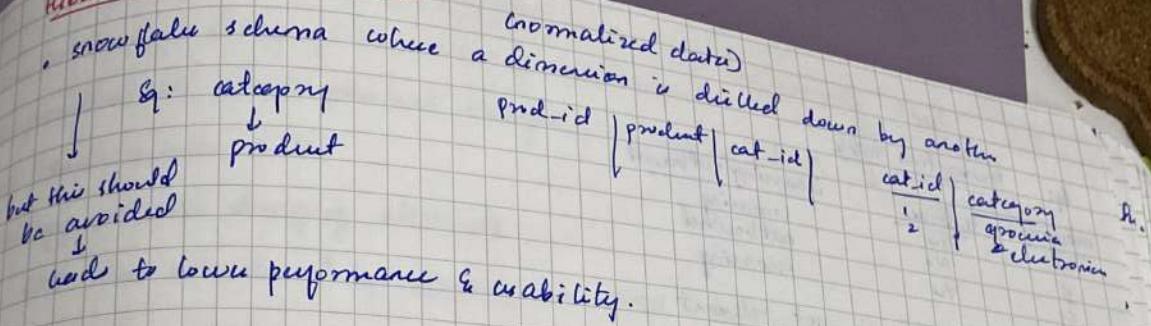
① one char
② combination

• Used
dimension

dimension

D

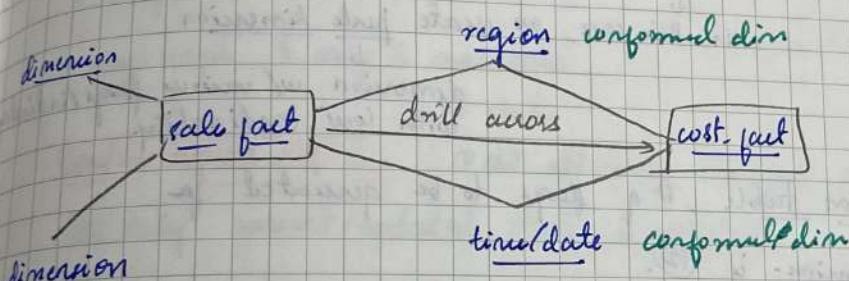
Hierarchies in Dimensions



- ① one should always denormalize / flatten the data for DW.
- ② combine the hierarchies & collapse

Conformed Dimensions

: a dimension shared by multiple fact tables/stars.
 used to compare facts across different tables



- same granularity is not necessary for conformed dim in 2 fact table

Eg: cost fact
 m1 | m2 | date-fact
 2022-01 | 2022-02

sale-fact-fle
 2022-01
 2022-02

or associated dimension

Degenerate Dimensions : A dimension that does not have separate dimension table, but still functions as dimension.

Eg: transaction sale fact

<u>transaction-fle</u>	<u>amount</u>	<u>payment-fle</u>
		263-023
		263-023

hence no separate table
 (dealt under one)

<u>payment-fle</u>	<u>header</u>
263-022	A
263-023	B
263-024	C

- All relevant information have already been extracted. Go other dimensions
- Attribute can still be useful (avg of payment amt)
- indicate that it is not a file & no dim table is associated w/ it.
 (-DD)

Eg: payment-DD.

- Occurs mostly in transactional fact

Eg: Invoice no., billing no., order id no

Junk Dimension

Eq: transaction table	
category	is-Bonus
A	No
B	No
C	Yes

X 2 X 2
incoming/outbound

- Eliminate them if not relevant
- If relevant, leave them in fact table.
- If the dim is long text value / long table size

↓
create flag for each dimension

↓
but this would be an issue if table has many columns / is wide

↓
in this case, we create junk dimension

dimension w/ various flags/indicators
with low cardinality.
(few unique values)

Eq: for above transaction table, # of flags to be generated for

junk dimension is 12.

But what if we have many dimensions?

9 indicators w/ 4 combinations
(take dimensions)

total rows/unique flags for junk dim = $4^9 = 262144$

① In this case, we only extract the available combinations of fact table.

② or split up into 2 or more junk dimension.

$$\text{Eq: } 4^{15} = 1024.$$

Role Playing Dimensions

A dimension that is referenced multiple times by diff fact table. Eq: date.

Eq: Prod table

id	order_date_FK	# products	product_FK	product_start_FK
role 1	date_FK	role 2		

date_FK / date / month / short month / year-quarter / year / weekday / is-weekend

date_FK

left
join

can analyse
① dt of product
② # of orders run
on that dt.

Slowly changing Dimensions

- How to deal w/ changing dimensions in DWD (Dimensions are usually static).
- Be proactive: Ask about potential changes.
- Business user & IT
- Strategy for each changing attribute
- SCD was introduced by Kimball in 1995 & distinguished b/w diff types (1, 2, 3, ...)
- Type 0 - Original
 - retain original data
 - Applicable if there are no changes occurring in our dimensions.
 - usually applicable for date tables (except company holidays).
 - 'Original'
 - very simple & easy to maintain
 - \rightarrow is static

- Type 1 - overwrite type :
- Old attributes are just overwritten w/ new values.
 - Only current/updated state is reflected.

g:	prod-key	name	category		\Rightarrow	prod-key	name	category
1	1	unglazed	acehn		updated	1	unglazed Te-7	acehn
2	2	choco-ri	cooki		version	2	chocolate bar- 20%	cooki
3	3	oats	Biscuite			3	<u>Delicious</u> oats	Rint

only this updated version is maintained.

- Very simple to implement
as we only update values in the dimension tables.

• No fact table needs to be modified. (only dimension table is modified)

Cons:

① History is lost

② Insignificant changes

③ might affect/break existing queries (e.g. if case when query are handled).

Type 2 - Additional row / New row (powerful type)

(difficult strategy).
to maintain history

• Powerfully partitions history

• changes are reflected w/ history.

• Instead of updating same row, create a new row w/ updated values!

Type 1:

prod-key	Name	category
1	sunglass	accessories
2	choco bar	sweet
3	oat meal bsc	sweet

↓ updated

prod-key	Name	category
1	sunglass	accessories
2	choco bar	sweet
3	delicious oat meal bsc	Biscuit

Type 2:

prod-key	Name	category
1	sunglass	accessories
2	choco bar	sweet
3	oat meal bsc	sweet

↓ add row

prod-key	Name	category
1	sunglass	accessories
2	choco bar	sweet
3	oat meal bsc	sweet
4	delicious oat meal bsc	Biscuit

↓ to get # of products?

do we want prod-ID natural key
(some generated
does not change)

In question only new foreign key
for 'oat meal biscuit' - i.e. it will
be used in fact table

hence we don't need to make any
changes in fact table.

hence respecting the history.

Ex: Prod-Pk

Prod-Pk	Prod-ID	Name	category
1	SG-Prod	sunglass	accessories
2	CH-Prod	choco bar	sweet
3	OT-BSC	oat meal	sweet
4	OT-BSC	del. oat meal	Biscuit

↓ count of distinct prod-IDs → gives # of grades.

Administering Type 2 SCD

how to get & identify current product name based on prod-ID

hence include effective date & exp date to show period in which values are valid.

current PK must be used in fact table.

requires surrogate key instead of natural key.

Current PK in fact table: ① Add new row in dimension (ef-date, exp-date)
② lookup in dimension w/ natural-key + ef/ex-date.
③ Is-current flag column can be used.

Q whether
Ex: prod-
But for

Type 3: Ad

+ in fact

. Instead

so: Prod

and instead

of null

default value

Choosing still valid

used to filter data

range.

18-current

Year

Year

Year

Year

Mixing Type 1 and Type 2

- Q whether to use Type1, Type2 or both - type1, type2 depends on dimensions.
- Eg: prod-name can be type1 (changing out meal to "delicious out meal")
But for category, we need to maintain both old & new category.
↓
this is not a technical decision.
↓
need to ask business user to then define

Type 3: Additional Attribute

- is in-between: switching back and forth between static previous state and static current state.
- Instead of adding rows, we add a new column.

prod-Pk	prod-ID	Name	category	prev-category
1	SC-TR7	sunglass	accessories	
2	CM-B70	choco bar	sweets	
3	OT-BSC	oat meal	biscuits	Biscuits

current state		prev state	
category	amount	prev. category	amount
accessories	\$25	accessories	\$25
sweets	\$6	sweets	\$14
biscuits	\$8		

- typically used for significant changes at a time
(eg: restructuring in organization)

region-Pk	region	prev-region
1	North	North
2	West	West
3	South	West
4	East	not applicable

↓
can add multiple additional rows

- instead of adding a row, we add a column
- typically used for significant changes @ a time
- enables switching b/w historic & current view.

Limitations

- not suitable for frequent unpredictable changes like → go for type 2
- minor change → type1

understanding ETL Process

Extract

- Data is part of DWH
- Understanding data
- data is transformed from staging layer to core layer.
- Data is truncated (most commonly)
 - ↓
 - i.e. once all the data is copied from staging to core layer, & all transformations are complete, the data is then deleted/truncated from staging layer.

some
alternatives

Extracting types

Initial load

- first (real) run
- all data

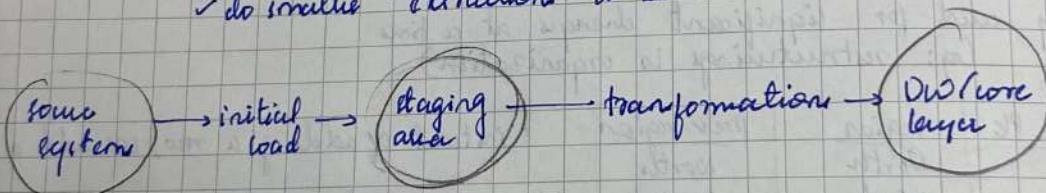
Delta load

- subsequent runs
- only additional data is loaded.

Initial load: ① extract all initial data from source systems

② After discussion w/ business user + IT:

- ✓ check what data is needed
- ✓ when is a good time to load data (nights or weekends?)
- ✓ do smallish extractions to test



• After all transformation steps have been designed

Delta load: Incremental periodic extraction / load.

- there is a delta column for every table.

sales_data	Name	amount
		→ a column to identify new data.

- transaction data, sales_dt etc can be used as delta columns.

- alternate columns such as primary-key can be used to identify new data

Ex:	sales_key	name	Amount
1		singlserve	\$25
2		chow bar	\$3
3		oat meal bisc	\$4
4		choco bar	\$3
5		oat meal bisc	\$4

⇒ ✓ MAX(sales_key)

↓

✓ MAX(sales_key) → variable X

✓ Next run: sales_key > X

Load

- Some tools can capture what data has been loaded & what has not been loaded
- alternative case: for dim table (not fact table).

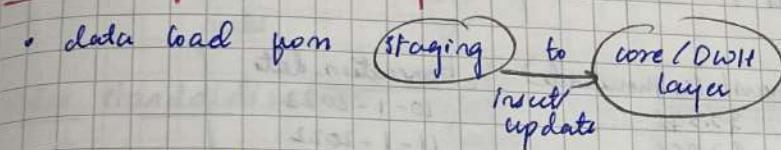
① Perform full load everytime

② compare data w/ data that has been already loaded

③ Depending on data volume → do performing of data load

More data load = longer ETL process.

Load Workflow (Insert/Update)



• insert is done when the data is new (e.g. was not loaded before)

• update: updating/replacing product name in product table
if name exist or

if name does not exist, then add product name in product table

• Typically we do not delete data

• DELETE: if row in source system gets deleted

↓
then do not remove row from DW/H

add a new column w/ flag whether the row is active/valid
in source system or not.

Transforming Data

Goal: Create a consolidated view of all data for analysis purpose

① Consolidate (from multiple systems)

② Reshape (for analytical purposes)

③ Consolidate:

table 1	transaction-ID	amount	date
	T1	\$5030	10/1/2022
	T2	\$5053	11/1/2022
	T3	\$5654	12/1/2022

table 2	transaction-ID	Amount (in thousands)	transaction-date
	T14	\$5.345	10-1-2022
	T15	\$7.953	11-1-2022
	T16	\$9.654	12-1-2022

- Making data compatible & consistent w/ each other.

④ Reshape data according to business requirement

e.g.: ① date format: 10/1/2022 to 2022-10-10

② pivot the data

Kinds of transformations

Basic : ① Deduplication

④ Value standardization (Categorization)

② Filtering (rows & columns)

⑤ Key generation.

③ cleaning & mapping (Integration)

Advanced : ① Joining ③ Aggregating

② Splitting ⑥ Deriving new values.

Basic transformation: ① Deduplication: Merge 2 store tables that have (remove duplicates) same product in both tables
↓
no distinct values

② Filtering rows: filter out irrelevant rows (e.g. when amt < 0 or type = 'refund')

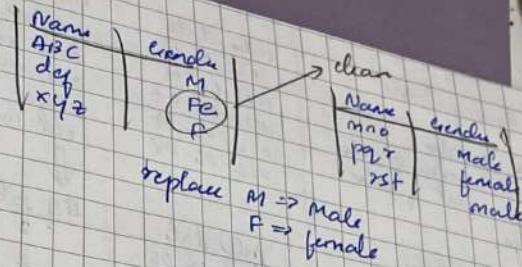
③ Filtering columns: If a column in a table contains same value in all rows
(has 1 unique value)
↓

then delete the column

① Cleaning & Mapping (Integration):
map different values: Eq:

map/replace null values
w/ 0's

Eq:	sale	sale
	300	300
	100	100
	null	0



Value standardization:

Month	sale in thousands
Jan 22	\$1.5
Feb 22	\$4.55
March 22	\$3.321

Month	sale
Jan 22	\$1500
Feb 22	\$4550
March 22	\$3321

Key generator: Add an auto-generated surrogate key

(generated in ETL system or DBMS)

Advanced Transformation:

① Joining

Product Dim		prod_id	name	category	of-date	exp-date
1	PS21	PS21	Almond	nuts	2021-01-01	2021-01-01
2	PS22	PS22	garlic	fruit & veg	2021-01-01	2021-01-01
3	PS23	PS23	banana	fruit & veg	2021-01-01	2021-01-01
4	PS24	PS24	choco	sweet & snack	2021-01-01	2021-01-01
5	PS25	PS25	chips	sweet & snack	2021-01-01	2021-01-01

Sales-fact

sale_pk	product_id	date	prod_FK
3	PS33	2022-01-01	6
4	PS21	2022-01-01	1
5	PS22	2022-01-02	2
6	PS29	2022-01-02	9
7	SS30	2022-01-02	10

product
surrogate key introduced
as FK in sales fact table

② Splitting a column into 2 columns:

store dim

store_id	location
1	New York, NY, 10011
2	Orland Park, IL, 60462
3	Houston, TX, 77002

store_id	city	state	zip
1	New York	NY	10011
2	Orland Park	IL	60462
3	Houston	TX	77002

split by length/position

③ Aggregations

- sum of sales
- count number of sales
- distinct count of sales
- average

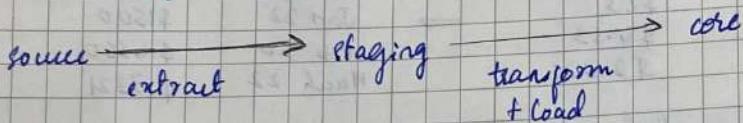
• subtract/multiply

- calc abs amount of tax using tax % & sale value.

sale_date	Name	Amount
2022-06-06	sunglass	\$25
2022-06-06	choco bar	\$8
2022-06-07	oat meal bisc	\$4
2022-06-07	choco bar	\$3
2022-06-08	oat meal bisc	\$4

sale_date	Qty sales	Amount
2022-06-06	2	\$28
2022-06-07	2	\$7
2022-06-08	1	\$4

Scheduling Jobs



- requires enterprise version in ETC tool (free version not possible)
 - ↓
 - use external tool

Guideline: ① How frequent to update data? 1x, 3x/day → ask business user every 30 min

② How long does it take? 5 min? 1 hr?

③ What is a good time to execute ETL? (∴ productive some data night slow down)

if we use resources when there is heavy load

- needs to be figured out what will not provide problems
- initial load vs delta load
 - effect on productive system.
 - short road access
 - night? early morning?

Demo: Plan of Attack

- ① Look at problem & plan
- ② Set up tables & schemas
- ③ Create staging table (truncate)
- ④ Transform + Load:
 - ① read data from staging
 - ② transform (clean + extract)
 - ③ update / insert

problems

No new values \Rightarrow staging will be empty afterwards & before the next run.
empty staging \Rightarrow max (product_id) = null \Rightarrow full load.

ETL Tools

Enterprise

commercial

most mature

graphical interface

✓ architectural needs offered ✓ support

support

e.g. Alteryx,
Informatica
Oracle data
integrator,
Microsoft SSIS.

open-source

source code

✓ open free

✓ graphical interface

✓ support

e.g. Talend open
studio,
Pentaho data
integration,
Hadoop

cloud-native

cloud-tech

data already in
cloud?

✓ efficiency.

✓ flexibility
(can we to
get data from
other service
provider)

e.g. Azure data
factory,
AWS Glue,
Google cloud
data flow.

custom

own development

✓ automated

internal resources

✓ maintainence?

✓ training?

Choosing right ETL tool:

① Evaluate current situation / needs in company

- What do you want to improve
- Data sources & other tools
- Define your requirements
- Define responsibilities
- Who are the users?

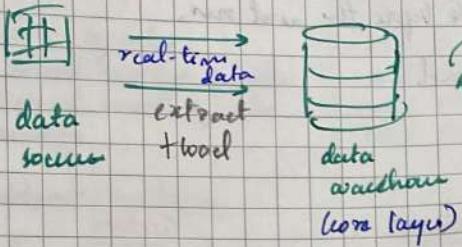
Table to evaluate while selecting ETL tool	Text	Must have	Weighted importance rating
cost			-5
connectors			5
capabilities			5
ease of work/use			5
reliability			5
support/contracts			5

② Evaluate Tools

total weighted score:

③ Contact providers to
get test/demo/free trial
for POC.

ELT (Extract Load Transform)



Leverage compute power
transform [since transformation is done in data warehouses this target data warehouse must have high compute power]

- we scale for transformation (e.g. Snowflake uses massive ML processing)

- use when streaming data sources are used.
(when real-time processing is used)

ETL

- more stable w/ defined transformations
(based on defined transformation flow stage → core)
- more generic use-cases
- security can be handled a lot better
(by agg send only anonymized / encrypted data to DW)

use-cases : ① reporting (no real-time req/ defined transformation)

② generic use-cases

③ can be w/ (default strategy)

ELT

- requires high performance target DB (but cloud platform can help)
- more flexible
- transformations can be changed quickly
- real-time requirements
- use cases : ① req from data scientist, ML
② real-time req
③ big data & high volume

Data Warehouse use cases:

DW: centralized data location for the purpose integrated data that is optimized for the purpose of data analysis

- ① used for reporting (to make data driven strategic decision) from diff data sources that are integrated.
- ② Performance for analysis by quickly & easily connecting to DW by Business user
- ③ Data quality, availability
- ④ continuous training of ML models for predictive analytics.
- ⑤ Use big data : Aggregate & filter

Issue: Business Using index
② Data in D

Optimizing A Data Warehouse

Issue: Business users facing bad query performance after setting DW
wait long for result to be retrieved

Using index

Data in DW disk space is not stored in a particular order

↓
it is stored whenever some free space is available

↓
this becomes a problem if we want to retrieve specific data.

e.g.: select product_id from sales where cust_id = 5;

table scan & read becomes inefficient.

fact_table

trans_id	prod_id	cust_id	payment	psn
4				
5				
5				
8				
5				

⇒ query: select prod_id from fact_table where cust_id = 5

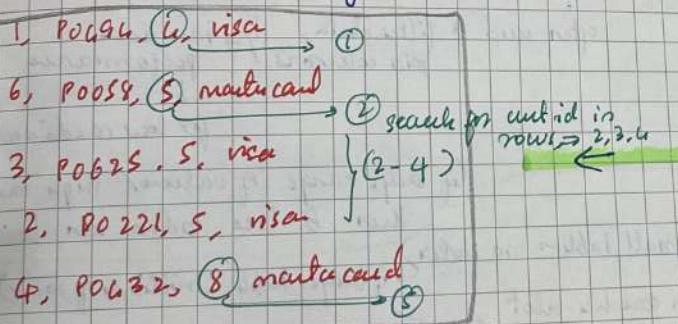
read is large
scan

∴ make index based on customer_id

location	value (cust_id)
1	4
2	5
5	8

↓
go to location using pointer assigned

disk-space ordering



- makes data retrieval faster
- ↓
- no need of full table scan

Cons: ① Writing/updating data is slower
as pointers need to be maintained

② needs additional storage
default/

① B-tree indexes: (std. ② index)

• multi-level tree structure

• pages break data down into pages/blks.

• used when we have unique values in column

(high cardinality - high range of diff values)

• Never put an index on all columns (as it is costly in terms of storage, slower updates/writes).

• automatically set up on PKs (primary keys)

(good for many replicating values)

② Bitmap indexes

• used when data is large + low cardinality

• very storage efficient.

• data is stored in so called bit

• hence storage efficient as

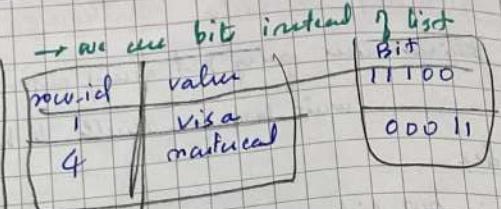
• not efficient for write/uploads

• very good for read operations

• few DML operations

Bitmap-Index

trans_id	prod_id	cust_id	payment	price
1	P004	4	visa	18.2
2	P0221	5	visa	1.5
3	P0625	5	visa	5.9
4	P0431	8	mastercard	11.6
5	P0058	5	mastercard	12.4



Bit
11100
00011

value	1	2	3	4	5	6	7	8
mastercard		x		x				
visa	x	x	x					

- ① Should we put index on every column?
- ② No, they come w/ a cost
- ③ storage + create/update time
- ④ should be used when necessary
 - ↓
 - (slow query performance due to full table scans)
- ⑤ small tables do not require indexes

- ⑥ On which columns should we place index?
- ⑦ To those columns & to those tables, that often need to have data filtered.

Eg: if filtering is just done on 10% table then full table scan is avoided → if index is used.

Fact Table

- ① B-tree on surrogate key (PK)
- ↓
subsumes automatically

- ② Bitmap index on foreign keys.

↓
often used to filter data, → increases join columns performance
for low cardinality

(if large range of values/ high cardinality then B-tree index can be used).

Dimension Table

- ① size of the table (small tables → no index) if large table (millions of rows)

- ② Are they used in queries a lot? → then?
- ③ choose based on cardinality

Syntax:

(postgres)
syntax

create index index-name on table-name [USING method]

(
 column-name [ASC/DESC],
); ...
)

if there are 100 columns & only 5 columns are needed → thus we don't have to process all 100 columns → 5% of data need to be processed.

Columnar Storage (used to improve query performance)

- ① traditionally data is stored in rows (based on above w/ of index)
↓
all rows need to be scanned, even if we need to check 1 row of data

1, 2, 3, 4, 5
P004, P001, P000, P0200, P0800

4, 5, 5, 8, 5

visa, visa, visa, mastercard, mastercard

18.2, 1.5, 5.9, 11.6, 12.4

Eg: select prod_id from sales
only this needs to be scanned.

↓
less data needs to be processed
storage is less ∵ data type stored is same → hence can be compressed efficiently)

• Need to have including:
①
②
③
Benefits: ✓ full

problems: ✓

Massive

this i
a

E c

In p

The Modern Data Warehouse

On-premises DW

- need to have own local hardware including:
 - ① storage layer
 - ② compute layer
 - ③ software layer

need to have a physical data center

Benefit: ✓ full control (own full infrastructure, hence own data governance)

Problems:

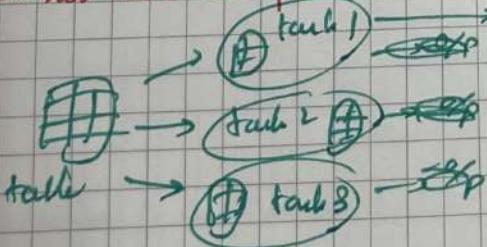
- ✓ full responsibility: (Buy if need for large workload, maintain)
- ✓ high cost for update & maintenance
- ✓ more internal resources
- ✓ less flexible
- ↓
if less workload, then there is wastage of resources which cannot be sold easily)

Massive Parallel Processing (MPP): Eg: if a query needs to be run: Select * from sales where cust_id = 5

this is traditional system,
a big task is broken
down to smaller tasks
& are executed in sequential
manner

↓
can be time consuming (need to wait until previous subtask are completed)

In MPP, all sub-tasks are processed [↓]! No need to wait for other sub-tasks to complete!
each nodes work independently shared disk - shared nothing architecture
compute resource are not shared but disk is shared



- solves performance issues
- millions of rows → process faster
- many people can run queries @ same time w/ good performance
- helpful w/ centralizing a massive amt of data.



Cloud DW

- software-as-a-service: Not owing infrastructure
but pay for using infrastructure
- pay for what you use
- managed security
- optimized for scalable analytics.

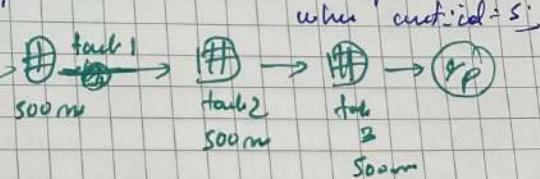
Benefit:

- ✓ fully managed scalable (can opt for additional subscription, resell)
- ✓ cost-efficient.
- ✓ managed security (better than on-prem)
- ✓ availability (w/ 99.9% uptime)
- ✓ implement quickly (link to market is quick)

Problems:

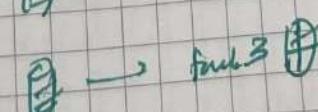
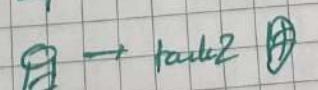
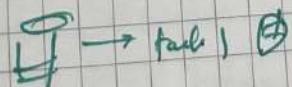
① regulations can be difficult to fulfill in case of special needs.

Most companies hence opt for cloud DW.
↓
they use technologies like MPP, columnar storage.



or data storage is centralized
but distributed to multiple data storage

buy & compute nodes also work independently



shared-nothing architecture
workload is split to process individually