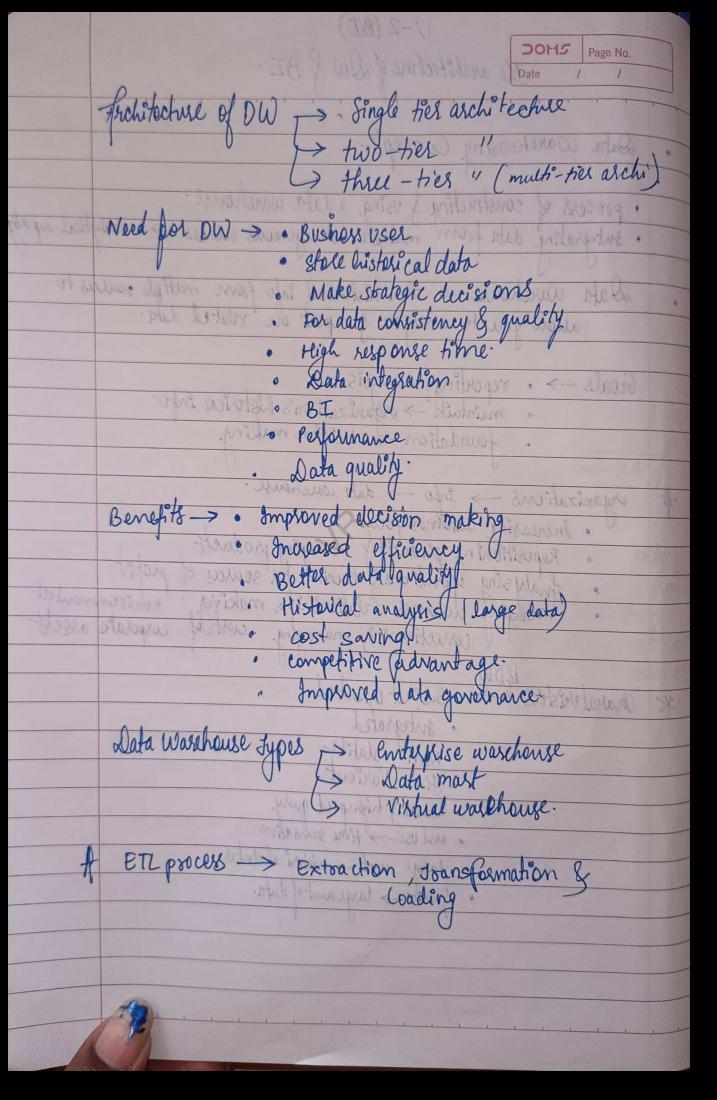
	U-2(BL)
	The architecture of DW & BI. Date / /
	in shotshower of DW -> single her architecture
	Data Warehousing Concept-
	1977 - 1974 - 1974 () 52 d - 1784 () - 5 - 5 - 5 - 5 - 5 - 5 - 5 - 5 - 5 -
	· process of constructing & using a data wantous.
	· process of constructing & using a data warehouse. · integrating data from mulipple heterogeneous sources. analytical reports
•	Data Warehouse > store historical into from muliple zoulcus to. allow you to analyze of report on selected data.
	· resp. resp once time
	Goals -> reporting analysis
	· maintain - organization's historical into. foundation - decision making.
	foundation - decision making.
M	Sala quality
#	organizations -> info -> data workhouse.
	Repositioning automer focus
	Repositioning products & managing products finally sing sporations & looking for sources of profit managing customer relationships, making environmental corrections & managing cost of corporate assets
	· managing customer relationships, making environmental
	Corrections & managing cost of corporate asset
11	Charaturisties > Subject oriented
*	
	· sutegrated Non-Volatile
	Jime Valient
	access & high speed query.
	· end user -> time sensative
	· large and - historical dato
	· guerres -> large and of data.
	U · · · · · · · · · · · · · · · · · · ·



		Date / /
4	Business Intellighel (BI)	Data Warchouse (Ow).
		· contral to cation -> store
- 10 - 10 - 10 - 10 - 10 - 10 - 10 - 10	set of strategies of technologic to another of visualize data to make divisions	contral to cation -> store consilidated data from multiple data sources
	data - vser's behaviour belongs to BI:	· customer commune tron > data
- Soul		· Car, bandage .
0100	set of technologies & strategies:	· storage
P Page	de sales	
	present data -> reports, charle & graphs.	· presents data in tables
via el	top executives & senior manager use BI.	· data enginers, data & busi'ness analystr use data warehouses example - Amazon Realshift
-41	require extensional et mes	analysis of the work
•	example - Datapine.	. example - Amazon Redshift
		. storing data from seweral somes.
•	create business insights.	, 31 of the form forms
	movidus user-kiendly foots	· Typically accessed via SQL queries.
- 93	provides user-friendly tools for data analysis & visulization.	· Typically accessed via SQL quelies.
-31-1	The state of the s	endingerit la office
•		can be added dynamically.
	· Small to redding.	sice. Medium to large
	· CHEST SEWEV	hordischus; diens Schuck
MARKE.	. Committed to prove defined in	Acers - support and-hor requests
	The of them sol sides.	Speak there there was . 6- Jung?
	a male speed but par .	applied muslem 20 pp.A.
		the state of the s
	0,03.	Just . Phillips
	1000	1 Miles A

			Date / /
(010)	ROLAP	MOLAP (38)	AAJOHiness chledingene
implementation	n · based on · ·	based on at ago	both relational mustiple mension technique
	relational DBMS	multidimensional DBMS's	multiple memoristichique
Adv >	· handle large data	· excellent performance	MOLAP tools
	· can devotage.	· perform complex	uhilize both pre-calculates
	· Functional ties.	calculations:	pre-calculate
	· inherit in the		ubes & retir
	relational database	th chart & graphs.	Libes & retir
Disadu->	. Performance →slow.	· 1Priviled dita handling	supports disali
estrature.	· Limited by SQL functionalties	· requires additional	of MOLAP.
The Johns	· Limited by SQL functionalities	investment in	of MOLAP.
		· C	
alul Sente	ROLAPE	Myhto.	OLA Picked shows
Sp. gravia	· dypically accused his	day plane	sor Data cubes
	· User Star Schema.	a harsulization	do" dimensione e mi se
	· Additional dimensions		dd" dimensions require resaction of data bost
	· Additional dimensions can be added dynamically		
atabase			and to all
Size	· Medium to large		Small to medium.
nchetochire	· client sever	, (lient/server.
ceus -	· support ad-hoc requests	• li	mited to pre-dyined dirent
poed ->	· crood with small data sets	·F	aster for small to medium
1	· Any for medium to large		my for leage datasets.
	data set.		110
Flexibility	· High		Low
Scalability			

OLTP

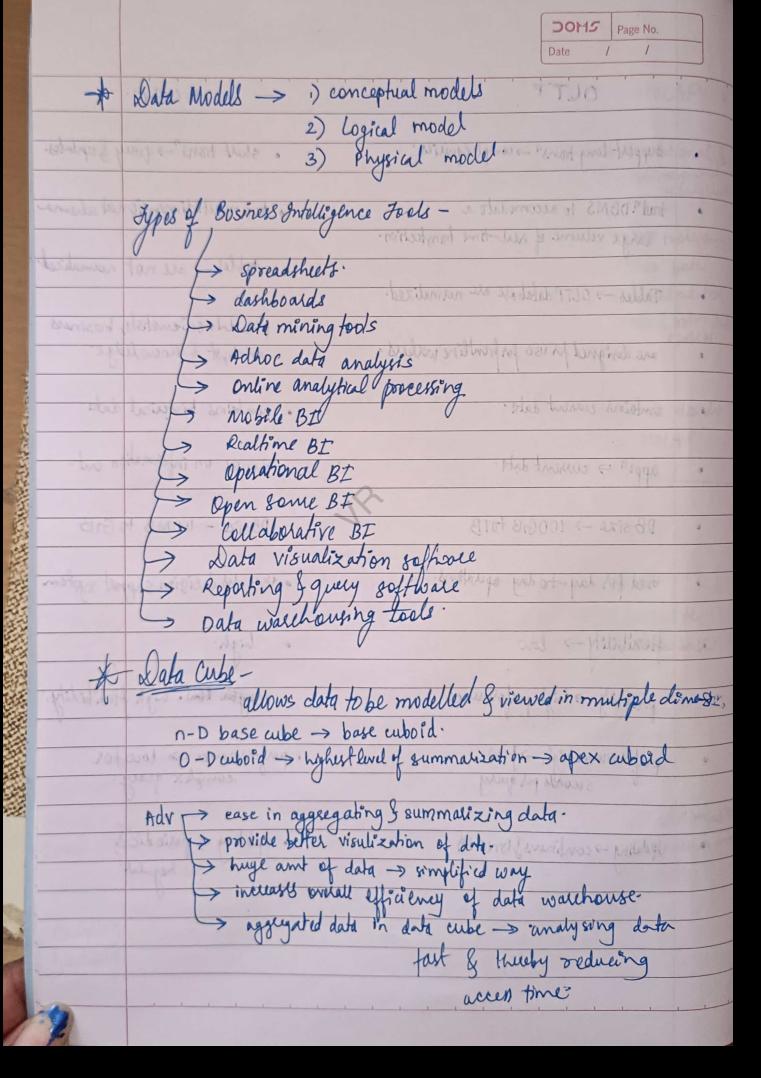
OLAP.

support long trans -> complex queies. · short trans" -> query & updates. trad" DBMS to accomodate a large volume of real-time transliction. · has multidimensional schema . tables - are not normalized. Tables -> OLTP database are normalized. for data scientists, business analyst & knowledge are designed for use for frontline workers. · contains historical data. contains current data. · fecuses on information out. appin -> current data. DBSIZE - 100MB to GIB DB size -> 100GB toTB used for day - to day equations. · used for decision support system flexibility -> low show and of plate on may tow. high flexibility. privily -> high performance performance -> loco for complex query. performance high > few seconds per query Updating -> continous & toregular . updating -> perio dic & segulati

sources brief the war of days warrander

appropried date in date cube - surabysing date

When historinas 6



Fact table

D'mension table

attributes of dimension title.

· contains attorbutes alongwhich fact table calculates metric

- contains less attoibuties of more records.
- contains more affir butes gless records.

· Fact table grows verticially

· dimension table grows

Fact table -> primary key which concatenation of primary keys -> dimension table

catains -> primary for

schema contains less no of fact tables

· Schema cartains more no of dimension tables -

fact table can have data innumers.

· always contains attributes

Star Schema

Snowflake Schema.

- simple frommon modelling paradigm where dots warehouse comprise of fact table.
- which includes hivatchial form of dimensional

ster schema -> not use normalization;

elimination reduces of

totais fact & dimension tables.

sub-dimersion tables including facts dimersional tables

simple to undertand y early.

high level of data redundancy.

· low level of data ordinary

uses more space.

· cube processing might be slow because of complex

. It uses less space

EE