

Date: \_\_\_\_\_

## ETL Processing :-

There are number of architecture deployed the data warehouse.

### Top-Down:-

#### Hub and Spoke:-

In this Datawarehouse contain the historical or daily data.

In this dependent data Mart is present.

Data warehouse mostly is normalized we can build it in the dimension model like snowflake or starflake.

Data Mart contains the summarized data if in the past frequent question is asked then we can add data mart.

Bottom's up:- In this independent data Mart more is present we cannot apply across analysis. Because across analysis are not related to maf.

Date: \_\_\_\_\_

## Data Mart Bus:-

We build data marts but if's requirement are according to enterprise level.

⇒ ETL

We have number of data sources for warehouse or Mart.

We have to identify which technique we use to access data and on which time and on what frequency.

First we create data extraction design

- 1) first identify the requirement
- 2) identify the data sources
- 3) identify the techniques to extract the data from sources which filter we have to apply.
- 4) first design is created then proper is created on ETL tool and on the

Date: \_\_\_\_\_

We define if collect or create programme.

There are three types of meta data

- 1) Source meta data
- 2) Target meta data
- 3) Intermediate meta data

Staging area

→ We can't directly convert data from source to target. We have to move it to the staging area.

→ In staging area we can apply different transformation according to the requirements of enterprise level.

Like for male we chose 1, 2 because meaning of 0, 1 is in sense of true, false

Date: \_\_\_\_\_

1) In source system there are multiple domains for first gender but in target we have one domain  
In this case we do the scalar transformation This is one to one mapping.

First m L

2) A<sub>1</sub>, A<sub>2</sub>, A<sub>3</sub>

↓  
MERGE

A

$1:N$  | is simple  
 $M:1$

3) A 20 1307

↓  
split

This one is easy because primary key ending is known

A<sub>1</sub>, A<sub>2</sub>, A<sub>3</sub>  
20 1307

Some time we have to convert or split Address if takes some time weeks.

$1:N$   
 $M:N$

is difficult

That's why ETL takes cost and time.

Date: \_\_\_\_\_

⇒ Data Quality is a continuous process if we deploy system with 50% accuracy then we can better the quality or refine the programme.

⇒ Business requirements must be clear so we have to move with the 100 percent quality or with less quality with multiple filter step. We verify it in the future. Verifying process must take less time.

⇒ Decision maker get data from sources and then give feedback to the operational system to take action.

Date: \_\_\_\_\_

## Loading Strategies:-

option 1 is that in the Staging area we do aggregated construction then move the detail data to ~~data warehouse~~ and aggregated to Maf.

Option 2:- is that in the warehouse we do aggregated then move the aggregated data to maf. this increase load on the

In logical there are three layers we get data from source and in the staging area transform if and then move to target. In physical there are some time staging area is in the warehouse we do the decision of Staging area separate or not on the basis of load.

Date: \_\_\_\_\_

## Data Extraction techniques-

Inmediate

Deferred.

If Data Extraction is Deferred and time frequency is on daily basis.

Current value vs periodic

Operation system store only current value if we extract data of midnight then if the change is done three times in a single day. Then we only get the last value. For this we use Inmediate extraction.

Every dimension has specific frequency technique and time.

If the operation system use the periodic means preserve the change. If the data is periodic then we can use the deferred exact method.

Date: \_\_\_\_\_

→ Immediate

1) Transaction log

if we perform any transaction in system then all the information is stored in the log file. It is the job function in the RDBMS if we change value then it is stored previous and next.

if the transaction failed then Database admins use log file

for recover the database.

if the log file full then if extend or create new log file.

if the transaction not failed then we not need log file.

Through log cache we can write in log file

log file means history of change.

Date:

.. In log file no activity perform for the warehouse, wife don't for the operation system and we write for the log. So performance high read.

First option is to extract from from the transaction log if available.

There is no performance issue or the extraction cost.

## 2) DB Trigger Programme.

These trigger are call explicitly.

We write trigger for the tabel.

In warehouse if we want to extract data from trigger then write the programme in the source system.

If the log file is not available or we need all data then

log. Then we write trigger in the source system from where

we extract data from source.

When trigger can and

it extract data from this

activity and move it in

the staging area.

Performance is high but low other

Date: \_\_\_\_\_

If log and cost is high to create programs then there is addition cost.

Minor impact of performance of trigger on the source system and development cost of the trigger.

### 3) Source Application

In this case the RDBMS is not available and we don't have log file and we are unable to trigger in the RDBMS in this case we write programme to extract files from the source Application and move it in the staging area. Write programme in the source Application.

There is a huge development cost of this programme and also huge performance issue.

In option we extract file from source Application according to target data requirement.

Date:

- We use option 3 in the worst case.

The best option is for log file if the data that we need is not available in the log then use option 2.

These are the option if we need data extraction immediately. We check what is the strategy of the operation like for

- 1) current value
- 2) periodic value

if current then

preserve history

we need to preserve them.

⇒ 2) Deffered =  
option 2 - Date Time stamp

we extract data on daily basis or on monthly or on hourly then how can we know is it in the today date or not.

for this we use timestamp then it is very easy to extract data.

Date: \_\_\_\_\_

## option 2 - Comparison

we have the copy of the yesterday  
and at mid night we extract  
today and subtract the  
Today from ~~the~~ yesterday and  
get the today

performance issue and we have  
to maintain large file cos. high

And file comparison programme create  
performance issue.

## Data loading Strategies

- Full Data Refresh
- Incremental Data Refresh
- Trickle feed - continuous update.

### 1) full Data Refresh

In this we  
empty the loaded data and  
then in the staging area  
we merge the previous and

Date: \_\_\_\_\_

new data and then load if  
in the load label.

In this we reload the previous  
data again.

if the existing/ previous data is  
95% and new data is  
5%. then if the new come  
we reload this 95% by merging  
it in the Staging area.

only insertion operation perform in the  
warehouse because we do merge in the  
staging we don't need to compare. It take  
less time if  
2) Incremental Data Refresh. The data is high.

In this we only load the  
new data in the loading  
Staging. By using ~~insert~~ destructive  
merge or constructive merge  
or append.

Destructive merge replace the  
previous.

Constructive merge store previous and  
new

Append only add the new.

In this we need to do the comparison and  
insertion so it takes more time

if the data is small

Date: \_\_\_\_\_

if the ratio of new data and previous data is low then the full DDL refresh is high  
if the ratio increases then the full DDL refresh is better.

if the ratio is low then we use incremental.

But this is not hard and fast rule.

- 3) if we use active warehouse then we use Trickle feed - because in this configuration update is needed.

Date: \_\_\_\_\_

In case of small dataset we use the nested loop join.

Nested loop join perform efficient.

In case of large dataset or very large dataset we use the hash join, sort merge join, merge join.

→ DSS vs OLTP.

DSS is for the different flavour we do analysis in DSS and there is a read operation only so there is no conflict. In OLTP there is a conflict of read and write. we must have concurrency control protocol.

→ Optimized is the component of DBMS Software.

Main role of query optimizer is to get the request and give the execution plan.

query optimizers have number of option then we select the best one from this.

if the optimizer select the top 4 option for the query

Date: \_\_\_\_\_

then we can say it as good  
if it select the average option  
then average optimizer.

Optimizer get the request or the  
input and generate the execution  
plan other task of this is  
to select the best path to  
execute this query.

path selected on the basis of  
temporary space communication channel in  
case of remote access it is important  
check the I/O operation.  
1) I/O operations: Almost every optimizer  
consider it. I/O operations are most  
costly. In I/O operation we read data  
from disk to cache and write from  
cache to disk.

Two category of Optimizer

1) Rule based OLTP

2) Cost based Data warehouse

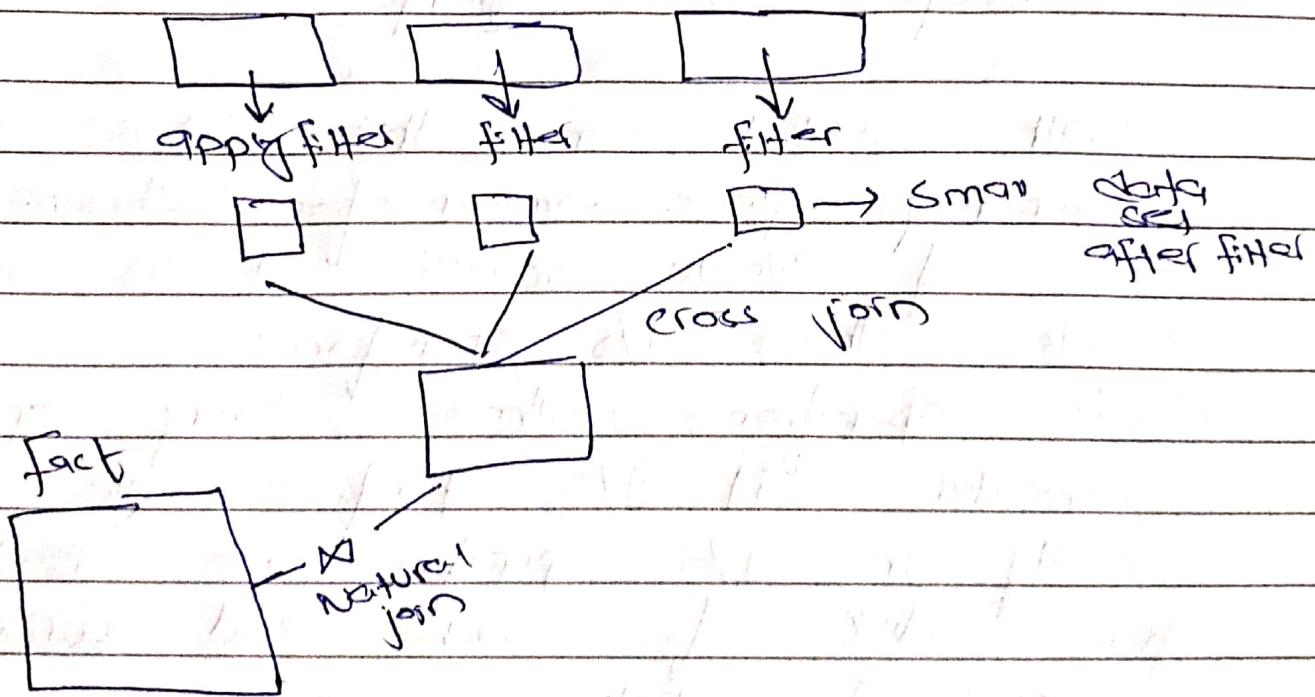
In OLTP Rule based is  
efficient. In this according to set

Date: \_\_\_\_\_

of define rule plan is executed.

In Data warehouse cost base is efficient. Cost base consider the other option like cost and plan can execute on the cost base of cost.

Note:- When we apply star join



This star join is the cost base optimizer. In Rule base we join the first with fact.

The work of optimizer is to get the query rewrite it

Date: \_\_\_\_\_

according to define paradigm. This functionality is build in the optimizer.

It is the work of designer to create the path and optimizer has to select the best path.

### DSS vs OLTP

DSS contain mostly the large dataset and it may consist of both OLTP and DSS. In this Merge or hash perform best.

OLTP contain the small dataset. In this mostly nested join perform better.

As the administrator or the user we can give hint to the optimizer to force optimizer to consider if and execute if

join with high selectivity is select first. Means when we join that join have the minimum No. of rows.

Date: \_\_\_\_\_

join with high selectivity is (select means low rows) join with low selectivity is not selected first (means high rows)

$\Rightarrow$  No. of rows and the distribution of the data table, distribution of the block and the no. of duplicated and null value in the column is predefine. All the statistics are precomputed in ETL we recollect the statistics when the major change done.

when we do the first join with high selectivity then second and third if poor then it does not much matter.

Block NEJ =  $O(R \times S)$   
for each qualifying block of outer table we have to read the inner table one time.  $O(R \times S)$

Basic NLJ =  $O(R \times S)$   
for each row we read

if we have index NLJ :-

$$R: S \\ 1: 100$$

(Index + Base)  
cost cost

$$3 \times (1 + 100) \quad \downarrow \text{average no. of rows}$$

if we have clustered NLJ :-

$$3 + (1 + 1)$$

If you place the file in order then there is  
a cost of maintenance when we remove  
or add file we have to maintain it.  
If you place the file randomly then there  
is not a cost of maintenance but  
difficulty in search.

In Block NLJ we cannot access or  
use index access path

Block factor in one is the number  
of rows by bfr. It is represented

Now we divide the

In clustered the index tabel give the index of the anchor block those have to block scan are correlated but we are all blocks.

If the index is on column and base of NLJ. the column then it is clustered.

It is not necessary that all the blocks are colocated. All blocks are correlated.

Example:-

index tabel is less than label of record.

DHW block size is greater than the OLTP block size

Combine selectivity :-

$$3.1 \times (6.1 \times (15-1) \times (128,000))$$

$$K = 10^6 \text{ blocks}$$

$$r_D = 128,000$$

(No. of record)

$$r_A = 256,000$$

$$R_{IS} = 16$$

$$B = 16 \text{ KB}$$

$$R_S = 256 \text{ bytes}$$

(Rec size)

$$R_A = 256 \text{ (Rec size)}$$

$$R_{IA} = 16 \text{ bytes}$$