

warehouse for sales and shipping)

DMQL (Data mining query language)  
Data warehouses can be defined using two  
language primitives — cube definition,  
dimension definition.

Syntax of cube definition

define cube <cube\_name> [<dimension\_list>] ! <measure\_list>

Syntax of dimension definition

define dimension <dimension\_name> as (<attribute\_or\_dimension\_list>)

Ex

## Star schema definition

define cube sales\_star [time, item, branch, location]:

dollars\_sold = sum(sales-in-dollars),  
units\_sold = count()

define dimension time as (time-key, day, day-of-week, month, quarter, year)

define dimension item as (item-key, item-name, brand, type, supplier-type)

define dimension branch as (branch-key, branch-name, branch-type)

define dimension location as (location-key, street, city, ~~area~~ state, country)

The define cube statement defines a datacube called sales\_star, which corresponds to the central sales fact table.

Define dimension statement is used to define each of the dimensions.

Ex Snowflake schema.

define cube as sales\_snowflake [time, item, branch, location]:  
dollars\_sold = sum(sales-in-dollars), units\_sold = count()

define dimension time as (time-key, day, day-of-week, month, quarter, year)

define dimension item as (item-key, item-name, brand, type, supplier-key, supplier-type)

define dimension branch as (branch-key, branch-name, branch-type)

define dimension location as (location-key, street, city, city-key, city, state, country)

### Example Fact constellation schema

define cube sales [time, item, branch, location]:  
dollars\_sold = sum(sales\_in\_dollars), units\_sold = count()

define dimension time as (time\_key, day, day-of-week,  
months, quarter, year)

define dimension item as (item\_key, item\_name, brand, type,  
supplier\_type)

define dimension branch as (branch\_key, branch\_name,  
branch\_type)

define dimension location as (location\_key, street, city,  
state, country)

define cube shipping [time, item, shipper, from\_location,  
to\_location]:  
dollars\_cost = sum(cost\_in\_dollars), units\_shipped = count()

define dimension time as time in cube sales

define dimension item as item in cube sales

define dimension shipper as (shipper\_key, shipper\_name,  
location as location in cube sales, shipper\_type)

define dimension from\_location as location in cube sales

define dimension to\_location as location in cube sales.

### Explanation

define cube statement is used to define data cubes for sales and shipping.

time, item and location dimensions of the sales cube are shared with the shipping cube.

so, define dimension time as time in cube sales.

## concept hierarchies

It defines a sequence of mappings from a set of low-level concepts to higher-level, more general concepts.

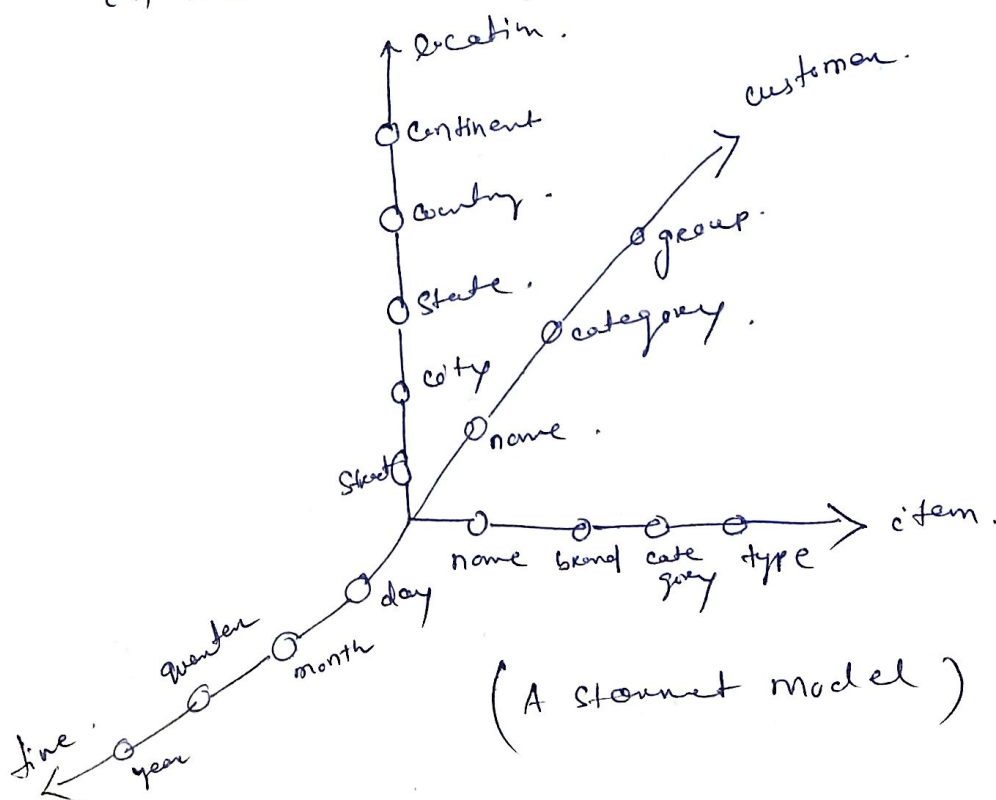
In the multidimensional model, data are organized into multiple dimensions, and each dimension contains multiple levels of abstraction defined by concept hierarchies.

This organization provides users with the flexibility to view data from different perspectives.

## Star and Snowflake query model

A star model consists of radial lines emanating from a central point, where each line represents a concept hierarchy for a dimension. Each abstraction level in the hierarchy is called a factoid.

customer



(A star model)

From multiple heterogeneous sources, analytics



Measures: Their categorization and computation

Measure  $\rightarrow$  is a numerical function that can be evaluated at each point in the data cube space.

A measure value is computed for a given point by aggregating the data corresponding to the respective dimension-value pairs defining the given points.

Measures can be divided into three categories

(1) Distributive

(2) Algebraic

(3) Holistic.

Distributive

An aggregate function is distributive if it can be computed in a distributed manner.

Suppose the data are partitioned into  $n$  sets. We apply the function to each partition, resulting in  $n$  aggregate values. If the result derived by applying the function to the  $n$  aggregate values is the same as that derived by applying the function to the entire dataset, the function can be computed in a distributed manner.

Ex  $\text{count}()$   
 $\text{sum}()$   
 $\text{min}()$   
 $\text{max}()$

For example, count() can be computed for a data cube by first partitioning the cube into a set of subcubes, computing  $\text{count}()$  for each subcube, and then summing up the counts obtained for each subcube.

## Algebraic

An aggregate function is algebraic if it can be computed by an algebraic function with  $M$  arguments, each of which is obtained by applying a distributive aggregate function.

$$\underline{\text{Ex}} \quad \text{avg}(C) = \frac{\text{sum}(C)}{\text{count}(C)}$$

where both  $\text{sum}(C)$  and  $\text{count}(C)$  are distributive aggregate functions.

Ex       $\text{standard\_deviation}(C)$ .  
          $\text{min}_N(C) \rightarrow$  find  $N$  minimum values.  
          $\text{max}_N(C) \rightarrow$  find  $N$  maximum values.

## Holistic

An aggregate function is holistic if there is no constant bound on the storage size needed to describe <sup>a self</sup> aggregate. That is there does not exist an algebraic function

Ex       $\text{median}(C)$       with  $M$  arguments that  
          $\text{mode}(C)$       do the computation.  
          $\text{rank}(C)$

## Concept Hierarchy

### Star schema model

The process of Data Warehouse Design

The design process consists of the following steps:

- ① Choose a business process to model, for example  
: orders, invoices, shipments, inventory etc.  
If the process is organizational and involve  
multiple complex object collections, a data  
warehouse model should be followed.  
If the process is departmental, a  
data mart model should be chosen.

② Choose the grain of the business process.  
The grain is the fundamental, atomic  
level of data to be represented in the  
fact table for this process.

③ Choose the dimensions that will apply  
to each fact table record. Ex: - time,  
item, customer, supplier etc.

④ Choose the measures that will populate  
each fact table record. Typical measures  
are numeric additive quantities like dollars-sold  
and units-sold.

A three-tier Data Warehouse Architecture.

- ① The bottom tier is a warehouse database  
server that is a relational database system.  
Back-end tools and utilities are used to feed  
data into the bottom tier from operational  
databases or other external sources.  
These tools and utilities perform data extraction,  
cleaning and transformation, as well as loading  
refresh functions to update the data warehouse.



if two relations R(RID, A) and S(SID, B) join on the attributes A and B, then the join index contains the pairs (RID, SID) where RID and SID are record identifiers from R and S respectively.

Quick searching in data cube  
 Bitmap indexing (Indexing OLAP Data).

It is an alternative representation of the record-ID list.

on the bitmap index for a given attribute, there is a distinct bit vector,  $B_v$ , for each value  $v$  in the domain of the attribute. If the domain of a given attribute consists of  $n$  values, then  $n$  bits are needed for each entry in the bitmap index.

if the attribute has the value  $v$  for a given row in the data table, then the bit representing that value is set to 1 in the corresponding row of the bitmap index. All other bits for that row are set to 0.

Base table

RID	item	city
R1	H	V
R2	C	V
R3	P	V
R4	S	V
R5	H	T
R6	C	T
R7	P	T
R8	S	T

item bitmap index table

RID	H	C	P	S
R1	1	0	0	0
R2	0	1	0	0
R3	0	0	1	0
R4	0	0	0	1
R5	1	0	0	0
R6	0	1	0	0
R7	0	0	1	0
R8	0	0	0	1

city bitmap index

RID	V	T
R1	1	0
R2	1	0
R3	1	0
R4	1	0
R5	0	1
R6	0	1
R7	0	1
R8	0	1

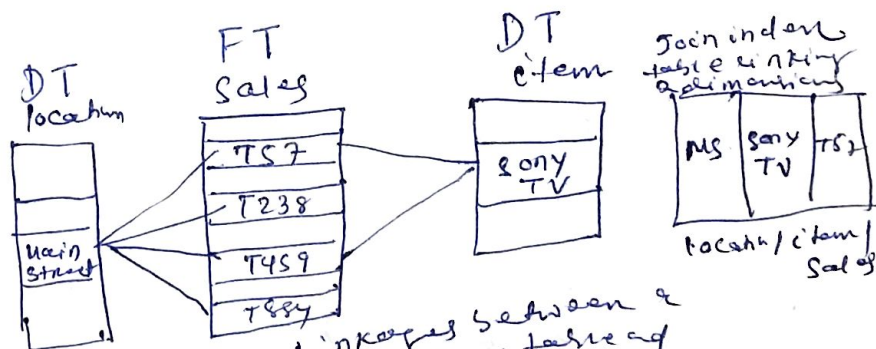
## Join indexing

join index table for location/sales

location	sales-key
Main Street	T57
-do-	T238
-do-	T459
-do-	T884

For item/sales

item	sales-key
Sony-TV	T57
-do-	T459



Linkages between a sales key and item in table.

The join index records can identify tuples without performing costly join operation. It is useful for many relationship between a foreign key and its matching primary keys.