

Module 1Data Mining & WarehousingNithin KumarAsst professorDept of CS&EVUCE, Mysuru* Data Warehouse :-

Data Warehousing provides Architecture & tools for business Executives to systematically organize, understand & use their data to make Strategic decisions.

→ A data Warehouse refers to a data repository that is maintained separately from an organization's operational data-bases.

Data Warehouse Systems allow for integration of variety of Application Systems.

→ William H. Inmon states "A data Warehouse is a Subject-oriented, Integrated, time-Variant & non-Volatile Collection of data in Support of Management's decision making process".

* Key features of Data Warehouse :- The key features of Data Warehouse are as follows.

→ "Subject - Oriented": Data Warehouse is organized around Major Subjects such as Customers, Suppliers, product & Sales. Rather than concentrating on day-to-day operations & Transaction processing of organization.

Data Warehouse focuses on Modeling & Analysis of Data for Decision Making.

→ "Integrated": A data Warehouse is usually constructed by Integrating Multiple Heterogeneous Sources such as Relational DB

Flat files & online Transaction Records.

"Data Cleaning" & "Data Integration" techniques are applied to Ensure Consistency in naming Conventions, Encoding Structures attributes, Measures & so on.

→ "Time - Variant": Data are Stored to provide Information from an historic perspective. for Ex 5-10 years.

Every key Structure in data Warehouse Contains, Either Implicitly or Explicitly a time Element.

→ "Non - volatile": A data Warehouse is always a physically Separate Store of data transformed from the Application data found in operational Environment.

Due to this separation, A data Warehouse does not require Transaction processing, Recovery & Concurrency Control Mechanism.

→ The Data Warehouse allows "Knowledge Workers" Such as Managers, Analysts & Executives to use the Warehouse to quickly & obtain an overview of the data, & make decisions based on Info in Warehouse.

The Construction of Data Warehouse Requires "Data Cleaning", "Data Integration" & "Data Consolidation".

→ Data Warehousing is very useful from the point of view of "Heterogeneous Database Integration".

In Traditional Database Approach to heterogeneous Database Integration is to build "Wrappers & Integrators" on top of Multiple Heterogeneous Databases.

→ When a Query is posed to a Client Site, a Metadata dictionary is used to Translate the Query into queries Appropriate for the Individual heterogeneous Sites Involved.

* "Query-Driven Approach" Requires complex information filtering & integration processes & Computers with local space for processing resources.

"Traditional Database" Make use of Query-Driven Approach.

* "Update-Driven Approach" In which information from multiple heterogeneous sources is integrated in advance & stored in a Warehouse for direct querying & Analysis.

"Data Warehousing" Make use of Update-Driven Approach.

* Differences between Operational Database System & Data Warehouse

The Major task of online operational database system is to perform online Transaction & Query processing. These systems are called "Online Transaction processing (OLTP)".

For Example, Most of day-to-day operations of an organization such as "purchasing", "Manufacturing", "Banking" etc.

→ "Online Analytical processing (OLAP)" serve users or knowledge workers in the role of data Analysis & Decision Making. Such systems can organize & present the data in various formats in order to accommodate the diverse needs of different users.

The Major distinguishing features of OLTP & OLAP are

* User & System Orientation :- An OLTP system is "Customer-oriented" and is used by clerks, clients, & IT professionals.

→ An OLAP system is "Market-oriented" and is used by knowledge workers such as Managers, Executives & Analysts.

* Data Contents :- OLTP system manages current data that, typically, are too detailed to be used.

→ OLAP System Manage large Amount historic data, provides facilities for Summarization & Aggregation & Store & Manage Information at different levels of granularity.

* Database Design :- OLTP usually adopts an Entity - Relationship (ER) data model & Appn oriented database design.

→ OLAP System adopts either "Star" or "Snowflake" data Models.

* View :- An OLTP System Focuses Mainly on Current data within an Enterprise or department, without referring to historic data or data in different organization.

→ OLAP Systems Span Multiple Version of database Schema, due to Evolutionary process of an organization.

OLAP S/m deal with Info that originates from different organization, Integrating Info from different sources.

* Access :- OLTP System consist mainly of Short atomic transaction so Concurrency Control & Recovery is required.

→ OLAP Systems are Mostly Read-only operations for any complex queries.

<u>Feature</u>	<u>OLTP</u>	<u>OLAP</u>
* <u>Characteristic</u>	* operational processing	* Informational processing
* <u>Orientation</u>	* Transaction	* Analysis
* <u>User</u>	* DBA, Clerk, Database professional	* Knowledge Workers
* <u>Function</u>	* Day-to-day operation	* long-term for Decision Support
* <u>DB-Design</u>	* ER-Model, Appn oriented	* Star/Snowflake, Subject oriented
* <u>View</u>	* Detailed, Flat Relational	* Summarized & Multi-dimensional
* <u>Access</u>	* Read / Write	* Read
* <u>Focus</u>	* Data-in	* Information-out
* <u>Unit of Work</u>	* Short, Simple Transaction	* Complex Query

* <u>No of Records Accessed</u>	* <u>Time</u>	* <u>MPI/Port</u>
* <u>No of Users</u>	* <u>Throughput</u>	* <u>Headreads</u>
* <u>DB Size</u>	* <u>OLB to high-order OLB</u>	* <u>≥ TB</u>
* <u>priority</u>	* <u>High performance & Availability</u>	* <u>High flexibility, End-user Autonomy</u>
* <u>Metric</u>	* <u>Transaction Throughput</u>	* <u>Query throughput & Response time</u>

* Data Warehousing : A Multitiered Architecture :-

Data warehouse often

adopts a 3-tier Architecture Such as

- * Bottom tier
- * Middle tier
- * Top tier

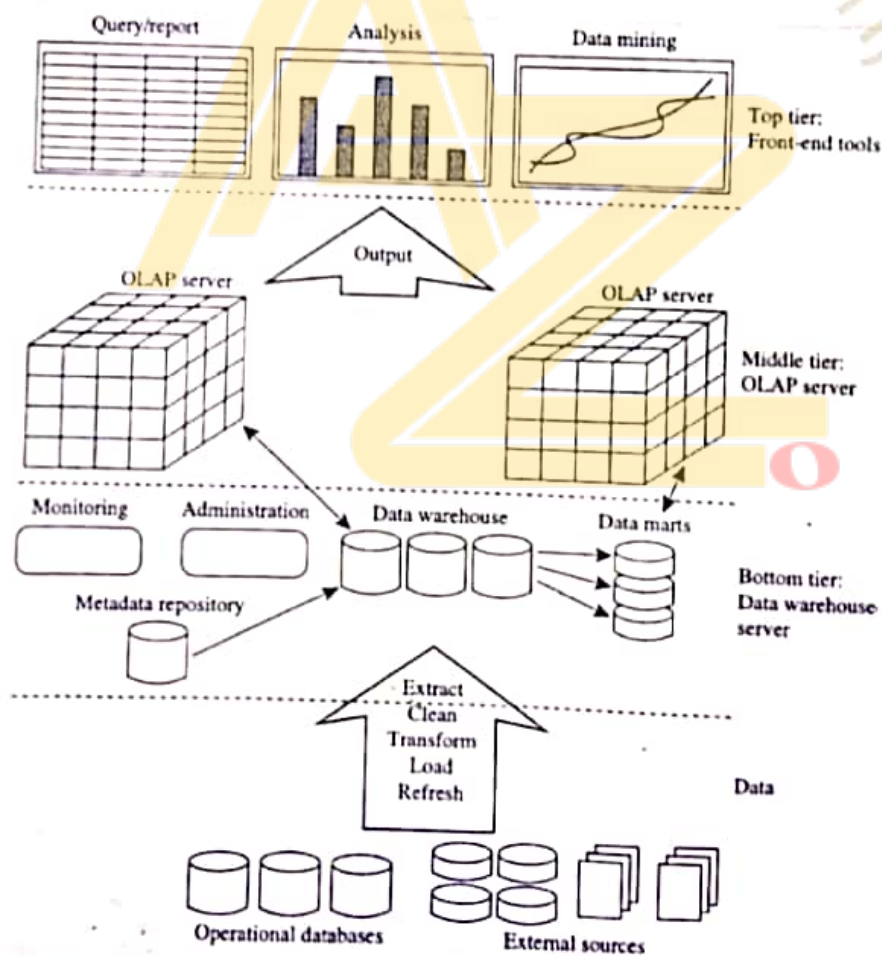


Fig:- Three - tier Datawarehouse Architecture

* Bottom tier :- Warehouse Database Server is almost a Relational database & other External Sources.
Back-End tools & utilities are used to perform Extraction, cleaning & Transformation (Ex:- to Merge Similar data from different Sources into a unified format).

→ The data are Extracted using Appln program Interface known as "gateways".
Example of gateway Include ODBC (Open Database Connection) & OLEDB (Object Linking & Embedding Database) by Microsoft & JDBC.
→ This tier also contains a Metadata Repository, which Stores Information about the Data Warehouse & Contents.

* Middle tier :- Consists of OLAP Server that is typically implemented using

→ "Relational OLAP (ROLAP)" :- Extended Relational DBMS that Maps Multidimensional data to Standard Relational-operation.

→ "Multidimensional OLAP (MOLAP)" :- Direct Implementation of Multidimensional data & operation.

* Top tier :- Consists of Front-End Client layer, which contains Query & reporting tools, Analysis tools & Data mining tools. Such as "Trend Analysis", "prediction" & so on.

* Data Warehouse Models :-

There are three Data Warehouse Models

* Enterprise Warehouse

* Data Mart

* Virtual Warehouse

* Enterprise Warehouse :- Collects all the Information about Subjects Spanning Entire organization.

It provides Corporate-Wide data with Corporate-Wide data - Integration.

→ Enterprise Warehouse Contains detailed data as well as Summarized data & can range in size from of few Gigabytes to hundreds of GB, Terabytes or beyond.

→ An Enterprise data Warehouse may be implemented on traditional Mainframes, Computer Super Servers or parallel Architectured platforms.

It requires Extensive Business Modeling & may take years to Design.

* Data Mart :- Contains a subset of Corporate-Wide data that is of Value to Specific Group of users.
For Ex: A Marketing Data Mart may confine its Subjects to Customers, Items & Sales.

→ Data Marts are usually implemented on "Low-Cost Department Servers" that are Unix/Linux or Windows based.

The Implementation Cycle of a Data Mart is more likely to be measured in Weeks rather than Months or Years.

→ Data Marts can be Categorized as

* Independent → Sourced from operational Systems or External Info providers.

* Dependent → Sourced Directly from Enterprise Data Warehouse

* Virtual Warehouse :- Is a set of Views over operational databases for Efficient Query processing.

A Virtual Warehouse is Easy to build but requires Extra Capacity on operational Database Servers.

→ For Ex: Air ticketing System, collects & displays business data relating to a Specific Moment in time, creating a Snapshot

of the Condition of business at that Moment.

* Recommended Approach for Warehouse development :-

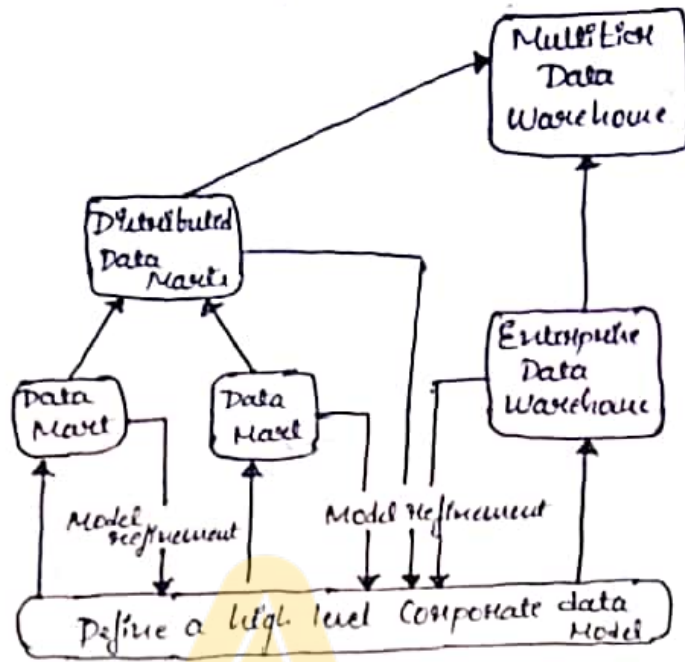


Fig: A Recommended Approach for Warehouse Development

A Recommended Method for the development of data Warehouse Systems is to implement the Warehouse in an Incremental & Evolutionary Manner as shown in above figure.

- Firstly, a high-level Corporate data Model is defined within a short period of time, that provides a Corporate-wide View of data among different Subjects & potential users.
- Second, Independent data Marts can be implemented in parallel with the Enterprise Warehouse based on the Corporate data Model.
- Third, Distributed data Marts can be constructed to integrate different data Marts via hub servers.

Finally a "Multi-Tier Data Warehouse" is constructed, where the Enterprise Warehouse is the sole Custodian of all Warehouse data, which is then distributed to various dependent Data Marts.

* Data Cube : A Multidimensional Data Model

A Data Cube allows data to be Modeled & Viewed In Multiple Dimensions. It is defined by "Facts" & "Dimensions".

* Dimension → are the perspectives on Entities with respect to which an organization wants to keep Records.

For Ex:- "All Electronics" May Create a "Sales" data Warehouse in order to keep Records of Stores Sales with respect to Dimensions "Time", "Item", "branch" & "Location".

* Facts → are Numeric Measures. Quantities by which we want to analyze Relationship between Dimensions.

For Ex:- Facts for a Sales data Warehouse Includes "Dollars-Sold", "Units-Sold", & "Amount-budgeted".

→ Consider a 2D Simple data Cube Example from "All Electronics" Sales data for Items Sold per Quarter in City Vancouver.

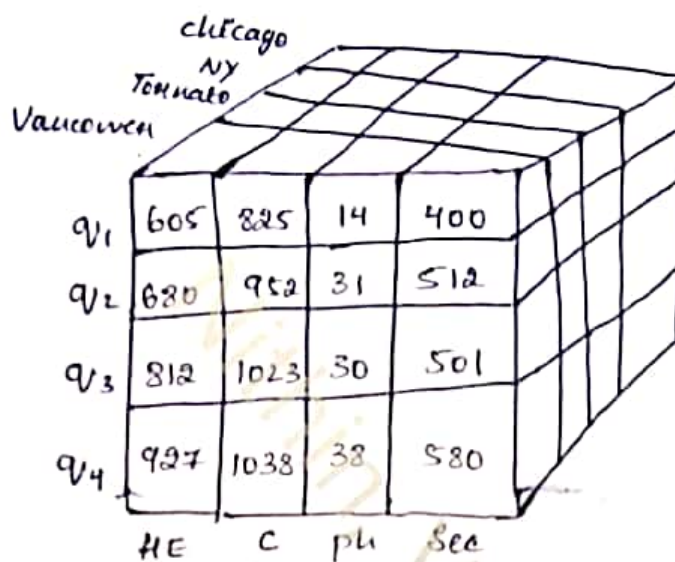
Location = "Vancouver"				
Time (Quarter)	Item (Type)			
	Home Ent	Computer	Phone	Security
Q ₁	605	825	14	400
Q ₂	630	952	31	512
Q ₃	812	1023	30	501
Q ₄	927	1038	38	580

In this 2-D Rep'n, the Sales for Vancouver are shown w.r.t time dimension & Item.

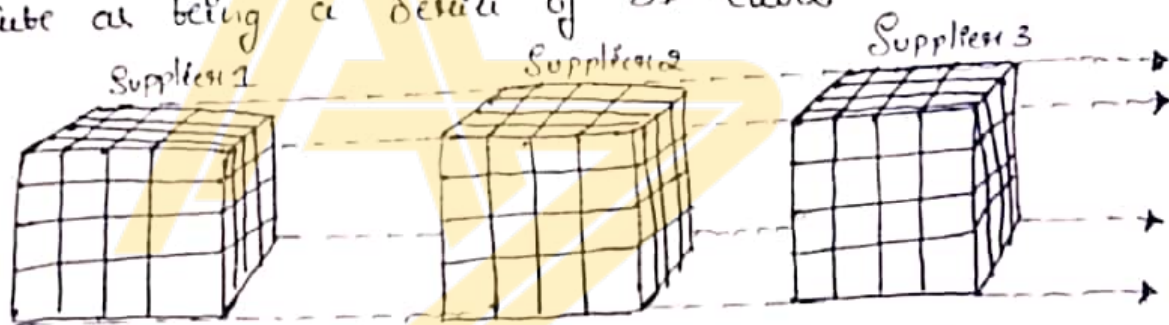
→ Suppose that we would like to view the Sales data with 3-D Rep'n the data according to Time, Item & Loc is as shown below

Time	Loc = "Chicago"				Loc = "New York"				Loc = "Toronto"				Loc = "Vancouver"			
	Item				Item				Item				Item			
	HE	C	Ph	Sec	HE	C	Ph	Sec	HE	C	Ph	Sec	HE	C	Ph	Sec
Q ₁	854	822	89	623	1037	720	64	620	1023	760	65	620	1000	620	65	600
Q ₂	943	890	64	698	1130	841	72	619	1024	820	66	619	1001	619	60	500
Q ₃	1032	924	59	789	1034	900	63	618	1025	823	67	720	1002	700	55	400
Q ₄	1129	992	63	876	1142	621	79	617	1026	824	68	723	1003	710	50	300

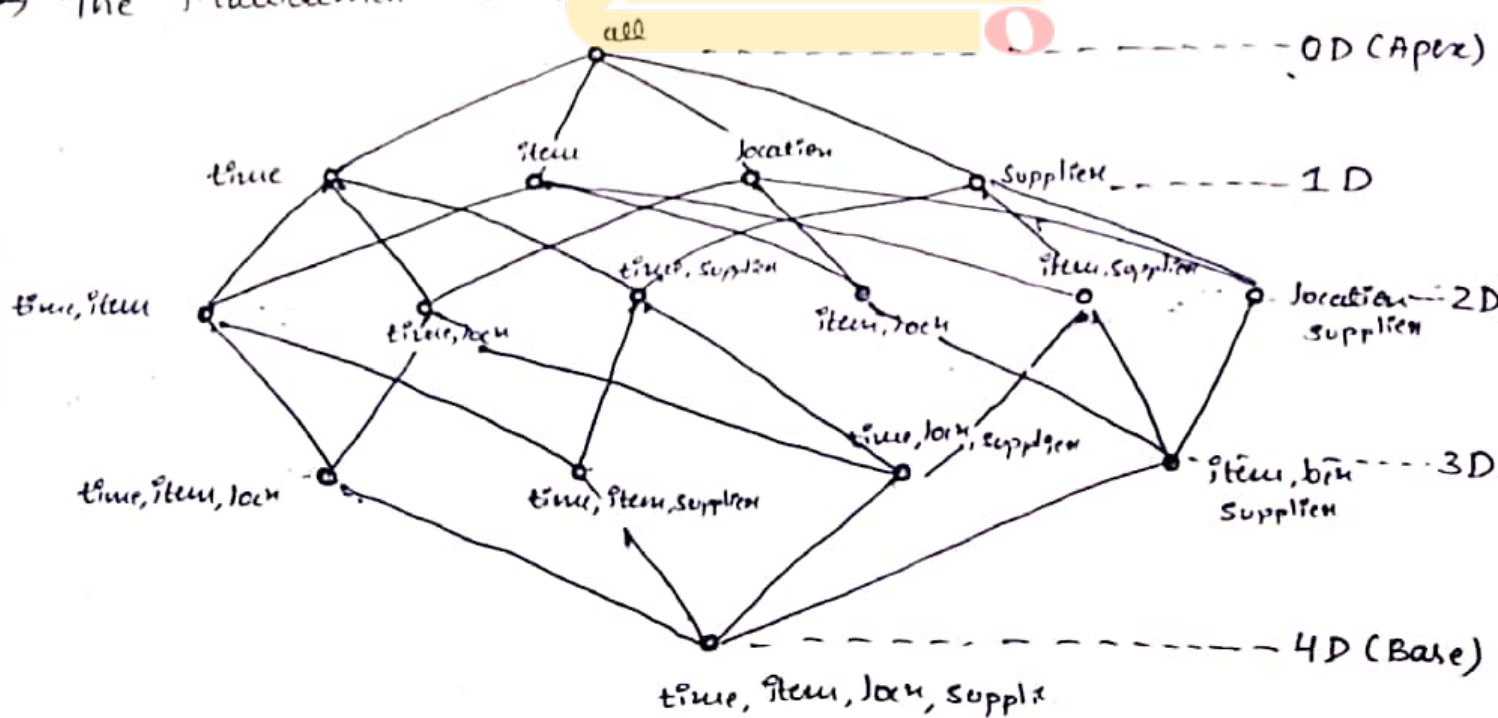
→ We May also represent the Same data in the form of 3D data Cube.



→ Suppose that we would like to view our Sales data with addition of 4th Dimension such that Supplier. Viewing things in 4-D cube becomes complex. However, we can think of a 4-D cube as being a series of 3D-Cubes.



→ The Multidimensional Data Model can be summarized as show



The OD Cuboid, which holds the highest level of Summary is called "Apex Cuboid".

In our Example, the "total Sales" or "dollar Sold", Summary-ed over all four dimensions.

* Star, Snowflake & Fact Constellations : Schema for Multidimensional Data Model

The Entity - Relationship data Model is commonly used to design the relational database, where database Schema consists of a set of Entities & the Relationships between them. Such data Model is appropriate for online Transaction processing.

→ The Most popular data Model for data Warehouse is a Multidimensional Model which can exist in the form of

- * Star Schema
- * Snowflake Schema
- * Fact Constellation Schema

* Star Schema :- Is the Most commonly used Multidimensional data Model, in which the data Warehouse contains

- * A large Central table (Fact table) containing the bulk data with no redundancy.
- * A set of smaller Attendant-tables (Dimension tables) one for each dimension.

The Schema Graph resembles a Star burst, with the Dimension tables displayed in radial pattern around the Central Fact table.

→ A Star Schema for "all Electronics" is shown in the following figure.

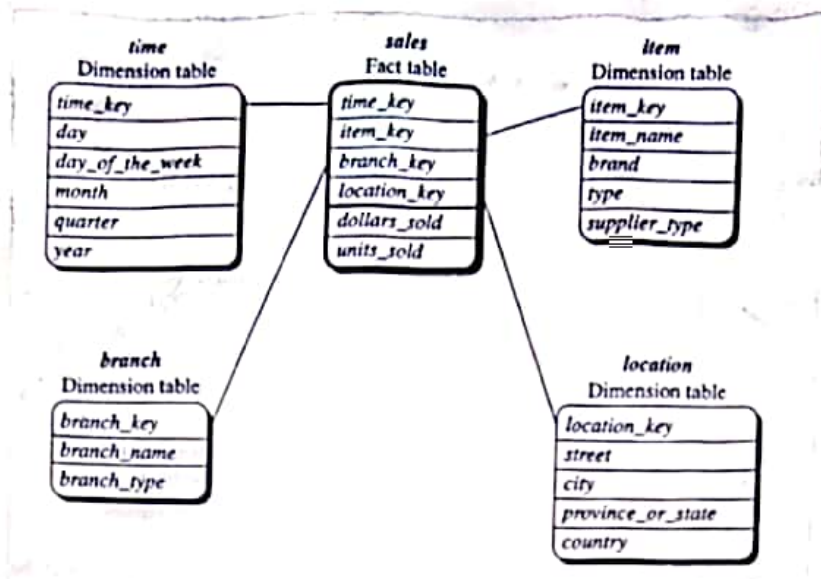


Fig:- Star Schema of Sales Data Warehouse

In the above Schema, Sales is considered along four Dimensions "Time, Item, branch & location". The Schema contains a Central Fact table for Sales that contains keys to each 4 Dimensions with two Measures "dollars_sold" & "units_sold".

To Minimize the Size of fact table, dimension-Identifiers (time_key, item_key etc) are System generated Identifiers.

→ In Star Schema, Each dimension is represented by only one table & Each table contains a Set of attributes.

For ex:-, the "Location" Dimension table contains the attribute Set [location_key, Street, City, State, Country]. This Constraint may introduce some Redundancy.

* Snowflake Schema:- Is a Variant of Star-Schema Model, where some dimension tables are Normalized, thereby splitting the data into additional tables.

→ The Major difference between the Star & Snowflake Schema Model is that the Dimension tables of the Snowflake Model may be kept in Normalized form to reduce redundancy. Such a table is Easy to Maintain & Save Storage Space.

→ The Snowflake Structure can reduce the Effectiveness of

Knowing, Since More Join will be needed to Execute a Query. The System performance May be adversely Impacted. Hence, although the Snowflake Schema Reduces Redundancy, It is not as popular as Star Schema in Data Warehouse Design.

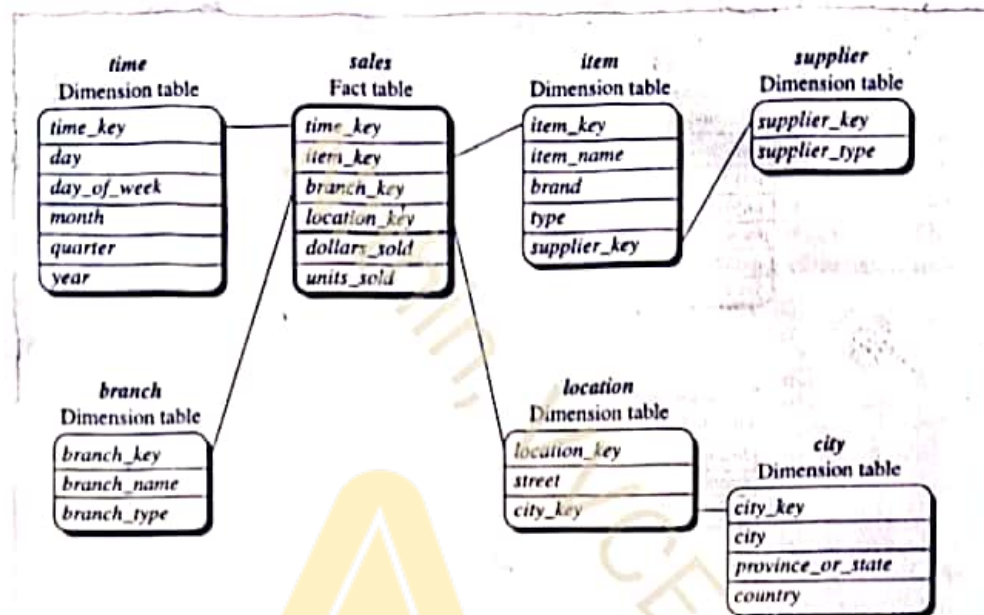


Fig:- Snowflake Schema of Sales

As Shown in above Snowflake Schema for "All Electronics", the Sales fact table is identical to the Star Schema.

→ The Main difference between Star & Snowflake Schema is Definition of Dimension table.

- * Dimension table of "Item" in Star Schema is Normalized in Snowflake Schema, resulting table form "Item" and "Supplier" Dimension tables.

- * Similarly Dimension table of "Location" in Star Schema is Normalized in Snowflake Schema, resulting table form "Location" & "City" Dimension tables.

* Fact Constellation :-

For Each Star Schema or Snowflake Schema, it is possible to Construct "Fact Constellation Schema".

→ The Fact Constellation Architecture Contains Multiple Fact tables that Share Many Dimensions.

→ Consider the following Example of Fact Constellation

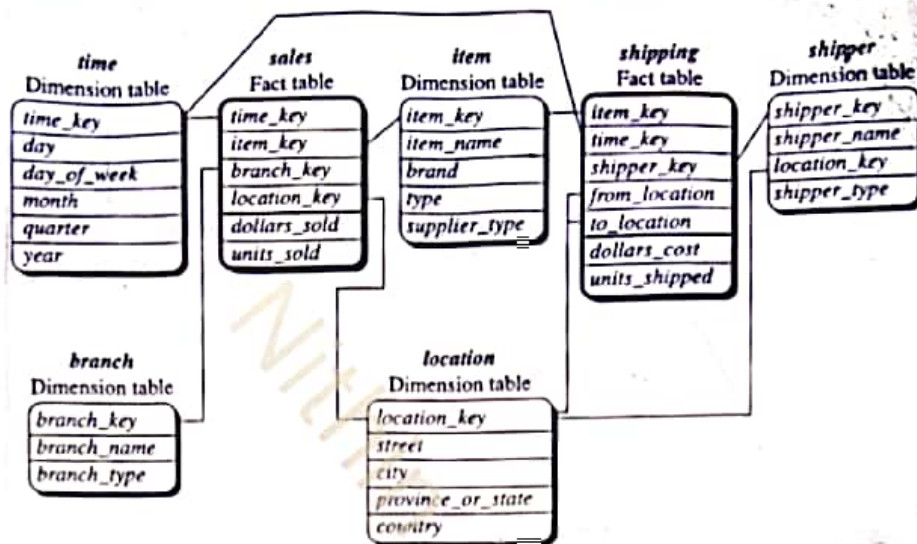


Fig. Fact Constellation Schema of Sales

The above Schema specifies two Fact tables "Sales" & "Shipping".

→ The Shipping table consists of 5 Dimensions on keys & Two Measures

* Dimensions → Item-key, Time-key, Shipper-key, to-location, From-Location

* Measures → Dollars-Cost, Units-Shipped.

A Fact Constellation Schema allows Dimension tables to be Shared between Fact tables.

→ A "Data Warehouse" Collects Information about Subjects that Span whole organization Such as Customers, Sales, Shipping Etc.

For data Warehouse, the fact Constellation Schema is commonly used to Model Multiple Interrelated Subjects.

→ A "Data Mart" is a Subset of Data Warehouse that Focus on Selected Subjects.

For Data Mart, the Star or Snowflake Schema is commonly used since both are used to Model Single Subjects.

* Role of Concept Hierarchy

A "Concept Hierarchy" defines a Sequence of Mapping from a Set of low-level concepts to High-level More general Concepts.

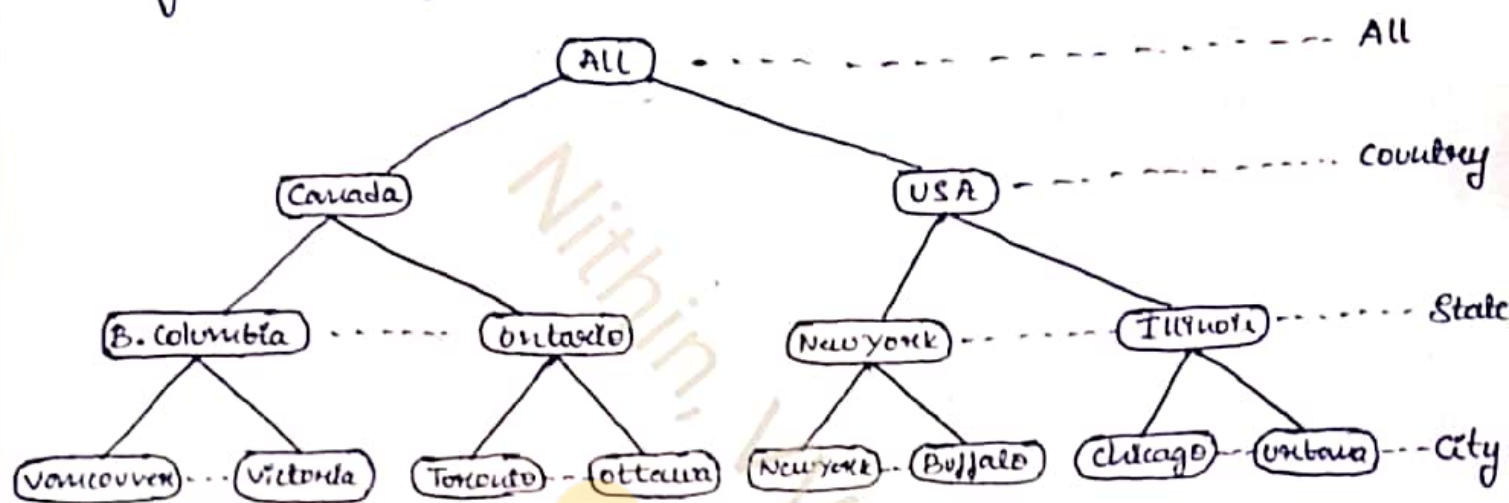


Fig:- Concept Hierarchy of Location

Mapping a Set of low-level concepts (Cities) to high-level More general Concepts (Countries).

- Many concept hierarchies are implicit within the database Schema
- Due to Space Limitation, not all of the hierarchy nodes are shown
- we can use lattice structure of attributes in Warehouse dimension.

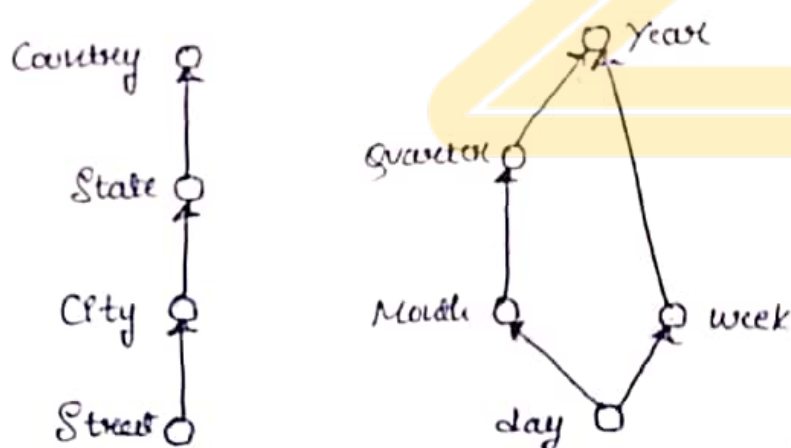


Fig:- Lattice Structure of Location & Time

→ Concept Hierarchy may also be defined by discretizing or grouping values for a given dimension or attribute, resulting in a "Set-Grouping Hierarchy".

For the dimension "price" with interval $(\$x \text{ --- } \$y)$ denote the range from $\$x$ (Exclusive) to $\$y$ (Inclusive)

→ There May be More than one Concept Hierarchy for a given attribute or Dimension based on different user View-points.

Concept Hierarchies May be provided Manually by System Users, Domain Experts or Knowledge Engineers or May be Automatically generated based on Statistical Analysis of the Data Distribution.

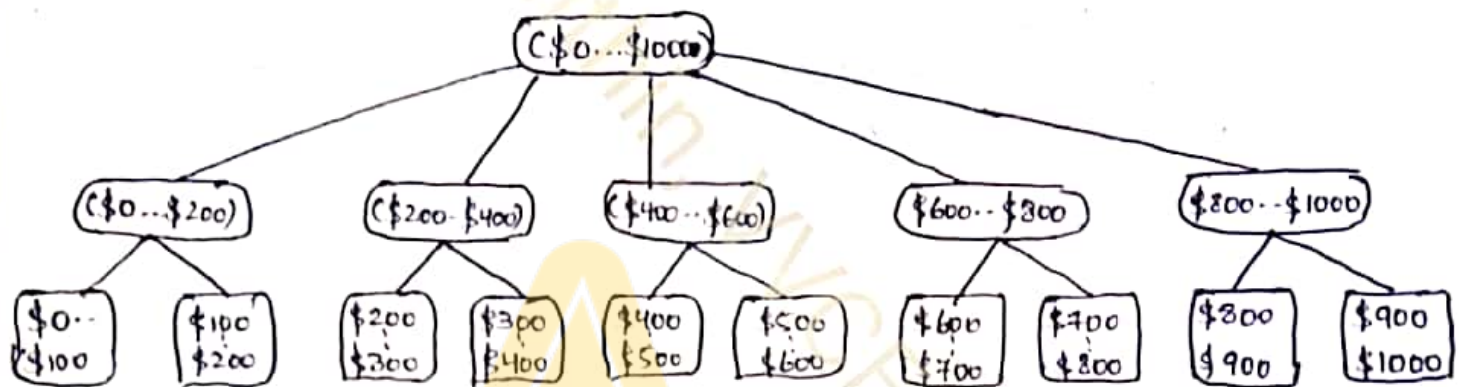


Fig:- Set Grouping Hierarchy of price

* Measures & Their Categorization & Computation:-

Measures can be

Organized into three categories Such as

* Distributive

* Algebraic

* Holistic

Based on the kind of Aggregate Functions used.

* 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 which results 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 Entire Data Set in Distributed Manner.

→ Some of good Examples of Distributive Aggregate Functions are

"Count()", "Min()", & "Max()".

* 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.

→ A Measure is "Algebraic" if it is obtained by Applying an Algebraic Aggregate Function.

For Ex:- "Avg()" can be computed by $\text{Sum}() / \text{Count}()$, where both $\text{Sum}()$ & $\text{Count}()$ are distributive Aggregate Functions.

* Holistic :- An Aggregate Function is Holistic, if there is no constant bound.

There does not exist any Algebraic function that characterizes the computation. Such Measures are called "Holistic".

→ Some of good Examples of Holistic Functions include "Median", "Mode()", & "Rank()".

More large data cube Applications require Efficient Computations of distributive & Algebraic Measures. Many Efficient techniques for this exist.

In contrast, it is difficult to compute Holistic Measures Efficiently.

* Typical OLAP operations :-

Number of operations may be applied to Data-Cubes are

* Roll-up

* Drill-Down

* Slice & Dice

* Pivot or Rotate

Consider the foll Example of "All Electronics" Sales Data Cube which

Contains the Dimensions [Location, Time & Item].

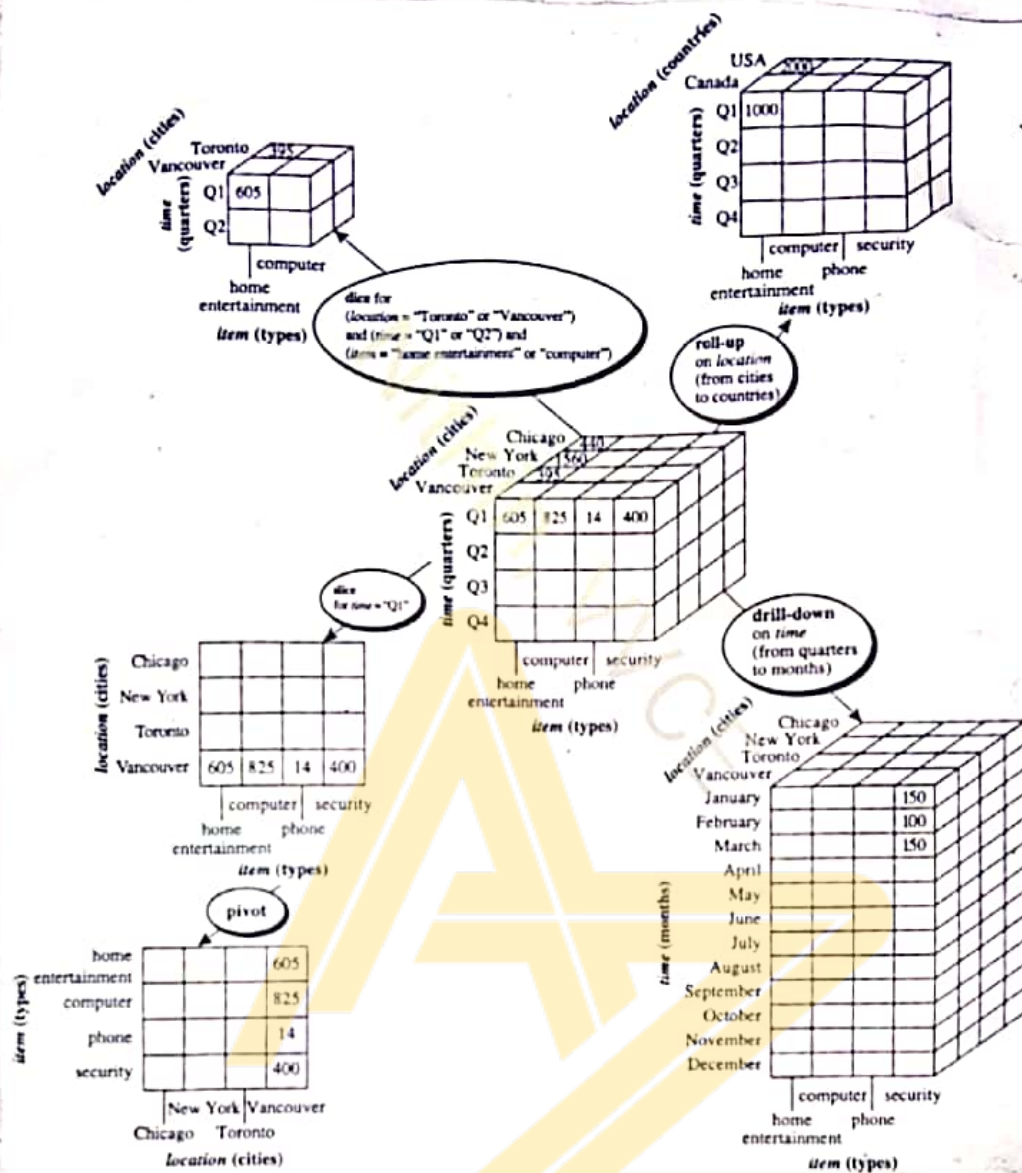


Fig: Typical OLAP operations

* Roll-up :- Is like "Zooming-out" on the Data Cube, It is required when the user needs further Abstraction or Less Detail.

→ This operation performs Further Aggregation on the Data.

For ex- Roll-up on Location from "cities" to "countries"

* Drill-Down :- Is like "Zooming-In" on the Data Cube, i.e. Reverse of Roll-up.

→ It is an Appropriate operation when the user needs further detail or when user wants to partition more finely on

Wants to focus on some particular values of certain dimension.
→ For ex, if the time in data cube is represented in quarter then it can be drilled down or zoomed-in the time to month wise.

* Slice & Dice :- Slice & Dice are operations for browsing the data in the cube. the term refers to the ability to look at information from different viewpoints.

* A "Slice" is a subset of cube corresponding to a single value for one or more members of dimension.

* A "Dice" operation is similar to slice but dicing does not involve reducing the no of dimension. A dice is obtained by performing a selection on two or more dimension.

* Pivot or Rotate :- Pivot operation is used when the user wishes to re-orient the view of data cube. It may involve swapping the row & column, or moving one of row dimension into column dimension.