# Unit 5 (DWDM)

By

Saurabh Saxena
Institute of Technology & Science,
Ghaziabad

# Data Aggregation

- Data aggregation is any process in which information is gathered and expressed in a summary form, for purposes such as statistical analysis.

- A common aggregation purpose is to get more information about particular groups based on specific variables such as age, profession, or income.

- The information about such groups can then be used for Web site personalization to choose content and advertising likely to appeal to an individual belonging to one or more groups for which data has been collected. For example, a site that sells music CDs might advertise certain CDs based on the age of the user and the data aggregate for their age group.

- Online analytic processing (OLAP) is a simple type of data aggregation in which the marketer uses an online reporting mechanism to process the information.

# Data Aggregation

- A common aggregation purpose is to get more information about particular groups based on specific variables such as age, profession, or income. The information about such groups can then be used for Web site personalization to choose content and advertising likely to appeal to an individual belonging to one or more groups for which data has been collected. For example, a site that sells music CDs might advertise certain CDs based on the age of the user and the data aggregate for their age group. Online analytic processing (OLAP) is a simple type of data aggregation in which the marketer uses an online reporting mechanism to process the information.

# OLAP Servers

- Online Analytical Processing Server (OLAP) is based on the multidimensional data model. It allows managers, and analysts to get an insight of the information through fast, consistent, and interactive access to information. This chapter cover the types of OLAP, operations on OLAP, difference between OLAP, and statistical databases and OLTP. Types of OLAP Servers

  We have four types of OLAP servers −

- Relational OLAP (ROLAP)

- Multidimensional OLAP (MOLAP)

- Hybrid OLAP (HOLAP)

- Specialized SQL Servers

# Relational OLAP

ROLAP servers are placed between relational back-end server and client front-end tools. To store and manage warehouse data, ROLAP uses relational or extended-relational DBMS.

ROLAP includes the following −

- Implementation of aggregation navigation logic.

- Optimization for each DBMS back end.

- Additional tools and services

# Multidimensional OLAP

- MOLAP uses array-based multidimensional storage engines for multidimensional views of data.

- With multidimensional data stores, the storage utilization may be low if the data set is sparse.

- Therefore, many MOLAP server use two levels of data storage representation to handle dense and sparse data sets.

# Hybrid OLAP

- Hybrid OLAP is a combination of both ROLAP and MOLAP.

- It offers higher scalability of ROLAP and faster computation of MOLAP.

- HOLAP servers allows to store the large data volumes of detailed information.

- The aggregations are stored separately in MOLAP store.

# OLAP Operations

Since OLAP servers are based on multidimensional view of data, we will discuss OLAP operations in multidimensional data. Here is the list of OLAP operations −

- Roll-up
- Drill-down
- Slice and dice
- Pivot (rotate)

# OLAP Operations

Since OLAP servers are based on multidimensional view of data, we will discuss OLAP operations in multidimensional data. Here is the list of OLAP operations −

- Roll-up
- Drill-down
- Slice and dice
- Pivot (rotate)

# Roll-up

Roll-up performs aggregation on a data cube in any of the following ways −

- By climbing up a concept hierarchy for a dimension

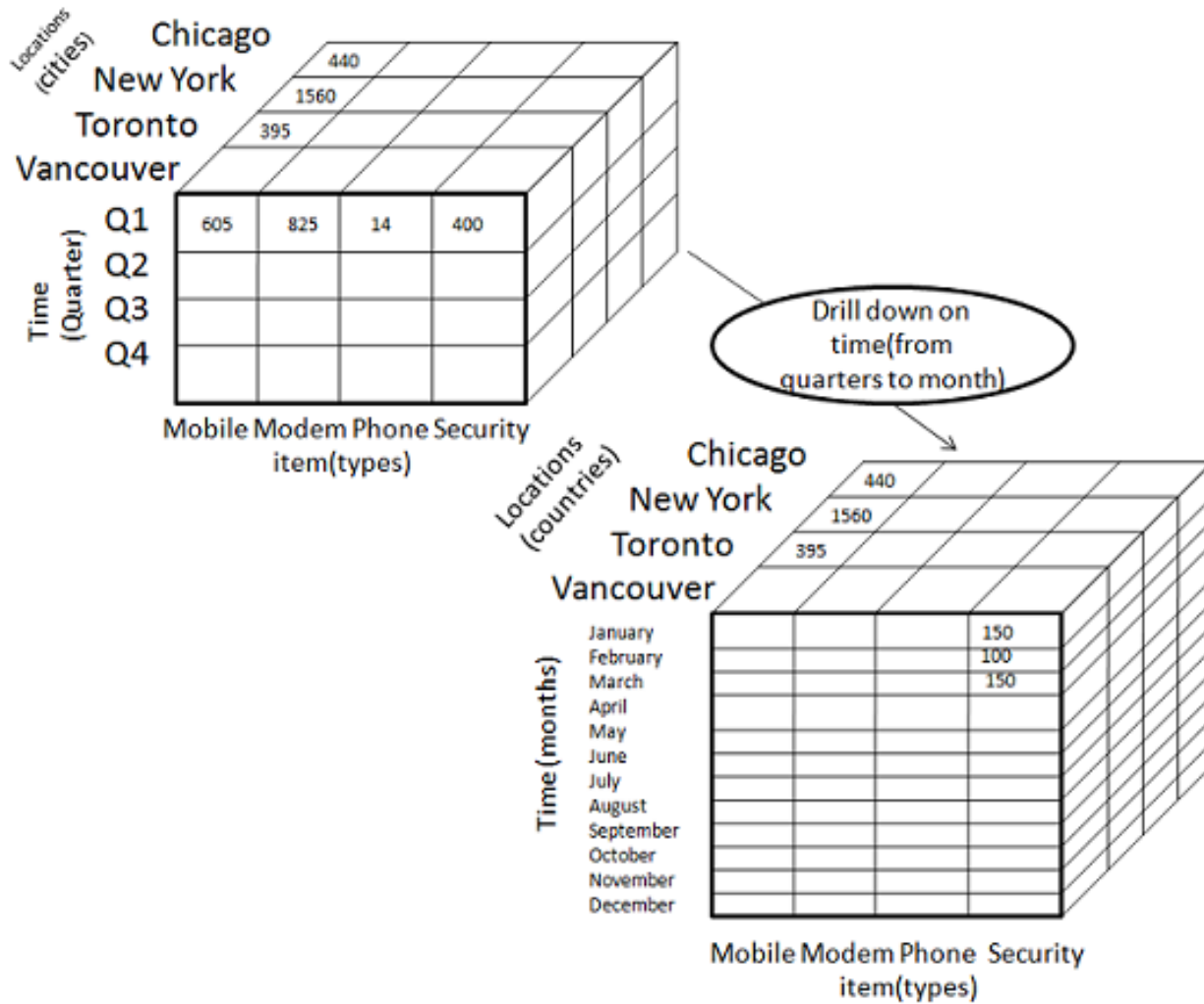- By dimension reduction

# Roll-up

# Roll-up

- Roll-up is performed by climbing up a concept hierarchy for the dimension location.
- Initially the concept hierarchy was "street < city < province < country".
- On rolling up, the data is aggregated by ascending the location hierarchy from the level of city to the level of country.
- The data is grouped into cities rather than countries.
- When roll-up is performed, one or more dimensions from the data cube are removed.

# Drill-down

Drill-down is the reverse operation of roll-up. It is performed by either of the following ways –

- By stepping down a concept hierarchy for a dimension

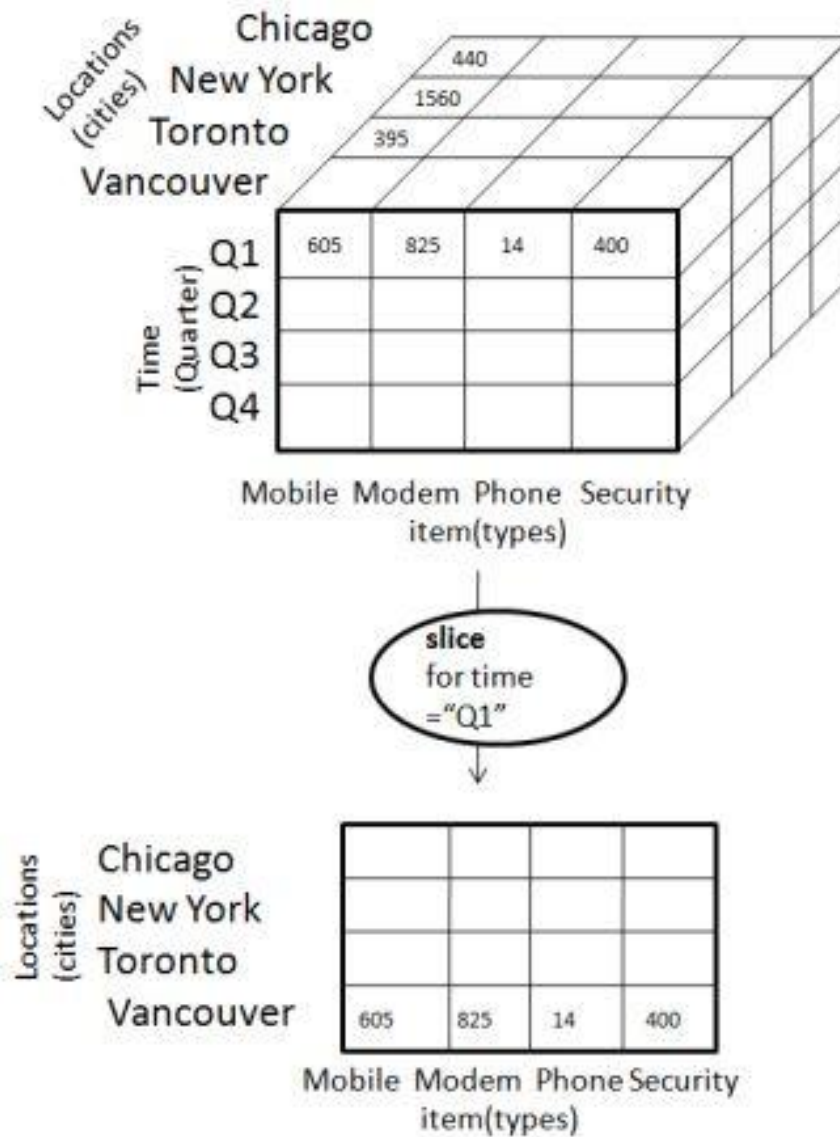- By introducing a new dimension.

# Drill-down

# Drill-down

- Drill-down is performed by stepping down a concept hierarchy for the dimension time.

- Initially the concept hierarchy was "day < month < quarter < year."

- On drilling down, the time dimension is descended from the level of quarter to the level of month.

- When drill-down is performed, one or more dimensions from the data cube are added.

- It navigates the data from less detailed data to highly detailed data.

# Slice

The slice operation selects one particular dimension from a given cube and provides a new sub-cube.

- Here Slice is performed for the dimension "time" using the criterion time = "Q1".
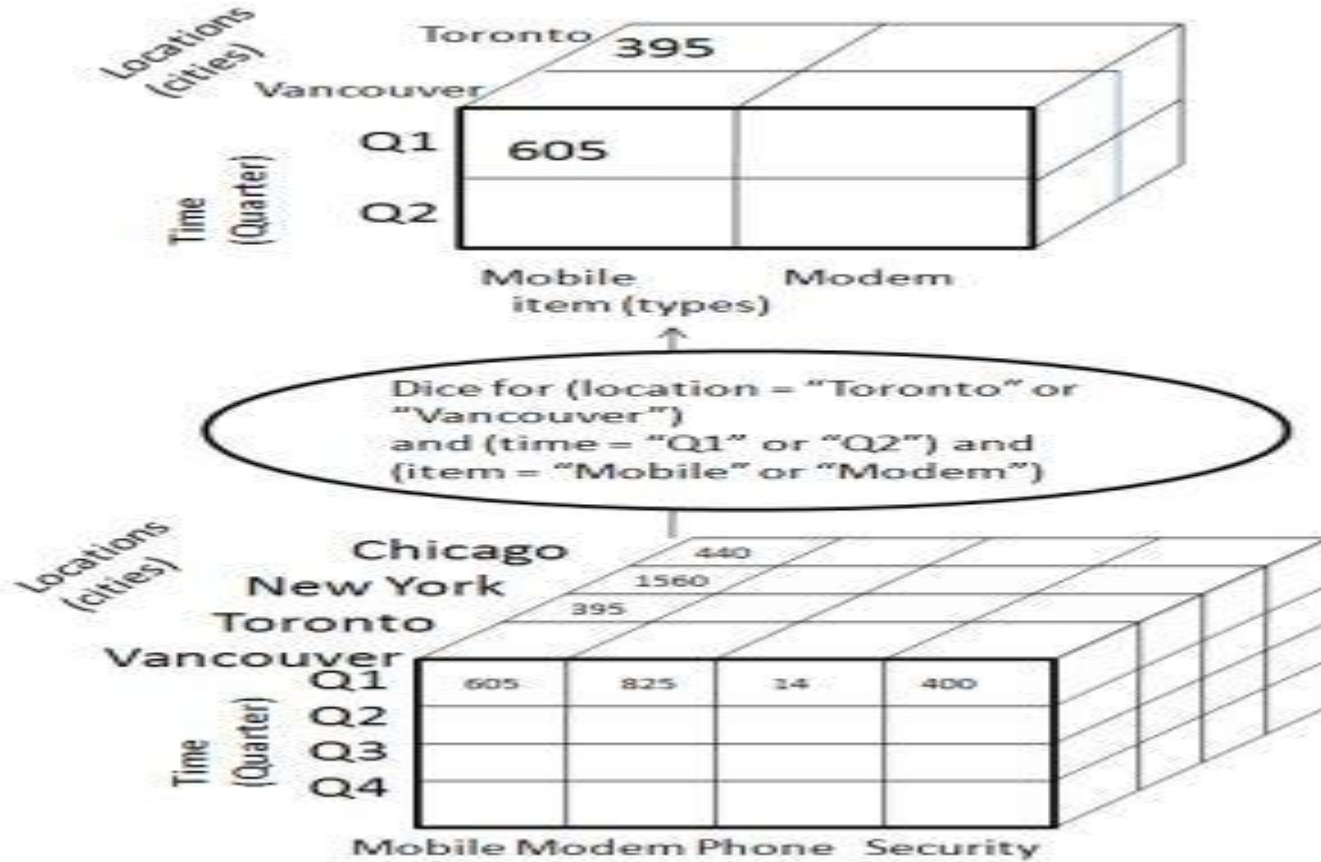- It will form a new sub-cube by selecting one or more dimensions

# Slice

# Dice

Dice selects two or more dimensions from a given cube and provides a new sub-cube:

# Dice

# Dice

- The dice operation on the cube based on the following selection criteria involves three dimensions.
- (location = "Toronto" or "Vancouver")
- (time = "Q1" or "Q2")
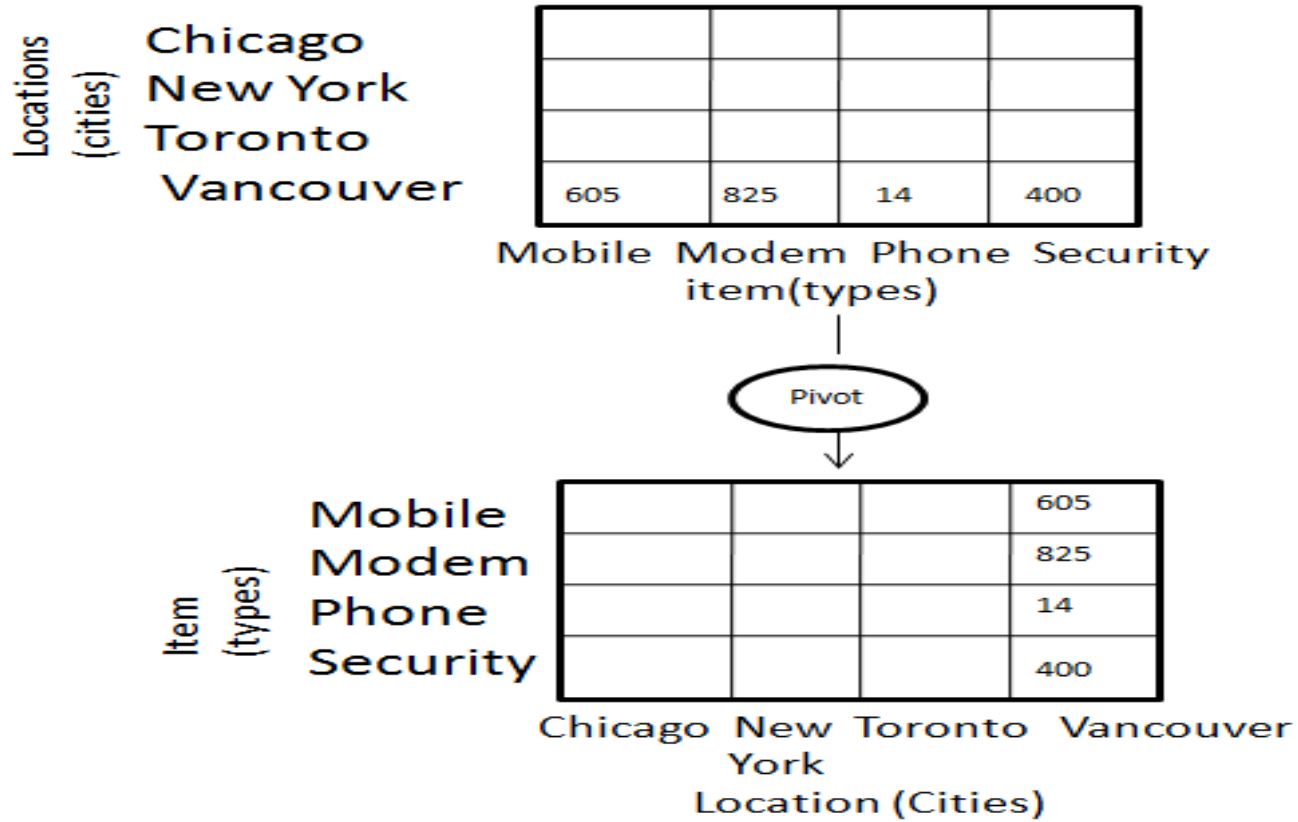- (item =" Mobile" or "Modem")

# Dice

- The dice operation on the cube based on the following selection criteria involves three dimensions.
- (location = "Toronto" or "Vancouver")
- (time = "Q1" or "Q2")
- (item =" Mobile" or "Modem")

# Pivot

- The pivot operation is also known as rotation. It rotates the data axes in view in order to provide an alternative presentation of data.

# Pivot

# Data Warehousing - Tuning

➤ A data warehouse keeps evolving and it is unpredictable what query the user is going to post in the future.

➤ Therefore it becomes more difficult to tune a data warehouse system.

➤ So it is needed to tune the different aspects of a data warehouse such as performance, data load, queries, etc.

# Difficulties in Data Warehouse Tuning

Tuning a data warehouse is a difficult procedure due to following reasons −

- Data warehouse is dynamic; it never remains constant.
- It is very difficult to predict what query the user is going to post in the future.
- Business requirements change with time.
- Users and their profiles keep changing.
- The user can switch from one group to another.
- The data load on the warehouse also changes with time.

# Performance Assessment

Here is a list of objective measures of performance −

- Average query response time
- Scan rates
- Time used per day query
- Memory usage per process
- I/O throughput rates

# Performance Assessment

Following are the points to remember.

- It is necessary to specify the measures in service level agreement (SLA).

- It is of no use trying to tune response time, if they are already better than those required.

- It is essential to have realistic expectations while making performance assessment.

- It is also essential that the users have feasible expectations.

- To hide the complexity of the system from the user, aggregations and views should be used.

- It is also possible that the user can write a query you had not tuned for.

# Data Warehousing Tuning : Data Load Tuning

Data load is a critical part of overnight processing. Nothing else can run until data load is complete. This is the entry point into the system.

- The very common approach is to insert data using the **SQL Layer**. In this approach, normal checks and constraints need to be performed. When the data is inserted into the table, the code will run to check for enough space to insert the data. If sufficient space is not available, then more space may have to be allocated to these tables. These checks take time to perform and are costly to CPU.

- The second approach is to bypass all these checks and constraints and place the data directly into the preformatted blocks. These blocks are later written to the database. It is faster than the first approach, but it can work only with whole blocks of data. This can lead to some space wastage.

# Data Warehousing Tuning : Data Load Tuning

- The third approach is that while loading the data into the table that already contains the table, we can maintain indexes.

- The fourth approach says that to load the data in tables that already contain data, **drop the indexes & recreate them** when the data load is complete. The choice between the third and the fourth approach depends on how much data is already loaded and how many indexes need to be rebuilt.

# Integrity Checks

Integrity checking highly affects the performance of the load. Following are the points to remember −

- Integrity checks need to be limited because they require heavy processing power.

- Integrity checks should be applied on the source system to avoid performance degrade of data load.

# Tuning Queries

We have two kinds of queries in data warehouse –

- Fixed queries
- Ad hoc queries

# Tuning Queries : Fixed queries

Fixed queries are well defined. Following are the examples of fixed queries −

- Regular reports
- Canned queries
- Common aggregations

Tuning the fixed queries in a data warehouse is same as in a relational database system. The only difference is that the amount of data to be queried may be different. It is good to store the most successful execution plan while testing fixed queries. Storing these executing plan will allow us to spot changing data size and data skew, as it will cause the execution plan to change.

# Tuning Queries : Ad hoc Queries

To understand ad hoc queries, it is important to know the ad hoc users of the data warehouse. For each user or group of users, you need to know the following −

- The number of users in the group
- Whether they use ad hoc queries at regular intervals of time
- Whether they use ad hoc queries frequently
- Whether they use ad hoc queries occasionally at unknown intervals.
- The maximum size of query they tend to run
- The average size of query they tend to run
- Whether they require drill-down access to the base data
- The elapsed login time per day
- The peak time of daily usage
- The number of queries they run per peak hour

# Tuning Queries : Ad hoc Queries

**Points to Note**

- It is important to track the user's profiles and identify the queries that are run on a regular basis.

- It is also important that the tuning performed does not affect the performance.

- Identify similar and ad hoc queries that are frequently run.

- If these queries are identified, then the database will change and new indexes can be added for those queries.

- If these queries are identified, then new aggregations can be created specifically for those queries that would result in their efficient execution.

# Data Warehousing - Testing

Testing is very important for data warehouse systems to make them work correctly and efficiently. There are three basic levels of testing performed on a data warehouse –

- Unit testing
- Integration testing
- System testing

# Data Warehousing - Testing

Unit testing:

- In unit testing, each component is separately tested.

- Each module, i.e., procedure, program, SQL Script, Unix shell is tested.

- This test is performed by the developer.

# Data Warehousing - Testing

Integration Testing

- In integration testing, the various modules of the application are brought together and then tested against the number of inputs.

- It is performed to test whether the various components do well after integration.

# Data Warehousing - Testing

System Testing

- In system testing, the whole data warehouse application is tested together.

- The purpose of system testing is to check whether the entire system works correctly together or not.

- System testing is performed by the testing team.

- Since the size of the whole data warehouse is very large, it is usually possible to perform minimal system testing before the test plan can be enacted.

# Data Warehousing - Testing

Test Schedule

- First of all, the test schedule is created in the process of developing the test plan. In this schedule, we predict the estimated time required for the testing of the entire data warehouse system.

- There are different methodologies available to create a test schedule, but none of them are perfect because the data warehouse is very complex and large. Also the data warehouse system is evolving in nature. One may face the following issues while creating a test schedule −

- A simple problem may have a large size of query that can take a day or more to complete, i.e., the query does not complete in a desired time scale.

- There may be hardware failures such as losing a disk or human errors such as accidentally deleting a table or overwriting a large table.

# Data Warehousing - Testing

Testing Backup Recovery

Testing the backup recovery strategy is extremely important. Here is the list of scenarios for which this testing is needed –

- Media failure
- Loss or damage of table space or data file
- Loss or damage of redo log file
- Loss or damage of control file
- Instance failure
- Loss or damage of archive file
- Loss or damage of table
- Failure during data failure
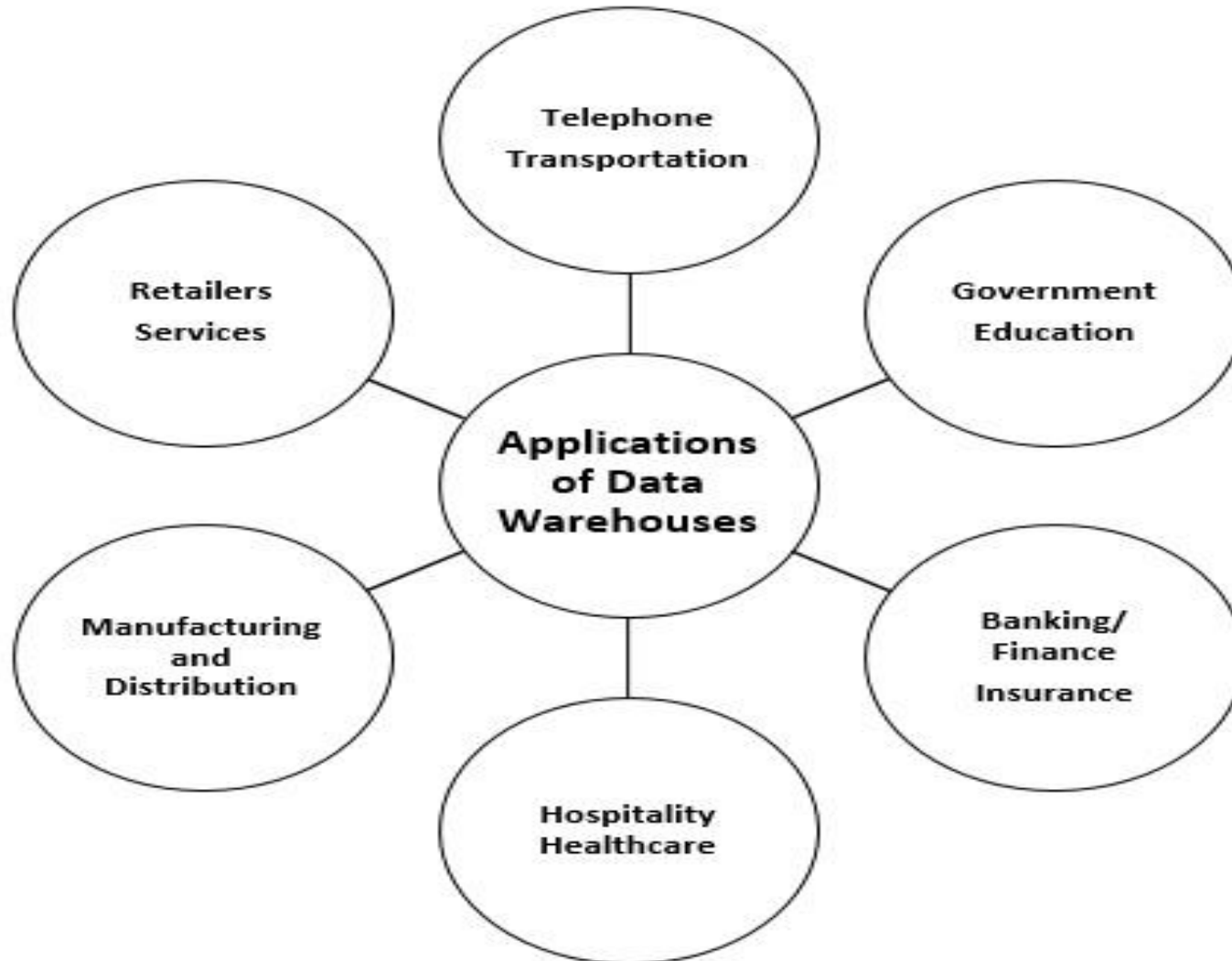
# Data Warehousing - Testing

Testing Operational Environment

- There are a number of aspects that need to be tested. These aspects are listed below.

- **Security** – A separate security document is required for security testing. This document contains a list of disallowed operations and devising tests for each.

- **Scheduler** – Scheduling software is required to control the daily operations of a data warehouse. It needs to be tested during system testing. The scheduling software requires an interface with the data warehouse, which will need the scheduler to control overnight processing and the management of aggregations.

# Data Warehousing - Testing

- **Disk Configuration.** – Disk configuration also needs to be tested to identify I/O bottlenecks. The test should be performed with multiple times with different settings.

- **Management Tools.** – It is required to test all the management tools during system testing. Here is the list of tools that need to be tested.
  - Event manager
  - System manager
  - Database manager
  - Configuration manager
  - Backup recovery manager

# Applications of Data Warehouse

# Applications of Data Warehouse

**1.Banking Industry**

- In the banking industry, concentration is given to risk management and policy reversal as well analyzing consumer data, market trends, government regulations and reports, and more importantly financial decision making.

**2. Finance Industry**

- Similar to the applications seen in banking, mainly revolve around evaluation and trends of customer expenses which aids in maximizing the profits earned by their clients.

**3. Consumer Goods Industry**

- They are used for prediction of consumer trends, inventory management, market and advertising research. In-depth analysis of sales and production is also carried out. Apart from these, information is exchanged business partners and clientele.

# Applications of Data Warehouse

**4. Government and Education**

- The federal government utilizes the warehouses for research in compliance, whereas the state government uses it for services related to human resources like recruitment, and accounting like payroll management.

**5. Healthcare**

- One of the most important sector which utilizes data warehouses is the Healthcare sector. All of their financial, clinical, and employee records are fed to warehouses as it helps them to strategize and predict outcomes, track and analyze their service feedback, generate patient reports, share data with tie-in insurance companies, medical aid services, etc.

# Applications of Data Warehouse

**6. Hospitality Industry**

- A major proportion of this industry is dominated by hotel and restaurant services, car rental services, and holiday home services. They utilize warehouse services to design and evaluate their advertising and promotion campaigns where they target customers based on their feedback and travel patterns.

**7. Insurance**

- As the saying goes in the insurance services sector, "Insurance can never be bought, it can be only be sold", the warehouses are primarily used to analyze data patterns and customer trends, apart from maintaining records of already existing participants. The design of tailor-made customer offers and promotions is also possible through warehouses.

# Applications of Data Warehouse

**8. Hospitality Industry**

- A major proportion of this industry is dominated by hotel and restaurant services, car rental services, and holiday home services. They utilize warehouse services to design and evaluate their advertising and promotion campaigns where they target customers based on their feedback and travel patterns.

**9. Insurance**

- As the saying goes in the insurance services sector, "Insurance can never be bought, it can be only be sold", the warehouses are primarily used to analyze data patterns and customer trends, apart from maintaining records of already existing participants. The design of tailor-made customer offers and promotions is also possible through warehouses.

# Applications of Data Warehouse

**10. Manufacturing and Distribution Industry**

- This industry is one of the most important sources of income for any state. A manufacturing organization has to take several make-or-buy decisions which can influence the future of the sector, which is why they utilize high-end OLAP tools as a part of data warehouses to predict market changes, analyze current business trends, detect warning conditions, view marketing developments, and ultimately take better decisions.

**11. The Retailers**

- Retailers serve as middlemen between producers and consumers. It is important for them to maintain records of both the parties to ensure their existence in the market.

- They use warehouses to track items, their advertising promotions, and the consumers buying trends. They also analyze sales to determine fast selling and slow selling product lines and determine their shelf space through a process of elimination.

# Applications of Data Warehouse

**12. Services Sector**

- Data warehouses find themselves to be of use in the service sector for maintenance of financial records, revenue patterns, customer profiling, resource management, and human resources.

**13. Telephone Industry**

- The telephone industry operates over both offline and online data burdening them with a lot of historical data which has to be consolidated and integrated.

# Spatial Data Mining

- Spatial data mining is the application of data mining to spatial models.

- In spatial data mining, analysts use geographical or spatial information to produce business intelligence or other results.

- This requires specific techniques and resources to get the geographical data into relevant and useful formats.

# Temporal Mining

- Temporal Data Mining is a single step in the process of Knowledge Discovery in Temporal Databases that enumerates structures (temporal patterns or models) over the temporal data, and any algorithm that enumerates temporal patterns from, or fits models to, temporal data is a Temporal Data Mining Algorithm.