

Data Warehousing Concepts



Copyright © 2011 IGATE Corporation. All rights reserved. No part of this publication shall be reproduced in any way, including but not limited to photocopy, photographic, magnetic, or other record, without the prior written permission of IGATE Corporation.

IGATE Corporation considers information included in this document to be Confidential and Proprietary.

## Data Warehousing Concepts

### Document History

Date	Course Version No.	Software Version No.	Developer / SME	Change Record Remarks
	0.1D	NA		Content Creation
Jan-2009	0.1	NA	BI CDI team	Review
16-Apr-2009	2.0	NA	Priya Rane	Material Revamp
04-Feb-2010		NA	CLS Team	Review
31-July-2012		NA	Coordinators	Change to Igate Format

June 15, 2016

Proprietary and Confidential

- 2 -



### Course Goals and Non Goals

#### ➤ Course Goals

- At the end of this program, participants gain an understanding of basic concepts in Data warehousing.



#### ➤ Course Non Goals

- Implementation of dimensional modeling is not the part of this course.

### Pre-requisites

- Fair knowledge of Database

### Intended Audience

- Software Engineers and Senior Software Engineers



### Day Wise Schedule

#### ➤ Day 1

- Lesson 1: Business Intelligence
- Lesson 2: General concept of Data Warehouse
- Lesson 3: Dimensional modeling
- Lesson 4: ETL and Metadata
- Lesson 5: Online Analytical Processing (OLAP)
- Lesson 6: Data Mining
- Lesson 7: Case Studies

### Table of Contents

- **Lesson 1: Business Intelligence**
  - 1.1: Business Intelligence
  - 1.2: Need for Business Intelligence
  - 1.3: Terms used in BI
  - 1.4: Components of BI
- **Lesson 2: General concept of Data Warehouse**
  - 2.1: Data Warehouse
  - 2.2: Characteristics of Data Warehouse

### Table of Contents

- 2.3: Need for Data Warehouse
- 2.4: Data Warehouse Architecture
- 2.5: Features of Data warehouse
- 2.6: Data Mart
- 2.7: Application Areas

#### ➤ Lesson 3: Dimension modeling basic concepts

- 3.1: Dimension modeling
- 3.2: Fact and Dimension tables
- 3.3: Database schema
- 3.4: Schema Design for Modeling

### Table of Contents

- **Lesson 4: ETL and Metadata**
  - 4.1: ETL process
  - 4.2: Metadata used in ETL
  - 4.3: Metadata in Data Warehousing
  - 4.4: Simple Data warehouse model
- **Lesson 5: Online Analytical Processing (OLAP)**
  - 5.1: Online Analytical Processing (OLAP)
  - 5.2: Nature of OLAP analysis
  - 5.3: Types of OLAP

### Table of Contents

- 5.4: OLAP Tools
- 5.5: OLTP and OLAP
- 5.6: Operational versus Informational System

#### ➤ Lesson 6: Data Mining

- 6.1: Data mining
- 6.2: The Knowledge Discovery process
- 6.3: Need of Data Mining
- 6.4: Use of Data mining
- 6.5: Data mining and Business Intelligence

### Table of Contents

- 6.6: Types of data used in Data mining
- 6.7: Data Mining applications
- 6.8: Data Mining products
- 6.9: Data Mining market

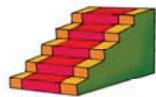
### References

- **Student material:**
  - Class Book (presentation slides with notes)
- **Book:**
  - The Data Warehousing Toolkit – Ralph Kimball
  - Introduction to Database Systems – C.J. Date
  - Advanced Data Warehouse – IBM
- **Web-site:**
  - <http://www.datawarehouse.org>
  - <http://etl-tools.info/>



### Next Step Courses (if applicable)

- BI related tool training



### Other Parallel Technology Areas

- NA

## Data Warehousing Concepts

### Lesson 1: Business Intelligence

June 12, 2014

Proprietary and Confidential

- 1 -



### Lesson Objectives

➤ In this lesson, you will learn:

- What is Business Intelligence?
- Need of Business Intelligence
- Terms used in Business Intelligence
- Components of Business Intelligence



1.1: Business Intelligence  
**What is Business Intelligence (BI)?**

- The term BI was coined by Gartner group in 1993.
- It is an important component in today's business information systems environment.
- It is the process of turning data into knowledge and knowledge into business gains.
- It collects and stores data into meaningful information in order to achieve better and timelier business decisions.
- It is an end user's activity supported by various analytical and collaborative tools.

June 12, 2014

Proprietary and Confidential

- 3 -

**IGATE**  
Spend Agilely Imagine

**Business Intelligence:**

- **Business Intelligence (BI)** is the process of getting useful information from data.

BI is an important component in today's business information systems environment.

- As the business environment has become increasingly competitive, the need to use corporate data as a strategic resource has intensified. However, most of the organizations in technology based businesses are **data rich** and are **information poor**. Much of the essential information that is needed to anticipate changing market conditions and customer preferences is locked in various transactional systems, spread sheets, and log files. So without the ability to deliver the right information to the right people at the right time, companies cannot stay competitive in this fast changing economy. So the **BI value proposition** is a term for the ability to navigate complex sales channels by maximizing knowledge about the customer base and developing strategies that leverage that knowledge from decision to action.
- **BI applications** are decision support tools that enable real-time, interactive access, analysis, and manipulation of mission-critical corporate information.

### What is Business Intelligence (BI)?

- BI is used for enhancement and optimization of organizational performance and operation.
- It delivers critical business information to end-users.
- It supports internal enterprise users in the assessment.
- It is applied across disciplines, namely Finance, CRM, and SCM
- It encompasses all types of data such as RDBMS, text, hierarchical, audio, and video.

June 12, 2014

Proprietary and Confidential

- 4 -



#### **Business Intelligence:**

- Business Intelligence gives answers to the questions such as given below:
  - Who are my top ten customers?
  - How effective was my last sales campaign?
  - Who is my best sales person by volume, and by dollar revenue, per region, during the last week of each month? How does that compare with last year?
  - How much more intelligent can you make your business processes?
  - How much more insight can you gain into your business?
  - How much more integrated can your business processes be?
  - How much more interactive can your business be with customers, partners, employees and managers?
  - BI solutions answer all these questions.

1.2: Need for Business Intelligence

## Why Business Intelligence?

➤ BI is required to meet the following business needs:

- To support the process of exploring data, relationships existing within data, and trends
- To make more accurate and more informed decision making
- To provide timely and accurate information to better understand your organization and to make more informed, real-time decisions

June 12, 2014 | Proprietary and Confidential | -5-

**IGATE**  
Spend Agilely Imagine

### Need for Business Intelligence:

➤ BI exhibits the following utility features:

- BI is a general term for applications, platforms, tools and technologies that support the process of exploring data, relationships existing within data, and trends.
- BI is important in helping organizations to stay ahead of the competition by providing the means for quicker, more accurate, and more informed decision making.
- BI provides timely and accurate information to better understand the organization and to make more informed, real-time decisions.

➤ But why do you need Business Intelligence?

- For many years, database vendors have focused on getting data into a database. The emphasis has led to great achievements in **online transaction processing** and capacity. Many companies have accumulated data that can be measured in gigabytes, terabytes, and even petabytes.
- **Transactional data**, which is the data that is used to run the business, is good for keeping track of what is happening in an organization. However, it is not well suited to finding out why things are happening or predicting future performance.
- Hence there arises a strong need for BI applications.

### Why Business Intelligence?

- Data Analysis is a huge and crucial part of Business Intelligence.
- Many organizations need to know the overall performance and the way its business is functioning.
- BI is used to gather past as well as present data.
- Modern BI systems are capable of managing large amount of unstructured data.

1.3: Terms used in BI

## Frequently used BI Terms

➤ Let us discuss some of the frequently used BI terms:

- Relational Database (RDB)
- Relational Database Management System (RDBMS)
  - Example: Informix, Microsoft SQL Server, Oracle.
- Online Transaction Processing (OLTP)
- Online Analytical Processing (OLAP)

June 12, 2014

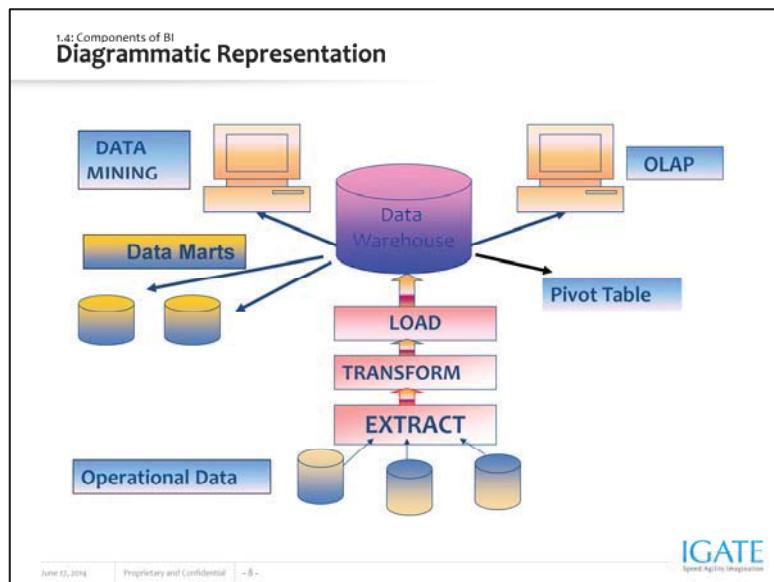
Proprietary and Confidential

-7-

**IGATE**  
Spend Agilely Imagine

### Terms used in BI:

- **Relational Database (RDB):**
  - It is a database that conforms to the relational model.
- **Relational Database Management System (RDBMS):**
  - It refers to the software used to create a RDB.
  - Example: Informix, Microsoft SQL Server, Oracle
- **Online Transaction Processing (OLTP):**
  - OLTP is a process which is used for day to day transaction processing.
  - **Example:** Operational systems, High volume data collection
- **Online Analytical Processing (OLAP):**
  - This processing method provides fast access to shared multidimensional data.
  - It is used to generically refer to software and applications that provide users with the ability to store and access data multi-dimensionally.



### Components of BI:

Following are the various components of BI:

- Operational Data:** Typically data is sourced from transaction processing systems. It is also called as Data Source. Typically data is sourced from transaction processing systems (Manufacturing, ERP, Sales). Example: Customer, Inventory, Credit, Sales, Operation and External are the data source.
- ETL Tools:**
  - Extract:** It is the process of pulling the data from external and operational data sources in order to source data for the data warehouse.
  - Transform:** It is the process that converts data to the format required by data warehouse. It cleanses data to ensure accuracy. It validates primary keys against defined owner. It converts to different numbering schema.
  - Load:** It is the process that loads data to data warehouse. It follows guidelines as outlined by the data warehouse.
- DWH:** Data Warehouse integrates and aggregates data from various operational and external database maintained by different Business Units.
- Data Mart:** Data mart is a repository of data collection from operational data source and other sources that are designed to serve a particular community of knowledge workers.
- Reports:** A report presents the data in a format understandable by the end user.

**Components of BI:**

Following are the various components of BI (contd.):

6. **OLAP:** OLAP is a category of software technology that enables the users to gain insight into data through fast, consistent, interactive access to a wide variety of possible views of information.
7. **Pivot Table:** A pivot table is the simplest tool to aggregate data by creating a dimension for each field and grouping the same values in a field. A pivot table is a data summarization tool found in data visualization programs such as spreadsheets. It allows you to reorganize and summarize selected columns and rows of data in a spreadsheet or database.

## Summary

➤ In this lesson, you have learnt:

- BI helps to extract information from data.
- BI helps organizations in making real time decisions.
- Components of BI are given below:
  - Data Warehouse
  - OLTP
  - OLAP
  - ETL tools
  - Data marts
  - Reports
  - Pivot table



### Review Questions

- **Question 1:** This is a huge and crucial part of Business Intelligence.
  - Option 1: Data collection
  - Option 2: Data analysis
  - Option 3: Data availability
- **Question 2:** OLAP Analysis is not the part of BI presentation.
  - True / False
- **Question 3:** \_\_\_ operation converts data to format required by data warehouse.



**Review Question: Match the Following**

1. pulling the data from external and operational data sources

2. part of BI presentation

3. Software for relational database

A. OLAP Analysis

B. RDB

C. Extract

D. OLTP

E. RDBMS



## Data Warehousing Concepts

### Lesson 2: General Concept of Data Warehouse

June 17, 2018

Proprietary and Confidential

+ 9 +



### Lesson Objectives

➤ In this lesson, you will learn:

- What is a Data Warehouse?
- History of Data Warehouse
- Need Of Data Warehouse
- Data Warehouse Architecture
- Data Warehouse Components
- Features of Data warehouse
- Data Mart
- Application areas



2.1: Data Warehouse

## What is a Data Warehouse?

➤ **Data Warehouse is a single, complete, and consistent store of data.**

- It is obtained from a variety of sources.
- It is made available to users in a way they understand and use in a business context.
- It is Central repository of information.
- It is a collection of key information.
- It contains read-only data.
- It contains historical data used for analysis purpose.
- It enables managers to make business decisions.

June 15, 2014 | Proprietary and Confidential | + 3 +



### Data Warehouse:

- A **Data Warehouse** is collection key of pieces of information to manage and direct the business for profitability.
- A Data warehouse is simply a single, complete, and consistent store of data obtained from a variety of sources and made available to end users in a way that they can understand and use it in a business context. It is nothing but a **database** or a **data store**. It is a database, so data has to be structured. The data is logically and physically transformed from multiple source applications to align with the business structure. It requires more historical data than that is generally maintained in operational database.
- Data is non-changing. It does not get updated. Data is never erased, so it is called **non-volatile**. Data Warehouse is designed for the analysis of non-volatile data.
- Data Warehouse integrates and aggregates data from various operational and external databases maintained by different Business Units.
- It is a process that managers use to load the warehouse query that makes information available. It enables people to make informed decisions. It is maintained for a long time period.
- A Data Warehouse is a **central repository** of information with appropriate tools.

### **Data Warehouse (contd.):**

- A Data Warehouse can also be defined as a **structured, extensible environment** designed for the analysis of non-volatile data, which is logically and physically transformed from multiple source applications to align with business structure, updated and maintained for a long time period, expressed in simple business terms.
- A Data Warehouse is used by different people in different fields. Companies use Data Warehouses to store information for marketing, sales, and manufacturing to help managers to get the feel of the data and run the business more effectively.
- A **database application** is a piece of software, which provides a user interface for users to add, delete, query, and update data, updates is called an on-line transaction processing (OLTP) application. An application that issues queries to the **read-only database** is called a **Decision Support System (DSS)**.

## 2.2 Characteristics of a Data Warehouse?

- A data warehouse is a subject-oriented, integrated, nonvolatile, time-variant collection of data in support of management's decisions.  
— WH Inmon

June 12, 2016

Proprietary and Confidential

→ 6 ←



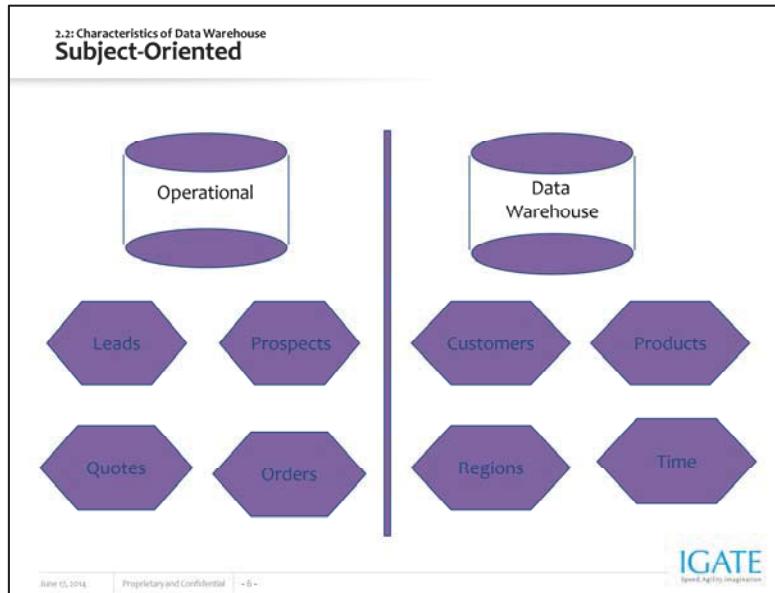
**Historical :** The data is continuously collected from sources and loaded in the warehouse. The previously loaded data is not deleted for long period of time. This results in building historical data in the warehouse.

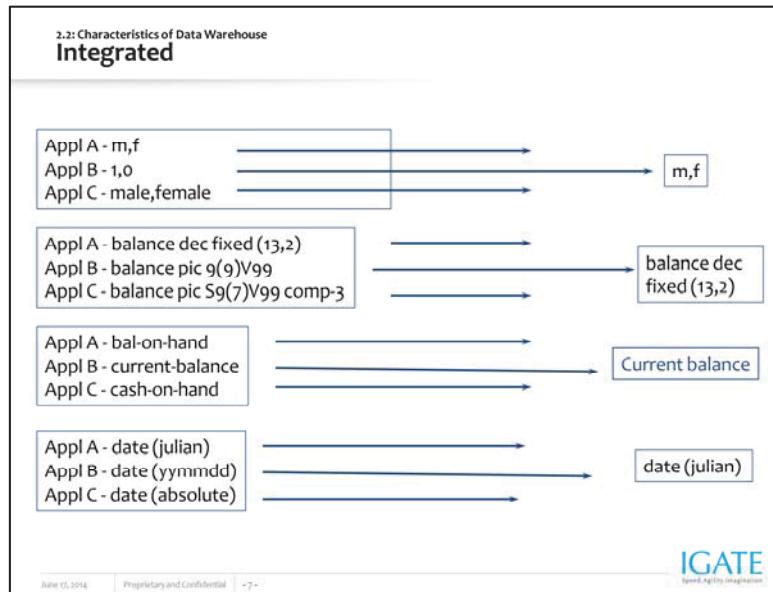
**Subject Oriented:** we mean data grouped into a particular business area instead of the business as a whole.

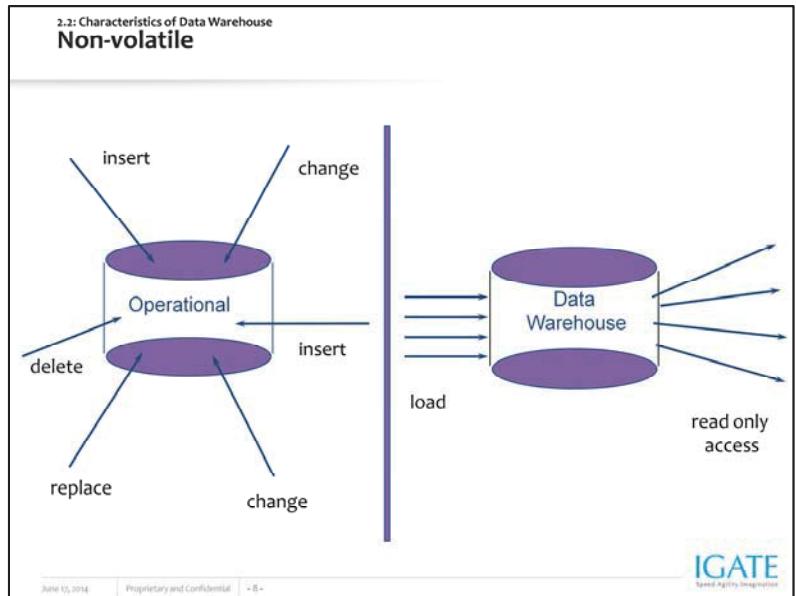
**Integrated:** It means, collecting and merging data from various sources. These sources could be disparate in nature.

**Time-variant:** It means that all data in the data warehouse is identified with a particular time period.

**Non-volatile:** It means, data that is loaded in the warehouse is based on business transactions in the past, hence it is not expected to change over time







2.2: Characteristics of Data Warehouse  
**Time Variant -**

The diagram illustrates two data storage models as cylinders. On the left, a cylinder labeled 'Operational' contains the text 'Current Value data'. Below it is a bulleted list: 'time horizon : 60-90 days' and 'key may not have element of time'. On the right, a cylinder labeled 'Data Warehouse' contains the text 'Snapshot data'. Below it is another bulleted list: 'time horizon : 5-10 years', 'key has an element of time', and 'data warehouse stores historical data'.

June 15, 2014 | Proprietary and Confidential | + 0 +

**IGATE**  
Speed. Agility. Integration.

2.3: Need for Data Warehouse

## Why Data Warehouse?

➤ **Data Warehouse is required to meet the following needs:**

- Companies want to tap on the vast potential of information to:
  - Have a separate informational system from operational systems
  - Improve quality of decision making
- Companies seek profitability through focused action.
- IT business requires an integrated, company-wide view of high quality information.
- Organizations want to analyze their activities in a balanced way.
- Organizations need to build on Customer Relationship Management.

June 15, 2014 | Proprietary and Confidential | + 10 +



### Need for Data Warehouse:

- The **Informational Systems** department must separate **informal systems** from **operational systems** in order to dramatically improve performance in managing company data. **Operational Data systems** are typically fragmented and are inconsistent. They are distributed over a variety of incompatible **hardware** and **software platforms**.
- IT professionals, in turn, must ensure that the enterprise's IT infrastructure properly supports a myriad set of requirements from different business users, each of whom has different and constantly changing needs.
- Example:** One file containing customer data may be located on a UNIX based server running an Oracle DBMS, while another is located on IBM main frame running the DB2 DBMS.
- Organizations want to analyze the activities in a balanced way.
- Customer Relationship Management is a building block of organizations. Organizations, in all sectors, are realizing that there is value in having a total picture of their interactions with customers across all touch points like for a bank, these touch points include ATM, electronic funds transfers, investment portfolio management, and loans.

### Why a separate Data Warehouse?

- A Data Warehouse helps in finding missing data.
- It provides consolidated data from multiple data sources.
- It helps in maintaining data quality coming from different sources.
- Special data organization is needed for vast volume of data.
- Complex OLAP queries degrade performance.

### Why a separate Data Warehouse?

#### Functions of a Data Warehouse:

- A Data Warehouse is typically used for data consolidation and enforcing uniform data quality.
  - **Data consolidation:** Decision support requires consolidation (aggregation, summarization) of data from many heterogeneous sources, namely operational databases, external sources.
  - **Data quality:** Different sources typically use inconsistent data representations, codes, and formats that have to be reconciled.

### 2.4: Data Warehouse Architecture What is Data Warehouse Architecture?

- **Data Warehouse Architecture is a description of the components and services of the Data Warehouse.**
  - It provides the mechanism to achieve enterprise integration to support business.
  - It provides an organizing framework that will improve data sharing.

### Data Warehouse Architecture Layers

- Data Warehouse Architecture consists of interrelated parts called as “layers” or “components”.
- Four layers of Data Warehouse Architecture are:
  - Operational: Functions as data storage
  - Informational: Stores business logic
  - Data access: Acts as a bridge between operational and informational layer
  - Meta data: Stores data dictionary

June 15, 2014

Proprietary and Confidential

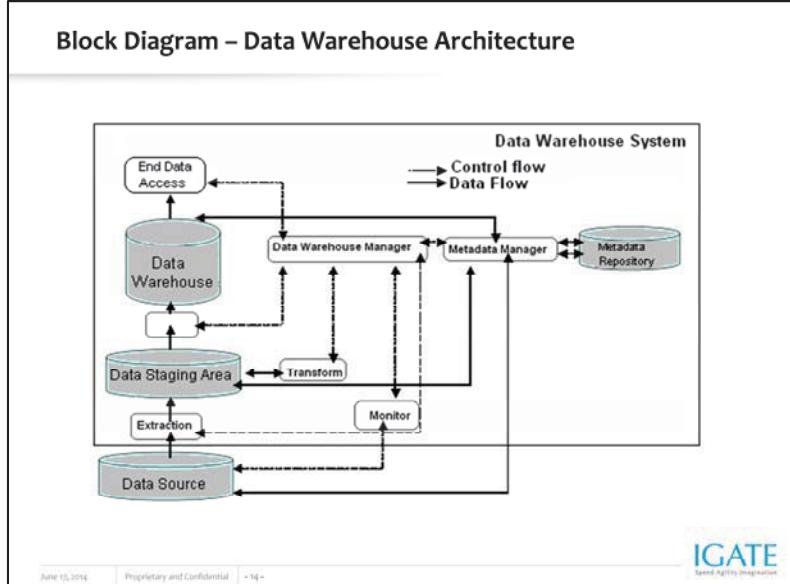
&gt; 13 &lt;



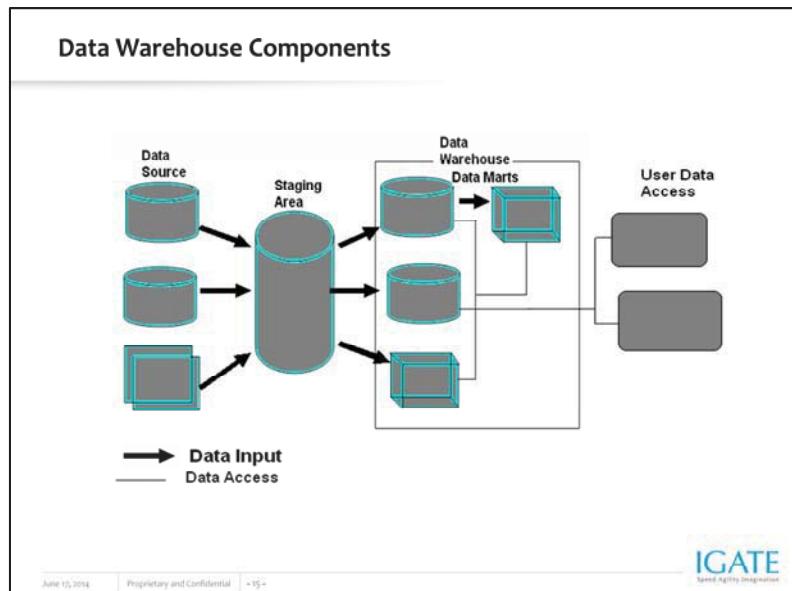
### Data Warehouse Architecture:

- It consists of four interrelated layers:

- **Operational:** It is the data source for Data Warehouse. It is also called as Internal / Physical layer. It takes care of how data is stored physically on disk.
- **Informational:** It performs data extraction for conducting analysis and reporting. It is also known as External / Logical layer. It is concerned with the way data is presented to the end user.
- **Data access:** It is an interface between Operational and Informational layer. It is also known as Conceptual layer.
- **Meta data:** It serves as a data dictionary for Data Warehouse.

**Data Warehouse Architecture:**

- Let us go through the different aspects of the Data Warehouse Architecture:
  - **Data Staging Area:** You need to clean and process your operational data before putting it into the warehouse. You can do this programmatically, although most data warehouses use a Staging Area instead.  
The Data Warehouse Staging Area is a **temporary location** where data from source systems is copied.  
A staging area is mainly required in a Data Warehousing Architecture for timing reasons. In short, all required data must be available before data can be integrated into the Data Warehouse.
  - **Metadata:** It provides a guide for warehouse users to understand DW.
  - **End User Access Tools:** High performance is achieved by pre-planning the requirement for joins, summations, and periodic reports by end users.
  - **Data Warehouse Manager:** It performs all operations associated with the management of the data in the warehouse.
- First, at the Data Access layer, the Data Source contains information.
- At the Operational layer, the data is extracted from Data Source and put into Data Staging Area.
- Metadata Repository stores the guidelines about Data Warehouse. With the help of transformation techniques, the Data Warehouse Manager and Metadata Manager load data into Data Warehouse.
- Finally, in the Informational layer, with the help of external view of database, the end user accesses the data.



### Data Warehouse Components:

- There are various components of Data Warehouse:
  - **Data Source:** Typically data is sourced from transaction processing systems (Manufacturing, ERP, Sales).
    - Data often resides in heterogeneous databases.
    - It comprises of different relational data (ORACLE, DB2, SQL Server, etc.).
    - Data could be on Mainframe (VSAM, IMS).
  - **Data Staging Area:** You need to clean and process your operational data before putting it into the warehouse. You can do this programmatically, although most data warehouses use a Staging Area instead.
  - **Data Marts:** You may want to customize your warehouse's architecture for different groups within your organization. You can do this by adding data marts, which are systems designed for a particular line of business.
  - **End User Access Tools:** High performance is achieved by pre-planning the requirement for joins, summations and periodic reports by end users.

2.5: Features of a Data Warehouse

## Salient Features

➤ Here are some of the features of a Data Warehouse:

- Time-variant data:
  - Data is meant for analysis and decision-making over the time.
  - Changes to the data are recorded against time dimension.
  - Data is stored as snapshots over past and current periods.
- Non-volatile data:
  - Data is not needed to run the daily business.
  - Data is primarily used for query and analysis.
  - Individual transactions are not updated in a Data warehouse.
  - Data is never over-written or deleted. It is non-updatable data.

June 15, 2014 | Proprietary and Confidential | +16+



### Features of a Data Warehouse:

Here are some of the features of a Data Warehouse:

➤ **Time-variant data:**

- Data in the Data warehouse contains a time dimension so that it may be used to study trends and changes.
- This nature of data:
  - Allows for analysis of the past.
  - Relates information to the present.
  - Enables forecast of the future.

➤ **Non-volatile data:**

- Data in the Data warehouse is loaded and refreshed from operational systems. However, it cannot be updated by end users.
  - Non-volatile data is not needed to run the daily business.
  - Non-volatile data is primarily used for query and analysis.
  - Individual transactions are not updated in a data warehouse.
  - Data is never over written or deleted.
  - Data warehouse consists of only non-updatable data.

### Salient Features

- Data granularity:
  - It refers to the level of detail.
  - It is inversely proportional to the amount of data stored.
  - Data is summarized at different levels.
  - Many Data warehouses have at least two levels of granularity.
  - Summarized data is stored.
  - It reduces storage costs.
  - It reduces CPU usage.
  - It increases performance since smaller number of records have to be processed.
  - Design is around traditional high level reporting needs.
  - Tradeoff with volume of data to be stored and detailed usage of data.

## Salient Features

— Subject oriented:

- Data is stored by subjects, not applications.
- Data is organized in the Data Warehouse such that it will infer the real world.
- Data is organized around major subjects, such as customer, product, sales.
- Focus is on the modeling and analysis of data for decision makers.
- DW provides a simple and concise view around a particular subject.
- DW is organized around the key subject of the enterprise.
- Major subjects may include customers, patients, students, products, and time.

June 15, 2014

Proprietary and Confidential

+ 18 +



### Features of Data Warehouse (contd.):

➤ **Data Warehouse is subject-oriented:**

- Focus is on the modeling and analysis of data for decision makers, not on daily operations or transaction processing.
- DW provides a simple and concise view around particular subject issues by excluding data that is not useful in the decision support process.

### Salient Features

#### Integrated data:

- Data is pulled from various databases from all applications.
- Operational platforms and operating systems for the data could be different.
- Data has to undergo a process of transformation, consolidation, and integration.
- Data inconsistencies are removed, standardization is achieved.

2.6: Data Mart

## What is a Data Mart?

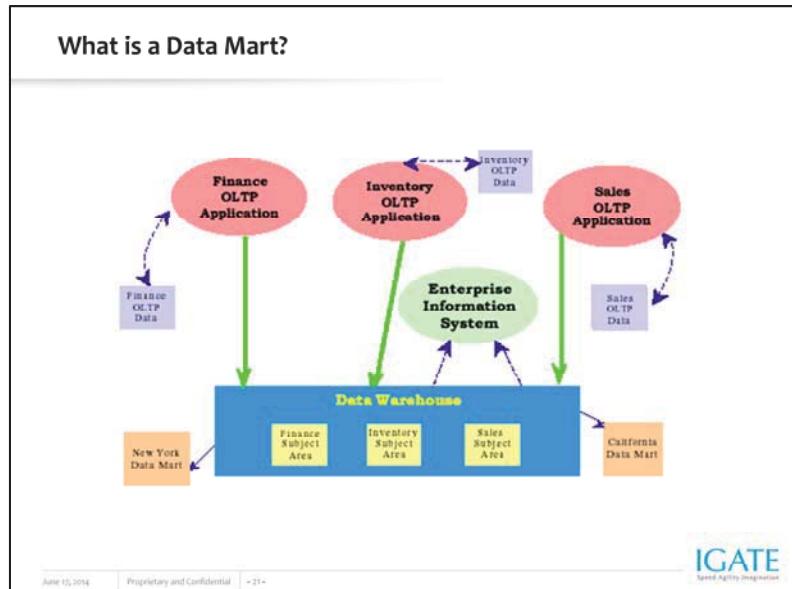
➤ **Data Mart is a subset of the Data warehouse.**

- It is typically fed from the Data warehouse.
- It is a Data warehouse that has limited scope.
- It is a repository of data gathered from operational data and other sources.
- It is used for decision making by a particular end-user group.
- Emphasis is on meeting the specific demands of a particular group of knowledge users.
- Maintain the ability to access the underlying base data.

June 15, 2014 | Proprietary and Confidential | +20 \*

**Data Mart:**

- **Data Mart** is a logical subset of a Data Warehouse that may make it simpler for users to access key corporate data. A Data Mart has a smaller data model, users only need a piece of data from the data warehouse.
- A Data Mart is a repository of data gathered from operational data and other sources. It is designed to serve a particular community of knowledge workers.
- In scope, the data may derive from an enterprise-wide database or data warehouse or be more specialized. The emphasis of a Data Mart is on meeting the specific demands of a particular group of knowledge users in terms of analysis, content, presentation, and ease-of-use. Users of a Data Mart can expect to have data presented in terms that are familiar.
- In practice, the terms **Data Mart** and **Data warehouse**, each tend to imply the presence of the other in some form. However, most writers using the terms seem to agree that the design of a **Data Mart** tends to start from the analysis of **user needs**. Similarly, the design of a **Data warehouse** tends to start from an analysis of the **data** that already exists and the manner in which it can be collected in such a way that the data can be used later.
- A **Data warehouse** is a central aggregation of data (which can be distributed physically). Whereas a **Data Mart** is a data repository that may or may not derive from a Data warehouse and emphasizes on the ease of access and usability for a particular design purpose. In general, a **Data warehouse** tends to be a strategic but somewhat unfinished concept. A **Data Mart** tends to be tactical and aimed at meeting an immediate need.
- In short, we need one large and complete Data warehouse that provides information to more focused, department-specific, and efficient Data Marts.
- Data Mart may derive from an enterprise-wide database or data warehouse or be more specialized.

**Data Mart (contd.):**

- Data Mart is a Data warehouse that is limited in scope.
- The emphasis of a data mart is on meeting the specific demands of a particular group of knowledge users in terms of analysis, content, presentation, and ease-of-use.
- Users of a Data Mart can expect to have data presented in terms that are familiar.
- It is important to maintain the ability to access the underlying base data to enable drilldown analysis as necessary. The only difference between a Data Warehouse and a Data Mart is the scope. One can define the Data Warehouse from various Data Marts. On the other hand, one can define Data Marts from the Data Warehouse.

## Types of Data Marts

### ➤ Dependent Data Mart

- A Data Mart whose source is the Data Warehouse
- All dependent Data Marts are loaded from the same source – the Data Warehouse

### ➤ Independent Data Mart

- A Data Mart whose source is the legacy application environment
- Each independent Data Mart is fed uniquely and separately by the individual source systems

The main difference between independent and dependent data marts is how you populate the data mart; that is, how you get data out of the sources and into the data mart. This step, called the Extraction-Transformation-and Loading (ETL) process, involves moving data from operational systems, filtering it, and loading it into the data mart.

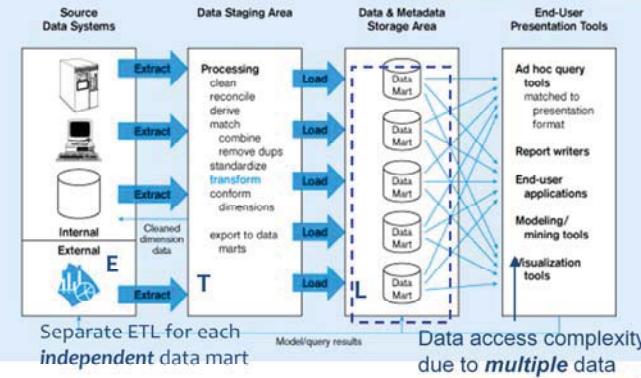
With dependent data marts, this process is somewhat simplified because formatted and summarized (clean) data has already been loaded into the central data warehouse. The ETL process for dependent data marts is mostly a process of identifying the right subset of data relevant to the chosen data mart subject and moving a copy of it, perhaps in a summarized form.

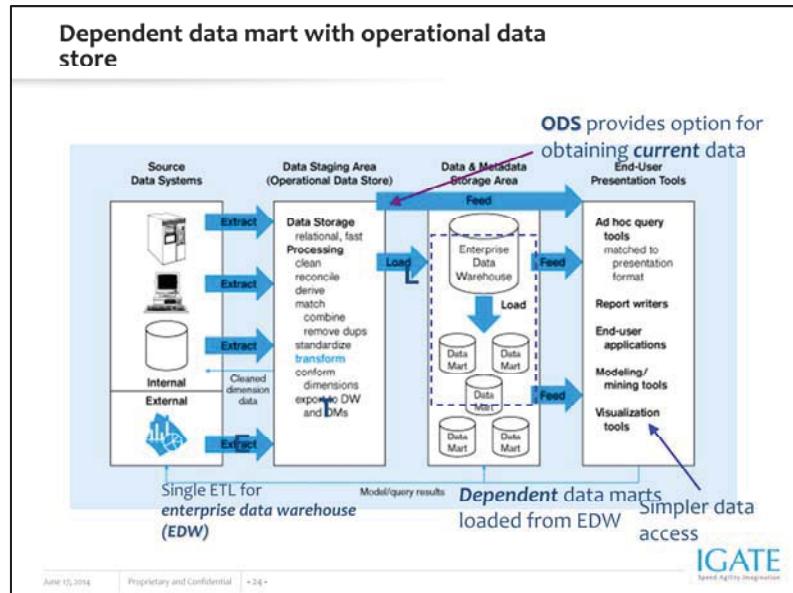
In Independent data mart, each data mart is sourced directly from the operational systems. One must deal with all aspects of the ETL process, much as you do with a central data warehouse. The number of sources is likely to be fewer and the amount of data associated with the data mart is less than the warehouse, given your focus on a single subject.

The motivations behind the creation of these two types of data marts are also typically different. Dependent data marts are usually built to achieve improved performance and availability, better control, and lower telecommunication costs resulting from local access of data relevant to a specific department. The creation of independent data marts is often driven by the need to have a solution within a shorter time.

### Independent Data Mart

**Data marts:**  
Mini-warehouses, limited in scope





2.7: Data Warehouse Application Areas  
**Industry-wise Application**

Industry	Application
Finance	Credit Card Analysis
Insurance	Claims, Fraud Analysis
Telecommunication	Call record analysis
Transport	Logistics management
Consumer goods	Promotion analysis
Data Service providers	Value added data
Utilities	Power usage analysis

## Summary

➤ In this lesson, you have learnt:

- Data Warehouse stores historical data.
- Data Mart emphasizes on meeting the specific demands of a particular group of knowledge users.
- Features of Data Warehouse are:
  - Time variant data
  - Non volatile data
  - Data granularity
  - Subject oriented
  - Integrated data



### Review Question

- Question 1: \_\_\_\_\_ is a subset of data warehouse.
- Question 2: Data Mart is a structure for corporate view of data.  
True/ False
- Question 3: \_\_\_ is used for decision making by a particular end-user group.



## Data Warehousing Concepts

### Lesson 3: Dimensional Modeling

June 12, 2014

Proprietary and Confidential

- 1 -



## Lesson Objectives

➤ In this lesson, you will learn:

- What is Dimensional modeling ?
- Facts and Dimension tables
- Database schema
- Schema Design for Modeling



3.1: Dimensional Modeling  
**What is Dimensional Modeling?**

- Dimensional Modeling (DM) is the name of a logical design technique often used for Data Warehouses.
- DM is the technique for databases that are designed to support end-user queries in a Data Warehouse.
- A Dimension Model is composed of dimension tables and fact tables.
- It provides a conceptual framework.
- It simplifies the business flow.
- It is structurally classified as fact or dimension.

June 12, 2014

Proprietary and Confidential

- 3 -



**Dimensional Modeling:**

- Dimensional Modeling has the characteristic for organizing data roughly into base facts and dimensions of those facts.
- Dimensional Modeling provides the Conceptual Framework. It is basically used for faster query performance for the business users. **Facts** are basically organization's business processes. They are usually numeric values. **Dimension** is a context that describes the fact.
- Every organization has Dimensional Modeling for its business processes, and it consists of **fact tables** and **dimensional tables**. It helps business users in easily understanding the typical system model.
- Dimensional Modeling represents the complexities of the business process in a simple manner. **Understandability** and **Query performance** are two major reasons for which dimensional modeling is accepted widely in the industry.
- Dimensional Modeling is a logical design technique that allows to retrieve the data with high-performance.

3.2: Fact and Dimension Tables

## Concepts of Fact and Dimension Tables

- Fact tables and Dimension tables are the two types of objects that are commonly used in designing database schemas.
- Fact table contains two columns, namely numeric facts and foreign keys of dimension tables.
- Dimension tables contain the attributes that describe fact records.

June 12, 2014 | Proprietary and Confidential | - 4 -

**IGATE**  
Spend Agilely Imagine

### Fact and Dimension Tables:

- A **fact table** has two types of columns.
  - The first column type contains numeric facts (often called measurements).
  - The other column type contains the foreign keys of dimension tables.
- A **fact table** contains multiple foreign keys.
- Each pair of primary key in dimension and foreign key of fact table contains the measurements.
- A **Dimension table** contains the attributes that describe fact records. Some dimension table attributes provide descriptive information and other attributes (primary key) are used to join with fact tables.  
**Example:** A customer dimension table contains two attributes, namely customer id (Primary key) and customer description. So we will use the primary key attribute customer id to join with fact tables.
- However, **dimensional** and **fact modeling** is not of the highest Normal Form, but makes use of a key of performance indicators. Dimensions can strive to be in Boyce Codd (BCNF) 3rd Normal Form. Whereas Fact tables may be in 1st Normal Form, having only a primary key being unique.

### Multidimensional Data

- Designed to resolve complex business queries
- Helps to analyze data from different dimensions
- Different dimensions form a cube
- Every edge represents a dimension

Dimensions: Product, Region, Time  
Hierarchical summarization paths

Product Industry  
Category  
Product Week

Region Country  
Region  
City  
Office

Time Year  
Quarter  
Month  
Day

June 12, 2016 Proprietary and Confidential - 5 -

**IGATE**  
Speed. Agility. Imagination.

#### Multidimensional Data :

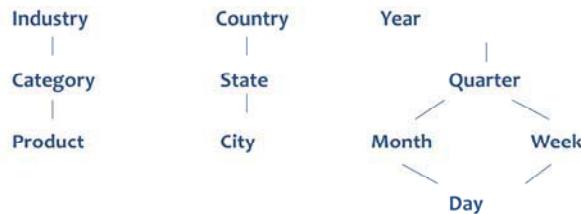
The multidimensional data model is the integral part of On line analytical Processing. The multidimensional data model is designed to resolve the complex queries.

In the logical multidimensional model, a cube represents the measures with same shape. In a cube every edge represents a dimension. Members of Dimension are aligned on the edges and divide the cube shape into cells in which stored the data values. It is basically used for developing data mart.

In above cube three edges represent the three dimension table Product, Region and time.

### Multidimensional Data Analysis

- Analysis based on multiple dimensions
- Result varies with the dimension change across analysis
- Customers (city, state, country)
- Time (day, week, month, quarter, year)
- Products (product, category, industry)
- Hierarchies on dimensions:



[June 12, 2016] Proprietary and Confidential - 6 -

**IGATE**  
Speed. Agility. Imagination.

Multidimensional Data Analysis is the analysis of data based on dimensions. It includes analysis of a particular data with respect to different and multiple dimensions. The value varies when there is a change in the dimensions across the analysis. It changes in terms of context one wishes to analyze data.

For E.g. Analysis of Product by City, Transactions for last 3 years.

3.3: Database Schema  
**Concept of Database Schema**

- **Database schema includes various elements to store data.**  
Example: facts, dimensions, attributes, hierarchy, cube
- **Facts are numeric values to be stored in the database.**
- **Dimensions are description about facts.**
- **Attributes are characteristics of dimensions.**
- **Hierarchy is a logical representation of the order of the entities.**

June 12, 2014

Proprietary and Confidential

-7-



**Fact and Dimension Tables:**

**Database Schema:**

- **Database schema** is a set of facts in multi-dimensional data. A fact has a measure dimension quantity that is analyzed, for example, number of visas.
- It has a set of dimensions on which data is analyzed, for example, country, consulate, date of issue for a visa. Each dimension has a set of attributes
- **Example:** “Visa” dimension has visa date, visa type, and visa category
- Attributes of a dimension may be related by partial order, or Hierarchy: for example, post > county > region.

3.4: Schema Design for Modeling  
**Schema Types**

➤ **Schema design is the Database organization for modeling.**

- It must look like business.
- It must be recognizable by business user.
- It must be approachable by business user.
- It must be simple.

➤ **Schema Types:**

- Star Schema
- Snowflake schema

June 12, 2014

Proprietary and Confidential

- 8 -



**Fact and Dimension Tables:**

**Schema Design for Modeling:**

- **Schema design** is the organization of database for modeling.
- The design shows how the model will be implemented in a system. It must be kept simple and familiar with the business context. It must be easily understood by business user. It should be designed in such a way that the business users can fully understand it in terms of facts, measures, dimensions, and hierarchies.

**Schema Types:**

- **Star Schema-Fact and Dimension tables:** Star schema has all multi-leveled dimensions that are flattened.
- **Snowflake Schema:** It has dimensional hierarchy directly by normalizing tables. In Snowflake schema, at least one multi-leveled dimension is kept separate.

### Star Schema

- Star Schema consists of a central fact table surrounded by dimension tables.
- The measures of interest for OLAP are stored in the fact table (for example: Dollar Amount, Units in the table SALES).

June 12, 2014

Proprietary and Confidential

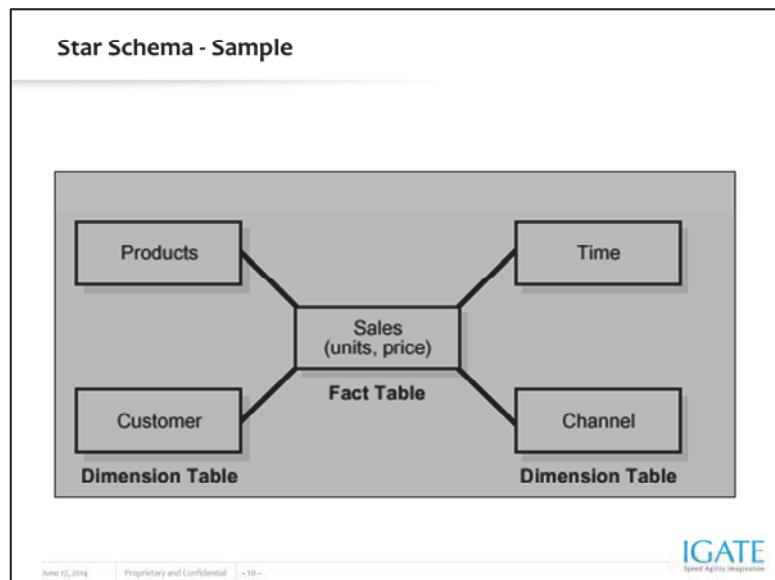
- 9 -

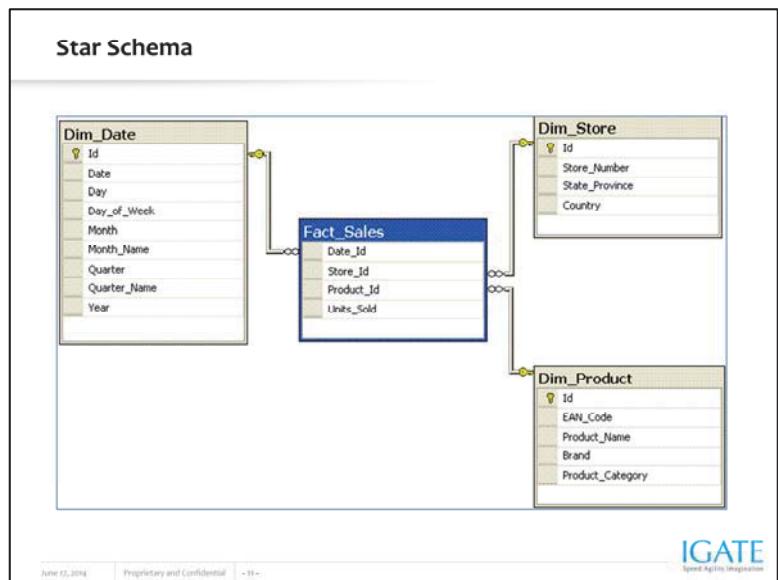


### Fact and Dimension Tables: Schema Types (contd.):

#### Star Schema (contd.):

- For each dimension of the multidimensional model there exists a dimension table (for example: Geography, Product, Time, Account) with all the levels of aggregation and the extra properties of these levels.
- It consists of Fact table.
- It consists Compound primary key.
- Star schema focuses on two major advantages, namely:
  - Ease of use
  - Efficient performance





### Fact and Dimension Tables:

#### Schema Types (contd.):

##### Star Schema:

➤ **Star schema** is commonly used by relational databases. The performance can be improved by using this design rather than traditional join operations. A **Star schema** is a database design that contains a central table, called a **fact table**, which is in relationship with many tables called **dimension tables**. This schema design resembles a star, thus the name is Star Schema. It is a very simple programmatic approach. It is very similar way in which a user thinks about a system, hence it is simple. It is easier to use. It is very efficient in the performance. It is best suited for **MOLAP application tools**. Typically, most of the fact tables in a star schema are in database Third Normal Form, while dimensional tables are de-normalized (Second Normal Form). Despite the fact that the Star schema is the simplest Data warehouse architecture, it is most commonly used in the Data warehouse implementations about 90-95%, across the world today.

##### Example:

- Fact Table: Fact\_Sales table
- Dimension table: Dim\_Date, Dim\_Store, Dim\_Product

### Snowflake Schema

- Snowflake Schema represents dimensional hierarchy directly by normalizing tables.
- It is a variation on the star schema.
- It is easy to maintain and saves storage, very large dimension tables.
- They have improved query performance.

June 12, 2014

Proprietary and Confidential

- 12 -



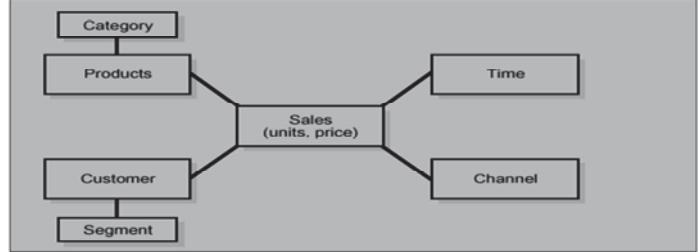
### Fact and Dimension Tables:

#### Schema Types (contd.):

##### **Snowflake Schema:**

- It is more complex than Star schema design. The main difference is that dimensional tables in a snowflake schema are normalized, so they have a typical relational database design.
- Snowflake schemas are generally used when a dimensional table becomes very big and when a Star schema cannot represent the complexity of a data structure. For example, if a PRODUCT dimension table contains millions of rows, then the use of Snowflake schemas should significantly improve performance by moving out some data to other table (with REGION for instance). The data redundancy is eliminated. The problem is that the more normalized the dimension table is, the more complicated SQL joins must be issued to query them. This is because in order for a query to be answered, many tables need to be joined.

### Snowflake - Sample

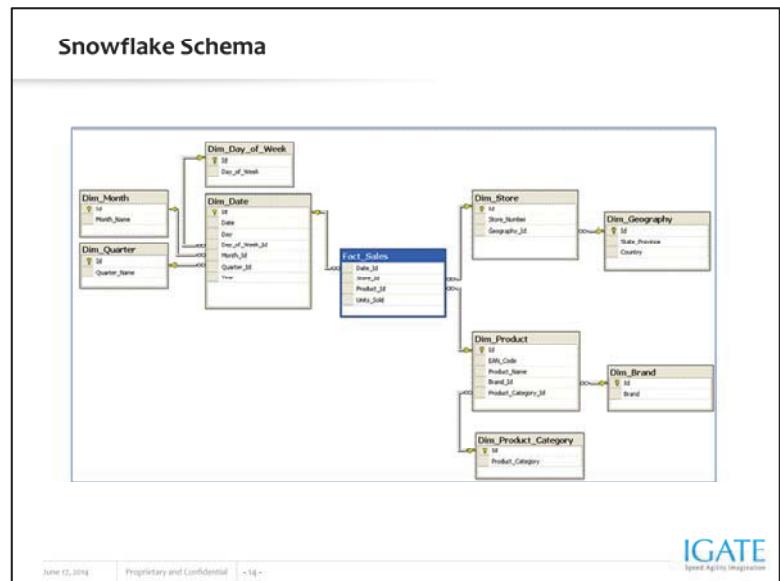


June 12, 2014

Proprietary and Confidential

-13-

**IGATE**  
Spend Agilely Imagine



## Summary

➤ In this lesson, you have learnt:

- Dimensional Modeling represents the complexities of the business process in a simple manner.
- The schema types are star schema and snowflake schema



## Summary

- Database schema has various elements, such as:
  - Fact
  - Dimension
  - Attributes
  - Hierarchy
  - Cube
- Schema design is the organization of database.



**Review Question**

- Question 1: \_\_\_\_\_ are description about facts.
- Question 2: \_\_\_\_\_ in Snowflake Schema are normalized into multiple tables.
- Question 3: \_\_\_ is the name of a logical design technique often used for Data Warehouses.

## Data Warehousing Concepts

### Lesson 4: ETL and Metadata

June 12, 2014

Proprietary and Confidential

- 1 -



### Lesson Objectives

➤ In this lesson, you will learn:

- ETL Process
- Metadata used in ETL
- Metadata in Data Warehousing
- Simple Warehouse Model



4.1: Extract Transform and Load (ETL) Process  
**Concept of ETL**

- The Data Warehouse always has enterprise data. Data comes from various sources, such as Spreadsheets, Mail lists, and Databases.
- The required data is extracted, transformed to suit information needs and finally loaded at a central location.
- This is done by ETL (Extract Transform and Load) process.
  - Extract: Data extraction and staging
  - Transform: Convert to format required by data warehouse
  - Load: Load data to data warehouse

## ETL Process

### ➤ Extraction

- Data extraction from various source (Heterogeneous systems)
- Different data representations, formats
  - e.g. RDBMS, Flat files, IMS, VSAM
- Data to be converted to a common format for transformation process
- Extracts the data from data source and keeps in staging.
- Data comes from an operational source or archive systems which are the primary sources of data for the Data warehouse.
- It minimizes impact on production data sources

## ETL Process

### ➤ Transformation

- Various sets of business rules and functions are applied on extracted data before the data gets loaded to Data warehouse
- One or more of the following steps may be involved in the transformation process
  - Selecting only certain columns to load
  - cleansing the data to remove duplicates and enforce consistency
  - Translating coded values (e.g., if the source system stores 1 for male and 2 for female, it may be translated as M for male and F for female in data warehouse)
  - Encoding free-form values (e.g., mapping "Male" and "I" and "Mr" into M)
  - Deriving a new calculated value
  - Joining together data from multiple sources (e.g., lookup, merge, etc.)

## ETL Process

### ➤ Loading

- Transformed data loaded to Data warehouse
- Load Dimensions and then Fact
- Indexes to be dropped before loading and recreated after loading the Data Warehouse
- Load cycle (Daily, weekly Monthly...)

4.2: Metadata

## Metadata used in ETL

➤ **Metadata in ETL contains data about Data:**

- Dimension
- Attribute
- Fact
- Measure

June 12, 2014 Proprietary and Confidential -7-

**IGATE**  
Spend Agilely Imagine

**Metadata:**

Metadata is the data about Data.

- **Dimension:** It is a perspective that can be used to analyze the data. Dimensions become more useful when there are many descriptive attributes that can be used for analyzing the data.
- **Attribute:** It is often used to describe the extended Dimension.  
Example: Customer, Item, Date, Fact
- **Fact:** It is the raw enumerable piece of information about the transaction. It is always a numeric value (usually aggregatable) about the transaction.  
Examples: Quantity, Unit Price, Count
- **Measure:** Measure can be the product of one or more fact tables. Measure can be the result of any formula which is derived from Relational Database or Business Intelligences Tool analytical engine. It is the product of one or more Facts.  
Examples: Quantity, Unit Price, Count, Quantity \* Unit Price, Average (Unit Price) and Minimum (Quantity).  
For example, if a customer is an attribute from your databases table, then customer metadata can give the information about the customer like name, address will be the dimensions, whereas telephone number can act as fact, etc. Age can be considered as measure, since it can be calculated from DOB-Current date.

4.3: Metadata in Data Warehousing

## Using Metadata in Data Warehousing

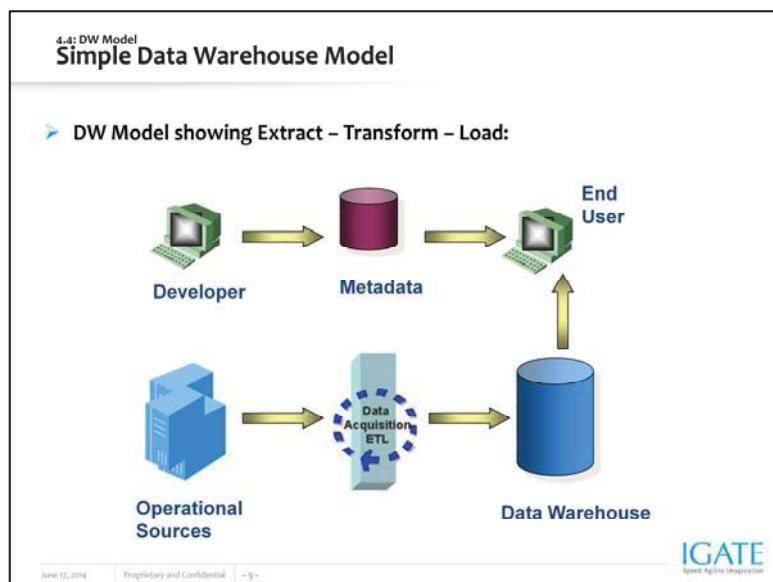
- Metadata plays vital role in Data Warehouse architecture.
- Metadata in Data Warehouse contains:
  - Data dictionary
  - Data flow
  - Data transformation
  - Version control
  - Data usage statistics
  - Alias information
  - Security

June 12, 2014 Proprietary and Confidential - 8 -

**IGATE**  
Spend. Align. Imagine.

### Metadata in Data Warehousing:

- Metadata is the blood of the Data Warehouse. It is the information that describes the system. Metadata plays a vital role in Data Warehouse architecture. It provides the information to the application to control warehouse activities. A single change in the metadata repository affects the entire architecture.
- Metadata in Data Warehousing:
  - **Data dictionary:** It contains definitions of the databases and relationship between data elements.
  - **Data flow:** It contains direction and frequency of data feed.
  - **Data transformation:** It contains transformations required when data is moved.
  - **Version control:** It records changes to stored metadata.
  - **Data usage statistics:** It is a profile of data in the warehouse.
  - **Alias information:** It contains alias names for a field.
  - **Security:** It contains the names of the data access authorized people.

**DW Model:**

- A Data Warehouse setup typically comprises of the following end points:
- **Developer:** The developer puts business rules for data transformation into the metadata repository.
  - **Metadata:** It indicates about the data is available in the warehouse and where the data is located.
  - **Data Warehouse:** Data Warehouse integrates and aggregates data from various operational and external databases maintained by different Business Units.
  - **Operational Sources:** It can comprise of Customer Database, Sales Database, and Product Database.
  - **End User:** High performance is achieved by pre-planning the requirement for joins, summations, and periodic reports by end users.

### Summary

- In this lesson, you have learnt:
- ETL Process
  - Metadata used in ETL
  - Metadata in Data Warehousing
  - Simple Warehouse Model



### Review Questions

- **Question 1:** Metadata contains the following:
  - Option 1: Data Dictionary
  - Option 2: Data Flow
  - Option 3: Data Mart
  
- **Question 2:** Multidimensional data represents business complexities
  - True/ False



**Review Question: Match the Following**

- |   |                    |
|---|--------------------|
| 1. Puts business rules                  | A. Data dictionary |
| 2. Product of one or more fact tables   | B. Measure         |
| 3. Direction and frequency of data feed | C. End user        |
|   | D. Data flow       |
|   | E. Developer       |



June 12, 2014 | Proprietary and Confidential | - 12 -

## Data Warehousing Concepts

### Lesson 5: Online Analytical Processing (OLAP)

June 12, 2014

Proprietary and Confidential

- 1 -



### Lesson Objectives

➤ In this lesson, you will learn about:

- The concept of Online Analytical Processing
- Need for Separate Operational and Informational Systems
- Nature of OLAP analysis
- Types Of OLAP
- OLAP Service Tools
- OLTP and OLAP
- Operational versus Informational Systems



S.1: Online Analytical Processing (OLAP)  
**Concept of OLAP**

- OLAP is a functionality available in Data Warehouse applications.
- It enables client applications to efficiently access data in a Data Warehouse or Data Mart.
- It is a multi-dimensional data model.
- It contains a variety of possible views of information.
- It simplifies evaluation of ad hoc complex queries.
- It provides a very fast response time to ad hoc queries.

June 12, 2014 | Proprietary and Confidential | - 3 -

**IGATE**  
Spend Agilely Imagine

**Online Analytical Processing (OLAP):**

- OLAP is a category of software technology. It enables the users to gain insight into data through fast, consistent, interactive access to a wide variety of possible views of information, which has been transformed from raw data, to reflect the real dimensionality of the business.
- It provides benefits like **pre-aggregation** of frequently required data, enabling a very fast response time to ad hoc queries. It gives a **multi-dimensional data model** that makes it easy to select, navigate, and explore the data.
- OLAP systems enable managers and analysts to rapidly and easily examine key performance data. OLAP systems allow comparison and trend analysis even on very large volumes. OLAP allows users to view data from various perspectives. It is fast and easy because some aggregations are computed in advance.

5.2: Nature of OLAP Analysis

## Use of OLAP

➤ **OLAP analysis is used for:**

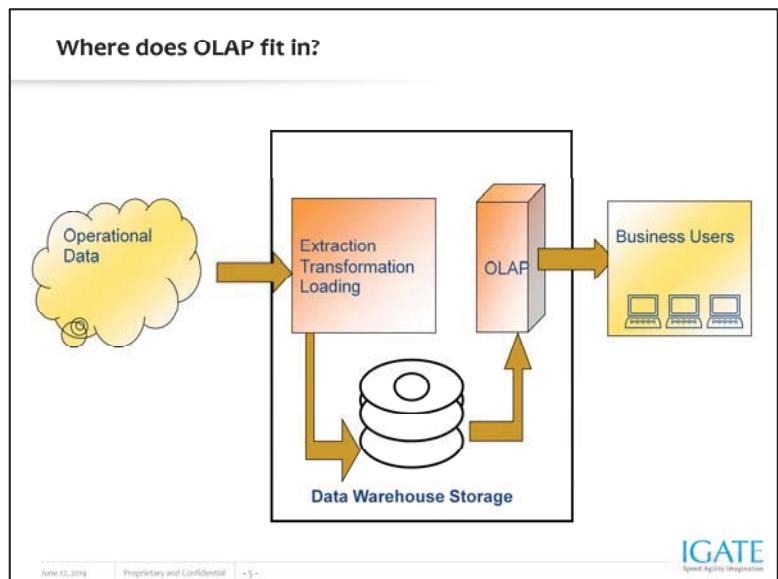
- Aggregation
- Comparison
- Ranking
- Access to data
- Complex criteria specification

June 12, 2014 | Proprietary and Confidential | - 4 -

**IGATE**  
Spend Agilely Imagine

### Nature of OLAP Analysis:

- The nature of OLAP analysis varies with the multiple ways of using it.
- It is used as the aggregation for summing up the data.
  - It provides the easy way of comparison.
  - It allows us to rank the data such that you will be able to find the top most and lower most values in the analysis.
  - It also helps in accessing the data in detailed way.
  - It allows to perform complex specification on the criteria.
  - It represents the data in a more simpler way such that it is easily visualized in terms of graphical presentation.
- OLAP Analysis is used for:
- Aggregation: (total sales, percent-to-total)
  - Comparison: Budget versus Expenses
  - Ranking: Top 10, quartile analysis
  - Access to detailed and aggregate data
  - Complex criteria specification

**Nature of OLAP Analysis:****Where does OLAP fit in?**

- Data from various sources goes through the ETL process and is integrated into a Data Warehouse. Subsequently, OLAP is used to analyze the data in the Data Warehouse. OLAP focuses on meeting end-user's analytical requirements.
  - **Operational Data:** It is the Customer Database, for example, Sales Database and Product Database.
  - **End User:** High performance is achieved by pre-planning the requirement for joins, summations and periodic reports by end users.
  - **Extract Transform and Load (ETL):**
    - **Extract:** It extracts data from data source and keeps it in staging.
    - **Transform:** It converts data into format required by Data Warehouse.
    - **Load:** It loads data to Data Warehouse.
  - **DWH:** Data Warehouse integrates and aggregates data from various operational and external data bases maintained by different Business Units.
  - **OLAP:** It has been in use to process and record transactions that create new data and update existing information in databases.
  - **Business User:** High performance is achieved by pre-planning the requirement and putting business rules by business users.

### OLAP Models

- OLAP models are of different types
- The processing in all these different types is the same:  
Online Analytical processing
- The storage methods are different in different models
- Different OLAP Models are
  - ROLAP
  - MOLAP
  - HOLAP

5.3: Types of OLAP  
**ROLAP**

➤ **Relational Online Analytical Processing (ROLAP):**

- It stores in a Relational form.
- It stores Data Mart (Star schema).

➤ **Advantages:**

- It has no data size limitation.
- It can leverage functions of RDB.

➤ **Disadvantages:**

- Each request must query the RDB.
- ROLAP itself is limited to RDB functionality.

June 12, 2014 | Proprietary and Confidential | -7-

**IGATE**  
Spend Agilely Imagine

**Types of OLAP:****Relational Online Analytical Processing (ROLAP):**

- The data stored in the relational database gives the appearance of traditional OLAP's slicing and dicing functionality.
- In Relational OLAP (ROLAP), Relational DBMS stores Data Mart (Star schema).

**Advantage:**

- ROLAP itself places no limitation on data amount.
- Relational database already comes with a host of functionalities. ROLAP technologies, can leverage these functionalities since they sit on top of the relational database.

**Disadvantage:**

- Each ROLAP report is essentially a SQL query (or multiple SQL queries) in the relational database. The query time can be long if the underlying data size is large.
- ROLAP itself is limited to RDB functionality.

S.3: Types of OLAP  
**MOLAP**

➤ **Multidimensional Online Analytical Processing (MOLAP):**

- Data is stored multi-dimensionally by using Multidimensional Databases (MDDB)
- MDDB's store data in the form of Multi dimensional cubes
- MOLAP cubes are built for fast data retrieval and are optimal for slicing and dicing operations.

➤ **Advantages:**

- It has excellent performance.
- It can return complex calculations.

➤ **Disadvantages:**

- It is limited in scope as definition of cube creates boundaries.
- Limited volume of data is churned.

June 12, 2014 | Proprietary and Confidential | -8-

**IGATE**  
Spend Agilely Imagine

**Types of OLAP:****Multidimensional Online Analytical Processing (MOLAP):**

- In MOLAP, the data is stored in a multi-dimensional cube. The storage is not in the relational database, but in proprietary formats.
- MOLAP storage structure is Array-based.

**Advantage:**

- MOLAP cubes are built for fast data retrieval, and is optimal for slicing and dicing operations.
- All calculations have been pre-generated when the cube is created. Hence, complex calculations are not only doable, they return quickly, as well.

**Disadvantage:**

- In case of MOLAP, scope is limited as definition of cube creates boundaries.
- It is not possible to include a large amount of data in the cube itself. This is because all calculations are performed when the cube is built. Only summary-level information will be included in the cube itself.

### HOLAP

- HOLAP is the product of the attempt to incorporate the best features of MOLAP and ROLAP into a single architecture.
- HOLAP systems stores larger quantities of detailed data in the relational tables while the aggregations are stored in the pre-calculated cubes.
- HOLAP also has the capacity to “drill through” from the cube down to the relational tables for delineated data.

June 12, 2014

Proprietary and Confidential

- 9 -



This tool tried to bridge the technology gap of both products by enabling access or use to both multidimensional database (MDDB) and Relational Database Management System (RDBMS) data stores.

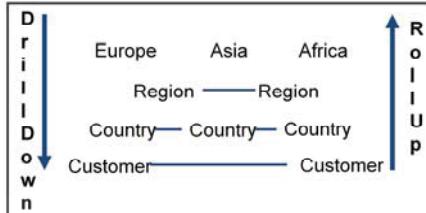
### Types of OLAP Operations

➤ Different OLAP operations are:

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

➤ Other operations:

- Drill across
- Drill through



### Types Of OLAP Operations:

- **Aggregation / Consolidation / Roll up (drill-up):** Roll Up operation is used to summarize data by climbing up hierarchy or by dimension reduction. This allows you to aggregate data from lower level details to the parent level. For example, the total revenue generated by a particular product type will be the rolled up value of the revenue generated by descendants, that is, products that belong to that particular product type.
- **Drill down (roll down):** Roll down operation is the reverse of roll-up, that is, a drill down from higher level summary to lower level summary or detailed data, or introducing new dimensions. It allows you to view data from a top-level to a detailed view by going down the hierarchy. For example, we can view the Sales data for the Year and drill down from the year level to a quarterly view and further down to a monthly view.
- **Slice and dice:** It is a general term for viewing data from any angle.
- **Pivot (rotate):** Rotate operation is used for reorienting the cube, visualization, 3D to series of 2D planes.

### **Other operations:**

- **Drill Across:** It involves (across) more than one fact table.
- **Drill Through:** It involves drilling through the bottom level of the cube to its back-end relational.

5.5: OLTP and OLAP

## Concepts of OLTP and OLAP

➤ **OLTP: Online Transaction Processing:**

- OLTP is used to process and record transactions that create new data.
- It updates existing information in databases.

➤ **OLAP: Online Analytical Processing:**

- Data is aggregated, warehoused, and then analyzed.
- Users query and generate reports without modifying any data.

June 12, 2014 | Proprietary and Confidential | - 11 -

**IGATE**  
Spend Agilely Imagine

### OLTP and OLAP:

- OLTP is basically operational data, wherein data is frequently changing.
- For example, you can consider an online Railway Reservation system. A passenger books a ticket for two people. This becomes an operational data. S/He can also change the number of passengers travelling online since OLTP data is frequently updated.
- On the other hand, OLAP data is non-operational data, wherein data is read-only data. It is used for analytical purpose.
- For example, suppose one wants to see, the number of trains running on the previous day. This becomes analysis of data that is not updatable.

5.6: Operational versus Informational system  
**Points of Difference**

	Operational (OLTP)	Informational (DW)
Typical User	Clerical	Management
System usage	Regular business	Analysis
Workload	Read/Write	Read only
Types of queries	Predefined	Ad-hoc
Unit of interaction	Transaction	Query
Level of isolation required	High	Low
No of records accessed	<100	>1,000
No of concurrent users	Thousands	Hundred
Focus	Data in and out	Information out

June 12, 2014 | Proprietary and Confidential | -12-



**Note:** The above table shows the comparison between **OLTP** and **Data Warehouse**.

### Points of Difference

	Operational (OLTP)	Analytical Systems
User	Clerk, IT Professional	Knowledge Worker
Function	Day to day operations	Decision support
DB Design	Application-oriented (E-R based)	Subject Oriented (Star, Snowflake)
Data	Current, Isolated	Historical, Consolidated
View	Detailed, Flat relational	Summarized, Multidimensional
Usage	Structured, Repetitive	Ad hoc
Unit of Work	Short, Simple transaction	Complex Query
Access	Read/Write	Read Mostly
Operation	Index/hash on prim. Key	Lots of scan
# Records accessed	Tens	Millions
#Users	Thousands	Hundreds
Db size	100 MB-GB	100GB-TB
Metric	Transaction throughput	Query throughput response

**Note:** Above table depicts the comparison between **Operational** and **Analytical systems**.

## Summary

➤ In this lesson, you have learnt:

- OLAP allows users to view data from various perspectives.
- The nature of OLAP analysis varies in multiple ways of using it.
  - Aggregation
  - Comparison
  - Ranking
  - Access to data
  - Complex criteria specification



## Summary

- OLAP is used to analyze the data in the Data Warehouse.
- Different types of OLAP are:
  - MOALP
  - HOLAP
  - ROLAP



### Review Question

- **Question 1:** OLAP analysis is used for:
  - Option 1: Retrieving data
  - Option 2: Updating data
  - Option 3: Summarizing data
- **Question 2:** OLAP makes use of multidimensional data model.
  - True/ False
- **Question 3:** \_\_\_ OLAP operation helps for viewing data from any angle.



## Data Warehousing Concepts

### Lesson 6: Data Mining

June 12, 2014

Proprietary and Confidential

- 1 -



## Lesson Objectives

➤ In this lesson, you will learn about:

- Online Analytical Processing
- Data Mining
- The Knowledge Discovery Process
- Why Use Data Mining Today?
- Data Mining Usage
- Data Mining and Business Intelligence
- Types of Data used in Data Mining
- Data Mining Applications



### Lesson Objectives

➤ In this lesson, you will learn about (contd.):

- Data Mining Products
- Mining market



6.1: Data Mining  
**What is Data Mining?**

➤ **Data Mining is:**

- Subset of BI.
- Extraction of necessary information from data in large databases.
- Process of analyzing large databases to find valid, novel, useful, and understandable patterns.
- Process of efficient discovery in large databases and Data warehouses.

June 12, 2014

Proprietary and Confidential

- 4 -



**Data Mining:**

- **Data mining** is the way of analyzing data by exploring large databases. It helps in understanding the business by extracting necessary information from the databases. It allows you to understand the pattern and helps in predicting the behavior of it.
- Data mining helps in increasing the business and forecasting the chunks related to it at early stages. It includes finding patterns that are suitable for the organization.

6.2: The Knowledge Discovery Process (KDD)  
**Concept of KDD**

➤ **Data Mining is also known as Knowledge Discovery in Databases (KDD).**

- It involves mining on different kinds of data.
- It is a process of using raw data.
- It is a collection of powerful techniques.
- It refers to the automated extraction of hidden information from databases.
- It helps customers to detect previously undetected facts present in their business critical data.

June 12, 2014

Proprietary and Confidential

- 5 -

**IGATE**  
Spend Agilely Imagine

**The Knowledge Discovery Process (KDD):**

- **Data Mining** involves mining on different kinds of data such as Relational databases, Data warehouses, Transactional databases, Advanced DB systems and information repositories, Object-oriented and object-based databases, Text databases and multimedia databases, Heterogeneous and legacy databases. Data mining is the process of using raw data to infer important business information. It is a collection of powerful techniques for analyzing large amounts of data. Data mining tools can access data directly in the Data Warehouse.
- The advantage of mining is that no separate copy of data is needed for data mining. Data may not be organized in a way that is efficient for the tool. Data Mining is done by running a software that examines a database and looks for patterns in the data. Data Mining will not tell users about patterns in data that users may not have thought about. Data mining is used to try and mine key information from a Data warehouse to find **patterns in data**. Data mining allows organizations to collect information and make themselves more productive and beat their competitors.

6.1: Need of Data Mining  
**Why Use Data Mining Today?**

- **Human analytical skills are inadequate when:**
  - Volume and dimensionality of the data increases.
  - Data growth rate is high.
- **Data Mining is used for availability of:**
  - Data
  - Data Storage
  - Computation of Data

June 12, 2014

Proprietary and Confidential

- 6 -



**Need for Data Mining:**

- Data mining is essential because of the following utilities:
  - Data mining helps to identify why customers buy certain products.
  - Data mining provides the ideas for very direct marketing.
  - Data mining provides the ideas for shelf placement.
  - It helps for training of employees versus employee retention.
  - It helps to identify employee benefits.

6.4: Use of Data Mining  
**Usage**

➤ Here are some instances where Data Mining is essential:

- The US Government needs to track fraudulent events.
- A Supermarket is aspiring to become an information broker.
- Basketball teams need it to track game strategy.
- Cross Selling
- Target Marketing
- Holding on to good customers
- Weeding out bad customers

### Usage Scenarios

➤ **Data warehouse mining is used in the following scenarios:**

- Assimilate data from operational sources
- Mine static data
- Mining log data
- Continuous mining in process control

➤ **Stages in mining:**

- Data selection -> Pre-processing-> cleaning -> Transformation -> Mining -> Result evaluation -> Visualization

June 12, 2014

Proprietary and Confidential

- 8 -



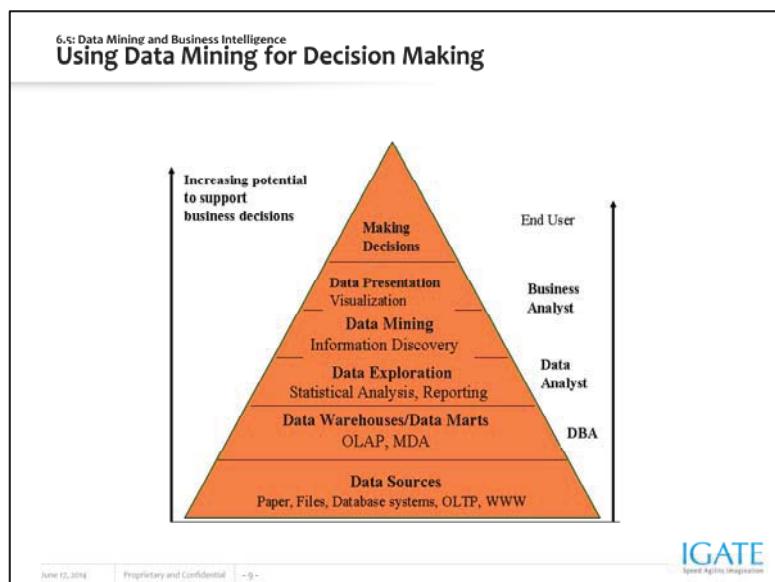
### Use of Data Mining:

#### **Usage scenarios:**

- Data warehouse mining assimilates data from operational sources.
- Data warehouse mining mines static data.
- Mining log data.
- Continuous mining in process control.

#### **Stages in mining:**

1. Data selection
2. Pre-processing: cleaning
3. Transformation
4. Mining
5. Result evaluation
6. Visualization



#### Data Mining and Business Intelligence:

- Data Mining has grown drastically in many businesses. Data Mining has become very popular since it helps in increasing organization's profit and achieving the target.
- When Data mining gets involved in Business Intelligence, it actually helps in understanding the functionality of the organization. It helps in increasing the potential for supporting the business decisions. It makes the data visible in a visual form to the business analysts. It helps in exploring data in terms of reporting and statistical analysis.
- Data Mining along with Business Intelligence takes the following steps in logical progression:
  - **Data Source:** Typically data is sourced from transaction processing systems (Manufacturing, ERP, Sales).
  - **Data Marts (OLAP & MDA)/DBA:** You may want to customize your warehouse's architecture for different groups within your organization. You can do this by adding data marts, which are systems designed for a particular line of business.
  - **End User/Making Decision:** The principle purpose of Data warehousing is to provide information to the business user for strategic decision making.

## 6.6: Types of Data

**Types of Data used in Data Mining****➤ The following types of data is drilled in Data Mining:**

- Relational data and transactional data
- Spatial and temporal data, spatio-temporal observations
- Time-series data
- Text
- Images, video
- Mixtures of data
- Sequence data
- Features from processing other data sources

6.7: Data Mining Applications

## Examples of Data Mining Applications

➤ **Here are some examples of Data Mining Applications:**

- Banking: Loan / Credit card approval
- Customer Relationship Management
- Targeted marketing
- Fraud detection: Telecommunications, Financial transactions
- Manufacturing and Production
- Web site/store design and promotion

June 12, 2014

Proprietary and Confidential

- 11 -



### **Data Mining Applications:**

Let us discuss some examples of Data Mining Applications:

- Banking: loan/credit card approval:
  - Predict good customers based on old customers
- Customer Relationship Management:
  - Identify those who are likely to leave for a competitor.
- Targeted marketing:
  - Identify likely responders to promotions.
- Fraud detection: Telecommunications, Financial transactions
  - From an online stream of events, identify fraudulent events.
- Manufacturing and production:
  - Automatically adjust knobs when process parameter changes.
- Web site/store design and promotion:
  - Find affinity of visitor to pages and modify layout.

### 6.8: Data Mining Products Examples of Data Mining Products

➤ Here are some examples of Data Mining Products:

- DataMind: neurOagent
- Information Discovery: IDIS
- SAS Institute: SAS/Neuronets

6.4: Data Mining Market  
**Mining Market and Vendors**

- There are around 20 to 30 mining tool vendors.
- Major tool players:
  - Clementine
  - IBM's Intelligent Miner
  - SGI's MineSet
  - SAS's Enterprise Miner
- Many embedded products:
  - Fraud detection
  - Electronic commerce applications
  - Health care
  - Customer Relationship Management: Epiphany

## Summary

➤ **In this lesson, you have learnt:**

- Data Mining is the way of analyzing data by exploring the large databases.
- Data Mining is used to mine key information from a data warehouse.
- It helps in exploring data in terms of reporting and statistical analysis.



**Review Question**

- **Question 1:** Data exploration for statistical analysis is done by:
  - Option 1: DBA
  - Option 2: Business analyst
  - Option 3: Data analyst
- **Question 2:** Data Mining is a subset of DW.
  - True/ False
- **Question 3:** Mining is also known as \_\_\_\_.



**Review - Match the Following**

1. End user

2. Business analyst

3. Data analyst

A. Data mining

B. Data warehouse

C. Data presentation

D. Making decisions

E. Data exploration

