# Lecture 1
## Data Warehouse and OLAP Technology

*CA4010: Data Warehousing and Data Mining*
2016/2017 Semester 1

Dr. Mark Roantree
Dublin City University

# Agenda

**1** **Basic Concepts**

**2** **Data Cubes and OLAP**

**3** **Warehouse Design and Usage**
- 4 Step Process
- Case Study

**4** **Warehouse Architecture**

**5** **Warehouse Implementation**
- Conforming
  - Conformed Dimensions
  - Conformed Facts

**Data Warehouse and OLAP Technology**

DCU

Basic Concepts

Data Cubes and OLAP

Warehouse Design and Usage

4 Step Process

Case Study

Warehouse Architecture

Warehouse Implementation

Conforming

Conformed Dimensions

Conformed Facts

# What is a Data Warehouse?

- A decision support database that is maintained separately from the organization's operational database.
- Support information processing by providing a solid platform of consolidated, historical data for analysis.
- Data Warehousing: The process of constructing and using data warehouses.

## Inmon's Definition

A data warehouse is a subject-oriented, integrated, time-variant, and nonvolatile collection of data in support of management's decision-making process.

# Data Warehouse: Subject-Oriented

- Organized around major subjects, such as customer, product, sales
- Focusing on the modeling and analysis of data for decision makers, not on daily operations or transaction processing
- Provide a simple and concise view around particular subject issues by excluding data that are not useful in the decision support process

### Inmon's Definition

A data warehouse is a **subject-oriented**, integrated, time-variant, and nonvolatile collection of data in support of management's decision-making process.

# Data Warehouse: Integrated

- Constructed by integrating multiple, heterogeneous data sources
  relational databases, flat files, on-line transaction records
- Data cleaning and data integration techniques are applied.
  Ensure consistency in naming conventions, encoding structures, attribute measures, among different data sources
  Hotel price: currency, tax, breakfast covered, etc.
  When data is moved to the warehouse, it is converted.

### Inmon's Definition

A data warehouse is a subject-oriented, **integrated**, time-variant, and nonvolatile collection of data in support of management's decision-making process.

# Data Warehouse: Time Variant

- The time horizon for the data warehouse is significantly longer than that of operational systems
  - Operational database: current value data
  - Data warehouse data: provide information from a historical perspective (e.g., past 5-10 years)
- Every key structure in the data warehouse
  - Contains an element of time, explicitly or implicitly
  - But the key of operational data may or may not contain *time element*

### Inmon's Definition

A data warehouse is a subject-oriented, integrated, **time-variant**, and nonvolatile collection of data in support of management's decision-making process.

# Data Warehouse: Non-volatile

- A physically separate store of data transformed from the operational environment
- Operational update of data does not occur in the data warehouse environment
    - Does not require transaction processing, recovery, and concurrency control mechanisms
    - Requires only two operations in data accessing: initial loading of data and access of data

### Inmon's Definition

A data warehouse is a subject-oriented, integrated, time-variant, and **non-volatile** collection of data in support of management's decision-making process.

# OLTP vs. OLAP

| Feature | OLTP | OLAP |
|---|---|---|
| Characteristic | operational processing | informational processing |
| Orientation | transaction | analysis |
| User | clerk, DBA, database professional | knowledge worker (e.g., manager, executive, analyst) |
| Function | day-to-day operations | long-term informational requirements, decision support |
| DB design | ER based, application-oriented | star/snowflake, subject-oriented |
| Data | current; guaranteed up-to-date | historical; accuracy maintained over time |
| Summarization | primitive, highly detailed | summarized, consolidated |
| View | detailed, flat relational | summarized, multidimensional |
| Unit of work | short, simple transaction | complex query |
| Access | read/write | mostly read |
| Focus | data in | information out |
| Operations | index/hash on primary key | lots of scans |
| Number of records accessed | tens | millions |
| Number of users | thousands | hundreds |
| DB size | 100 MB to GB | 100 GB to TB |
| Priority | high performance, high availability | high flexibility, end-user autonomy |
| Metric | transaction throughput | query throughput, response time |

**Figure 1:** OLTP vs. OLAP: A Comparison

# Properties of a Warehouse: Data Engineering

- Integrated data
- Detailed and summarised data
- Historical data
- Metadata

# Integrated Data

- Integrated data allows the miner to easily examine very high volumes of data.
- Without integrated data, too much time is spent cleaning and transforming data before data mining can commence effectively.
- Keys must be reconstituted; encoded values reconciled; data structures standardized; and many more tasks, before data mining.
- All of these tasks have been completed *before* the data warehouse.

# Detailed and Summarised Data

- Detailed data is necessary when the miner wishes to examine data in its most *granular* form.
- Very low levels of detail hide important patterns that require the study of very low level (detailed) data.
- By contrast, summarised data ensures the reuse of a prior analysis (using a pre-existing *cube*).
- Summarized data ensures the miner can build on the work of others rather than build everything from the ground up.
- This saves significant amounts of unnecessary work for the miner.

# Historical Data

- Historical data is important to the miner because important knowledge can be hidden there.
- A miner who has to work with only very current information can never detect trends and long-term patterns of behavior.
- Historical information is crucial to understanding the seasonality of business and the larger cycles of business to which every corporation is subject.

# Metadata

- Metadata serves as a *road map* to the miner, who uses metadata to describe the *context* of information.
- When examining information *over time*, context becomes as relevant as content.
- In other words, raw data becomes very difficult for the data miner to work with when there is no explanation for the *meaning* of the data.
- The data warehouse provides a platform for successful and efficient exploration of the world of data, with the miner exploiting metadata for more powerful analyses.

# Why a Separate Data Warehouse?

- High performance for both systems
  - DBMS tuned for OLTP: access methods, indexing, concurrency control, recovery
  - Warehouse tuned for OLAP: complex OLAP queries, multidimensional view, consolidation
- Different functions and different data
  - missing data: Decision support requires historical data which operational DBs do not typically maintain
  - data consolidation: DS requires consolidation (aggregation, summarization) of data from heterogeneous sources
  - data quality: different sources typically use inconsistent data representations, codes and formats which have to be reconciled

# Data Warehouse: A Multi-Tiered Architecture

**Figure 2:** ETL Architecture

# Three Data Warehouse Models

- **Enterprise Warehouse**: collects all of the information about subjects spanning the entire organization
- **Data Mart**: a subset of corporate-wide data that is of value to a specific groups of users.
  Its scope is confined to specific, selected groups, such as marketing data mart.
- **Virtual warehouse**: A set of views over operational databases.
  Only some of the possible summary views may be materialized

# Extraction, Transformation, and Loading (ETL)

- Data extraction: get data from multiple, heterogeneous, and external sources
- Data cleaning: detect errors in the data and rectify them when possible
- Data transformation: convert data from legacy or host format to warehouse format
- Load: sort, summarize, consolidate, compute views, check integrity, and build indicies and partitions
- Refresh: propagate the updates from the data sources to the warehouse

# Metadata Repository

Meta data is the data defining warehouse objects. It stores:

- Description of the data warehouse schema, view, dimensions, hierarchies, derived data definitions, data mart locations and contents
- Operational meta-data: data lineage (history of migrated data and transformation path), currency of data (active, archived, or purged), monitoring information (warehouse usage statistics, error reports, audit trails)
- The algorithms used for summarization
- The mapping from operational environment to the data warehouse
- Data related to system performance: warehouse schema, view and derived data definitions
- Business data: business terms and definitions, ownership of data, charging policies

# From Tables and Spreadsheets to Data Cubes

- A data warehouse is based on a **multidimensional** data model which views data in the form of a **data cube**.
- A data cube, such as sales, allows data to be modeled and viewed in multiple dimensions
    - Dimension tables, such as item (item_name, brand, type), or time(day, week, month, quarter, year)
    - Fact table contains measures (such as dollars_sold) and keys to each of the related dimension tables
- In data warehousing, an *n*-D base cube is called a **base cuboid**. The top most *0*-D cuboid, which holds the highest-level of summarization, is called the **apex cuboid**.
- The *lattice* of cuboids forms a data cube.

# Cube: A Lattice of Cuboids

**Data Warehouse and OLAP Technology**

**DCU**

Basic Concepts

Data Cubes and OLAP

Warehouse Design and Usage
4 Step Process
Case Study

Warehouse Architecture

Warehouse Implementation
Conforming
Conformed Dimensions
Conformed Facts

1.20

**Figure 3:** 4-D Data Cube Representation

# Conceptual Modeling of Data Warehouses

Modeling data warehouses: **dimensions** & **measures**

- **Star schema**: A fact table in the middle connected to a set of dimension tables.

- **Snowflake schema**: A refinement of star schema where some dimensional hierarchy is normalized into a set of smaller dimension tables, forming a shape similar to snowflake.

- **Fact constellations**: Multiple fact tables share dimension tables, viewed as a collection of stars, called *galaxy* schema or *fact constellation*.

DCU

# Star Schema

- Star schema: The most common modeling paradigm is the star schema, in which the data warehouse contains
  1. a large central table (fact table) containing the bulk of the data, with no redundancy, and
  2. a set of smaller attendant tables (dimension tables), one for each dimension.
- The schema graph resembles a starburst, with the dimension tables displayed in a radial pattern around the central fact table.

# Star Schema for a Sales Data Warehouse

Data Warehouse
and OLAP
Technology

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage
4 Step Process
Case Study

Warehouse
Architecture

Warehouse
Implementation
Conforming
Conformed Dimensions
Conformed Facts

**Figure 4:** AllElectronics Star Schema

- Sales are considered along four dimensions: *time*, *item*, *branch*, and *location*.
- The schema contains a central **fact table** for *sales* that contains keys to each of the four dimensions.
- Contains two measures: *dollars sold* and *units sold*.

1.23

# Star Schema Model

- Each dimension is represented by a single table, containing a set of attributes.

- The *location* dimension table contains the attribute set *{location key, street, city, province or state, country}*.

- This constraint may introduce some redundancy: Vancouver and Victoria are both cities in the Canadian province of British Columbia.

- This will create redundancy among the attributes province or state and country.

- Moreover, the attributes within a dimension table may form either a hierarchy (total order) or a lattice (partial order).

# Snowflake Schema for a Sales Data Warehouse

The snowflake schema is a variant of the star schema model, where some dimension tables are normalized, thereby further splitting the data into additional tables.

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage
 4 Step Process
 Case Study

Warehouse
Architecture

Warehouse
Implementation
 Conforming
 Conformed Dimensions
 Conformed Facts

**time**
dimension table

| time_key |
| day |
| day_of_week |
| month |
| quarter |
| year |

**sales**
fact table

| time_key |
| item_key |
| branch_key |
| location_key |
| dollars_sold |
| units_sold |

**item**
dimension table

| item_key |
| item_name |
| brand |
| type |
| supplier_key |

**supplier**
dimension table

| supplier_key |
| supplier_type |

**branch**
dimension table

| branch_key |
| branch_name |
| branch_type |

**location**
dimension table

| location_key |
| street |
| city_key |

**city**
dimension table

| city_key |
| city |
| province_or_state |
| country |

**Figure 5:** Sample Snowflake Schema

# Snowflake Schema

- The major difference between the *snowflake* and *star* schema models is that snowflake dimensions may use normalization to reduce redundancies.

- However, this saving of space is negligible in comparison to the typical magnitude of the fact table.

- Furthermore, the *snowflake* structure can reduce the effectiveness of browsing, since more joins will be needed to execute a query.

- Consequently, the system performance may be adversely impacted.

- Hence, although the *snowflake* schema reduces redundancy, it is not as popular as the *star* schema in data warehouse design.

# Fact Constellation

- Sophisticated applications may require multiple fact tables to share dimension tables.
- This kind of schema can be viewed as a collection of stars, and hence is called a galaxy schema or a fact constellation.

DCU

Basic Concepts

Data Cubes and OLAP

Warehouse Design and Usage
4 Step Process
Case Study

Warehouse Architecture

Warehouse Implementation
Conforming
Conformed Dimensions
Conformed Facts

1.27

# Fact Constellation Schema for Sales and Shipping

Data Warehouse
and OLAP
Technology

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage
4 Step Process
Case Study

Warehouse
Architecture

Warehouse
Implementation
Conforming
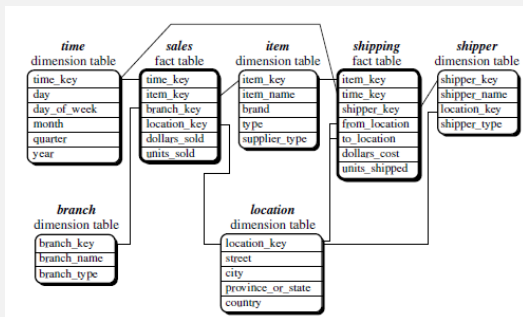Conformed Dimensions
Conformed Facts

1.28

**Figure 6:** Sample Fact Constellation Schema

- This schema specifies two fact tables: *sales* and *shipping*.
- Shipping has 5 dimensions (keys): *item key, time key, shipper key, from location*, and *to location*, and two measures: *dollars cost* and *units shipped*.
- A *fact constellation* schema allows dimension tables to be shared between fact tables.

# Concept Hierarchy

- A concept hierarchy defines a sequence of mappings from a set of low-level concepts to higher-level, more general concepts.
- Consider a concept hierarchy for the dimension location: City values for location include Vancouver, Toronto, New York, and Chicago.
- Each city can be mapped to the province or state to which it belongs.
- Provinces and states can in turn be mapped to the country to which they belong.
- These mappings form a concept hierarchy for the dimension location, mapping a set of low-level concepts (cities) to higher-level, more general concepts (countries).

# Location Concept Hierarchy

**Figure 7:** Location Concept Hierarchy

# OLAP Operations

- Roll up (drill-up): summarize data by climbing up hierarchy or by dimension reduction
- Drill down (roll down): reverse of roll-up from higher level summary to lower level summary or detailed data, or introducing new dimensions
- Slice and dice: project and select
- Pivot (rotate): reorient the cube, visualization, 3D to series of 2D planes
- drill across: involving (across) more than one fact table
- drill through: through the bottom level of the cube to its back-end relational tables (using SQL)

# Typical OLAP Operations

Data Warehouse
and OLAP
Technology

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage
4 Step Process
Case Study

Warehouse
Architecture

Warehouse
Implementation
Conforming
Conformed Dimensions
Conformed Facts

# DW Design: A Business Analysis Framework

- Top-down view: allows selection of the relevant information necessary for the data warehouse
- Data source view: exposes the information being captured, stored, and managed by operational systems
- Data warehouse view: consists of fact tables and dimension tables
- Business query view: sees the perspectives of data in the warehouse from the view of end-user

# Four-Step Dimensional Design Process

1. Select the business process to model.
2. Declare the grain of the business process
3. Choose the dimensions that apply to each fact table row
4. Identify the numeric facts that will populate each fact table row.

# Step 1. Select the business process to model.

- A process is a natural business activity performed in an organization that typically is supported by a source data-collection system.
- **Business measurement** processes provide the performance measurements users need to analyze in the data warehouse.
- Example business processes include raw materials purchasing, orders, shipments, invoicing, inventory, and general ledger.
- Note: we do not refer to an organizational business department or function when we talk about business processes.

- For example, build a single dimensional model to handle orders data rather than building separate models for the sales and marketing departments, which both want to access orders data.
- Focus on business processes (not departments) to deliver *consistent* information more economically throughout the organization.

# Step 2. Declare the grain of the business process

- This means specifying exactly what an individual fact table row represents.
- The grain conveys the level of detail associated with the fact table measurements.
- It provides the answer to the question: *How do you describe a single row in the fact table?*
- It is virtually impossible to reach closure in step 3 without declaring the grain.

**Example grain declarations include:**

- An individual line item on a customer's retail sales ticket as measured by a scanner device
- A line item on a bill received from a doctor
- An individual boarding pass to get on a flight
- A daily snapshot of the inventory levels for each product in a warehouse
- A monthly snapshot for each bank account

# Step 3. Choose the dimensions that apply to each fact table row

Data Warehouse
and OLAP
Technology

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage

4 Step Process

Case Study

Warehouse
Architecture

Warehouse
Implementation

Conforming

Conformed Dimensions

Conformed Facts

- Dimensions emerge from asking: *How do business people describe the data that results from the business process?*
- Boost fact tables with a robust set of dimensions representing all possible descriptions.
- Once clear about the grain, dimensions can be identified easily.
- For each dimension, list all discrete, text-like attributes that will flesh out each dimension table.

# Step 4. Identify the numeric facts that will populate each fact table row.

- Facts are determined by answering the question: *What are we measuring?*
- Business users are keenly interested in analyzing these business process performance measures.
- All candidate facts in a design must be true to the grain defined in step 2.
- Facts that clearly belong to a different grain *must be in a separate fact table*.
- Typical facts are numeric additive figures such as quantity ordered or dollar cost amount.

# Inputs to 4-step Design

Consider both business users' requirements *and* realities of source data, to make decisions regarding the 4 steps.

*Business Requirements*

Dimensional Model
1. Business Process
2. Grain
3. Dimensions
4. Facts

*Data Realities*

**Figure 9:** Key input to the four-step dimensional design process

# Large Grocery Chain: Case Study

- Business has 100 grocery stores spread over multiple regions.
- Each of the stores has a full complement of departments, including grocery, frozen foods, dairy, meat, produce, bakery, floral, and health/beauty aids.
- Each store has roughly 60,000 individual products, called stock keeping units (SKUs).
- About 55,000 of the SKUs come from outside manufacturers and have bar codes imprinted on the product package.

# Case Study: Grain

Data Warehouse and OLAP Technology

DCU

Basic Concepts

Data Cubes and OLAP

Warehouse Design and Usage
4 Step Process
Case Study

Warehouse Architecture

Warehouse Implementation
Conforming
Conformed Dimensions
Conformed Facts

- These bar codes are called universal product codes (UPCs).
- UPCs are at the same grain as individual SKUs.
- Each different package variation of a product has a separate UPC and thus, a separate SKU.
- The remaining 5,000 SKUs come from departments such as meat, produce, bakery, or floral.
- While these products don't have nationally recognized UPCs, the grocery chain assigns SKU numbers to them.

# Step 1. Select the Business Process

- The first dimensional model built should be that with most impact: it should answer the most pressing business questions and be readily accessible for data extraction.

- In our retail case study, management wants to better understand customer purchases as captured by the POS system.

- Thus, POS retail sales is selected.

- This data allows them to analyze what products are selling, in which stores, on what days, under what promotional conditions.

# Step 2. Declare the Grain

- *What level of data detail should be made available in the dimensional model?*
- Tackling data at its lowest, most atomic grain makes sense on multiple fronts.

# Atomic Data: Advantages

1. Atomic data is highly dimensional.
2. The more detailed and atomic the fact measurement, the more things we know for sure (correct, certain).
3. All those things we know for certain *translate into dimensions*.
4. Atomic data provides maximum analytic flexibility because it can be *constrained* and *rolled up* in every way possible.

# Higher Grains: Disadvantages

- Declare higher-level grains for a business process that represent an aggregation of the most atomic data.
- However, as soon as we select a higher-level grain, users are limited to fewer and/or potentially less detailed dimensions.
- The less granular model is immediately vulnerable to unexpected user requests to drill down into the details.
- Users hit an *analytic wall* when not given access to the atomic data.

# Step 3. Choose the Dimensions

**Data Warehouse and OLAP Technology**

DCU

Basic Concepts

Data Cubes and OLAP

Warehouse Design and Usage
4 Step Process
Case Study

Warehouse Architecture

Warehouse Implementation
Conforming
Conformed Dimensions
Conformed Facts

- Once the grain of the fact table has been chosen, the date, product, and store dimensions fall out immediately.

- We assume that the calendar date is the date value delivered to us by the POS system.

- Within the framework of the primary dimensions, we can ask whether other dimensions can be attributed to the data, such as the promotion under which the product is sold.

**Data Warehouse and OLAP Technology**

DCU

Basic Concepts

Data Cubes and OLAP

Warehouse Design and Usage

4 Step Process

Case Study

Warehouse Architecture

Warehouse Implementation

Conforming

Conformed Dimensions

Conformed Facts

- For the case study, assume we decide on the following descriptive dimensions: `date`, `product`, `store`, and `promotion`.
- In addition, include the POS transaction ticket number as a special dimension.
- We now have a preliminary schema as illustrated in Figure 10.
- Do not begin populating the dimension tables with descriptive attributes: instead complete the final step of the process (keep it simple for now!).

# Preliminary retail sales schema.

Data Warehouse
and OLAP
Technology

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage

4 Step Process

Case Study

Warehouse
Architecture

Warehouse
Implementation

Conforming

Conformed Dimensions

Conformed Facts

| Date Dimension |
| --- |
| Date Key (PK) |
| Date Attributes TBD |

| POS Retail Sales Transaction Fact |
| --- |
| Date Key (FK) |
| Product Key (FK) |
| Store Key (FK) |
| Promotion Key (FK) |
| POS Transaction Number |
| Facts TBD |

| Product Dimension |
| --- |
| Product Key (PK) |
| Product Attributes TBD |

| Store Dimension |
| --- |
| Store Key (PK) |
| Store Attributes TBD |

| Promotion Dimension |
| --- |
| Promotion Key (PK) |
| Promotion Attributes TBD |

**Figure 10:** Preliminary retail sales schema.

# Step 4. Identify the Facts

- There is a firm rule for this step: **the facts must be true to the grain**.
- Specifically, this is the individual line item on the POS transaction.
- At this point, it may be necessary to revise grain assumptions or our choice of dimensions.
- The facts collected by the POS system include the sales quantity, per unit sales price, and the sales amount.
- The sales amount is: *quantity by unit price*.

- The current fact table is shown in Figure 11.
- Three of the facts, *sales quantity*, *sales amount*, and *cost amount*, are perfectly additive across all the dimensions.
- We can slice and dice the fact table at will, and every sum of these three facts is valid and correct.

# Measured facts in the retail sales schema

| Date Dimension |
| --- |
| Date Key (PK)<br>Date Attributes TBD |

| Store Dimension |
| --- |
| Store Key (PK)<br>Store Attributes TBD |

| POS Retail Sales Transaction Fact |
| --- |
| Date Key (FK)<br>Product Key (FK)<br>Store Key (FK)<br>Promotion Key (FK)<br>POS Transaction Number<br>Sales Quantity<br>Sales Dollar Amount<br>Cost Dollar Amount<br>Gross Profit Dollar Amount |

| Product Dimension |
| --- |
| Product Key (PK)<br>Product Attributes TBD |

| Promotion Dimension |
| --- |
| Promotion Key (PK)<br>Promotion Attributes TBD |

**Figure 11:** Measured facts in the retail sales schema

# Computed Facts

- Compute gross profit by subtracting the *cost amount* from the *sales amount*.
- Although *computed*, this gross profit is also *perfectly additive* across *all* the dimensions.
- One can calculate the gross profit of any combination of products sold in any set of stores on any set of days.
- *Should a computed fact be stored physically in the database?*

# Computed Facts: Storage Benefits

# DCU

Basic Concepts

Data Cubes and OLAP

Warehouse Design and Usage
4 Step Process
Case Study

Warehouse Architecture

Warehouse Implementation
Conforming
Conformed Dimensions
Conformed Facts

- Computing and storage of facts eliminates potential user error.
- Storing them also ensures that all users and their reporting applications refer to gross profit consistently.
- Pre-computation is a form of optimisation.

# Non Additive Facts

- The gross margin can be calculated by dividing the gross profit by revenue (sales amount - cost).
- Gross margin is a non-additive fact because it cannot be summarized along any dimension.
- One can calculate the gross margin of any set of products, stores, or days by remembering to add the revenues and costs before dividing.

# Non Additive Facts: General Rule

- Percentages and ratios, such as gross margin, are non-additive.
- The *numerator* and *denominator* should be stored in the fact table.
- The ratio can be calculated in a data access tool for any slice of the fact table by remembering to calculate the ratio of the sums and *not* the sum of the ratios.
- Other attributes (eg. unit price in this case study) will also be non-additive (by their meaning).

# 3-tier DW Architecture
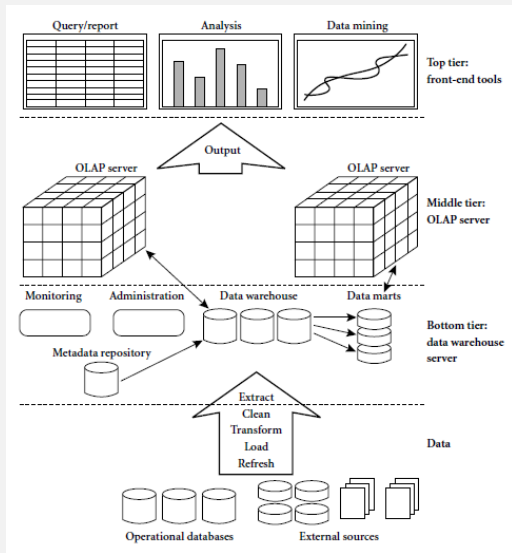
Data warehouses often adopt a three-tier architecture.

DCU

**Figure 12:** 3-tier DW Architecture

# Bottom Tier: DW Server

- The bottom tier is a warehouse database server that is almost always a relational database system.
- Back-end tools and utilities are used to feed data into the bottom tier from operational databases or other external sources (such as customer profile information provided by external consultants).
- These tools and utilities perform data extraction, cleaning, and transformation (to merge/standardize similar data from different sources), as well as load and refresh functions to update the data warehouse.
- The data are extracted using application program interfaces known as gateways.
- A gateway is supported by the underlying DBMS and allows client programs to generate SQL code to be executed at a server.
- This tier also contains a metadata repository, which stores information about the data warehouse and its contents.

# Middle Tier: OLAP Server

- The middle tier is an OLAP server that is typically implemented using either:
  1. a relational OLAP (ROLAP) model, that is, an extended relational DBMS that maps operations on multidimensional data to standard relational operations or
  2. a multidimensional OLAP (MOLAP) model, that is, a special-purpose server that directly implements multidimensional data and operations.

# Top Tier: Front End

- The top tier is a front-end client layer, which contains query and reporting tools, analysis tools, and/or data mining tools (e.g., trend analysis, prediction, and so on).

- From the architecture point of view, there are three data warehouse models: the enterprise warehouse, the data mart, and the virtual warehouse.

# Enterprise warehouse

- An enterprise warehouse collects all of the information about subjects spanning the entire organization.

- It provides corporate-wide data integration, usually from one or more operational systems or external information providers, and is cross-functional in scope.

- It typically contains detailed data as well as summarized data, and can range in size from a few gigabytes to hundreds of gigabytes, terabytes, or beyond.

- An enterprise data warehouse may be implemented on traditional mainframes, computer super servers, or parallel architecture platforms.

- It requires extensive business modeling and may take years to design and build.

# Data Mart

- A data mart contains a subset of corporate-wide data that is of value to a specific group of users.
- The scope is confined to specific selected subjects.
- For example, a marketing data mart may confine its subjects to customer, item, and sales.
- The data contained in data marts tend to be summarized.
- The implementation cycle of a data mart is more likely to be measured in weeks rather than months or years.
- However, it may involve complex integration in the long run if its design and planning were not enterprise-wide.

# Data Mart (2)

Data Warehouse
and OLAP
Technology

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage
4 Step Process
Case Study

Warehouse
Architecture

Warehouse
Implementation
Conforming
Conformed Dimensions
Conformed Facts

1.64

- Depending on the source of data, data marts can be categorized as *independent* or *dependent*.
- **Independent data marts** are sourced from one or more operational systems or external information providers, or from data generated locally within a particular department or geographic area.
- **Dependent data marts** are sourced directly from enterprise data warehouses.

# Virtual Warehouse

- A virtual warehouse is a set of views over operational databases.
- For efficient query processing, only some of the possible summary views may be materialized.
- A virtual warehouse is easy to build but requires excess capacity on operational database servers.

# Data Warehouse Bus Architecture

- By defining a *standard bus interface* for the data warehouse environment, separate data marts can be implemented by *different* groups at *different* times.
- The separate data marts can be integrated if they adhere to the standard.
- For long-term data warehouse success, use an architected, incremental approach to build the enterprise warehouse.
- This approach adopts the *data warehouse bus architecture*.

# Design Methodology

**The data warehouse bus architecture provides a rational approach to decomposing the enterprise data warehouse planning task.**

1. The DW team design a master suite of standardized dimensions and facts that have uniform interpretation across the enterprise.

2. This establishes the *data architecture framework*.

3. Build the first (next) data mart that closely adheres to the architecture.

4. Repeat Step 3 until enough data marts exist to make good on the promise of an integrated enterprise data warehouse.

# Benefits: reduce the Size of the Problem

**The data warehouse bus architecture provides a rational approach to decomposing the enterprise data warehouse planning task.**

- The bus architecture allows data warehouse managers to get the best of both worlds.

- They have an architectural framework that guides the overall design, but the problem has been divided into manageable data mart segments that can be implemented in realistic time frames.

- Separate data mart development teams follow the architecture guidelines while working fairly independently and asynchronously.

# Benefits: Easier Integration

- All flavors of relational and online analytical processing (OLAP)-based data marts can be full participants in the data warehouse bus if they are designed around *conformed dimensions and facts*.
- Data warehouses will inevitably consist of numerous separate machines with different operating systems and database management systems (DBMSs).
- If designed coherently, they will share a *uniform architecture* of conformed dimensions and facts that will allow them to be fused into an integrated whole.

# Data Warehouse Bus Matrix

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage
4 Step Process
Case Study

Warehouse
Architecture

Warehouse
Implementation
Conforming
Conformed Dimensions
Conformed Facts

- The rows of the bus matrix correspond to data marts.
- Create rows where sources/processes are different or where a process is too big.

**COMMON DIMENSIONS**

| BUSINESS PROCESSES | Date | Product | Store | Promotion | Warehouse | Vendor | Contract | Shipper |
|---|---|---|---|---|---|---|---|---|
| Retail Sales | X | X | X | X | | | | |
| Retail Inventory | X | X | X | | | | | |
| Retail Deliveries | X | X | X | | | | | |
| Warehouse Inventory | X | X | | | X | X | | |
| Warehouse Deliveries | X | X | | | X | X | | |
| Purchase Orders | X | X | | | X | X | X | X |

**Figure 13:** Sample Data Warehouse Bus Matrix

1.70

# Constructing the Bus Matrix

- The tool used to *create*, *document*, and *communicate* the bus architecture is the data warehouse bus matrix.
- Working in tabular fashion, lay out the business processes of the organization as matrix rows (should be sources of data).
- Matrix rows translate into **data marts** based on the organization's primary activities.
- Begin with data marts derived from a single primary source, known as *first-level data marts*.

# First Level Data Marts

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage
4 Step Process
Case Study

Warehouse
Architecture

Warehouse
Implementation
Conforming
Conformed Dimensions
Conformed Facts

- *Step 1*. Begin with *first-level data marts* as they minimize the risk of failure (overly ambitious implementation).
- Most of the overall risk of failure comes from taking on too much of the extract transformation-load (ETL) data staging design and development effort.
- In many cases, first-level data marts provide users with enough interesting data and ROI while the data mart teams focus on more difficult issues.

# Multisource Marts

- *Step 2*. Identify more complex multi-source marts
- Referred to as *consolidated data marts* because they typically cross business processes.
- It is prudent to focus on the first-level data marts as dimensional building blocks before tackling the task of consolidating.
- In some cases, the consolidated data mart is a simple union of data sets from the first-level data marts.
- Thus, it can often be easier to rebuild the consolidation *within* the data warehouse.

# Using the Matrix

- The columns of the matrix represent the common *dimensions* used across the enterprise.
- *Maybe create a comprehensive list of dimensions before filling in the matrix*?
- **X** indicates that the dimension column is related to the business process row.
- Looking across the rows is revealing as it is possible to see the *dimensionality* of each data mart at a glance.
- However, the real power of the matrix comes from looking at the columns as they depict the *interaction* between the *data marts* and *common dimensions*.

# Matrix Benefits

**Data Warehouse and OLAP Technology**

Basic Concepts

Data Cubes and OLAP

Warehouse Design and Usage
4 Step Process
Case Study

Warehouse Architecture

Warehouse Implementation
Conforming
Conformed Dimensions
Conformed Facts

1.75

Very powerful device for both planning and communication

- Defines the overall data architecture for the warehouse (by laying out rows & columns)
- Immediately recognise which dimensions warrant special attention given their participation in multiple data marts
- Helps prioritize which dimensions should be tackled first for *conformity* given their prominent roles.
- Visually: communicate effectively within and across data mart teams; convey the entire plan at once; directly communicate with senior IT and business management

# A Framework: Integration and Interoperability

- It is unacceptable to build separate data marts that ignore a framework to tie the data together.

- Isolated, independent data marts are worse than simply a lost opportunity for analysis.

- They deliver access to irreconcilable views of the organization and further enshrine the reports that cannot be compared with one another.

- Independent data marts become legacy implementations in their own right; by their very existence, they block the development of a coherent warehouse environment.

- Conformed dimensions are either identical or strict mathematical subsets of the most granular, detailed dimension.
- Conformed dimensions have:
  - consistent dimension keys,
  - consistent attribute column names,
  - consistent attribute definitions, and
  - consistent attribute values (which translates into consistent report labels and groupings).

# Non Conforming

- Dimension tables are not conformed if the attributes are labeled differently or contain different values.

- If a *customer* or *product* dimension is deployed in a non-conformed manner, then either the separate data marts cannot be used together or, worse, attempts to use them together will produce **invalid** results.

- Conformed dimensions come in several different flavors.
- At the most basic level, conformed dimensions have *identical meaning* with every possible fact table to which they are joined.
- The *date* dimension table connected to the *sales facts* is identical to the *date* dimension table connected to the *inventory facts*.
- In fact, the conformed dimension may be the *same physical table* within the database.
- In general, the date dimensions in both data marts will have the same number of rows, same key values, same attribute labels, same attribute definitions, and same attribute values.

# Conforming and Grain

- Most conformed dimensions are defined naturally at the *most granular level possible*.

- The grain of the *customer* dimension will be the individual customer.

- The grain of the *product* dimension will be the lowest level at which products are tracked in the source systems.

- The grain of the *date* dimension will be the individual day.

# Different Grains

- Sometimes dimensions are needed at a rolled-up level of granularity.
- Why? Perhaps the roll-up dimension is required because the fact table *represents aggregated facts* that are associated with *aggregated dimensions*.
- This would be the case if we had a weekly inventory snapshot *in addition* to our daily snapshot.
- In other situations, the facts simply may be generated by another business process at a higher level of granularity.

# Conforming by Roll-up

Data Warehouse
and OLAP
Technology

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage

4 Step Process

Case Study

Warehouse
Architecture

Warehouse
Implementation

Conforming

Conformed Dimensions

Conformed Facts

1.82

- Assume a sales process captures data at the **atomic** *product level*, while forecasting generates data at the *brand level*.

- It is not possible to share a single product dimension table across the two business process schemas because the granularity is different.

- Product and brand dimensions still conform if the brand table were a *strict subset* of the atomic product table.

- Attributes that are common to both the detailed and rolled-up dimension tables, such as the *brand* and *category* descriptions, should be **labeled**, **defined**, and **valued** identically in both tables, as Figure 14.

- Roll-up dimensions conform to atomic dimension if they are a strict subset of that atomic dimension.

| **Product Dimensions** |
| --- |
| Product Key (PK) |
| Product Description |
| SKU Number (Natural Key) |
| Brand Description |
| Subcategory Description |
| Category Description |
| Department Description |
| Package Type Description |
| Package Size |
| Fat Content Description |
| Diet Type Description |
| Weight |
| Weight Units of Measure |
| Storage Type |
| Shelf Life Type |
| Shelf Width |
| Shelf Height |
| Shelf Depth |
| … and more |

Conforms →

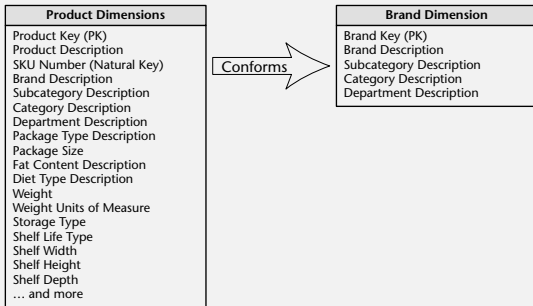| **Brand Dimension** |
| --- |
| Brand Key (PK) |
| Brand Description |
| Subcategory Description |
| Category Description |
| Department Description |

**Figure 14:** Conforming Rollup Dimension Subsets

1.83

- Another case of conformed dimension subsetting occurs when two dimensions are at the same level of detail but one represents only a subset of rows.

- For example, a corporate product dimension contains data for a *full* portfolio of products across multiple disparate lines of business, as illustrated in Figure 15.

- Analysts in the separate businesses may want to view only their subset of the corporate dimension, restricted to the product rows for their business.

- By using a subset of rows, they are not obliged to use the entire product set for the organization.
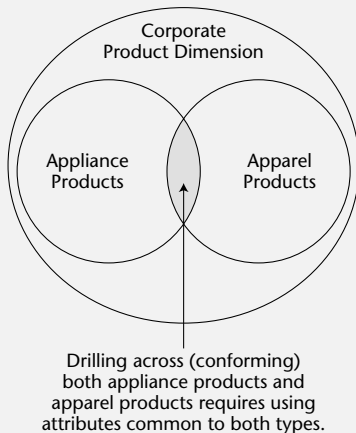
**Figure 15:** Conforming Rollup Dimension Subsets

# Issues with Dimension Row Subsetting

- The fact table joined to this subsetted dimension *must be limited* to the same subset of products.
- An attempt to use a subset dimension while accessing a fact table consisting of the complete product set, may result in unexpected query results.
- Technically, referential integrity would be violated.
- It is necessary to be aware of the potential for user confusion or error with dimension row subsetting.

# Fact Table Comparison

- There are three fundamental types of fact tables: **transaction**, **periodic snapshot**, and **accumulating snapshot**.
- These three fact table variations are not totally dissimilar because they share *conformed dimensions*.
- These are the keys to building separate fact tables that can be used together with common, consistent filters and labels.
- While the dimensions are shared, the administration and rhythm of the three fact tables are quite different.

# Transaction Fact Tables

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage
4 Step Process
Case Study

Warehouse
Architecture

Warehouse
Implementation
Conforming
Conformed Dimensions
Conformed Facts

- Most fundamental view of business's operations is at the *individual transaction level*.

- These fact tables represent an event that occurred at an instantaneous point in time.

- A row exists in the fact table for a given customer or product *only* if a transaction event occurred.

- On the other hand, a specific customer or product should be linked to *multiple* rows in the fact table as the customer or product is involved in multiple transactions.

- Transaction data often is structured quite easily into a dimensional framework.

- The lowest-level data is the most naturally dimensional data, supporting analyses that cannot be done on summarized data, and facilitating analysis of behaviour in extreme details.

- Even with transaction-level data, there is an entire class of urgent business questions that are impractical to answer using only transaction detail.

- Dimensional analysts cannot survive on transactions alone!

# Periodic Snapshot Fact Tables

- Periodic snapshots are needed to see the *cumulative performance* of the business at regular, predictable time intervals.

- Here, we take a snapshot of the activity at the end of a day, week, or month, then another picture at the end of the next period, and so on.

- Each periodic snapshots is stacked consecutively into the fact table.

- The periodic snapshot fact table is generally the sole mechanism for easy retrieval of a regular, predictable, trendable view of the key business performance metrics.

# Transaction and Periodic can provide the same data

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage
4 Step Process
Case Study

Warehouse
Architecture

Warehouse
Implementation
Conforming
Conformed Dimensions
Conformed Facts

- While transactions equate to little pieces of revenue, we can move easily from individual transactions to a daily snapshot by *summing the transactions*.
- In this situation, the periodic snapshot represents an *aggregation* of the transactional activity that occurred during that time period.
- We build the summary snapshot for performance reasons.
- The design of the snapshot table is closely related to the design of its companion transaction table in this case.

# Where to use

DCU

- In many cases, transactions are not components of revenue.
- Credit cards create transactions but the bank's source of revenue occurs when fees are assessed.
- In this situation, it is not possible to rely on transactions alone to analyze revenue performance.
- Not only would crawling through the transactions be time-consuming, but also the logic required to interpret the effect of different kinds of transactions on revenue or profit can have significant complications.

# Data Source

- The periodic snapshot benefits users with a quick, flexible view of revenue.

- It is quite possible the data for this snapshot schema is sourced directly from an operational system.

- If not, the warehouse staging area must incorporate very complex logic to interpret the financial impact of each transaction type correctly at data load time.

# Accumulating Snapshot Fact Tables

DCU

Basic Concepts

Data Cubes and OLAP

Warehouse Design and Usage
4 Step Process
Case Study

Warehouse Architecture

Warehouse Implementation
Conforming
Conformed Dimensions
Conformed Facts

- An accumulating snapshot is a type of fact table that records a single row for something the enterprise tracks closely, such as a problem product or a single customer's credit.

- They represent an *indeterminate* time span, covering the complete life of a transaction or *discrete* product (or customer).

- Accumulating snapshots almost always have multiple date stamps, representing the predictable major events or phases that take place during the course of a lifetime.

- Generally records an additional date column to indicate when the snapshot row was last updated.

# Updating the Data Warehouse

- Unlike other kinds of fact tables, the accumulating snapshot is intended to be *updated*.
- For example, dates are adjusted each time one of the milestones is reached.
- There may also be facts that track the number of days (or minutes) spent between each milestone.
- These *lags* in time are a convenience: they can be computed from the dates.
- Building them into the fact table makes analysis much easier, but does require that the ETL process revisit rows on a regular basis, rather than when status changes.

DCU

Basic Concepts

Data Cubes and OLAP

Warehouse Design and Usage
4 Step Process
Case Study

Warehouse Architecture

Warehouse Implementation
Conforming
Conformed Dimensions
Conformed Facts

- Sometimes accumulating and periodic snapshots work in conjunction.
- Such is the case when we build the monthly snapshot incrementally by adding the *effect* of each day's transactions to an accumulating snapshot.
- If we normally think of the data warehouse as storing 36 months of historical data in the periodic snapshot, then the current rolling month would be month 37.
- Ideally, when the last day of the month has been reached, the accumulating snapshot simply becomes the new regular month in the time series, and a new accumulating snapshot is started the next day.

# Transactions and Snapshots: the Holistic View

- Transactions and snapshots complement nicely in dimensional data warehouses.
- Used together, companion transaction and snapshot fact tables provide a complete view of the business.
- They exist separately as it is very difficult to combine these two contrasting perspectives.
- Data redundancy is not an issue: the goal is to publish data so that the organization can analyze it effectively.

# Summary

Data Warehouse
and OLAP
Technology

DCU

Basic Concepts

Data Cubes and
OLAP

Warehouse Design
and Usage
4 Step Process
Case Study

Warehouse
Architecture

Warehouse
Implementation
Conforming
Conformed Dimensions
Conformed Facts

| CHARACTERISTIC | TRANSACTION GRAIN | PERIODIC SNAPSHOT GRAIN | ACCUMULATING SNAPSHOT GRAIN |
|---|---|---|---|
| Time period represented | Point in time | Regular, predictable intervals | Indeterminate time span, typically short-lived |
| Grain | One row per transaction event | One row per period | One row per life |
| Fact table loads | Insert | Insert | Insert and update |
| Fact row updates | Not revisited | Not revisited | Revisited whenever activity |
| Date dimension | Transaction date | End-of-period date | Multiple dates for standard milestones |
| Facts | Transaction activity | Performance for predefined time interval | Performance over finite lifetime |

**Figure 16:** Fact Table Comparison