

1/30 4:33:07 ***

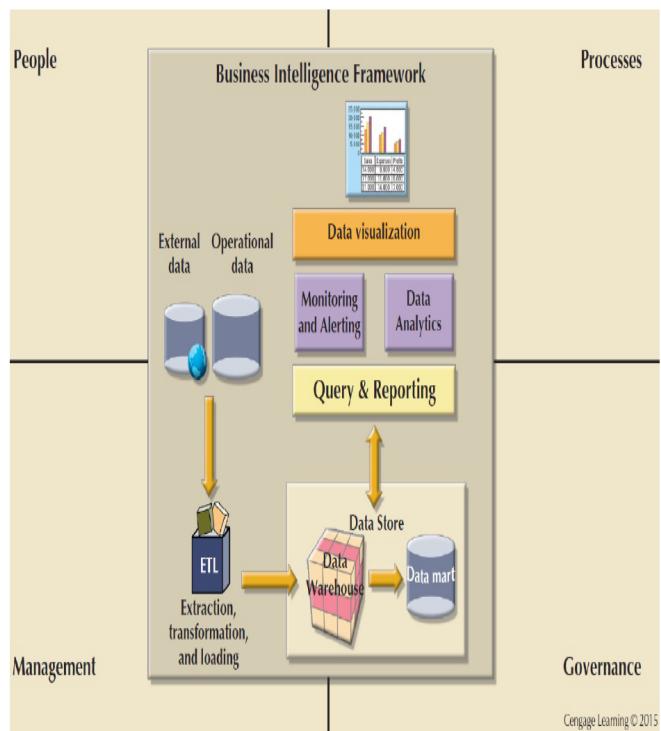


Business Intelligence ("BI")

Business Intelligence (BI)

- Comprehensive, cohesive, integrated set of tools and processes
 - Captures, collects, integrates, stores, and analyzes data
- Purpose - Generate and present information to support business decision making
- Allows a business to transform:
 - Data into information
 - Information into knowledge
 - Knowledge into wisdom

Figure 13.1 - Business Intelligence Framework



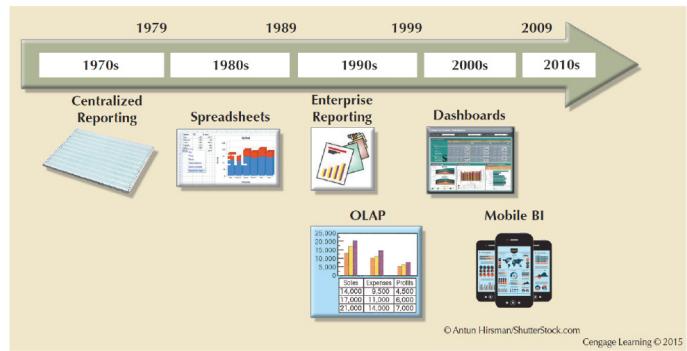
©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

Business Intelligence Benefits

- Improved decision making
- Integrating architecture
- Common user interface for data reporting and analysis
- Common data repository fosters single version of company data
- Improved organizational performance

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

Figure 13.3 - Evolution of BI Information Dissemination Formats



©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated; or posted to a publicly accessible website, in whole or in part.

14

| Decision Support Data Operational Data | Decision Support Data |
|---|---|
| <ul style="list-style-type: none">▪ Effectiveness of BI depends on quality of data gathered at operational level▪ Operational data▪ Seldom well-suited for decision support tasks▪ Stored in relational database with highly normalized structures▪ Optimized to support transactions representing daily operations | <ul style="list-style-type: none">▪ Differ from operational data in:<ul style="list-style-type: none">▪ Time span▪ Granularity<ul style="list-style-type: none">▪ Drill down: Decomposing a data to a lower level▪ Roll up: Aggregating a data into a higher level▪ Dimensionality |

Transactional, operational data: individual units.

'Analytical', decision-support data: aggregated.

Table 13.5 - Contrasting Operational and Decision Support Data Characteristics

| CHARACTERISTIC | OPERATIONAL DATA | DECISION SUPPORT DATA |
|---------------------|--|--|
| Data currency | Current operations Real-time data | Historic data Snapshot of company data Time component (week/month/year) |
| Granularity | Atomic-detailed data | Summarized data |
| Summarization level | Low; some aggregate yields | High; many aggregation levels |
| Data model | Highly normalized Mostly relational DBMSs | Non-normalized Complex structures Some relational, but mostly multidimensional DBMSs |
| Transaction type | Mostly updates | Mostly query |
| Transaction volumes | High-update volumes | Periodic loads and summary calculations |
| Transaction speed | Updates are critical | Retrievals are critical |
| Query activity | Low to medium | High |
| Query scope | Narrow range | Broad range |
| Query complexity | Simple to medium | Very complex |
| Data volumes | Hundreds of gigabytes | Terabytes to petabytes |

Cengage Learning © 2015

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

| Decision Support Database Requirements | Decision Support Database Requirements |
|--|---|
| <ul style="list-style-type: none">▪ Database schema<ul style="list-style-type: none">▪ Must support complex, non-normalized data representations▪ Data must be aggregated and summarized▪ Queries must be able to extract multidimensional time slices | <ul style="list-style-type: none">▪ Data extraction and loading<ul style="list-style-type: none">▪ Allow batch and scheduled data extraction▪ Support different data sources and check for inconsistent data or data validation rules▪ Support advanced integration, aggregation, and classification▪ Database size should support:<ul style="list-style-type: none">▪ Very large databases (VLDBs)▪ Advanced storage technologies▪ Multiple-processor technologies |

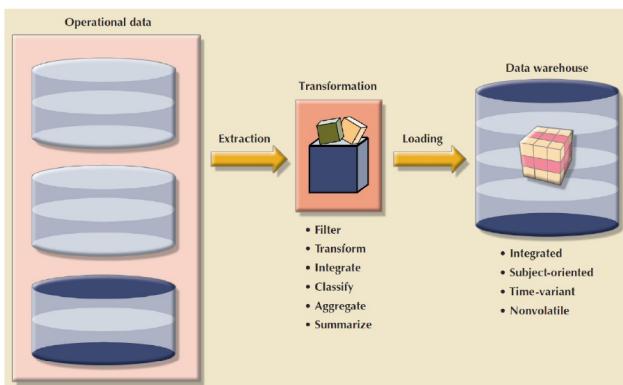
Table 13.8 - Characteristics of Data Warehouse Data and Operational Database Data

| CHARACTERISTIC | OPERATIONAL DATABASE DATA | DATA WAREHOUSE DATA |
|------------------|--|--|
| Integrated | Similar data can have different representations or meanings. For example, Social Security numbers may be stored as ####-##-#### or as #####-##, and a given condition may be labeled as T/F or 0/1 or Y/N. A sales value may be shown in thousands or in millions. | Provide a unified view of all data elements with a common definition and representation for all business units. |
| Subject-oriented | Data are stored with a functional, or process, orientation. For example, data may be stored for invoices, payments, and credit amounts. | Data are stored with a subject orientation that facilitates multiple views of the data and decision making. For example, sales may be recorded by product, division, manager, or region. |
| Time-variant | Data are recorded as current transactions. For example, the sales data may be the sale of a product on a given date, such as \$342.78 on 12-MAY-2014. | Data are recorded with a historical perspective in mind. Therefore, a time dimension is added to facilitate data analysis and various time comparisons. |
| Nonvolatile | Data updates are frequent and common. For example, an inventory amount changes with each sale. Therefore, the data environment is fluid. | Data cannot be changed. Data are added only periodically from historical systems. Once the data are properly stored, no changes are allowed. Therefore, the data environment is relatively static. |

Cengage Learning © 2015

21

Figure 13.5 - The ETL Process



©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

22

An ETL-related job ad:

DATA WAREHOUSE AND ANALYTICS DEVELOPERS (ETL/INFORMATICA)

Ascension Health-IS, Inc. is seeking two Data Warehouse and Analytics Developers (ETL/Informatica) in St. Louis, Missouri to code design and development on the data warehouse/analytics Extract Transform Load (ETL) toolset, Informatica PowerCenter; support Informatica toolset; integrate and develop other technologies. Research solutions and technology; participate in testing (e.g. user acceptance testing, unit, system, regression, integration testing); develop test plans and documentation; debug code. Contact Jenna Mihm, Vice President Legal Services & Associate General Counsel, Ascension Health, 4600 Edmundson Road, St. Louis, MO 63134, 314-733-8692, Jenna.Mihm@ascensionhealth.org To apply for this position, please reference Job Number 03.

Data Marts

- Small, single-subject data warehouse subset
- Provide decision support to a small group of people
- Benefits over data warehouses
 - Lower cost and shorter implementation time
 - Technologically advanced
 - Inevitable people issues

Table 13.9 - Twelve Rules for a Data Warehouse

| RULE NO. | DESCRIPTION |
|----------|--|
| 1 | The data warehouse and operational environments are separated. |
| 2 | The data warehouse data are integrated. |
| 3 | The data warehouse contains historical data over a long time. |
| 4 | The data warehouse data are snapshot data captured at a given point in time. |
| 5 | The data warehouse data are subject oriented. |
| 6 | The data warehouse data are mainly read-only with periodic batch updates from operational data. No online updates are allowed. |

Cengage Learning © 2015

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

24

Table 13.9 - Twelve Rules for a Data Warehouse

| RULE NO. | DESCRIPTION |
|----------|--|
| 7 | The data warehouse development life cycle differs from classical systems development. Data warehouse development is data-driven; the classical approach is process-driven. |
| 8 | The data warehouse contains data with several levels of detail: current detail data, old detail data, lightly summarized data, and highly summarized data. |
| 9 | The data warehouse environment is characterized by read-only transactions to very large data sets. The operational environment is characterized by numerous update transactions to a few data entities at a time. |
| 10 | The data warehouse environment has a system that traces data sources, transformations, and storage. |
| 11 | The data warehouse's metadata are a critical component of this environment. The metadata identify and define all data elements. The metadata provide the source, transformation, integration, storage, usage, relationships, and history of each data element. |
| 12 | The data warehouse contains a chargeback mechanism for resource usage that enforces optimal use of the data by end users. |

Cengage Learning © 2015

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

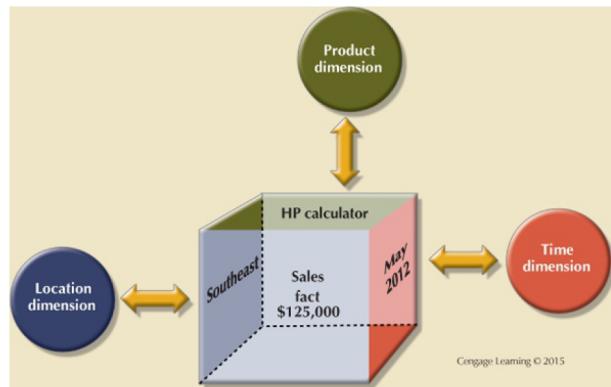
25

Star Schema

- Data-modeling technique
- Maps multidimensional decision support data into a relational database
- Creates the near equivalent of multidimensional database schema from existing relational database
- Yields an easily implemented model for multidimensional data analysis

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

26



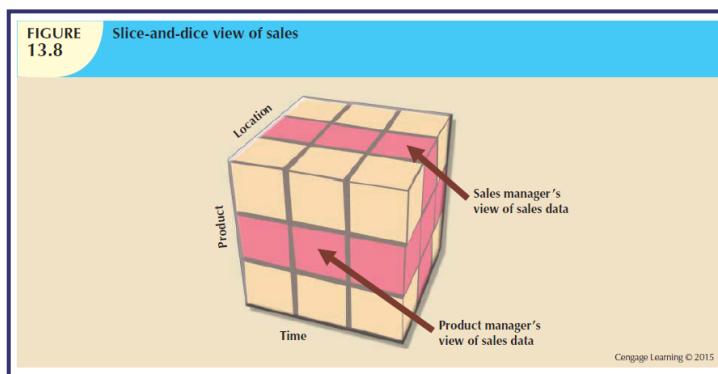
©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

Components of Star Schemas

- Facts**
 - Numeric values that represent a specific business aspect
- Dimensions**
 - Qualifying characteristics that provide additional perspectives to a given fact
- Attributes**
 - Used to search, filter, and classify facts
 - **Slice and dice:** Ability to focus on slices of the data cube for more detailed analysis
- Attribute hierarchy**
 - Provides a top-down data organization

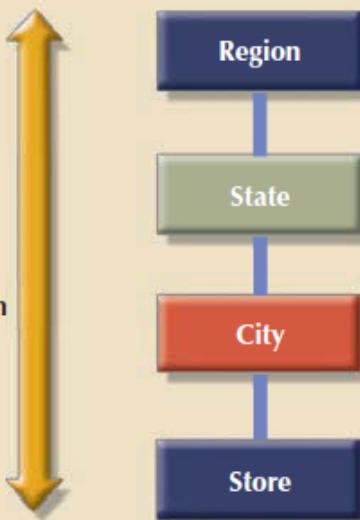
©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated; or posted to a publicly accessible website, in whole or in part.

28

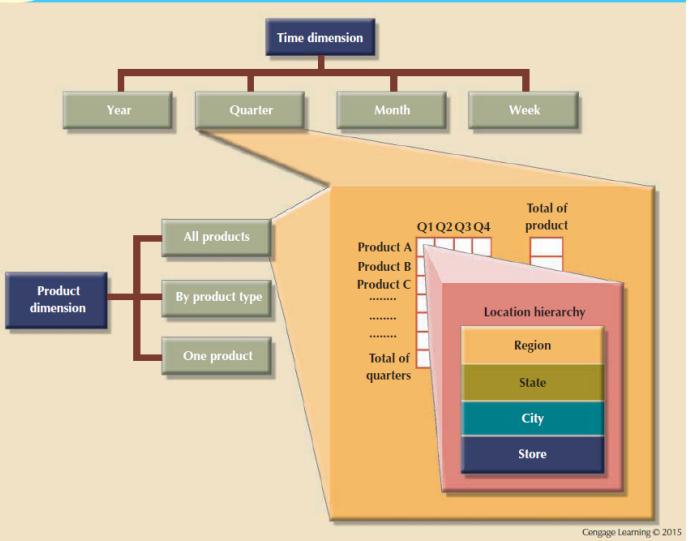


**FIGURE
13.9****Location attribute hierarchy**

The attribute hierarchy allows the end user to perform drill-down and roll-up searches.



Cengage Learning © 2015

**FIGURE
13.10****Attribute hierarchies in multidimensional analysis**

Cengage Learning © 2015

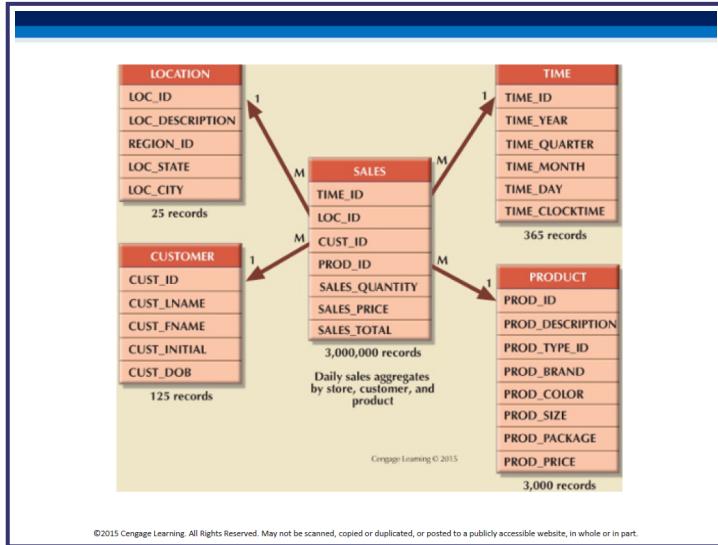
Star Schema Representation

- Facts and dimensions represented by physical tables in data warehouse database
- Many-to-one (M:1) relationship between fact table and each dimension table
- Fact and dimension tables
 - Related by foreign keys
 - Subject to primary and foreign key constraints

Star Schema Representation

- Primary key of a fact table
 - Is a composite primary key because the fact table is related to many dimension tables
 - Always formed by combining the foreign keys pointing to the related dimension tables

A sample star schema

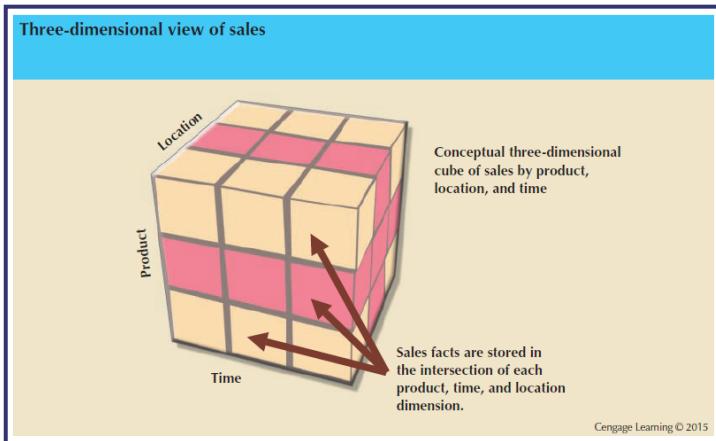


We have the fact table (of transactions) at the center, and denormalized (all-in-one) dimension tables all around.

Here is another representation.

Note: "dimensions are qualifying characteristics that provide additional perspectives to a given fact; dimensions provide descriptive characteristics about the facts through their attributes."

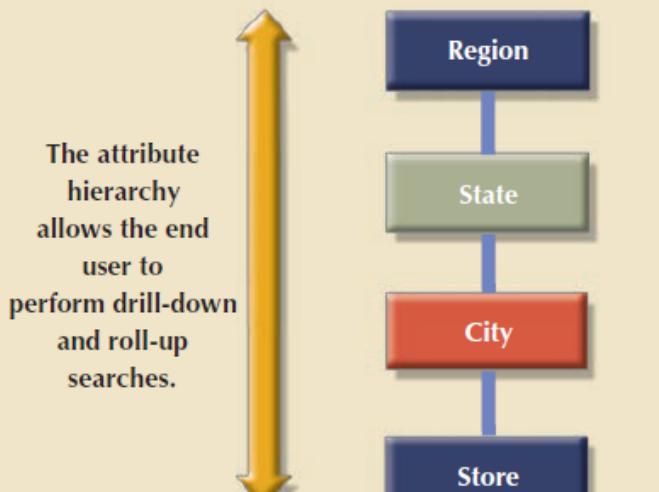
Each fact (transaction) can now pictured to be located in a multi-dimensional cube where the axes are dimensions. Eg. a 3D representation of our data for the above schema would look like this:



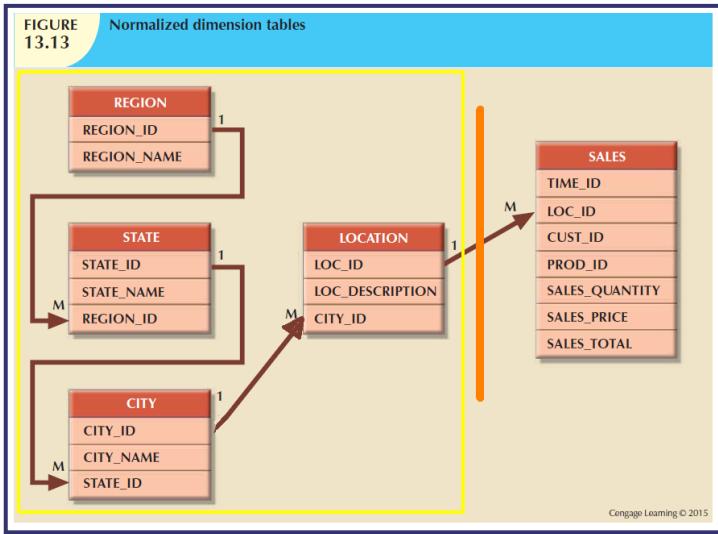
Slicing and dicing the cube provides specific insights..

Additionally, an attribute hierarchy would provide drill-down/roll-up capability as well, eg.

FIGURE 13.9 Location attribute hierarchy



Snowflake schema



Dimensional tables can be normalized so that they have their own dimensional tables - this is done to simplify the design, but requiring navigation across the normalized chains.

Here is another representation.

Techniques Used to Optimize Data Warehouse Design

- Normalizing dimensional tables
 - **Snowflake schema:** Dimension tables can have their own dimension tables
- Maintaining multiple fact tables to represent different aggregation levels
- Denormalizing fact tables

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

32

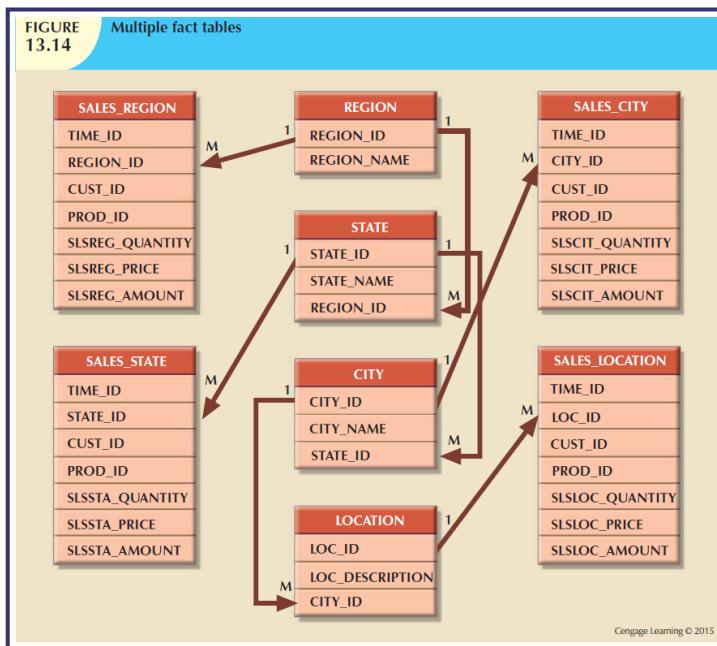
There are four (five!) different ways in which we can organize (structure) a data warehouse:

- 1. star schema: fact table FKs point to a single level of dimension tables (points of a star)
- 2. snowflake schema: each dimension table can be normalized to create a 1:M chain
- 3. the fact table can be supplanted with the columns in the star or snowflake dimensions
- 4. a separate fact table can be created for each attribute in a dimension hierarchy
- 5. ?

To denormalize a fact table (#3 above), we simply add extra 'dimension' columns to it, and fill them with redundant data - this permits fast queries (no joins needed) at the expense of disk space (and cleanliness of design).

Redundant fact tables!

Instead of a denormalized fact table (#3), or a fact table pointing to denormalized star dimensions (#1), or a fact table with lowest attrs pointing to a chain of rolled-up attrs, ie. snowflake schema (#2), we can create multiple fact tables, one for each level in an attr hierarchy (#4) - it is a different form of denormalization, where the redundant data is stored in physically separate tables.



There is also a fifth design option - horizontal fragmentation of tables.

Data Analytics

- Encompasses a wide range of mathematical, statistical, and modeling techniques to extract knowledge from data
 - Subset of BI functionality
- Classification of tools
 - **Explanatory analytics:** Focuses on discovering and explaining data characteristics and relationships based on existing data
 - **Predictive analytics:** Focuses on predicting future outcomes with a high degree of accuracy

Online Analytical Processing

- Advanced data analysis environment that supports decision making, business modeling, and operations research
- Characteristics
 - Multidimensional data analysis techniques
 - Advanced database support
 - Easy-to-use end-user interfaces

Figure 13.19 - OLAP Architecture

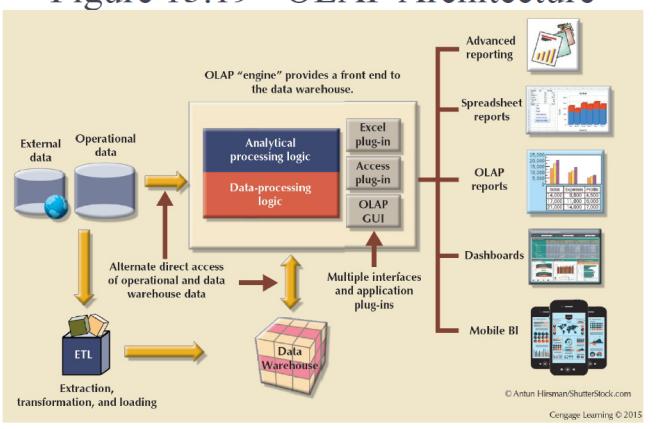
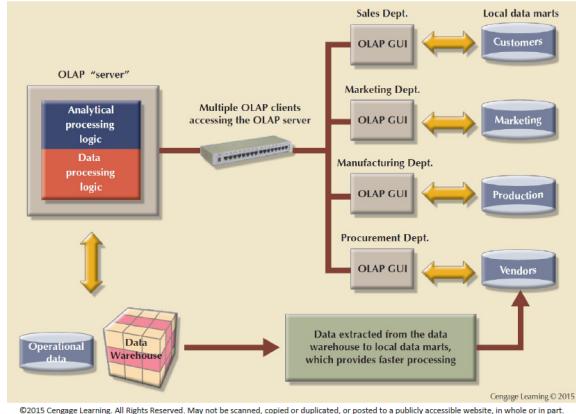


Figure 13.20 - OLAP Server with Local Miniature Data Marts



40

Relational Online Analytical Processing (ROLAP)

- Provides OLAP functionality using relational databases and familiar relational tools to store and analyze multidimensional data
- Extensions added to traditional RDBMS technology
 - Multidimensional data schema support within the RDBMS
 - Data access language and query performance optimized for multidimensional data
 - Support for very large databases (VLDBs)

Multidimensional Online Analytical Processing (MOLAP)

- Extends OLAP functionality to multidimensional database management systems (MDBMSs)
 - **MDBMS:** Uses proprietary techniques store data in matrix-like n-dimensional arrays
 - End users visualize stored data as a 3D **data cube**
 - Grow to n dimensions, becoming hypercubes
 - Held in memory in a **cube cache** to speed access
- **Sparsity:** Measures the density of the data held in the data cube

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

42

Table 13.12 - Relational vs.
Multidimensional OLAP

| CHARACTERISTIC | ROLAP | MOLAP |
|----------------|--|---|
| Schema | Uses star schema Additional dimensions can be added dynamically | Uses data cubes Multidimensional arrays, row stores, column stores Additional dimensions require re-creation of the data cube |
| Database size | Medium to large | Large |
| Architecture | Client/server Standards-based | Client/server Open or proprietary, depending on vendor |
| Access | Supports ad hoc requests Unlimited dimensions | Limited to predefined dimensions Proprietary access languages |
| Speed | Good with small data sets; average for medium-sized to large data sets | Faster for large data sets with predefined dimensions |

Cengage Learning © 2015

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

43

BI-oriented SQL extensions

ROLLUP and CUBE are GROUP BY modifiers - they help generate subtotals for a list of specified columns (see examples that follow). Depending on granularity of the columns (eg. US_REGION vs STORE_NUMBER), these subtotals help provide a rolled-up (aggregated) or drilled-down (detailed) analysis of data.

SQL Extensions for OLAP

The ROLLUP extension

- Used with GROUP BY clause to generate aggregates by different dimensions
- Enables subtotal for each column listed except for the last one, which gets a grand total
- Order of column list important

The CUBE extension

- Used with GROUP BY clause to generate aggregates by the listed columns
- Includes the last column

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

ROLLUP, CUBE: usage examples

ROLLUP extension

```
SELECT column1 [, column2, ...],  
       aggregate_function(expression)  
FROM table1 [, table2, ...]  
[WHERE condition]  
GROUP BY ROLLUP (column1 [, column2, ...])  
[HAVING condition]  
[ORDER BY column1 [, column2, ...]]
```

| TOTALES | |
|-----------------|----------|
| 21225 | 23189.90 |
| 21225 | 104.90 |
| 21225 | 41.90 |
| 21225 | 341.90 |
| 21225 13-02/2 | 231.90 |
| 21225 54778-21 | 52.90 |
| 21225 15M-B02 | 70.50 |
| 21225 15M-B02 | 70.50 |
| 21225 22826/QIV | 31.80 |
| 21225 89-ME1-Q | 512.90 |
| 21225 22827/PB | 72.50 |
| 21225 MB3/1TB | 219.50 |
| 21225 | 219.50 |
| | 21225.90 |

45

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

"Subtotal for each vendor - all products; sum of (the only set of) subtotals".

CUBE extension

```
SELECT column1 [, column2, ...],  
       aggregate_function(expression)  
FROM table1 [, table2, ...]  
[WHERE condition]  
GROUP BY CUBE (column1 [, column2, ...])  
[HAVING condition]  
[ORDER BY column1 [, column2, ...]]
```

| TOTALES | |
|-------------|----------|
| 9 13-02/2 | 156.91 |
| 9 2232/QIV | 189.92 |
| 9 23189.90 | 17.90 |
| 9 23189-HB | 59.7 |
| 9 89-ME1-Q | 39.8 |
| 9 89-ME1-Q | 39.8 |
| 9 89-ME1-Q | 256.99 |
| 9 89-ME1-Q | 512.90 |
| 9 SM-16277 | 28.97 |
| 9 SM-16277 | 259.85 |
| 9 SM-16277 | 129.90 |
| 18 13-02/2 | 186.92 |
| 18 2232/QIV | 189.92 |
| 18 23189-HB | 39.8 |
| 18 23189-HB | 99.5 |
| 18 89-ME1-Q | 39.8 |
| 18 89-ME1-Q | 256.99 |
| 18 89-ME1-Q | 512.90 |
| 18 SM-16277 | 28.97 |
| 18 SM-16277 | 259.85 |
| 18 SM-16277 | 129.90 |
| 13-02/2 | 229.85 |
| 2232/QIV | 219.80 |
| 23189-HB | 99.5 |
| 89-ME1-Q | 512.90 |
| SM-16277 | 41.90 |
| SM-16277 | 219.70 |
| | 21225.90 |

46

©2015 Cengage Learning. All Rights Reserved. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part.

"Subtotal for each month - all products; subtotal for each product - all months; sum of (either set of) subtotals".

Data Lakes

A 'traditional' data warehouse is an ETL-based, historical record of transactions - very RDB-like.

A 'modern' alternative is a 'data lake', which offers a more continuous form of analytics, driven by the rise of unstructured data, streaming, cloud storage, etc. In a data lake, data is NOT ETLD, rather, it is stored in its 'raw' ("natural") form [even incomplete, untransformed...].