## Data Warehousing - Server Issues

Chaudhuri and Dayal paper

CISC432/832                    DW Servers                    1

---

## Outline

- Performance and functionality requirements
- Server architectures
- Index structures
- Query processing
- Materialized views

CISC432/832                    DW Servers                    2

---

## Server Requirements

- Key issue is dealing with **large** volumes of data
- Update processing
  - updates typically batched and performed within a refresh window
  - load complex since it takes data from raw external sources

CISC432/832                    DW Servers                    3

## Server Requirements (cont.)

- Data quality management
  - data must be cleaned of error
  - data must be checked for local consistency, global consistency and referential integrity
- Query performance
  - queries complex and access large volumes of data
- Scalability
  - users, data sources and amount of data

CISC432/832          DW Servers          4

## DW Server Architectures

- Specialized SQL servers
  - provide advanced query language and query processing support
  - targeted to DW schemas and processing
- ROLAP servers
  - intermediate server on top of RDBMS back-end
  - add support for multidimensional OLAP queries

CISC432/832          DW Servers          5

## DW Server Architectures (cont.)

- MOLAP servers
  - provide multidimensional storage engine
  - direct mapping of OLAP queries to storage layer
- note: exploiting parallelism is a key feature of all architectures

CISC432/832          DW Servers          6
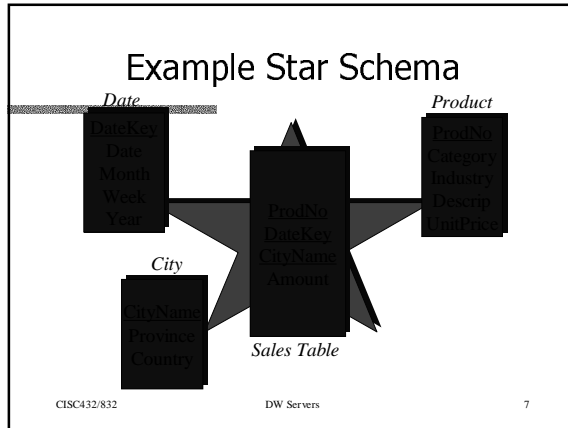
## Example Star Schema

*Date*

DateKey
Date
Month
Week
Year

*Product*

ProdNo
Category
Industry
Descrip
UnitPrice

ProdNo
DateKey
CityName
Amount

*City*

CityName
Province
Country

*Sales Table*

CISC432/832　　　　　　　DW Servers　　　　　　　7

## Typical OLAP Queries

- Give total sales for each product in each quarter of 1995.
- In 1996, for each city give the products with the top 5 sales.
- Give the average monthly sales of computer instruction books in each of the last 3 years for each city in Canada.

CISC432/832　　　　　　　DW Servers　　　　　　　8

## Index Structures (cont.)

- Indexes may be very large
  - need both space and time efficient organizations
- Queries may often involve multiple indexed attributes
  - need organizations that support efficient processing

CISC432/832　　　　　　　DW Servers　　　　　　　9

## Index Structures (cont.)

- **Bit-mapped index**
  - good when the multidimensional cube is sparse
  - each dimension has a bit-mapped index

| R(A) | | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 1 | 3 | 0 | 0 | 1 | 0 | 0 |
| 2 | 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 4 | 0 | 0 | 0 | 1 | 0 |
| 4 | 2 | 0 | 1 | 0 | 0 | 0 |
| 5 | 5 | 0 | 0 | 0 | 0 | 1 |
| 6 | 1 | 1 | 0 | 0 | 0 | 0 |

CISC432/832     DW Servers     10

## Index Structures (cont.)

SELECT * FROM Sales
WHERE prodNo = 3 AND City = "Toronto"

**Sales - ProdNo**

| 1 | 0 | 1 | 0 | 0 |
|---|---|---|---|---|
| 3 | 1 | 0 | 1 | 1 |
| 5 | 0 | 0 | 0 | 0 |

**Sales - City**

| Toronto | 0 | 1 | 1 | 0 |
|---|---|---|---|---|
| Beijing | 1 | 0 | 0 | 0 |
| Chicago | 0 | 0 | 0 | 1 |

**AND** → 0 0 1 0

CISC432/832     DW Servers     11

## Index Structures (cont.)

- **Join Index**
  - speeds up join by indicating which rows in dimension table a row in fact table joins with

Product.ProdNo    Product ⋈ Sales    Sales.ProdNo

| 1 | | | 1 | 1 | 0 | 1 | 0 | 0 | | 1 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | | | 2 | 0 | 1 | 0 | 0 | 1 | | 2 | 2 |
| 3 | | | 3 | 0 | 0 | 0 | 1 | 0 | | 3 | 1 |
| 4 | | | 4 | 0 | 0 | 0 | 0 | 0 | | 4 | 3 |
| 5 | | | 5 | 0 | 0 | 0 | 0 | 0 | | 5 | 2 |
| 6 | | | 6 | 0 | 0 | 0 | 0 | 0 | | | |

CISC432/832     DW Servers     12

## Query Processing

- Optimizing complex queries
  - nested queries
  - multiple joins
  - group-by and aggregation
  - sorts
- Exploiting materialized views

CISC432/832                    DW Servers                    13

## Example Query

Select ProdNo, CityName, SUM(Amount)
  from Sales, Date, City, Product
  where Year = 1998
  and Country = "Canada"
  and Category = "pencils"
  and Sales.ProdNo = Product.ProdNo
  and Sales.CityName = City.CityName
  and Sales.DateKey = Date.DateKey
  group-by ProdNo, CityName
  order by ProdNo

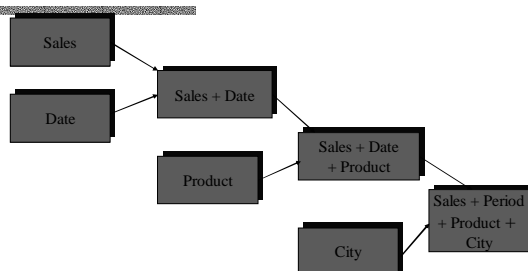CISC432/832                    DW Servers                    14

## Conventional STAR Join



CISC432/832                    DW Servers                    15

## Conventional STAR Join (cont.)

- Potential performance problems
  - limitations of pair-wise join
    - potentially large intermediate tables
  - selecting join order
    - number of possible orders is N!
    - only join "related" tables
  - generating reasonable cost estimates

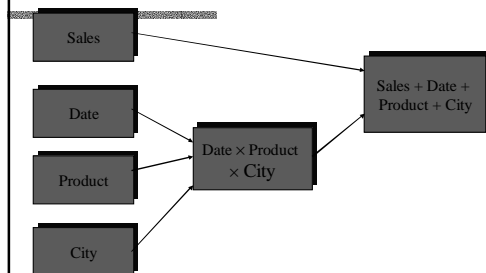CISC432/832                    DW Servers                    16

## Alternative 1: STAR Join with Cross-Product



CISC432/832                    DW Servers                    17

## Alternative 2: STAR Join with Join Indices



CISC432/832                    DW Servers                    18

## Materialized Views

- Precompute and store summary data
  - improves query performance but adds to storage and maintenance costs
- Issues
  - what views to materialize
  - how to update them
  - how to exploit them

CISC432/832          DW Servers          19

## Materialized Views(cont.)

MV = **select** Category, Year, SUM(Amount)
    **from** Sales, Product, Date
    **where** Sales.ProdNo = Product.ProdNo
    **and** Sales.DateKey = Date.DateKey
    **group by** Category, Year

**select** *
    **from** MV
    **where** Year = 1997

CISC432/832          DW Servers          20
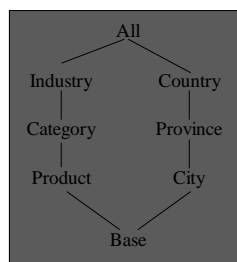
## Materialized Views(cont,)

- Choosing which view to materialize
  - materialize everything
  - precompute most frequently asked queries
  - methods to find optimal set of views

All
Industry          Country
Category          Province
Product          City
Base

CISC432/832          DW Servers          21

## Materialized Views(cont.)

- How to update them?
  - Periodically recompute (snapshot view)
  - incrementally update views
- How to exploit them?
  - Optimizer must recognize where a materialized view can be substituted in a query and then query rewritten

CISC432/832                DW Servers                22

## MV Selection Problem

Given a set of queries

$Q = \{Q_1, Q_2, ..., Q_n\}$

choose a set of MVs

$V = \{V_1, V_2, ..., V_k\}$

that minimize one or more of

storage costs

query processing costs

maintenance costs

CISC432/832                DW Servers                23

## An Example: TPC-D

partNo

Part

custKey
nation
region

Customer

custKey
partNo
ordKey
suppKey
quantity   LineItem

Select partNo, custKey,
SUM (quantity)
from LineItem
group by partNo, custKey

ordKey

Order

suppKey
nation
region

Supplier

CISC432/832                DW Servers                24

## An Example (cont.)

partNo

Part

custKey
partNo
ordKey
suppKey
quantity — LineItem

custKey
nation
region

Customer

ordKey

Order

suppKey
nation
region

Supplier

**Possible Views (group-bys)**
(part, supplier, customer) (6M)
(part, customer) (6M)
(part, supplier) (0.8M)
(supplier, customer) (6M)
(part) (0.2M)
(supplier) (0.01M)
(customer) (0.1M)
(none) (1)

CISC432/832                         DW Servers                                    25

## View Lattice Approach

- Dependence relation on views
  - $V_1 \leq V_2$ iff $V_1$ can be answered using only result of $V_2$
  - e.g. (part) $\leq$ (part, customer) but (part) $\not\leq$ (customer)

**psc** (6M)

**pc** (6M)    **ps** (0.8M)    **sc** (6M)

**p (**0.2M)    **s** (0.01M)    **c** (0.1M)

**none** (1)

CISC432/832                         DW Servers                                    26

## View Lattice (cont.)

- To answer a query $Q_i$
  - choose ancestor of $Q_i$ in lattice, $Q_A$ , that is materialized
  - cost of processing Q is number of rows in $Q_A$
- benefit of materializing view V is amount by which it improves costs of queries in Q

CISC432/832                         DW Servers                                    27

## View Lattice (cont.)

- Greedy MV Selection Algorithm

  S = {top view}

  for (i = 1 to k ) do

     select $V_i \notin$ S s.t. V has maximum benefit

     with        respect to S

     S = S $\cup$ $V_i$

  end

  return S

  CISC432/832                    DW Servers                    28