

# Week 13: Data Warehousing

## Warehousing

- Growing industry: \$8 billion in 1998
- Range from desktop to huge:
  - ◆ Walmart: 900-CPU, 2,700 disk, 23TB Teradata system
- Lots of buzzwords, hype
  - ◆ slice & dice, rollup, MOLAP, pivot, ...

# Outline

- What is a data warehouse?
- Why a warehouse?
- Models & operations
- Implementing a warehouse
- Future directions

3

## What is a Warehouse?

- Collection of diverse data
  - ◆ subject oriented
  - ◆ aimed at executive, decision maker
  - ◆ often a copy of operational data
  - ◆ with value-added data (e.g., summaries, history)
  - ◆ integrated
  - ◆ time-varying
  - ◆ non-volatile



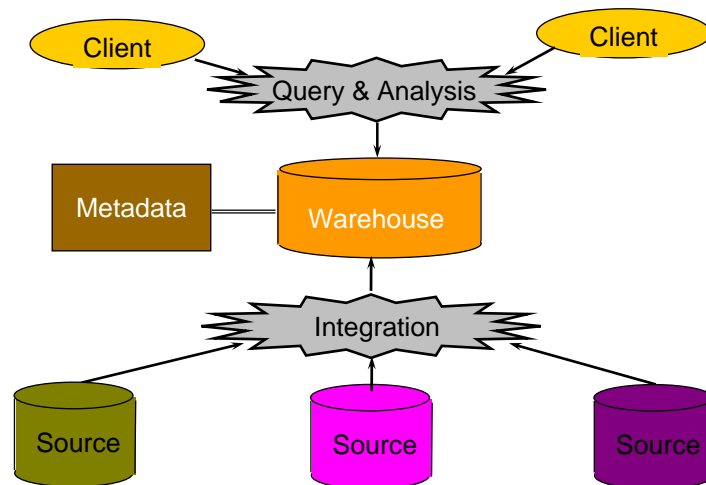
4

# What is a Warehouse?

- Collection of tools
  - ◆ gathering data
  - ◆ cleansing, integrating, ...
  - ◆ querying, reporting, analysis
  - ◆ data mining
  - ◆ monitoring, administering warehouse

5

## Warehouse Architecture



6

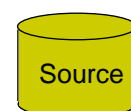
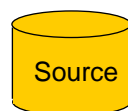
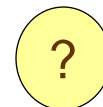
## Motivating Examples

- Forecasting
- Comparing performance of units
- Monitoring, detecting fraud
- Visualization

7

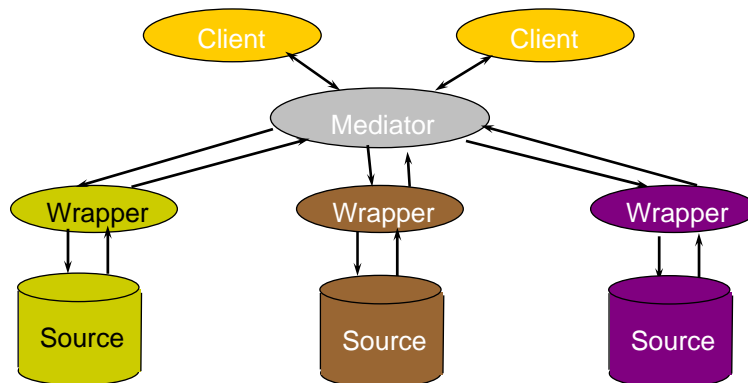
## Why a Warehouse?

- Two Approaches:
  - ◆ Query-Driven (Lazy)
  - ◆ Warehouse (Eager)



8

## Query-Driven Approach



9

## Advantages of Warehousing

- High query performance
- Queries not visible outside warehouse
- Local processing at sources unaffected
- Can operate when sources unavailable
- Can query data not stored in a DBMS
- Extra information at warehouse
  - ◆ Modify, summarize (store aggregates)
  - ◆ Add historical information

10

## Advantages of Query-Driven

- No need to copy data
  - ◆ less storage
  - ◆ no need to purchase data
- More up-to-date data
- Query needs can be unknown
- Only query interface needed at sources
- May be less draining on sources

11

## OLTP vs. OLAP

- OLTP: On Line Transaction Processing
  - ◆ Describes processing at operational sites
- OLAP: On Line Analytical Processing
  - ◆ Describes processing at warehouse

12

# OLTP vs. OLAP

## OLTP

- Mostly updates
- Many small transactions
- Mb-Tb of data
- Raw data
- Clerical users
- Up-to-date data
- Consistency, recoverability critical

## OLAP

- Mostly reads
- Queries long, complex
- Gb-Tb of data
- Summarized, consolidated data
- Decision-makers, analysts as users

13

# Data Marts

- Smaller warehouses
- Spans part of organization
  - ◆ e.g., marketing (customers, products, sales)
- Do not require enterprise-wide consensus
  - ◆ but long term integration problems?

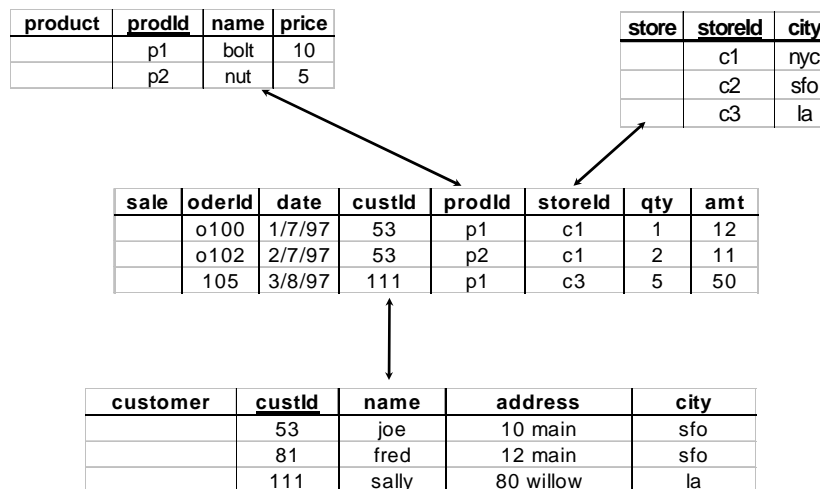
14

# Warehouse Models & Operators

- Data Models
  - ◆ relations
  - ◆ stars & snowflakes
  - ◆ cubes
- Operators
  - ◆ slice & dice
  - ◆ roll-up, drill down
  - ◆ pivoting
  - ◆ other

15

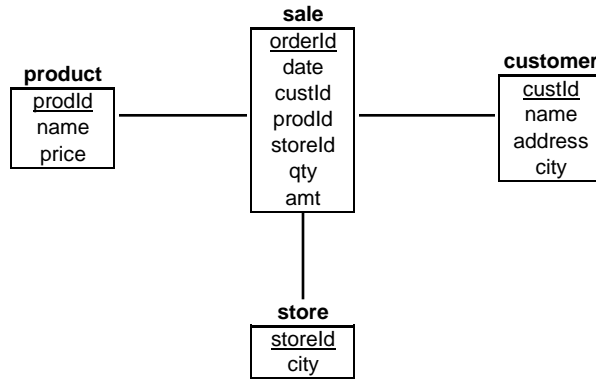
## Star



16



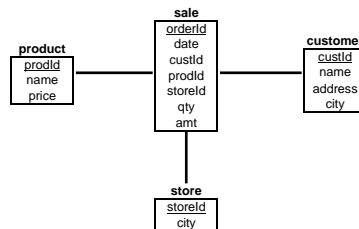
# Star Schema



17

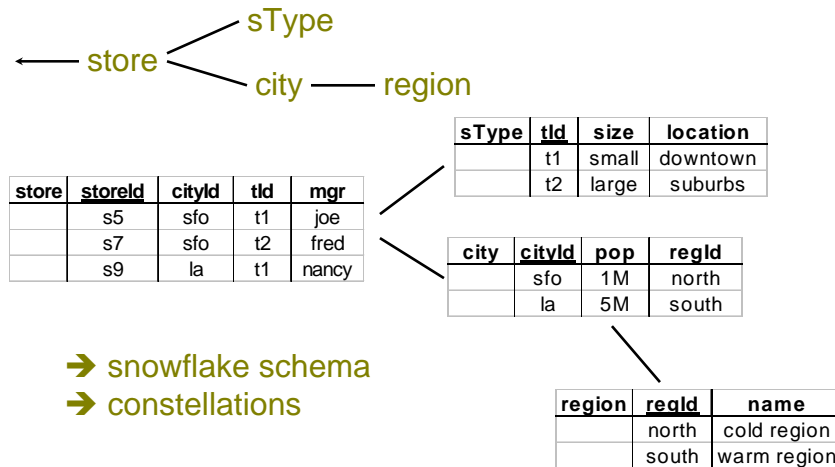
# Terms

- Fact table
- Dimension tables
- Measures



18

# Dimension Hierarchies



19

# Cube

Fact table view:

sale	prodId	storeId	amt
	p1	c1	12
	p2	c1	11
	p1	c3	50
	p2	c2	8

Multi-dimensional cube:

	c1	c2	c3
p1	12		50
p2	11	8	

dimensions = 2

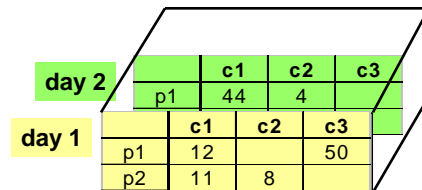
20

# 3-D Cube

Fact table view:

sale	prodId	storeId	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

Multi-dimensional cube:



dimensions = 3

21

## ROLAP vs. MOLAP

- ROLAP:  
Relational On-Line Analytical Processing
- MOLAP:  
Multi-Dimensional On-Line Analytical Processing

22

# Aggregates

- Add up amounts for day 1
- In SQL: `SELECT sum(amt) FROM SALE WHERE date = 1`

sale	prodlid	storeld	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

➔ 81

23

# Aggregates

- Add up amounts by day
- In SQL: `SELECT date, sum(amt) FROM SALE GROUP BY date`

sale	prodlid	storeld	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4



ans	date	sum
	1	81
	2	48

24

## Another Example

- Add up amounts by day, product
- In SQL: `SELECT date, sum(amt) FROM SALE GROUP BY date, prodl`

sale	prodl	storeld	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4



sale	prodl	date	amt
	p1	1	62
	p2	1	19
	p1	2	48

— rollup —→

← drill-down —

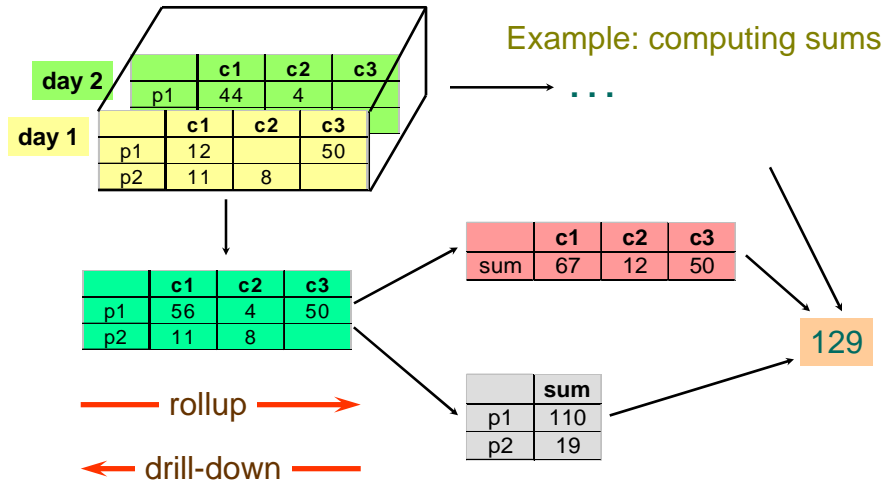
25

## Aggregates

- Operators: sum, count, max, min, median, ave
- “Having” clause
- Using dimension hierarchy
  - ◆ average by region (within store)
  - ◆ maximum by month (within date)

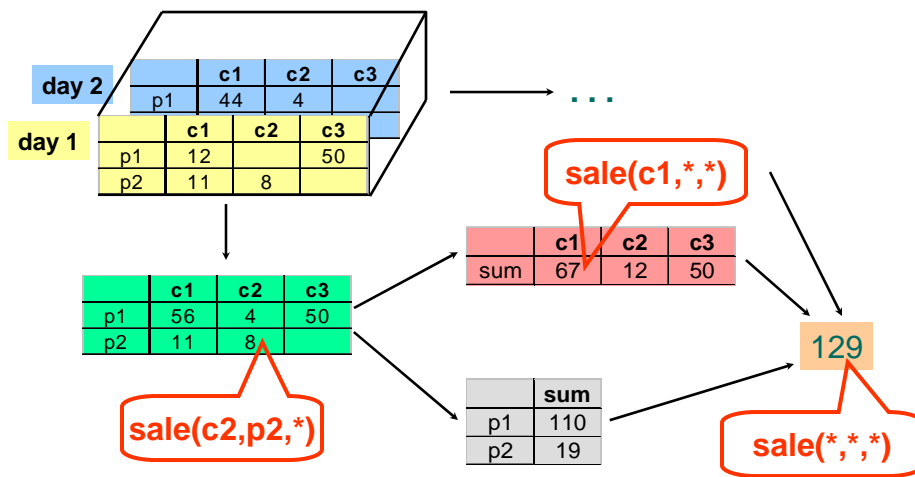
26

# Cube Aggregation



27

# Cube Operators



28

# Extended Cube

		day 2			
		c1	c2	c3	*
p1		56	4	50	110
p2		11	8		19
*					129

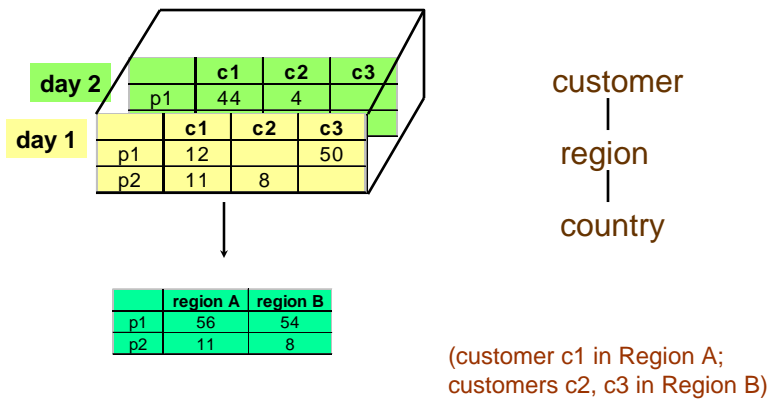
  

		day 1			
		c1	c2	c3	*
p1		12		50	62
p2		11	8		19
*		23	8	50	81

sale(\*,p2,\*)

29

# Aggregation Using Hierarchies



30

# Pivoting

Fact table view:

sale	prodlid	storeld	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

Multi-dimensional cube:

	c1	c2	c3	
day 2	p1	44	4	50
day 1	p1	12	8	50
day 1	p2	11	8	

	c1	c2	c3
p1	56	4	50
p2	11	8	

31

## Query & Analysis Tools

- Query Building
- Report Writers (comparisons, growth, graphs,...)
- Spreadsheet Systems
- Web Interfaces
- Data Mining

32



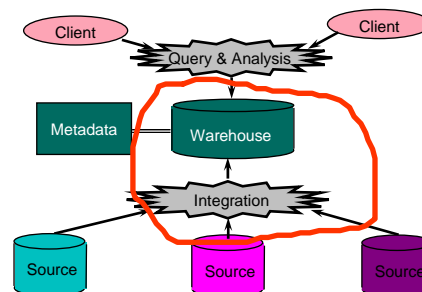
## Other Operations

- Time functions
  - ◆ e.g., time average
- Computed Attributes
  - ◆ e.g.,  $\text{commission} = \text{sales} * \text{rate}$
- Text Queries
  - ◆ e.g., find documents with words X AND B
  - ◆ e.g., rank documents by frequency of words X, Y, Z

33

## Integration

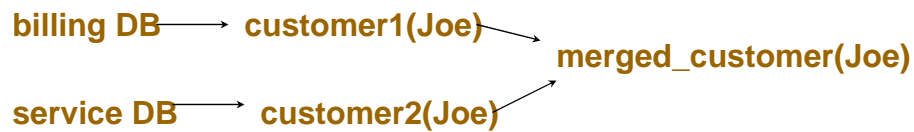
- Data Cleaning
- Data Loading
- Derived Data



34

## Data Cleaning

- Migration (e.g., yen  $\Rightarrow$  dollars)
- Scrubbing: use domain-specific knowledge (e.g., social security numbers)
- Fusion (e.g., mail list, customer merging)



- Auditing: discover rules & relationships (like data mining)

35

## Loading Data

- Incremental vs. refresh
- Off-line vs. on-line
- Frequency of loading
  - ◆ At night, 1x a week/month, continuously
- Parallel/Partitioned load

36

# Derived Data

- Derived Warehouse Data
  - ◆ indexes
  - ◆ aggregates
  - ◆ materialized views (next slide)
- When to update derived data?
- Incremental vs. refresh

37

# Materialized Views

- Define new warehouse relations using SQL expressions

sale	prodlid	storeld	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

product	id	name	price
	p1	bolt	10
	p2	nut	5

---

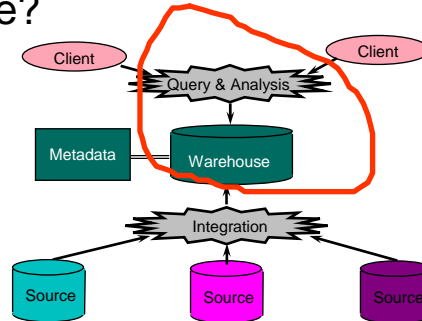
joinTb	prodlid	name	price	storeld	date	amt
	p1	bolt	10	c1	1	12
	p2	nut	5	c1	1	11
	p1	bolt	10	c3	1	50
	p2	nut	5	c2	1	8
	p1	bolt	10	c1	2	44
	p1	bolt	10	c2	2	4

does not exist  
at any source

38

# Processing

- ROLAP servers vs. MOLAP servers
- Index Structures
- What to Materialize?
- Algorithms

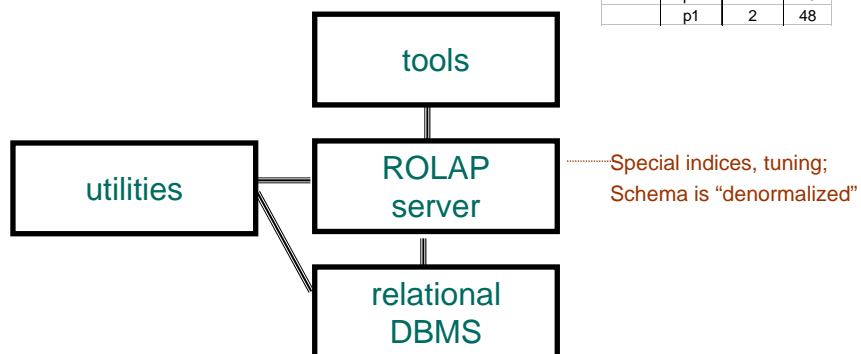


39

# ROLAP Server

- Relational OLAP Server

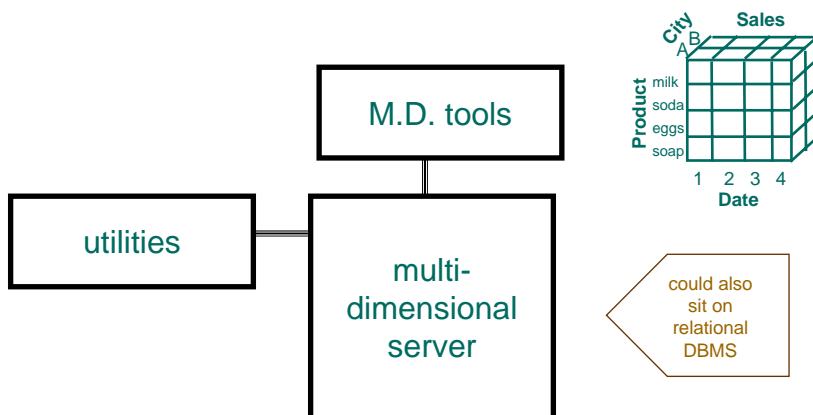
sale	prodId	date	sum
	p1	1	62
	p2	1	19
	p1	2	48



40

# MOLAP Server

- Multi-Dimensional OLAP Server



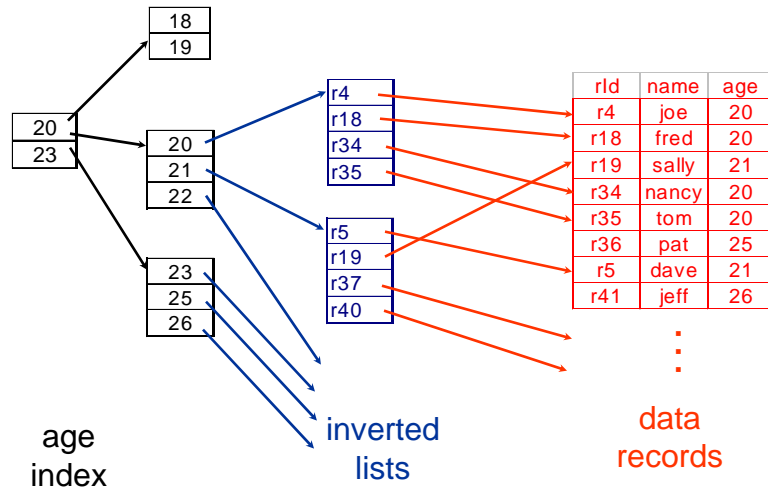
41

# Index Structures

- Traditional Access Methods
  - ◆ B-trees, hash tables, R-trees, grids, ...
- Popular in Warehouses
  - ◆ inverted lists
  - ◆ bit map indexes
  - ◆ join indexes
  - ◆ text indexes

42

## Inverted Lists



43

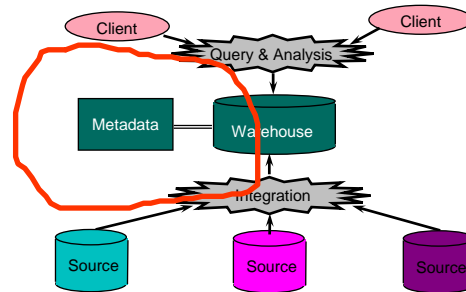
## Using Inverted Lists

- Query:
  - ◆ Get people with age = 20 and name = "fred"
- List for age = 20: r4, r18, r34, r35
- List for name = "fred": r18, r52
- Answer is intersection: r18

44

# Managing

- Metadata
- Warehouse Design
- Tools



45

# Metadata

- Administrative
  - ◆ definition of sources, tools, ...
  - ◆ schemas, dimension hierarchies, ...
  - ◆ rules for extraction, cleaning, ...
  - ◆ refresh, purging policies
  - ◆ user profiles, access control, ...

46

# Metadata

- Business
  - ◆ business terms & definition
  - ◆ data ownership, charging
- Operational
  - ◆ data lineage
  - ◆ data currency (e.g., active, archived, purged)
  - ◆ use stats, error reports, audit trails

47

# Design

- What data is needed?
- Where does it come from?
- How to clean data?
- How to represent in warehouse (schema)?
- What to summarize?
- What to materialize?
- What to index?

48



# Tools

- Development
  - ◆ design & edit: schemas, views, scripts, rules, queries, reports
- Planning & Analysis
  - ◆ what-if scenarios (schema changes, refresh rates), capacity planning
- Warehouse Management
  - ◆ performance monitoring, usage patterns, exception reporting
- System & Network Management
  - ◆ measure traffic (sources, warehouse, clients)
- Workflow Management
  - ◆ “reliable scripts” for cleaning & analyzing data

49

# Current State of Industry

- Extraction and integration done off-line
  - ◆ Usually in large, time-consuming, batches
- Everything copied at warehouse
  - ◆ Not selective about what is stored
  - ◆ Query benefit vs storage & update cost
- Query optimization aimed at OLTP
  - ◆ High throughput instead of fast response
  - ◆ Process whole query before displaying anything

50

## Future Directions

- Better performance
- Larger warehouses
- Easier to use
- What are companies & research labs working on?

51

## Research (1)

- Incremental Maintenance
- Data Consistency
- Data Expiration
- Recovery
- Data Quality
- Error Handling

52

## Research (2)

- Rapid Monitor Construction
- Temporal Warehouses
- Materialization & Index Selection
- Data Fusion
- Data Mining
- Integration of Text & Relational Data

53

## Conclusions

- Massive amounts of data and complexity of queries will push limits of current warehouses
- Need better systems:
  - ◆ easier to use
  - ◆ provide quality information

54