

# Data Warehousing and OLAP

*Hector Garcia-Molina  
Stanford University*

## 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

Hector Garcia Molina: Data Warehousing and OLAP

## 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

Hector Garcia Molina: Data Warehousing and OLAP

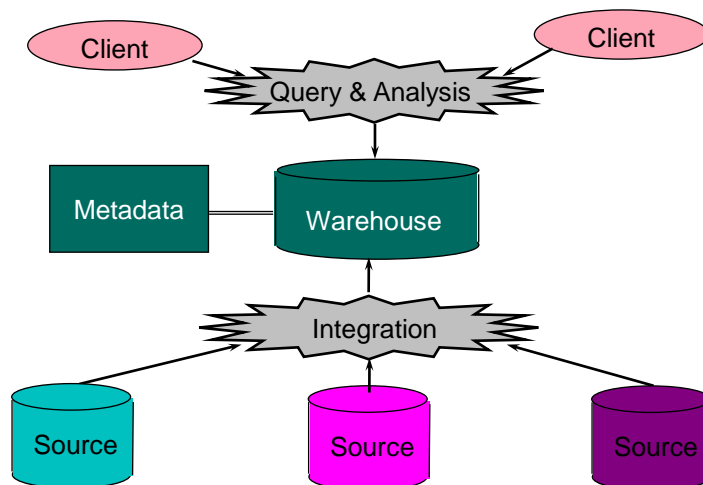
# What is a Warehouse?

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

5

Hector Garcia Molina: Data Warehousing and OLAP

# Warehouse Architecture

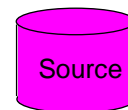
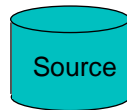


6

Hector Garcia Molina: Data Warehousing and OLAP

# Why a Warehouse?

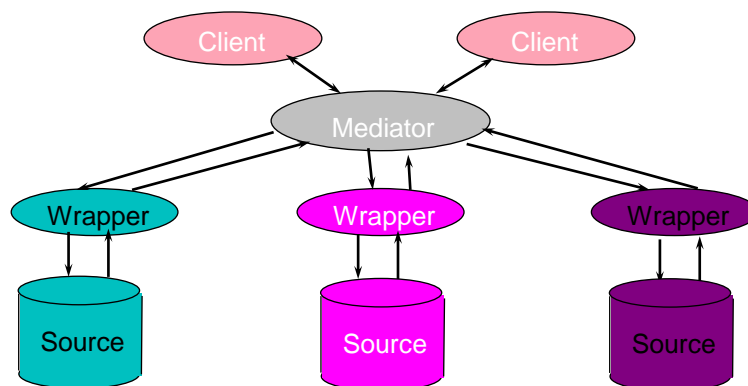
- Two Approaches:
  - Query-Driven (Lazy)
  - Warehouse (Eager)



7

Hector Garcia Molina: Data Warehousing and OLAP

# Query-Driven Approach



8

Hector Garcia Molina: Data Warehousing and OLAP

## 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

9

Hector Garcia Molina: Data Warehousing and OLAP

## 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

10

Hector Garcia Molina: Data Warehousing and OLAP

# OLTP vs. OLAP

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

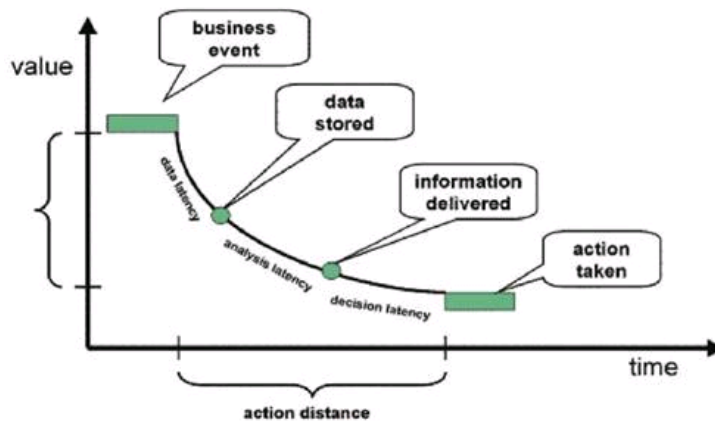
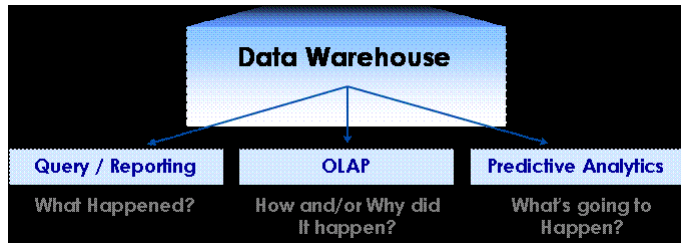
# OLTP vs. OLAP

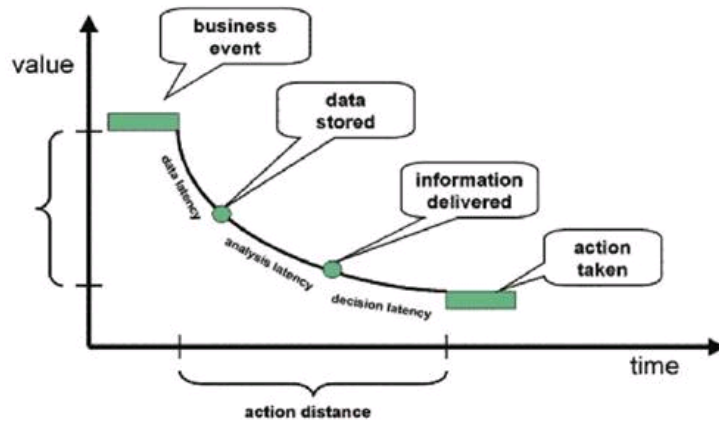
## OLTP

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

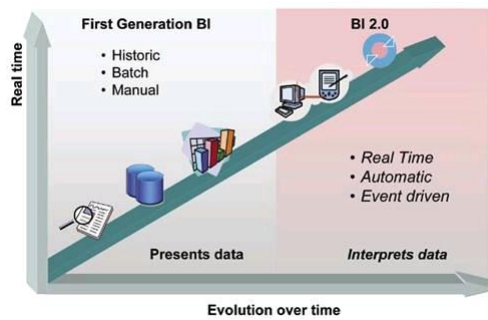
## OLAP

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





# Business Intelligence (BI)





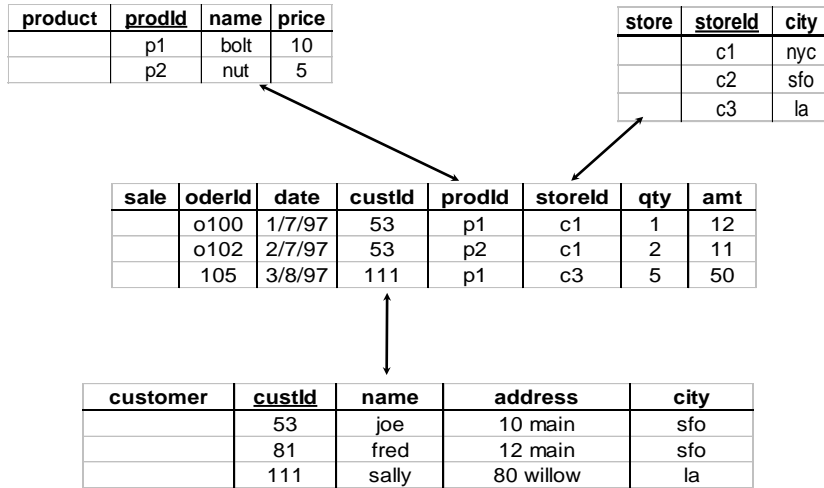
## Data Marts

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

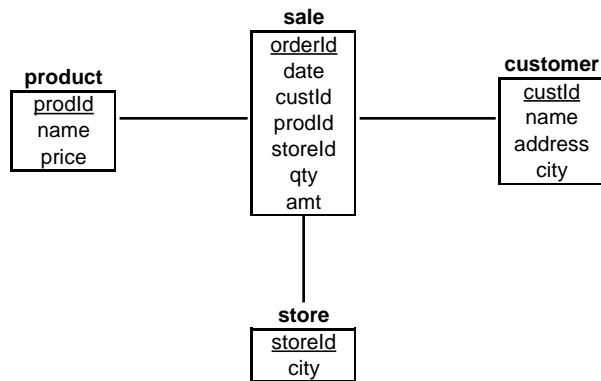
## Warehouse Models & Operators

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

# Star

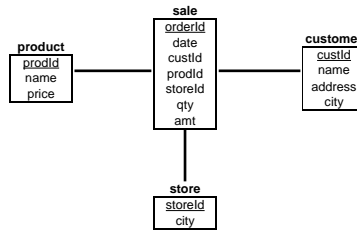


# Star Schema



# Terms

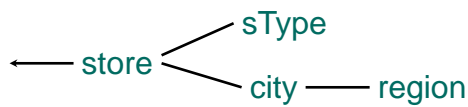
- Fact table
- Dimension tables
- Measures



21

Hector Garcia Molina: Data Warehousing and OLAP

# Dimension Hierarchies



store	storeid	cityid	tld	mgr
	s5	sfo	t1	joe
	s7	sfo	t2	fred
	s9	la	t1	nancy

sType	tld	size	location
	t1	small	downtown
	t2	large	suburbs

city	cityid	pop	regld
	sfo	1M	north
	la	5M	south

region	regld	name
	north	cold region
	south	warm region

- ➔ snowflake schema
- ➔ constellations

22

Hector Garcia Molina: Data Warehousing and OLAP

# Cube

Fact table view:

sale	prold	storeld	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

# 3-D Cube

Fact table view:

sale	prold	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 1	p1	12	50
day 1	p2	11	8
day 2	p1	44	4

dimensions = 3

## ROLAP vs. MOLAP

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

## Aggregates

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

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



81

# Aggregates

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

sale	prold	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

27

Hector Garcia Molina: Data Warehousing and OLAP

# Another Example

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

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

→ rollup →

← drill-down ←

28

Hector Garcia Molina: Data Warehousing and OLAP

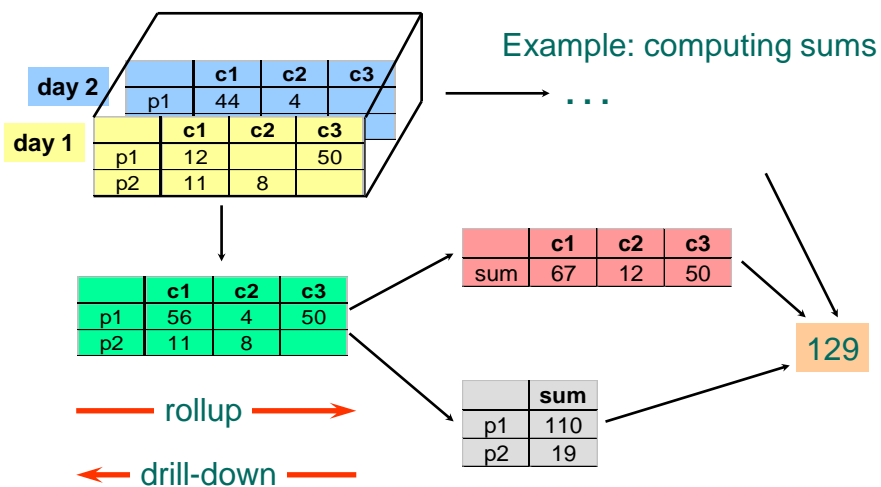
# Aggregates

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

29

Hector Garcia Molina: Data Warehousing and OLAP

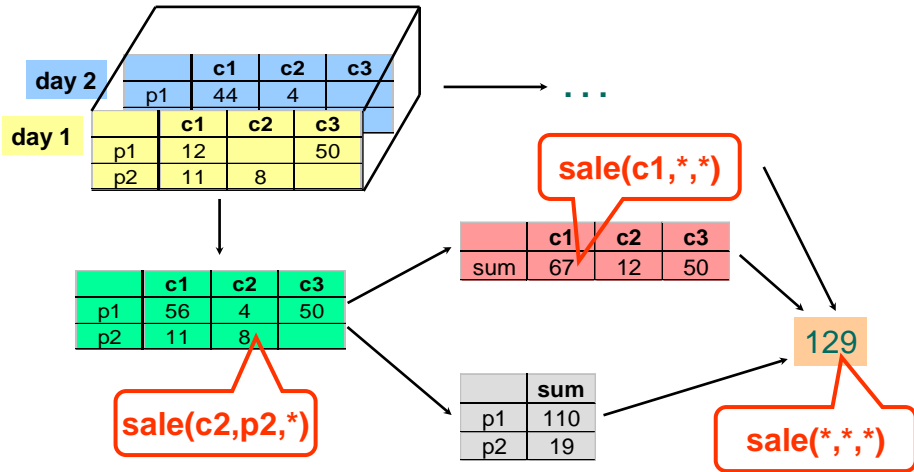
# Cube Aggregation



30

Hector Garcia Molina: Data Warehousing and OLAP

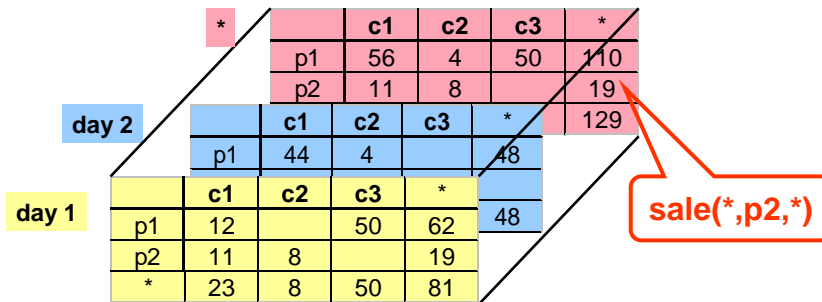
# Cube Operators



31

Hector Garcia Molina: Data Warehousing and OLAP

# Extended Cube

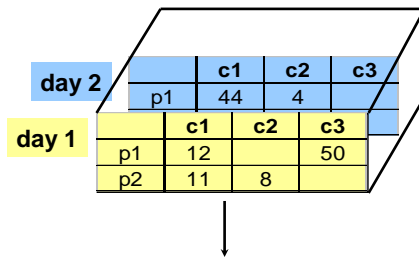


32

Hector Garcia Molina: Data Warehousing and OLAP



# Aggregation Using Hierarchies



	region A	region B
p1	56	54
p2	11	8

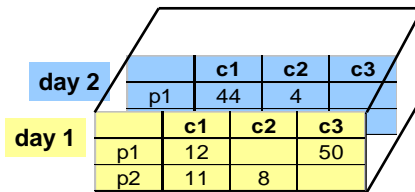
(customer c1 in Region A;  
customers c2, c3 in Region B)

# Pivoting

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:



	c1	c2	c3
p1	56	4	50
p2	11	8	

# Implementing a Warehouse

- *Monitoring*: Sending data from sources
- *Integrating*: Loading, cleansing,...
- *Processing*: Query processing, indexing, ...
- *Managing*: Metadata, Design, ...

35

Hector Garcia Molina: Data Warehousing and OLAP

# Monitoring

- Source Types: relational, flat file, IMS, VSAM, IDMS, WWW, news-wire, ...
- Incremental vs. Refresh

customer	id	name	address	city
	53	joe	10 main	sfo
	81	fred	12 main	sfo
	111	sally	80 willow	la



36

Hector Garcia Molina: Data Warehousing and OLAP

## Monitoring Techniques

- Periodic snapshots
- Database triggers
- Log shipping
- Data shipping (replication service)
- Transaction shipping
- Polling (queries to source)
- Screen scraping
- Application level monitoring

↑ Advantages & Disadvantages!!

37

Hector Garcia Molina: Data Warehousing and OLAP

## Monitoring Issues

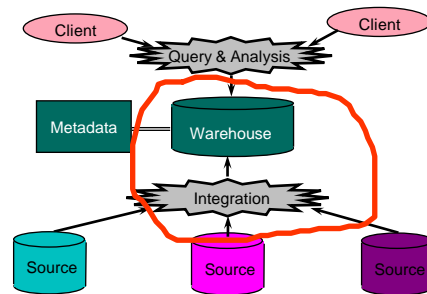
- Frequency
  - ◆ periodic: daily, weekly, ...
  - ◆ triggered: on “big” change, lots of changes, ...
- Data transformation
  - ◆ convert data to uniform format
  - ◆ remove & add fields (e.g., add date to get history)
- Standards (e.g., ODBC)
- Gateways

38

Hector Garcia Molina: Data Warehousing and OLAP

# Integration

- Data Cleaning
- Data Loading
- Derived Data

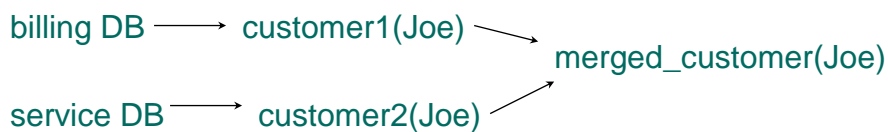


39

Hector Garcia Molina: Data Warehousing and OLAP

# 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)

40

Hector Garcia Molina: Data Warehousing and OLAP

## Loading Data

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

41

Hector Garcia Molina: Data Warehousing and OLAP

## Derived Data

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

42

Hector Garcia Molina: Data Warehousing and OLAP

# Materialized Views

- Define new warehouse relations using SQL expressions

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

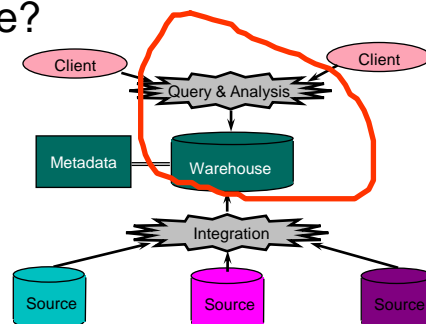
product	id	name	price
	p1	bolt	10
	p2	nut	5

joinTb	prodId	name	price	storeId	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

# Processing

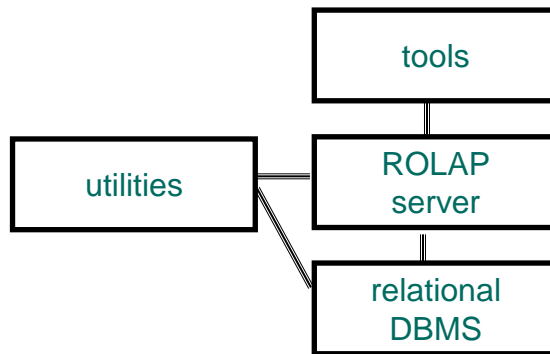
- ROLAP servers vs. MOLAP servers
- Index Structures
- What to Materialize?
- Algorithms



# ROLAP Server

- Relational OLAP Server

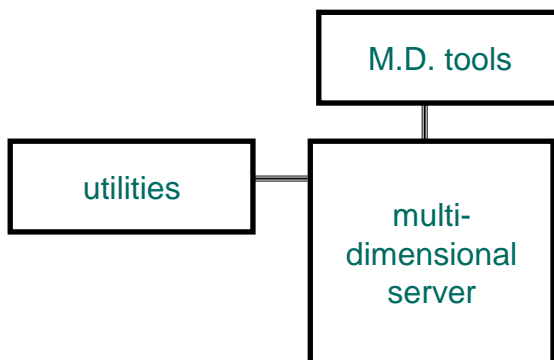
sale	prold	date	sum
	p1	1	62
	p2	1	19
	p1	2	48



Special indices, tuning;  
Schema is "denormalized"

# MOLAP Server

- Multi-Dimensional OLAP Server



could also  
sit on  
relational  
DBMS

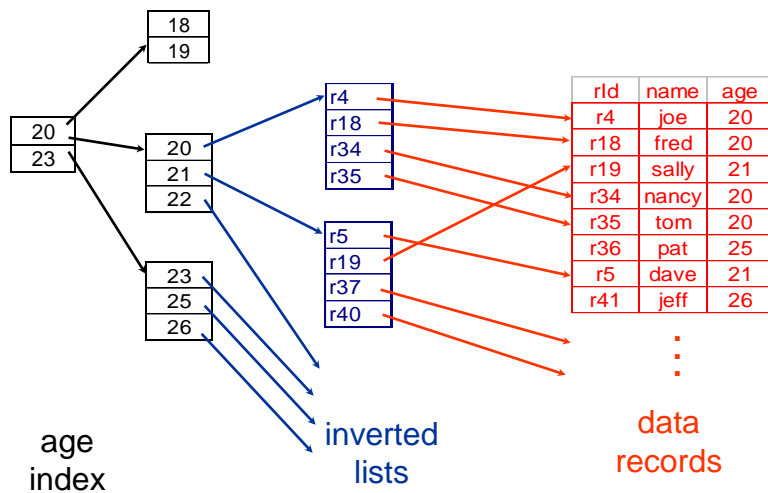
# Index Structures

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

47

Hector Garcia Molina: Data Warehousing and OLAP

# Inverted Lists



48

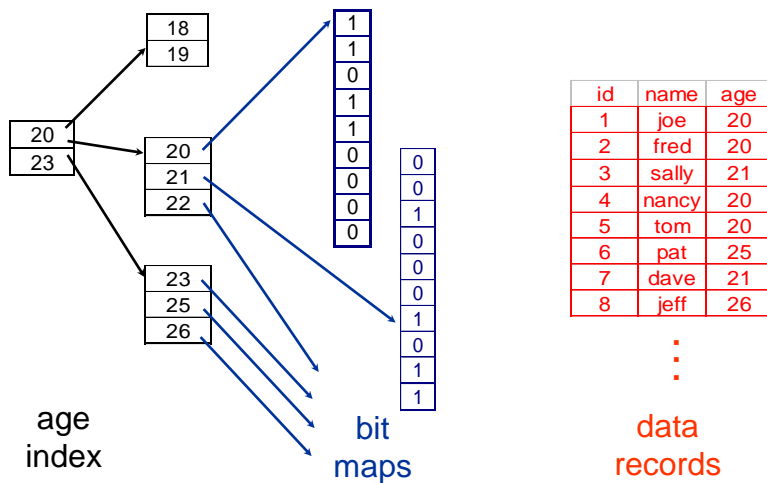
Hector Garcia Molina: Data Warehousing and OLAP



# 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

# Bit Maps



# Using Bit Maps

- Query:
  - ◆ Get people with age = 20 and name = "fred"
- List for age = 20: 1101100000
- List for name = "fred": 0100000001
- Answer is intersection: 010000000000
  
- Good if domain cardinality small
- Bit vectors can be compressed

51

Hector Garcia Molina: Data Warehousing and OLAP

# Join

- "Combine" SALE, PRODUCT relations
- In SQL: SELECT \* FROM SALE, PRODUCT

sale	prold	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

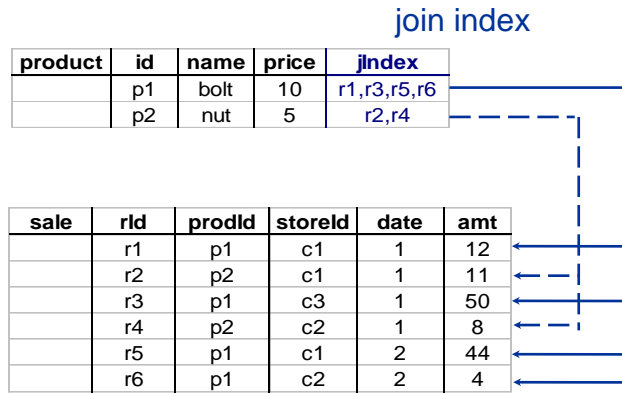
product	id	name	price
	p1	bolt	10
	p2	nut	5

joinTb	prold	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

52

Hector Garcia Molina: Data Warehousing and OLAP

# Join Indexes

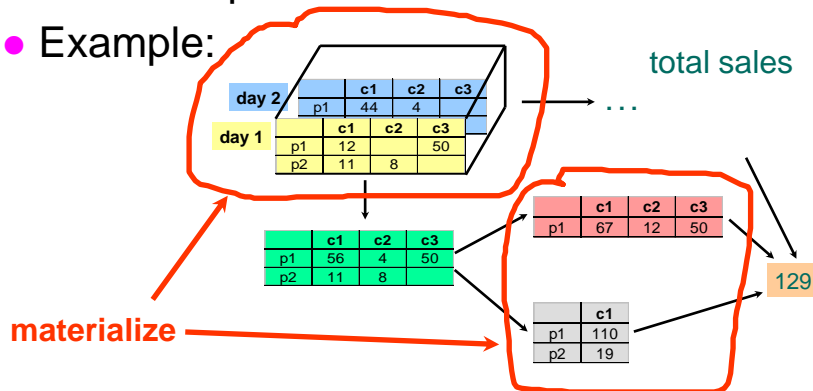


53

Hector Garcia Molina: Data Warehousing and OLAP

# What to Materialize?

- Store in warehouse results useful for common queries
- Example:



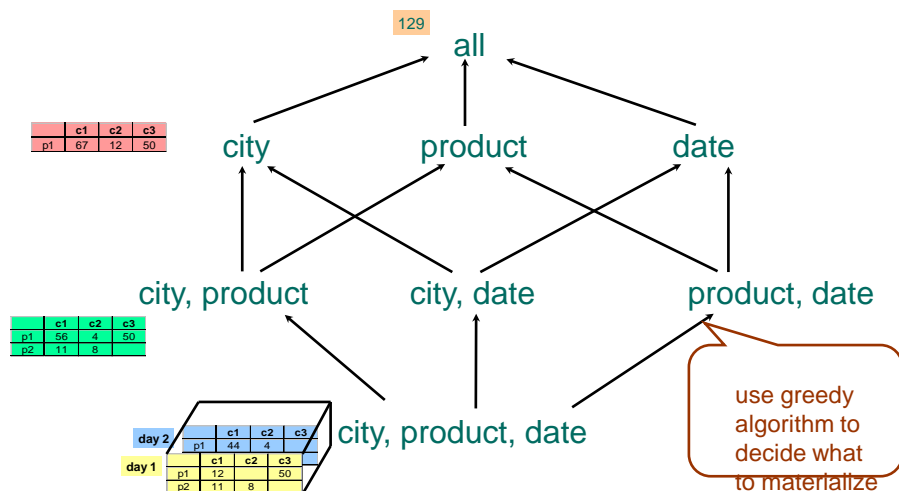
54

Hector Garcia Molina: Data Warehousing and OLAP

# Materialization Factors

- Type/frequency of queries
- Query response time
- Storage cost
- Update cost

# Cube Aggregates Lattice



# Dimension Hierarchies

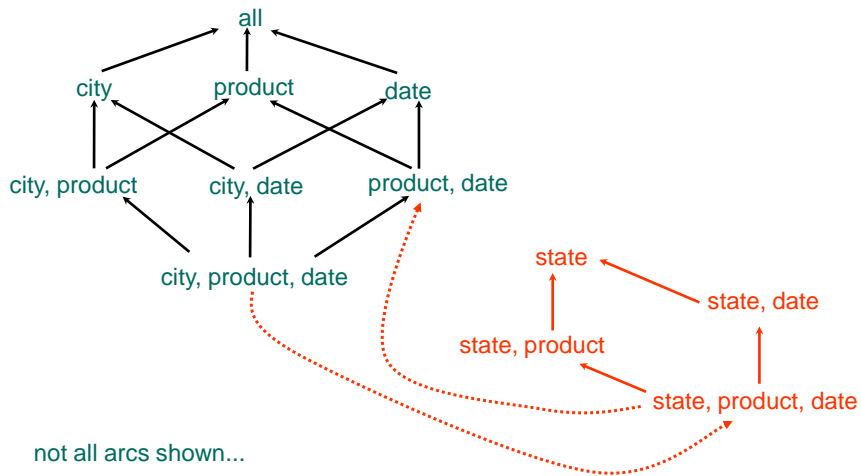


cities	city	state
	c1	CA
	c2	NY

57

Hector Garcia Molina: Data Warehousing and OLAP

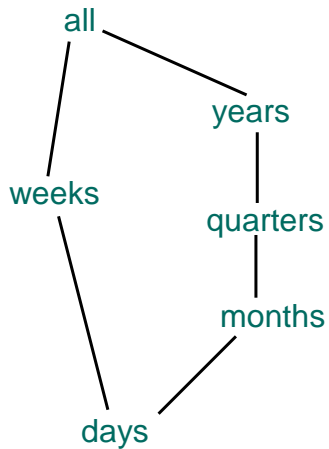
# Dimension Hierarchies



58

Hector Garcia Molina: Data Warehousing and OLAP

# Interesting Hierarchy



time	day	week	month	quarter	year
1	1	1	1	1	2000
2	1	1	1	1	2000
3	1	1	1	1	2000
4	1	1	1	1	2000
5	1	1	1	1	2000
6	1	1	1	1	2000
7	1	1	1	1	2000
8	2	1	1	1	2000

conceptual  
dimension table

59

Hector Garcia Molina: Data Warehousing and OLAP

# 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?

60

Hector Garcia Molina: Data Warehousing and OLAP

## 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

61

Hector Garcia Molina: Data Warehousing and OLAP

## 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

62

Hector Garcia Molina: Data Warehousing and OLAP