

Data Warehousing

INF 551

Wensheng Wu

Warehousing

- Growing industry:
 - \$8 billion in 1998
 - \$16.9 billion in 2018
- Range from desktop to high-end system:
 - Wal-Mart in 1990's
 - 900-CPU, 2,700 disk, 23TB, Teradata system
 - What about 2007?
 - Wal-Mart Plans for Its 4PB Data Warehouse
 - 100 billion rows in tables
 - 276 million records/day from POS systems
- Lots of buzzwords, hype
 - Data cube, roll-up, drill-down, slice/dice, ...

Outline

- What is a data warehouse?
- Why a warehouse?
- Data models & operations
- Implementing a warehouse

What is a Warehouse?

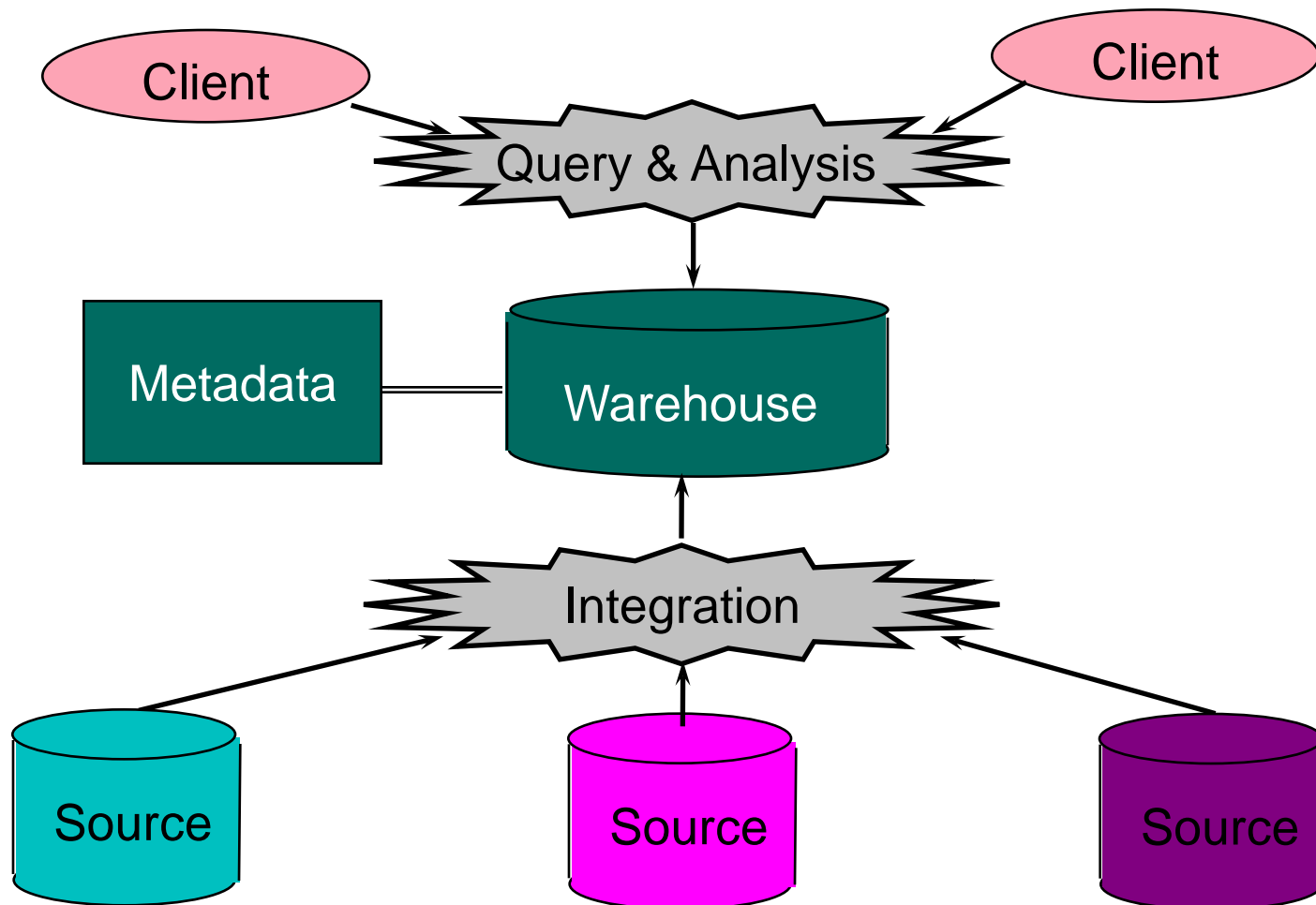
- Collection of diverse data
 - subject (area) oriented, e.g., sales
 - aimed at executive, decision maker
 - often a copy of operational data
 - with value-added data (e.g., summaries, history)
 - integrated
 - time-varying: historical data, discovering trend
 - non-volatile: once in warehouse, data do not change



What is a Warehouse?

- Collection of tools
 - gathering data
 - cleansing & integrating data
 - querying, reporting, analyzing data
 - mining data
 - monitoring, administering warehouse

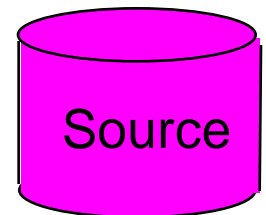
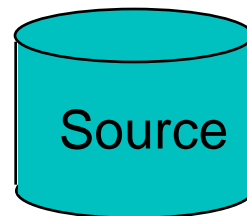
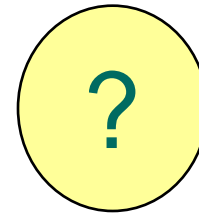
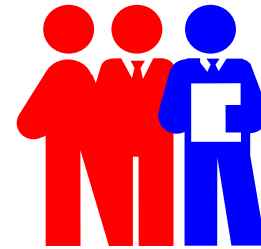
Warehouse Architecture



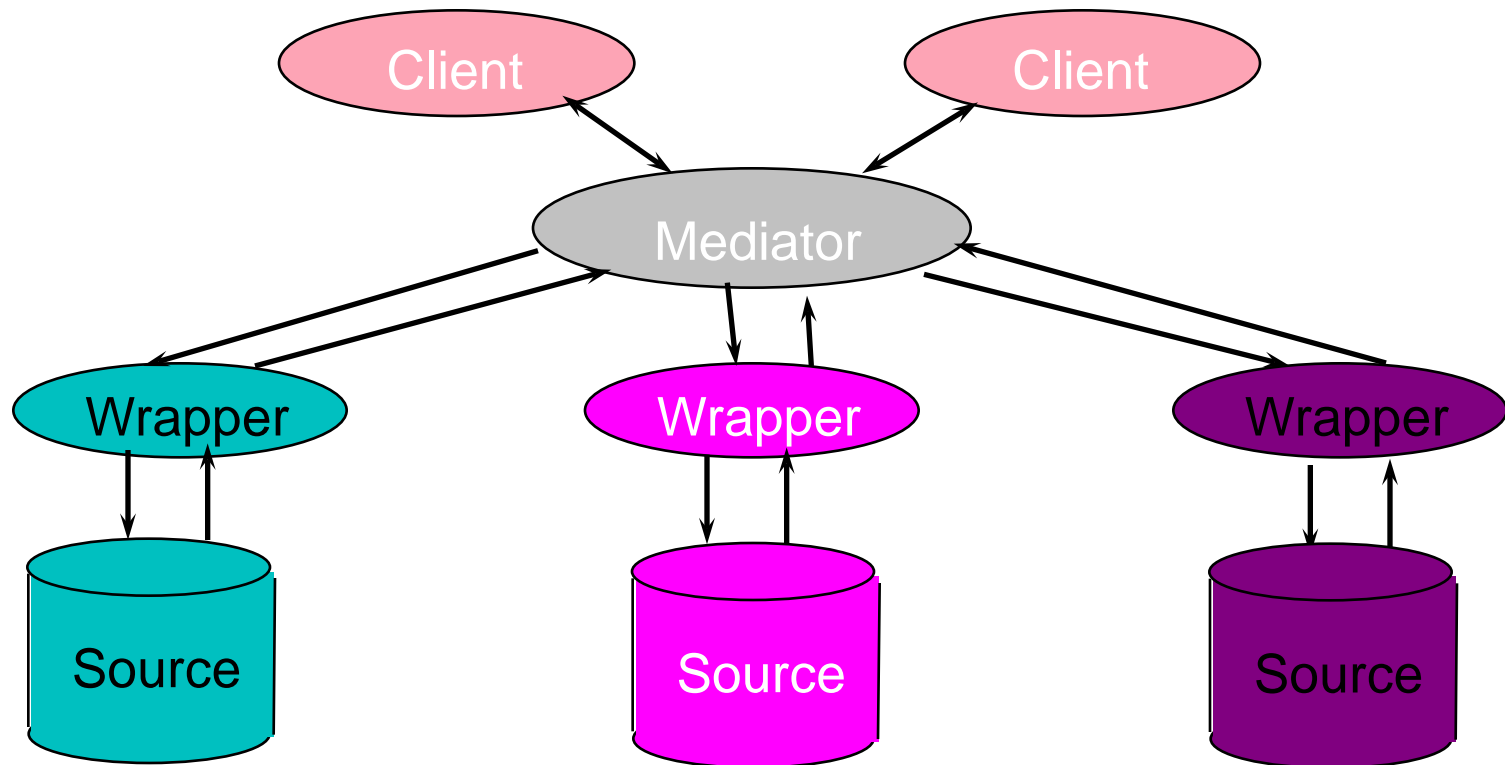
Why a Warehouse?

- Two Approaches to Information Integration:

- Warehouse (Eager)
- Query-Driven (Lazy)



Query-Driven Approach



Advantages of Warehousing

- High query performance
- Local processing at sources unaffected
- Can operate when sources unavailable
- Extra information (e.g., pre-computed summary) at warehouse

Disadvantages of Warehousing

- Decide what to store in advance
- Can only query data stored in warehouse
- Data get stale
- Must detect source changes & update warehouse

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

Disadvantages of Query-Driven

- Inefficient/delay in query processing
 - source unreliable
 - slow network
 - expensive translation, filtering, merging
- Sources might not permit ad-hoc queries
 - Examples?

OLTP vs. OLAP

- OLTP: On Line Transaction Processing
 - ◆ Describes processing at operational sites (order entry in POS/online, banking transactions, etc.)
- OLAP: On Line Analytical Processing
 - ◆ Describes processing (answering analytical queries: aggregation, rollup/drilldown, slice/dice, etc.) at warehouse

OLTP vs. OLAP

OLTP

- Mostly updates
- Many small transactions
- Mb-Tb of data
- Raw data
- Clerical users/clients/customers
- 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
- Historical data
- Query performance critical

Data Marts

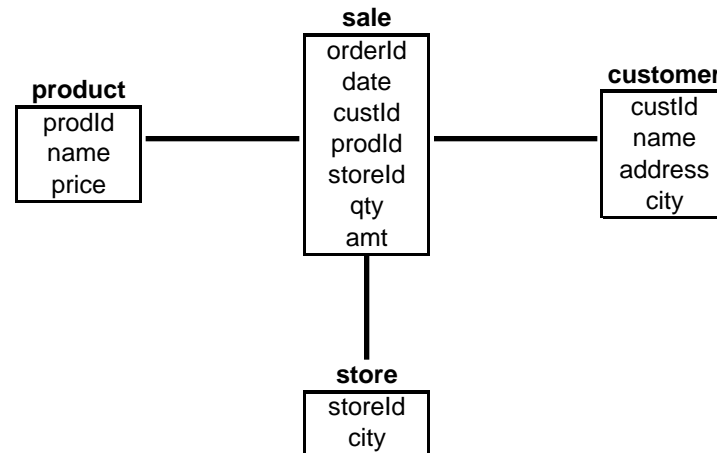
- Smaller warehouses
- Spans part of organization
 - e.g., marketing, sales, financial data marts
- Do not require enterprise-wide consensus
 - but long term integration problems?

Warehouse Models & Operators

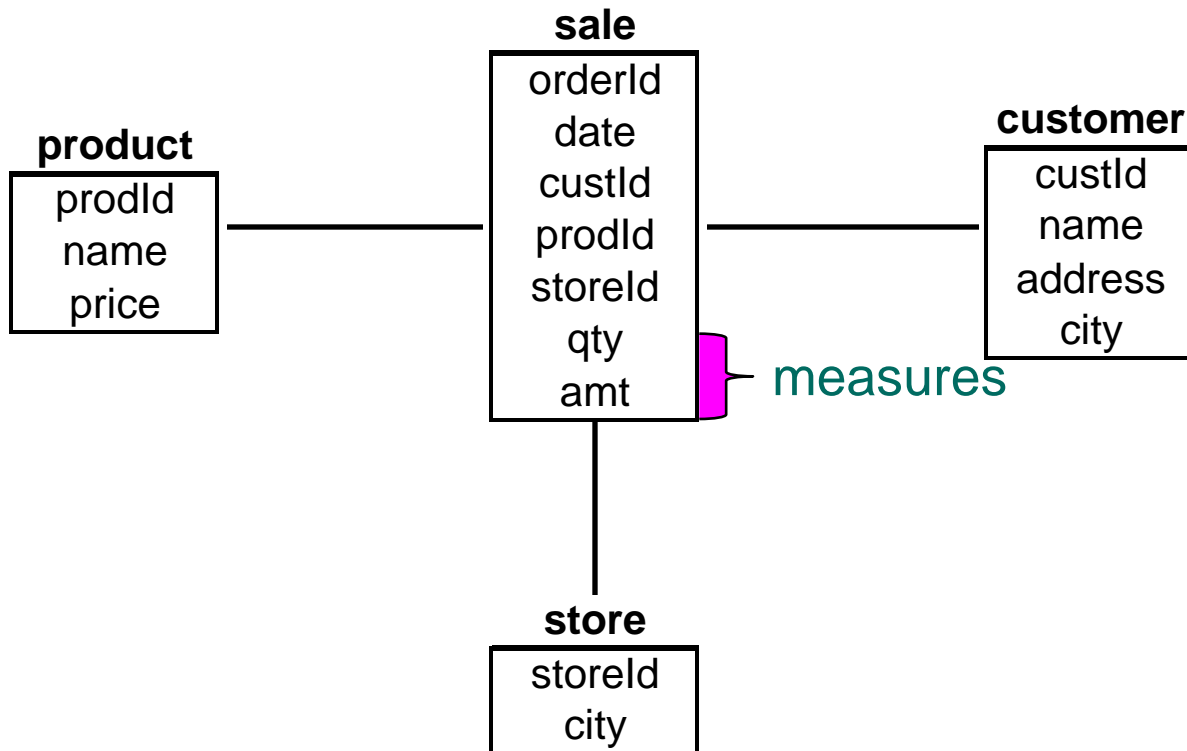
- Multi-dimensional data model (data cube)
 - Organizing measures/facts around different dimensions
- Relational implementation
 - Special schema design: star, snowflake, etc.
- Operators
 - roll-up, drill down
 - slice & dice
 - ...

Terms

- Fact table
- Dimension tables
- Measures



Star Schema



Star

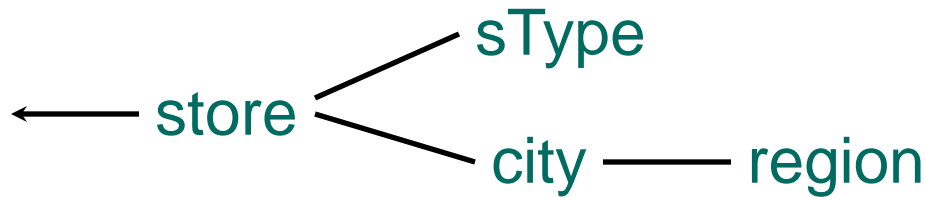
product	prodId	name	price
	p1	bolt	10
	p2	nut	5

store	storeId	city
	c1	nyc
	c2	sfo
	c3	la

sale	oderId	date	custId	prodId	storeId	qty	amt
	o100	1/7/97	53	p1	c1	1	12
	o102	2/7/97	53	p2	c1	2	11
	105	3/8/97	111	p1	c3	5	50

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

Dimension Hierarchies



store	storeld	cityld	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	cityld	pop	regld
	sfo	1M	north
	la	5M	south

→ snowflake schema

region	regld	name
	north	cold region
	south	warm region

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	

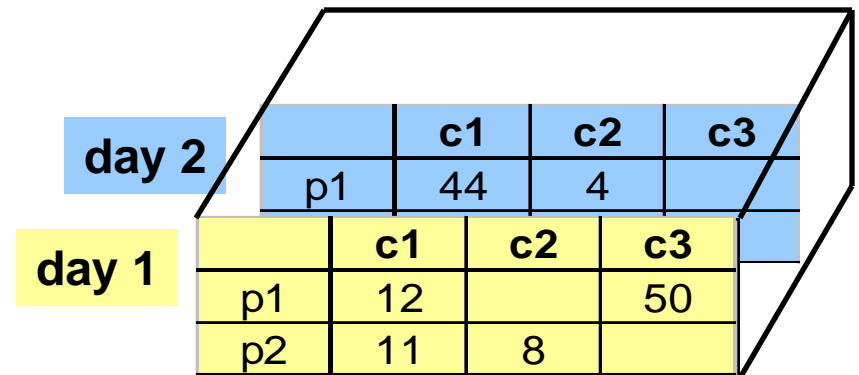
dimensions = 2

3-D Cube

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:

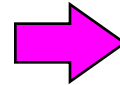


dimensions = 3

Aggregates

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

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



ans	date	sum
	1	81
	2	48

—— rollup —→

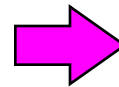
←—— drill-down

Another Example

- Add up amounts by product, day
- In SQL:

```
SELECT prodid, date, sum(amt)  
FROM SALE  
GROUP BY prodid, date
```

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

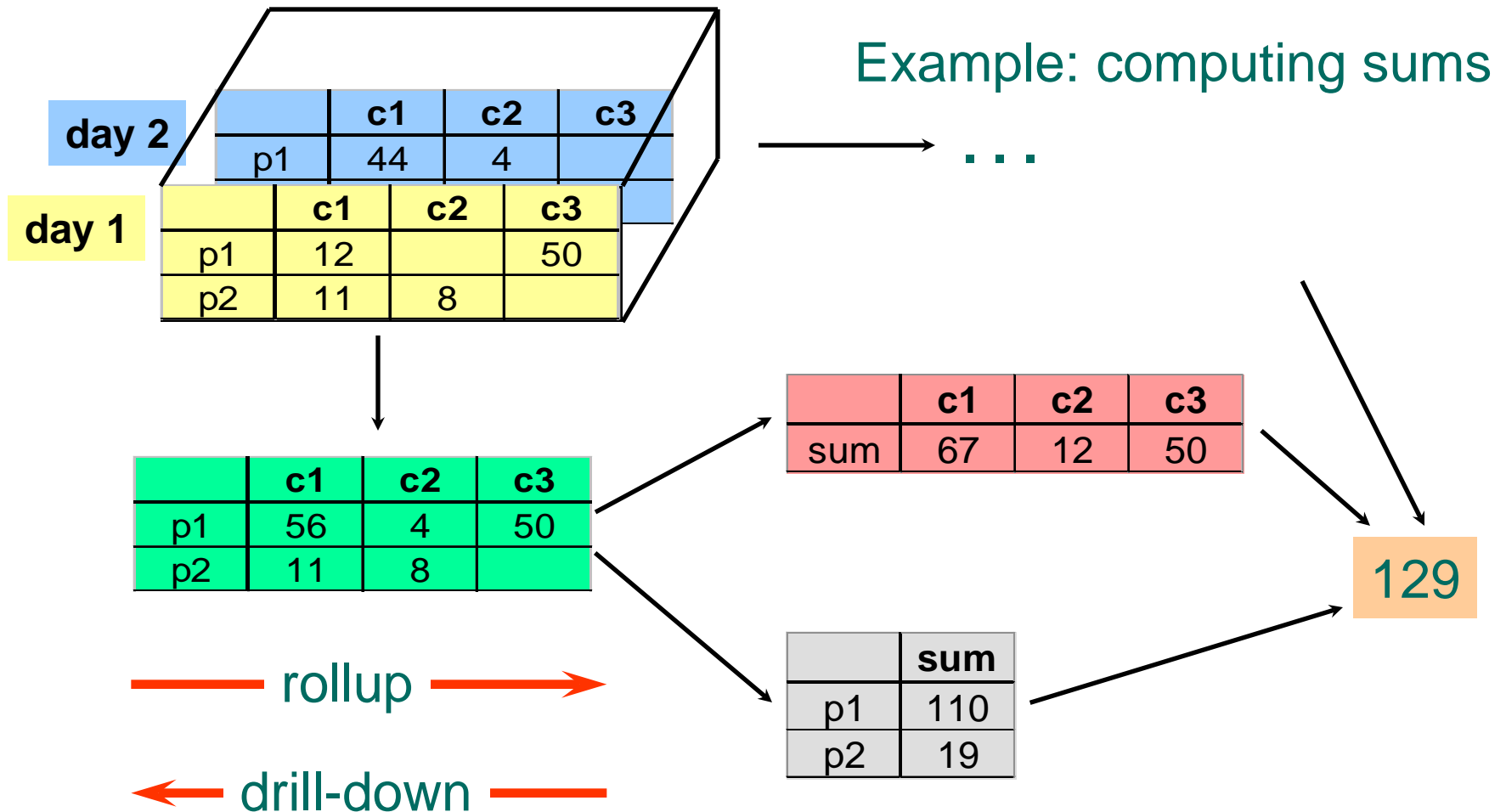
—— rollup —→

←—— drill-down

Aggregates

- Operators: sum, count, max, min, avg
- May have “having” clause
- May also use dimension hierarchy, e.g.,
 - average by region (store dimension)
 - maximum by month (date dimension)

Cube Aggregation



Aggregation Using Hierarchies

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

	region A	region B
p1	56	54
p2	11	8

country
|
region
|
store

(store c1 in Region A;
stores c2, c3 in Region B)

Implementing a Warehouse

- Extracting/monitoring: getting data from sources
- Integrating: cleaning, transforming, loading
- Processing: query processing, indexing, materialized view, ...
- Managing: metadata, performance monitoring, ...

Extracting & Monitoring

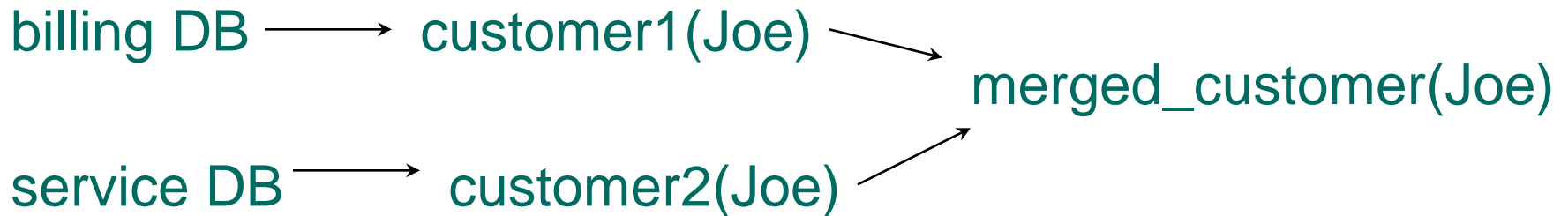
- Source Types: relational, flat file, XML, web sites, blogs, news feeds, ...
- Incremental vs. complete refresh

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



Data Cleaning

- Transformation (e.g., yen \Rightarrow dollars)
- Remove errors or inconsistency
 - E.g., 10 digit ssn, J. Smith vs Smith, John
- Fusion & entity resolution



- Auditing: discover rules & relationships (e.g., discover FK/FK)

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**
 - e.g., manage “reliable scripts” for cleaning & analyzing data

SQL Server 2005 (Analysis Service)

Analysis Services Tutorial - Microsoft Visual Studio

File Edit View Project Build Debug Database Cube Tools Window Community Help

Analysis Service...al.cube [Design] Start Page Solution Explorer

Cube Struct... Dimension Usage Calculations KPIs Actions Partitions Perspectives Translations Browser

Perspective: Analysis Se Language: Default

Dimension Hierarchy Operator Filter Expression

<Select dimension>

Drop Filter Fields Here

		Category			
		Accessories	Bikes	Clothing	Grand Total
Country-Region	State-Province	Internet Sales-Sales Amount	Internet Sales-Sales Amount	Internet Sales-Sales Amount	Internet Sales-Sales Amount
Australia	New South Wales	\$59,838.56	\$3,843,861.53	\$30,785.64	\$3,934,485.13
	Queensland	\$31,555.51	\$1,942,682.65	\$14,176.87	\$1,988,415.03
	South Australia	\$8,838.42	\$605,055.60	\$4,361.84	\$618,255.86
	Tasmania	\$4,040.85	\$233,197.18	\$2,699.87	\$239,937.90
	Victoria	\$34,417.29	\$2,227,253.04	\$18,235.73	\$2,279,906.06
	Total	\$138,690.63	\$8,852,050.00	\$70,259.95	\$9,061,000.58
Canada		\$103,377.85	\$1,821,302.39	\$53,164.62	\$1,977,844.86
France		\$63,005.65	\$2,523,523.04	\$26,793.77	\$2,613,322.46
Germany		\$62,232.59	\$2,808,514.35	\$23,565.40	\$2,894,312.34
United Kingdom		\$75,155.42	\$3,231,209.26	\$31,788.65	\$3,338,153.33
United States		\$258,297.82	\$9,081,545.60	\$134,200.22	\$9,474,043.64
Grand Total		\$700,759.96	\$28,318,144.65	\$339,772.61	\$29,358,677.22

Output Error List

ETL Example

- Consider integrating data from two dealers

- Dealer 1 schema:

- Cars(serialNo, model, color, autoTrans, navi)
- Attributes starting from autoTrans are binary attributes for car options
- E.g., autoTrans = 'Yes', if particular car has auto. transmission

123, 'accord', 'silver', 'Yes', 'No'

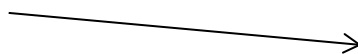


- Dealer 2 schema:

- Autos(serial, model, color)
- Options(serial, option)

123, 'autoTrans'

123, 'navi'

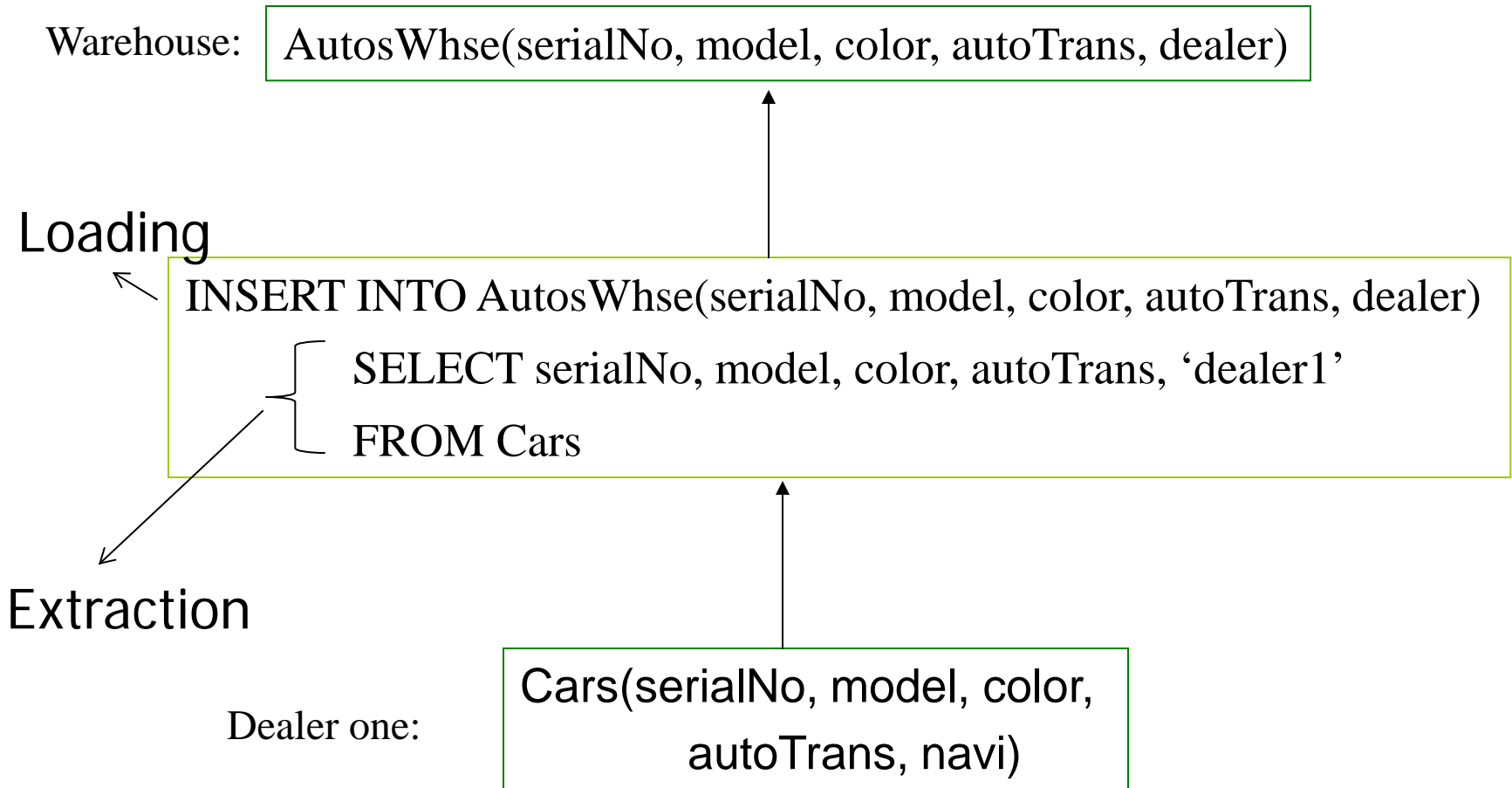


...

- Warehouse schema:

- AutosWhse(serialNo, model, color, autoTrans, dealer)

Extract Data from Dealer One



Exercise: Extracting from Dealer Two

Warehouse:

AutosWhse(serialNo, model, color, autoTrans, dealer)

INSERT INTO AutosWhse(serialNo, model, color, autoTrans, dealer)
... ?

Dealer two:

Autos(serial, model, color)
Options(serial, option)