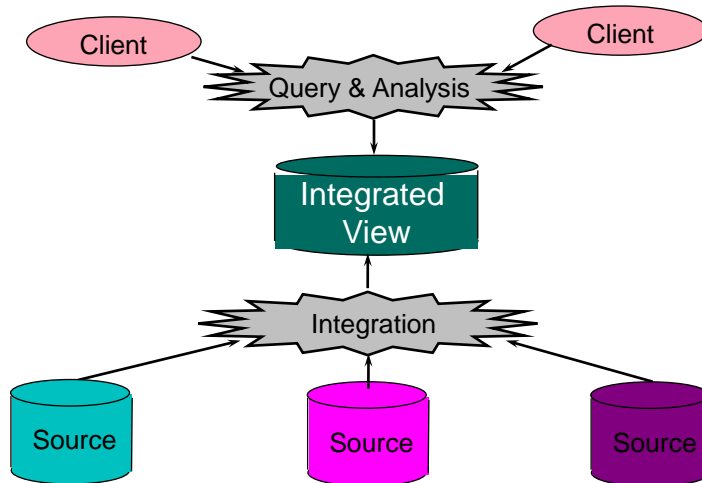
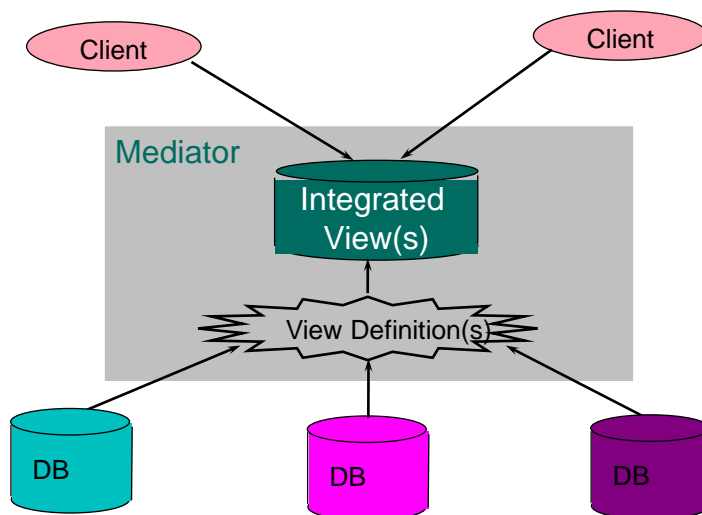


Warehouses & Virtual Databases offer Integrated & Added-Value Views

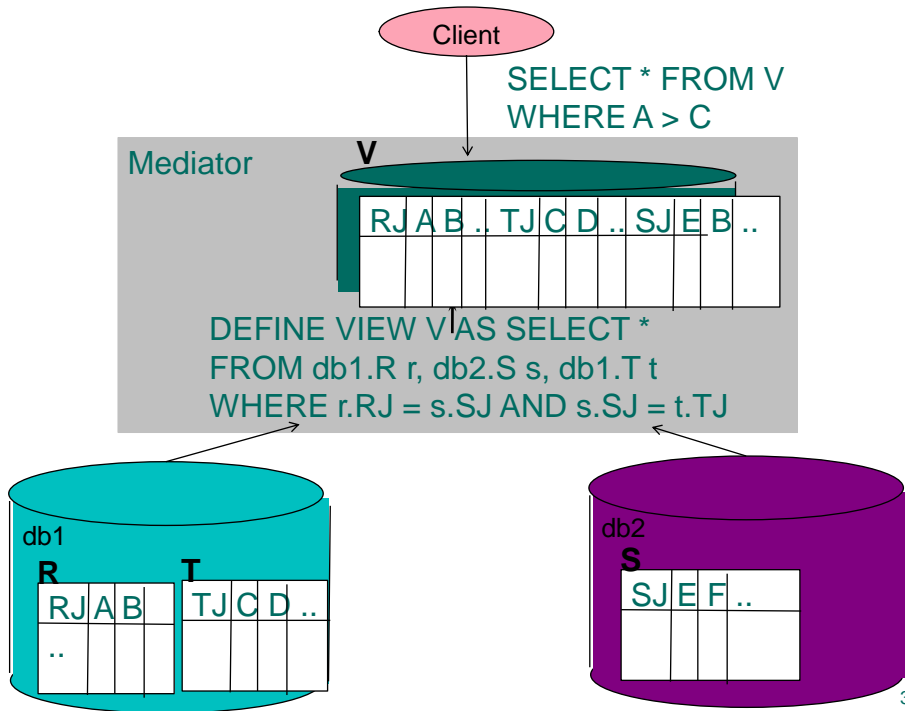


1

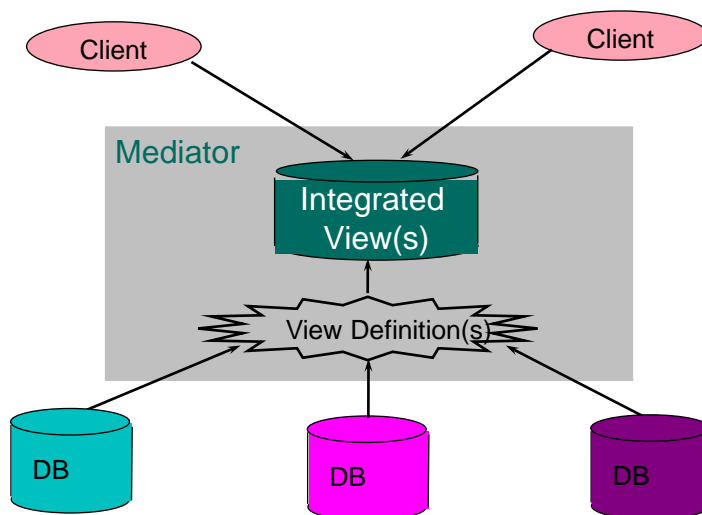
Focus: Sources are relational DBs. Integration specified by distributed view definition(s). Clients issue queries on views.



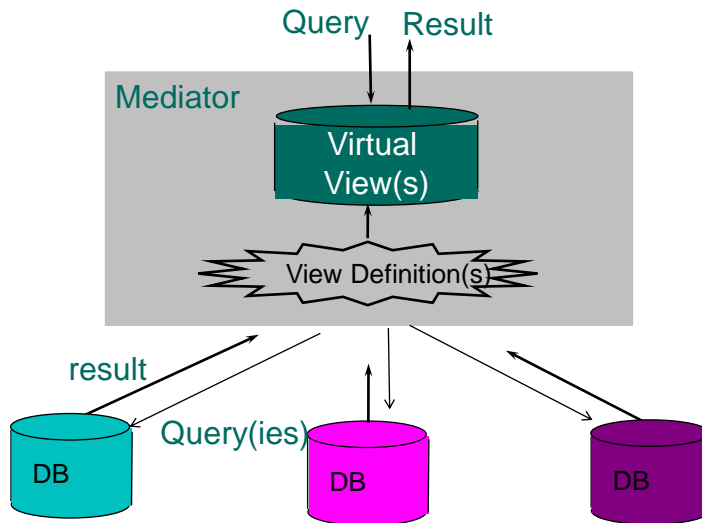
2



Virtual View -> Mediator is a Distributed Query Processor
 Materialized view (warehouse) -> Mediator actually stores
 integrated view

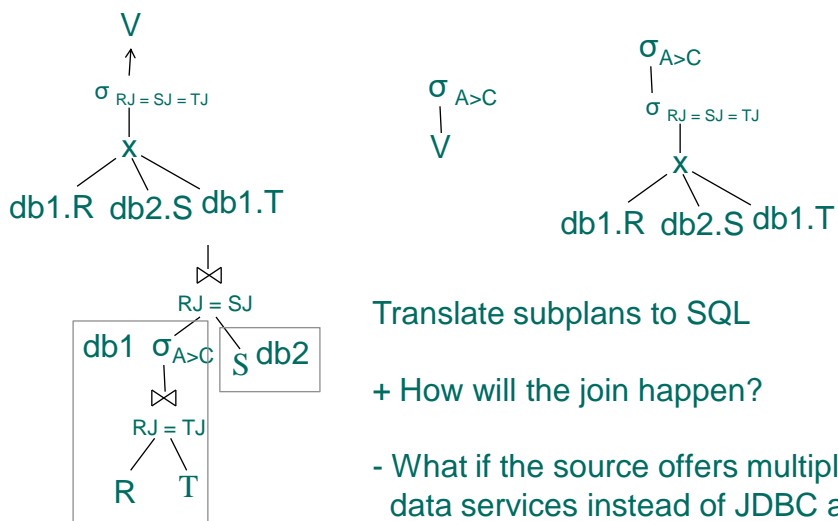


Distributed Query Processing in Mediators



5

Distributed Query processing



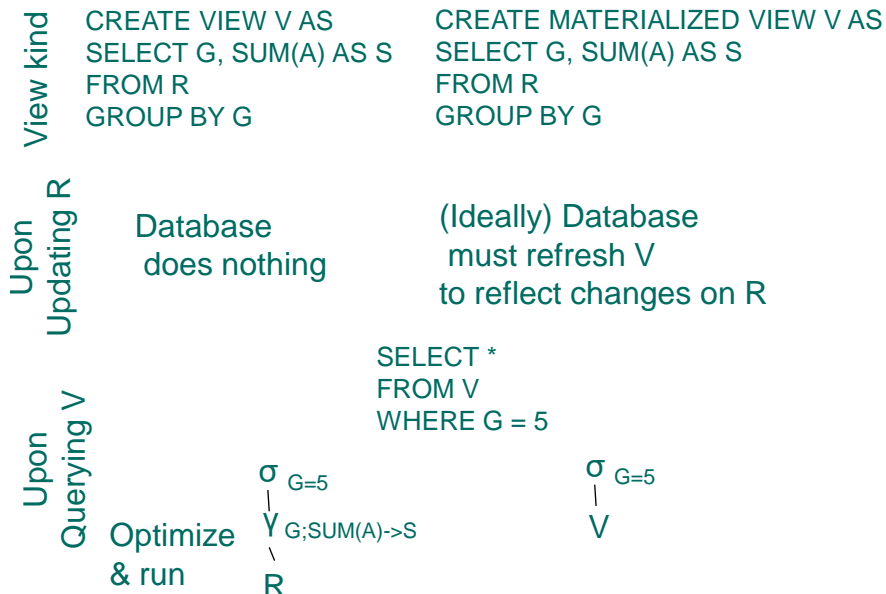
6

Distributed Join Types

- Mediator-based Join
 - ◆ Ship results of queries at mediator
- Parameterized Join
 - ◆ Right subquery is enhanced with selection on join attribute
 - ◆ For each join value of left hand side, execute another right subquery
- Data Ship Join
 - ◆ Insert the result of left hand side (lhs) in the db of right hand side (rhs).
 - ◆ Execute join at db of right hand side
- Semijoin Reduction Join
 - ◆ Send rhs parameters to lhs
 - ◆ (Data ship alike variation) Lhs sends to rhs the semijoin of its subquery with the parameters set.
 - ◆ Execute join at db of rhs
 - ◆ Also, variation that looks like mediator-based join

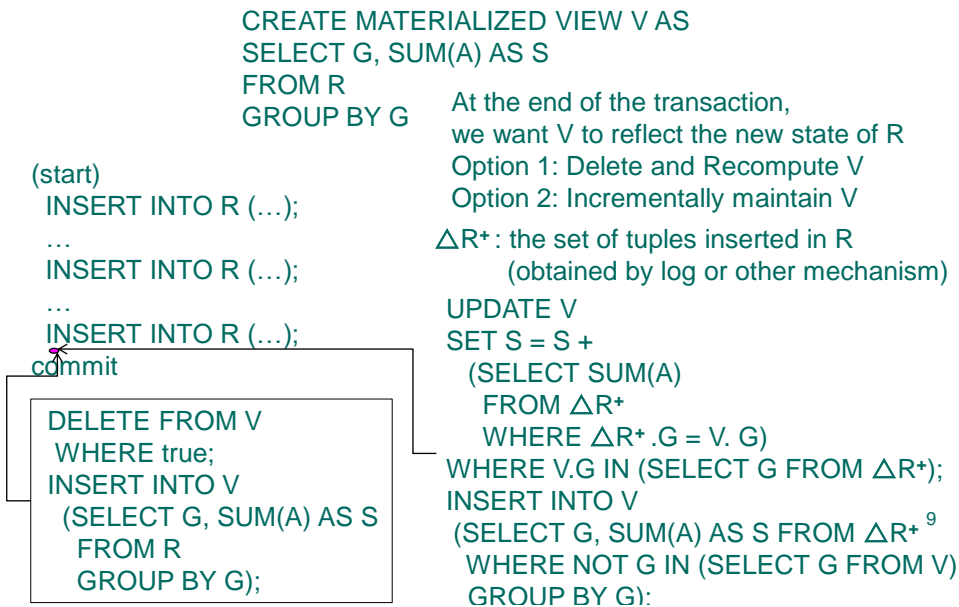
7

Virtual Views Vs Materialized Views

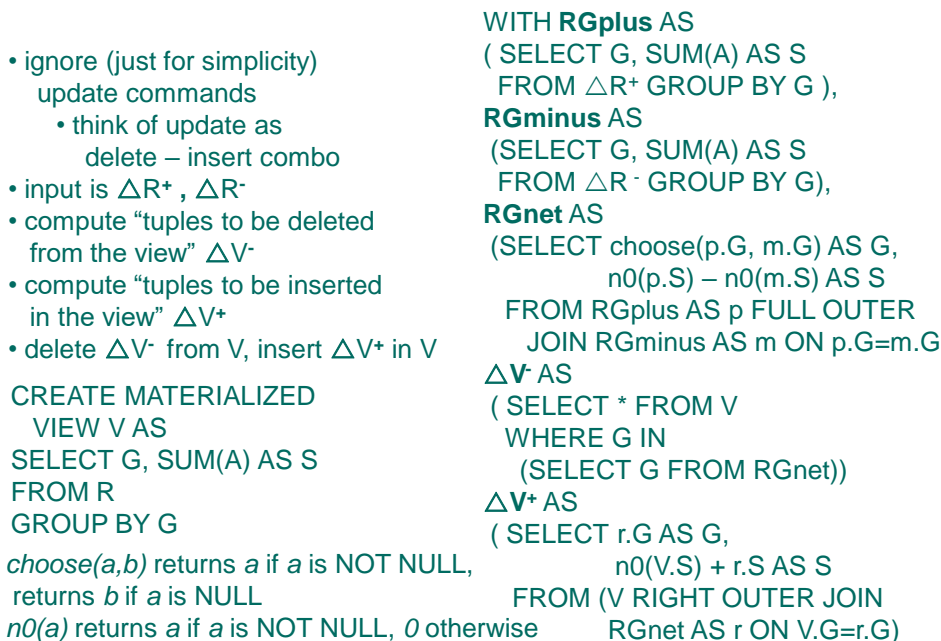


8

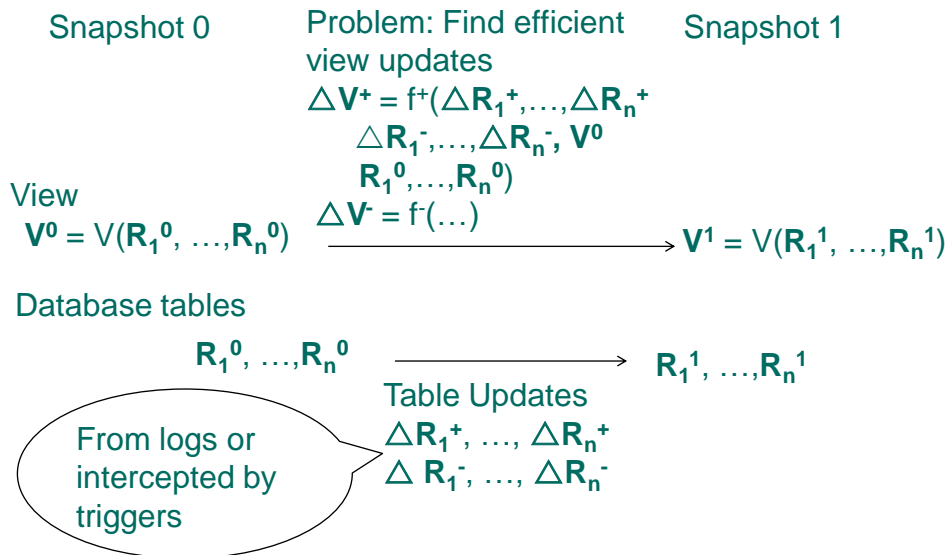
Recompute Vs Incremental (Materialized View) Maintenance – Informal Example



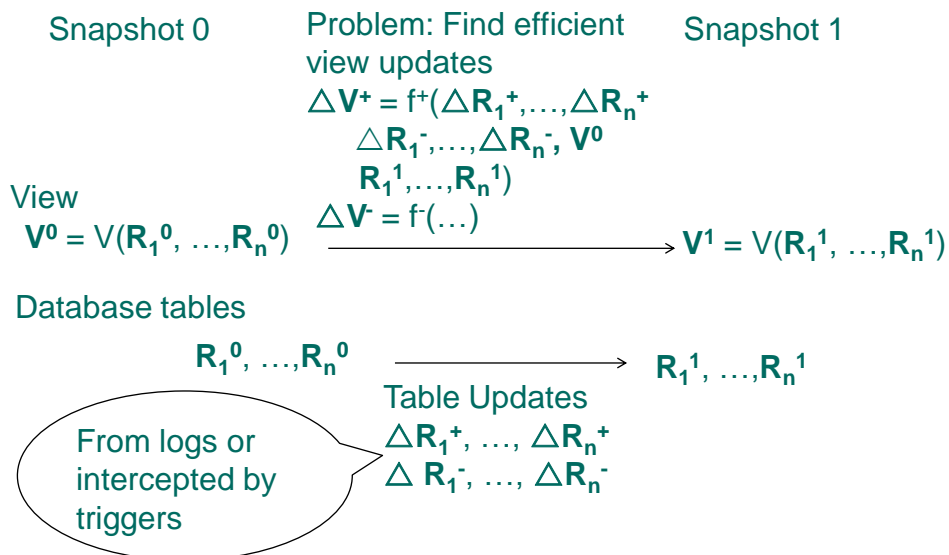
Capturing IVM as computation of ΔV^+ , ΔV^-



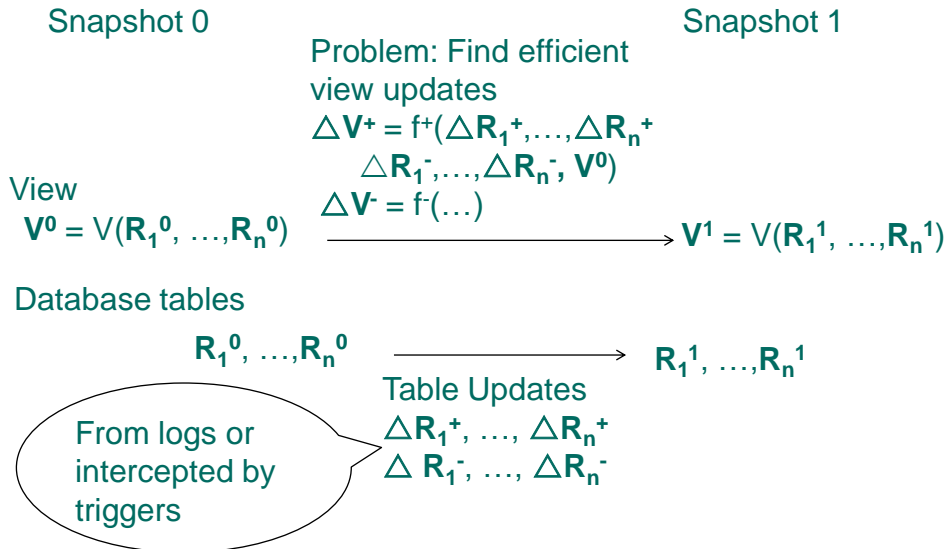
IVM: Incremental (Materialized) View Maintenance. Eager version.



IVM: Deferred version



IVM: Self-maintaining version (not always possible)



Basic IVM Algorithm: Compose operator IVM rules

Example (wlog deferred, i.e., R means R^1 and S means S^1)

- Rule for $V = R \bowtie S$
 - ◆ $\Delta \mathbf{V}^+ = ((\Delta \mathbf{R}^+ \bowtie S) \cup (R \bowtie \Delta \mathbf{S}^+)) - (\Delta \mathbf{R}^+ \bowtie \Delta \mathbf{S}^+)$
 - ◆ $\Delta \mathbf{V}^- = ???$
- Rule for $V = \sigma_c R$
 - ◆ $\Delta \mathbf{V}^+ = \sigma_c \Delta \mathbf{R}^+$
 - ◆ $\Delta \mathbf{V}^- = ???$

- Composition of rules leads to solutions for

$$V = T \bowtie \sigma_{A > 5} W$$

$$\Delta \mathbf{V}^+ = \dots$$

$$\Delta \mathbf{V}^- = \dots$$

- May rewrite initial expression

IVM with Caching

- May associate intermediate views (caches) with subexpressions
- Bottom-up: From updating caches to reaching the materialized view
- Caches will typically need indices
- Caches may or may not pay off as they incur cost for maintaining them (and their indices)

15

Generalizations

- Multiple views
 - ◆ self maintenance may involve a view utilizing the other views in its computation
- Genuine updates
 - ◆ Not simulated via insertions/deletions
- Insertions, deletions, updates on tables and views expressed as DML statements

16

Comparisons

Materialized View

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

Virtual View

- No need for yet another database
 - More up-to-date data
 - ◆ Depending on specifics of IVM
 - Query needs can be unknown
 - Only query interface needed at sources
- => Lower Total Cost of Ownership

17

Performance revisited: What if indices are not enough for decent online performance?

- Buy RAM
- Use a column database
 - ◆ In analytics queries can give a 10x easily
- Scalable, parallel processing
 - ◆ Mostly via no SQL
- **Precompute**
 - ◆ **Fast answers!**
 - ◆ **Penalty: Cost of maintaining precomputed results**
 - ◆ **Applicability depends on schema and queries**
 - ◆ **Star schemas and summation are a good (but not the only) target of precomputation**

Precomputation problems

Steps:

1. Choose what data to precompute
2. Use the precomputed data smartly in your queries
3. Update smartly the precomputed data as the database changes (IVM)

Tradeoff:

- Precomputed data accelerate analytics => faster queries
- But need to be updated => cost

Example: Precomputation and its Use

Database has huge table `Sales(product, store, date, amt)`

Application issues often this slow query and displays the results

```
SELECT product, SUM(amt) AS sumamt
FROM Sales
GROUP BY product
```

To improve performance we precompute table

```
ProductSales(product, sumamt)
```

and insert in it the precomputed data by

```
INSERT INTO ProductSales (
SELECT product, SUM(amt) AS sumamt
FROM Sales
GROUP BY product )
```

Now the application issues instead this fast query below

```
SELECT *
FROM ProductSales
```

Example (cont'd)

Now we have to keep up to date the

```
ProductSales(product, sumamt)
```

as new sales happen. E.g., if another \$10 of product 23 were just sold

```
UPDATE ProductSales
```

```
SET sumamt = sumamt + 10
```

```
WHERE product = 23
```

(in actual code it will use prepared queries)

You do not need the “exact” view

- Consider V1(Product, Customer, Sales) and V2(Product, Customer, Date, Sales) are precomputed

- ProductSales is not precomputed

- You need to answer the query

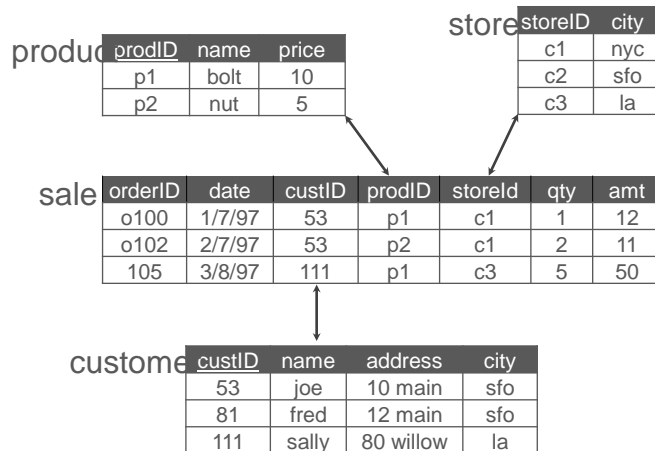
```
SELECT product, SUM(amt)
```

```
FROM SALES
```

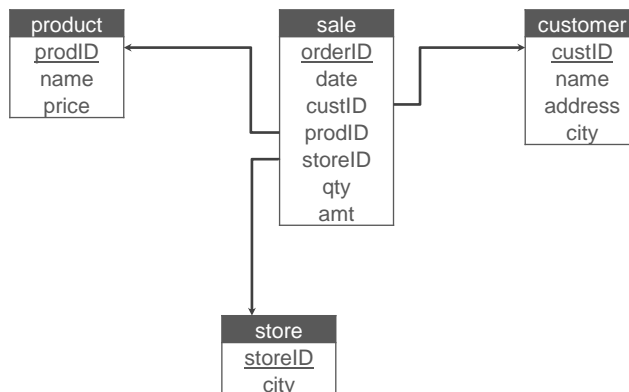
```
GROUPBY product
```

- Write it in an alternate way, using one of the views in the most efficient way

Star Schemas

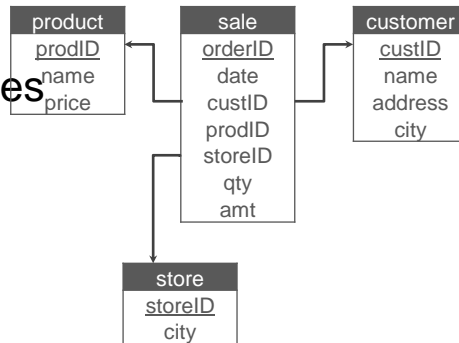


Star Schema

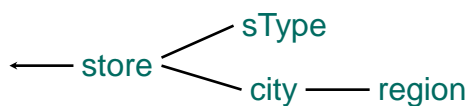


Terms

- Fact table
- Dimension tables
- Measures



Dimension Hierarchies



store	<u>storeID</u>	cityID	tID	mgr
	s5	sfo	t1	joe
	s7	sfo	t2	fred
	s9	la	t1	nancy

sType	<u>tID</u>	size	location
	t1	small	downtown
	t2	large	suburbs

city	<u>cityID</u>	pop	regID
	sfo	1M	north
	la	5M	south

region	<u>regID</u>	name
	north	cold region
	south	warm region

- Snowflake Schema
- Constellations

Cube

Fact table view

sale	prodlid	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

for

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

	c1	c2	c3
p1	12		50
p2	11	8	

dimensions = 3

Aggregates on Slices

- Add up amounts for day 1
 - ◆ 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

Aggregates

- Add up amounts by day
 - ◆ 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

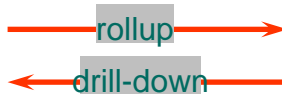
Another Example

- Add up amounts by day, product
 - ◆ SELECT date, sum(amt) FROM SALE
 - ◆ GROUP BY date, prodid

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



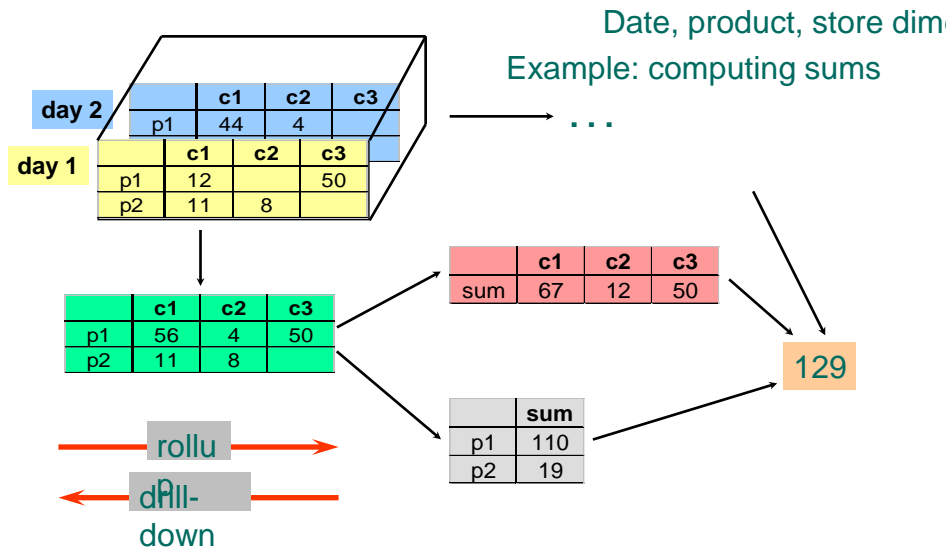
sale	prodid	date	amt
	p1	1	62
	p2	1	19
	p1	2	48



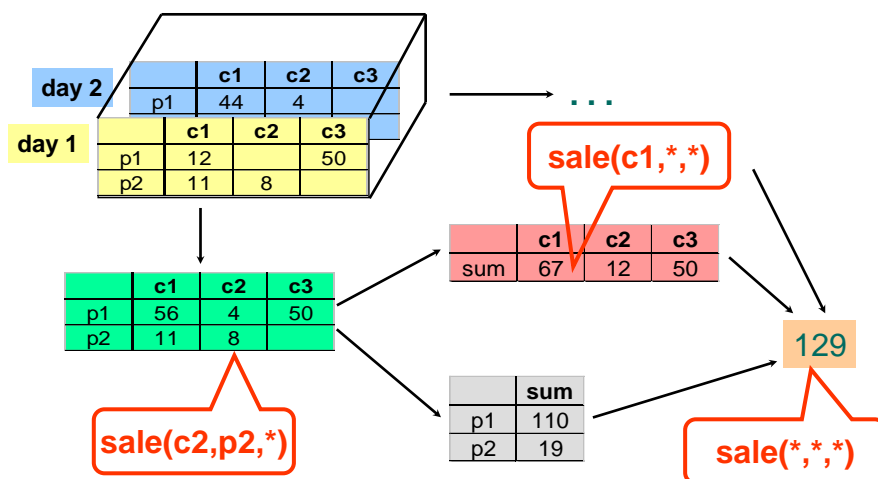
Aggregates

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

Cube Aggregation



Cube Operators

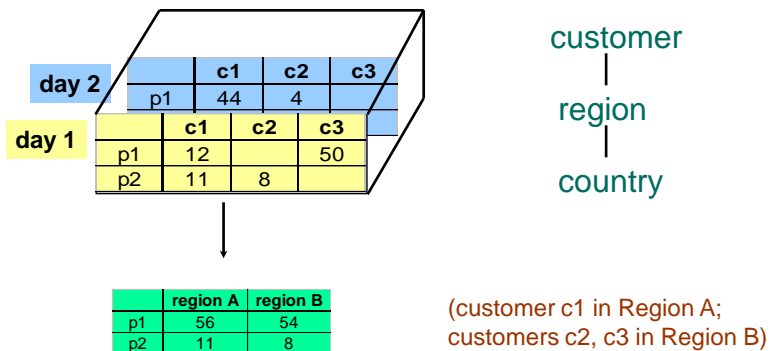


Extended Cube

		c1	c2	c3	*
day 2	*				
	p1	56	4	50	110
	p2	11	8		19
day 1					
	p1	44	4		48
	p2	12		50	62
	*	23	8	50	81

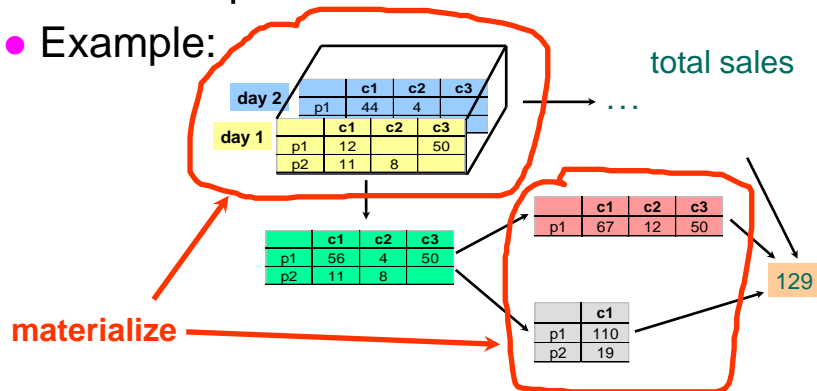
sale(*,p2,*)

Aggregation Using Hierarchies



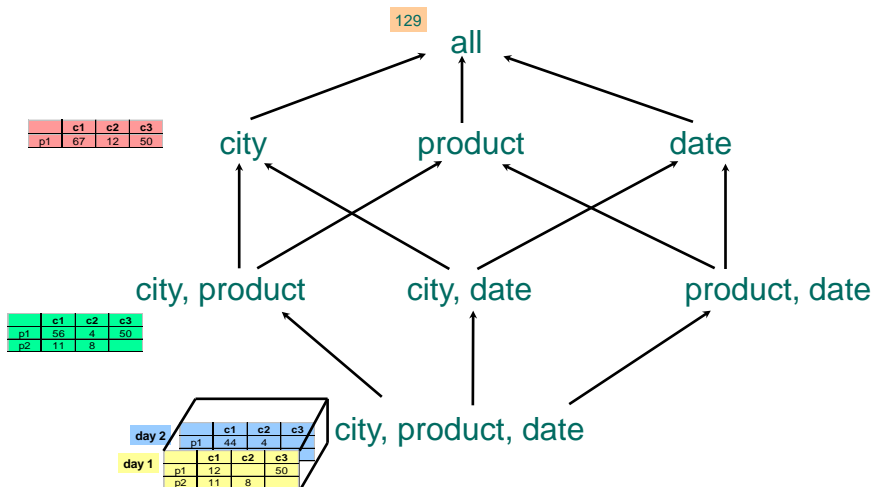
What to Materialize?

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



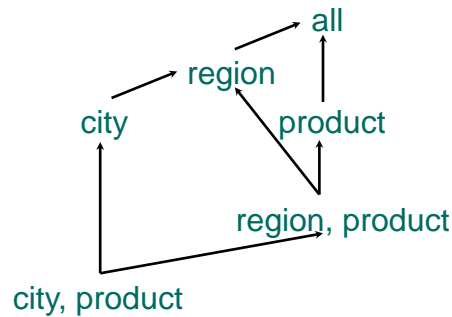
Cube Aggregates Lattice

Example assumes fact table is sales(city, product, date, amt)



Cube Aggregates Lattice

Example assumes fact table is sales(city, product, amt) and cities c



Should one precompute joins?

- Notice that we have featured foreign keys, not printable values. Why?
- Why (city product) and not (city region product)?
- Minor penalty to find the cities of a particular region
- Probably larger penalty by having a larger table
 - ◆ Think space in storage and time to scan it