Data Warehousing and Decision Support

Ashraf Aboulnaga

David R. Cheriton School of Computer Science University of Waterloo

CS 348 Introduction to Database Management Winter 2013

	CS 348	Warehousing	Winter 2013	1 / 24
Votes				

Outline

1 Introduction to Decision Support

On-Line Analytical Processing Multidimensional Data Multidimensional Queries

3 Data Warehousing

Creating and Maintaining a Warehouse Views Materializing Views

	CS 348	Warehousing	Winter 2013	2 / 24
Ιo	tes			
•				

Transaction Processing

The most common use of relational databases is for operational data.

- Examples:
 - Students enrolling in courses
 - Customers purchasing products
 - Passengers purchasing airline tickets

On-	-Line Transactional Processing (OLTP)
	abases that support the basic operations of a business are generally
	sified as OLTP systems. Workload characteristics:
	1 simple queries
•	2 many short transactions making small changes Systems tuned to maximize throughput of concurrent transactions

	CS 348	Warehousing	Winter 2013	3 / 24
lotes				

More recent uses of operational data:

Decision Support Summarizing data to support high-level decision making

• Complex queries with much aggregation

Data Mining Searching for trends or patterns in data for a business to exploit

• Simple queries, but very data-intensive

Data Warehousing

A data warehouse is a separate copy of the operational data used for executing decision support and/or data mining queries.

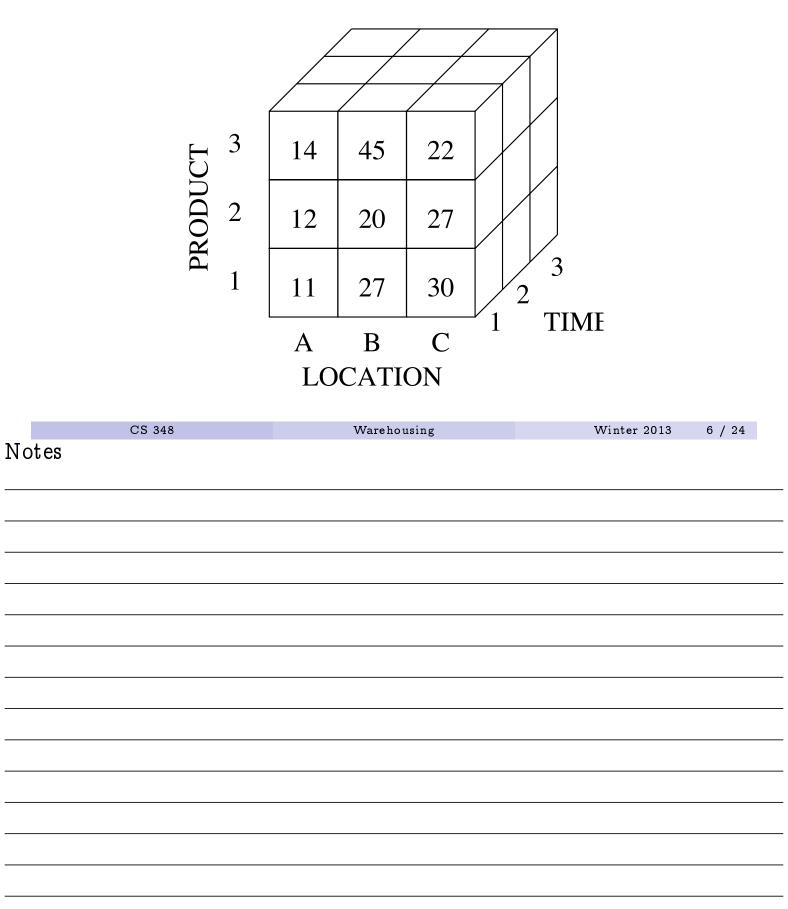
	CS 348	Warehousing	Winter 2013	4 / 24
Notes				

On-Line Analytical Processing

	On-Line Analytical Processing (OLAP)
	OLAP is a particular type of decision support
	• Data is modeled as multidimensional array
	• Queries are usually ad hoc
	• Queries select and aggregate cells of the array
	 OLAP systems are divided into two categories: Special-purpose OLAP systems store data as multidimensional arrays ("MOLAP") provide an OLAP-specific query language Relational databases
	 store data in relations ("ROLAP") queries written in SQL
	CS 348 Warehousing Winter 2013 5 / 24
Note	es

Multidimensional Data

• Example: Number of Sales



Star Schemas

Location

\underline{lid}	store	city	province	country
Α	Weber	Waterloo	ON	CA
В	F-H	Kitchener	ON	CA
С	Park	Kitchener	ON	CA

Product

\underline{pid}	pname	category	price
1	Bolt	Hardware	.10
2	Nut	Hardware	.05
3	Wrench	Tools	1.99

Time

tid	date	week	month	quarter	y ear
		virtı	ual relatio	n	

Sales

Sales				
lid	pid	tid	sales	
Α	1	1	11	
A	2	1	12	
A	3	1	14	
В	1	1	27	
В	2	1	20	
В	3	1	45	
C	1	1	30	
С	2	1	27	
С	3	1	22	
A	1	2	16	
A	2	2	20	
A	3	2	55	
÷				

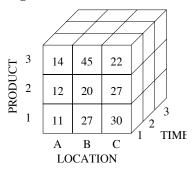
	CS 348	Warehousing	Winter 2013	7 / 24
Notes				

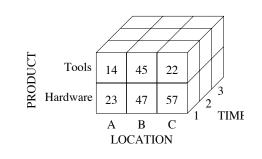
- OLAP queries typically aggregate over one or more dimensions. Examples:
 - Total sales
 - Total sales this year for each product category
 - Total sales for each store per quarter
- OLAP is a tool for *ad hoc* data exploration/visualization
 - Ad hoc queries tend to be iterative
 - Desirable to express queries using operations over previous result

	CS 348	Warehousing	Winter 2013	8 / 24
No	ces			

OLAP Query Operations

- Slicing and Dicing PRODUCT PRODUCT TIME В С В С А A LOCATION / LOCATION
- Roll-up and Drill-down



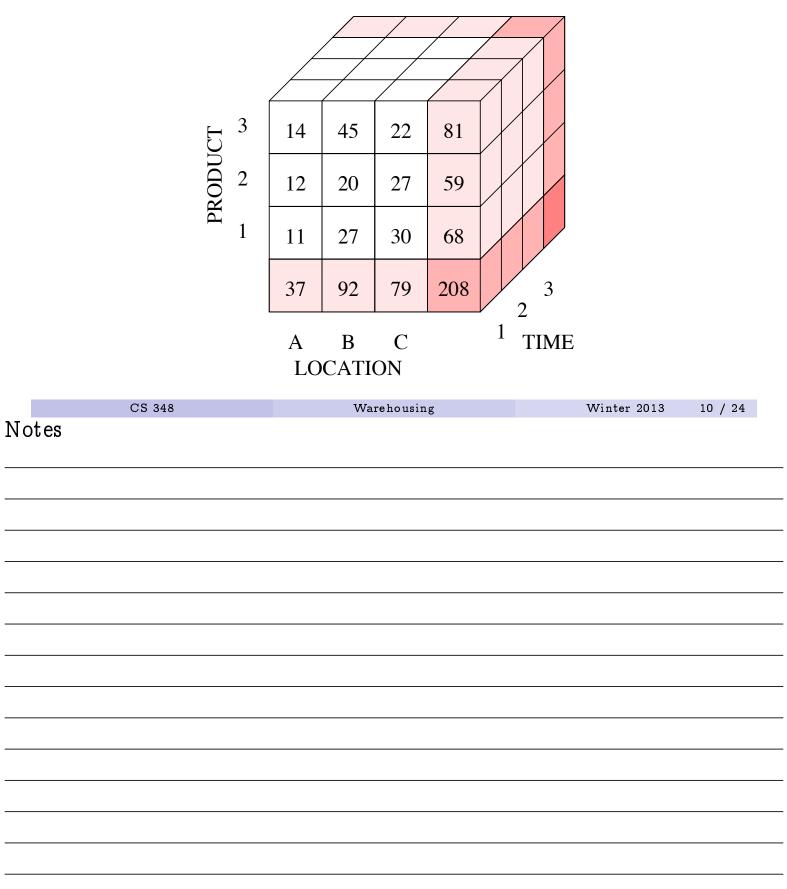


TIME

	CS 348	Warehousing	Winter 2013	9 / 24
otes				
0000				

Data Cube

• A *data cube* extends a multidimensional array of data to include all possible aggregated totals



Data Cubes as Relations

Sale	s												
lid	pid	tid	sales										
A	1	1	11										
Α	2	1	12										
A	3	1	14										
A	-	1	37										
В	1	1	27										
B	2	1	20										
В	3	1	45										
B	I	1	92										
C	1	1	30										
C	2	1	27										
C	3	1	22										
C	I	1	79										
-	1	1	68										
-	2	1	59										
-	3	1	81										
-	-	1	208										
A	1	2	16										
		•											
	Ware	housing				Wint	Winter	Winter 2	Winter 202	Winter 2013	Winter 2013	Winter 2013	Winter 2013

11 / 24

Notes

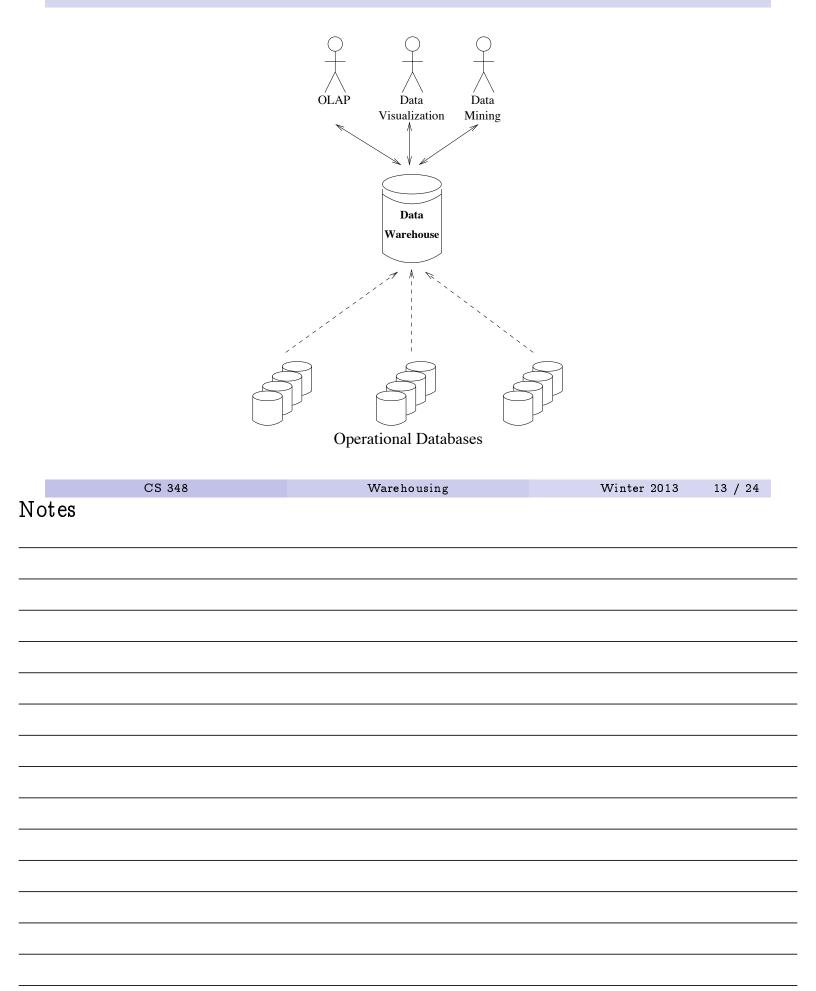
CS 348

- Generating the data cube:
 - **1** SUM(sales) GROUP BY location, product, time (raw cells)
 - 2 SUM(sales) GROUP BY location, time
 - **3** SUM(sales) GROUP BY product, time
 - 4 SUM(sales) GROUP BY product, location
 - **5** SUM(sales) GROUP BY product
 - 6 SUM(sales) GROUP BY location
 - **7** SUM(sales) GROUP BY time
 - 8 SUM(sales)
- CUBE operator in SQL:1999 groups by all combinations

SELECT lid, pid, tid, SUM(sales) FROM Sales GROUP BY CUBE(lid, pid, tid)

	CS 348	Warehousing	Winter 2013	12 / 24
Notes				

Data Warehousing



Creating and Maintaining a Warehouse

Necessary steps when creating a warehouse:

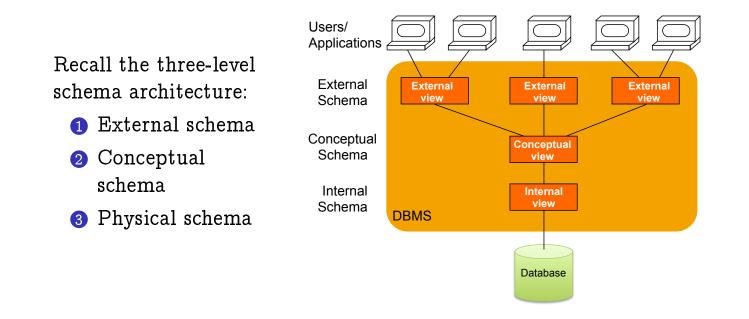
- Extract Run queries against the operational databases to retrieve necessary data
 - Clean Delete or repair tuples with missing or invalid information
- Transform Reorganize the data to fit the conceptual schema of the warehouse
 - Load Populate the warehouse tables; build indexes and/or materialized views

Note

The data in the warehouse needs to be refreshed periodically (typically nightly or weekly). To make this process efficient, the above steps need to be executed *incrementally*.

	CS 348	Warehousing	Winter 2013	14 / 24
Notes				

Views



	CS 348	Warehousing	Winter 2013	15 / 24
Notes				

Definition (View)

A view is a relation in the external schema whose instance is determined by the instances of the relations in the conceptual schema.

A view has many of the same properties as a base relation in the conceptual schema:

- its schema information appears in the database schema
- access controls can be applied to it
- other views can be defined in terms of it

	CS 348	Warehousing	Winter 2013	16 / 24
lotes				

Types of Views

- Virtual: Views are used only for querying; they are not stored in the database
- Materialized: The query that makes up the view is executed, the view constructed and stored in the database.

CS 348	Warehousing	Winter 2013	17 / 24
Notes			
110105			

SQL DDL: Views

• General form:

create [materialized] view <name>
 as <query>

• Example

create view ManufacturingProjects as

(select projno, projname, firstnme, lastname
 from project, employee
 where respemp = empno and deptno = 'D21')

	CS 348	Warehousing	Winter 2013	18 / 24
Notes				

Accessing a View

Query a view as if it were a base relation.

```
select projname
from ManufacturingProjects
```

What happens when you query a virtual view?

- At compile time, the view definition is found
- The query over the view is modified with the query definition
- The resulting query is optimized and executed

	CS 348	Warehousing	Winter 2013	19 / 24
No	tes			

Updating Views

- Modifications to a view's instance must be propagated back to instances of relations in conceptual schema.
- Some views cannot be updated unambiguously. Conceptual Schema

	Persons					Exte	ernal Schema	ι
	NAME		ENSHIP			Personal	lPastimes	
	Ed	Canad				NAME	PASTIME	
	Dave	Canad				Ed	Hockey	
	Wes	Ameri		\sim	> ∣	Ed	Curling	
	Nationa					Dave	Hockey	
	CITIZE		PASTIM	[E		Dave	Curling	
	Canadia		Hockey			Wes	Hockey	
	Canadia		Curling			Wes	Baseball	
	America		Hockey		<u> </u>			
	America	.n	Baseball					
	1 What	does it	mean to	insert (Darr	vl, Ho	ckev)?		
				delete (Dave		- ,		
					o, o d <u>i</u>	6/*		
	CS 348			Warehousing			Winter 2013	20 / 24
Notes								

According to SQL-92, a view is updatable only if its definition satisfies a variety of conditions:

- The query references exactly one table
- The query only outputs simple attributes (no expressions)
- There is no grouping/aggregation/distinct
- There are no nested queries
- There are no set operations

These rules are more restrictive than necessary.

	CS 348	Warehousing	Winter 2013	21 / 24
Notes				

Problem

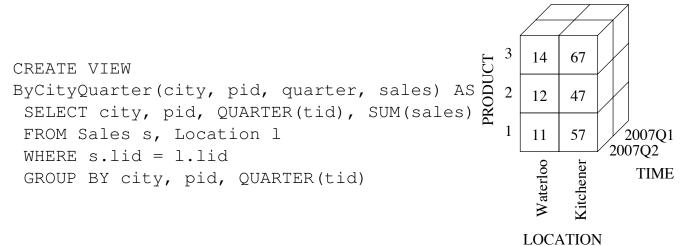
When a base table changes, the materialized view may also change.

• Solution?

- Periodically reconstruct the materialized view.
- Incrementally update the materialized view.
- Example: Data warehouses

	CS 348	Warehousing	Winter 2013	22 / 24
Notes				

• Consider the following view of the Sales data:



LOCATION

- View ByCityQuarter is useful for any query that
 - 1 Rolls-up the Location dimension to at least City; and
 - 2 Rolls-up the Time dimension to at least Quarter

	CS 348	Warehousing	Winter 2013	23 / 24
Jotes				

- Issues related to using materialized views:
 - 1 Which views to materialize (*view selection*)
 - 2 Which views are useful to answer a query (view matching)
 - **3** Which indexes to build on the views
 - 4 How to refresh the data in the view. Options:
 - Synchronous incremental maintenance
 - Asynchronous incremental maintenance
 - No synchronization (periodic re-creation)

Observation

These are the very same issues that apply to the entire data warehouse, relative to the data in the operational databases.

	CS 348	Warehousing	Winter 2013	24 / 24
Ιc	tes			