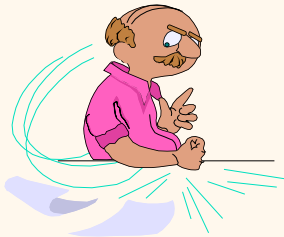
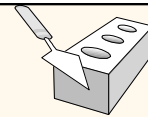


# Introduction to Data Management

## Lecture #6 (Relational Algebra)



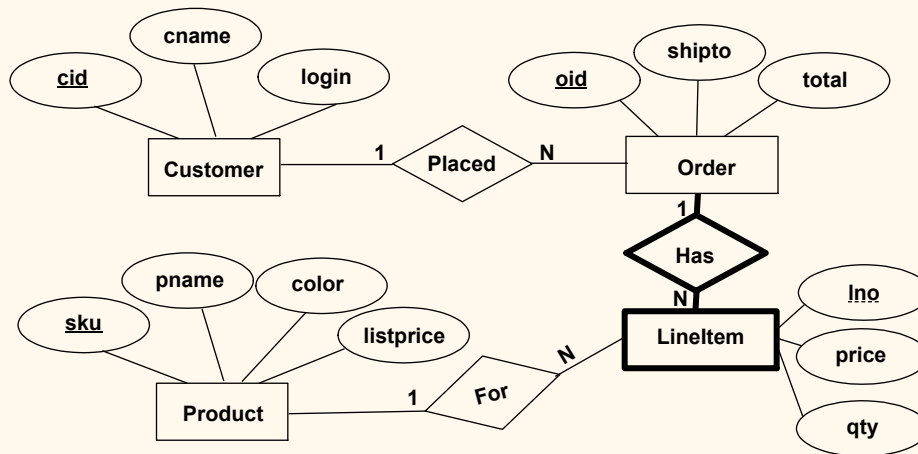
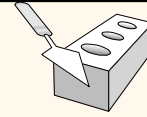
Instructor: Mike Carey  
mjcarey@ics.uci.edu



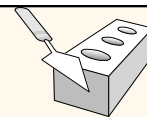
## Announcements

- ❖ HW#1 is due *right now!*
  - <https://grape.ics.uci.edu/wiki/asterix/wiki/cs122a-2014-spring#HomeworkAssignments>
  - Project Part 1 – again, please register your two-person teams by following Young-Seok's instructions (if you haven't already).
- ❖ Today's plan:
  - Quick peek at what I added to last time's PPT
  - Time to move on – next up are relational query languages!

## An Example: Putting It Together



## Putting It Together (Cont'd.)



```

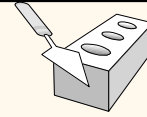
CREATE TABLE Customer (
    cid INTEGER,
    cname VARCHAR(50),
    login VARCHAR(20)
    NOT NULL,
    PRIMARY KEY (cid),
    UNIQUE (login))

CREATE TABLE Order(
    oid INTEGER,
    custid INTEGER,
    shipto VARCHAR(200),
    total DECIMAL(8,2),
    PRIMARY KEY (oid),
    FOREIGN KEY (custid) REFERENCES Customer))

CREATE TABLE Product (
    sku INTEGER,
    pname VARCHAR(100),
    color VARCHAR(20),
    listprice DECIMAL(8,2),
    PRIMARY KEY (sku))

CREATE TABLE LineItem(
    oid INTEGER,
    lno INTEGER,
    price DECIMAL(8,2),
    qty INTEGER,
    sku INTEGER,
    PRIMARY KEY (oid, lno),
    FOREIGN KEY (oid) REFERENCES Order
    ON DELETE CASCADE),
    FOREIGN KEY (sku) REFERENCES Product))
    
```

## Putting It Together (Cont'd.)



Customer

cid	cname	login
1	Smith, James	jsmith@aol.com
2	White, Susan	<u>suzie@gmail.com</u>
3	Smith, James	js@hotmail.com

Product

sku	pname	color	listprice
123	Frozen DVD	null	24.95
456	Graco Twin Stroller	green	199.99
789	Moen Kitchen Sink	black	350.00

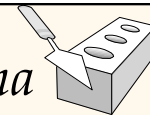
Order

oid	custid	shipto	total
1	3	J. Smith, 1 Main St., USA	199.95
2	1	Mrs. Smith, 3 State St., USA	300.00

LineItem

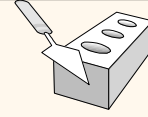
oid	lno	price	qty	item
1	1	169.95	1	456
1	2	15.00	2	123
2	1	300.00	1	789

## Relational Model and E-R Schema Translation: Summary



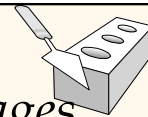
- ❖ Relations: a tabular representation of data.
- ❖ Simple and intuitive, also (very!) widely used.
- ❖ Integrity constraints can be specified by the DBA, based on application semantics. DBMS then checks for violations.
  - Important ICs: primary key, foreign key, unique.
  - In addition, we *always* have domain constraints.
- ❖ Powerful and natural query languages exist (next!)
- ❖ Rules to translate E-R schemas to relational schemas
  - Can be done by a human or automatically (using a tool)

## Relational Query Languages



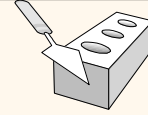
- ❖ Query languages: Allow manipulation and **retrieval of data** from a database.
- ❖ Relational model supports simple, powerful QLs:
  - Strong formal foundation based on logic.
  - Allows for much optimization.
- ❖ Query Languages **≠** programming languages!
  - QLs not expected to be “Turing complete”.
  - QLs not intended to be used for complex calculations.
  - QLs support easy, efficient access to large data sets.

## Formal Relational Query Languages




- ❖ Two mathematical Query Languages form the basis for “real” languages (e.g., SQL), and for implementation:
  - Relational Algebra: More **operational**, very useful for representing execution plans.
  - Relational Calculus: Lets users describe what they want, rather than how to compute it.  
(**Non-operational**, declarative.)

## Preliminaries



- ❖ A query is applied to *relation instances*, and the result of a query is also a relation instance.
  - *Schemas of input* relations for a query are *fixed* (but query will run regardless of instance!)
  - The *schema for the result* of a given query is also *fixed*! Determined by definition of query language constructs.
- ❖ Positional vs. named-field notation:
  - Positional notation easier for formal definitions, named-field notation more readable.
  - Both used in SQL (but try to avoid positional stuff!)

## Example Instances



<i>R1</i>	<u>sid</u>	<u>bid</u>	<u>day</u>
	22	101	10/10/96
	58	103	11/12/96

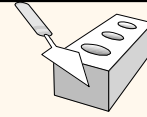
- ❖ “Sailors” and “Reserves” relations for our examples.

- ❖ We’ll use positional or named field notation, and assume that names of fields in query results are “inherited” from names of fields in query input relations.

<i>S1</i>	<u>sid</u>	sname	rating	age
	22	dustin	7	45.0
	31	lubber	8	55.5
	58	rusty	10	35.0

<i>S2</i>	<u>sid</u>	sname	rating	age
	28	yuppy	9	35.0
	31	lubber	8	55.5
	44	guppy	5	35.0
	58	rusty	10	35.0

# Relational Algebra



## ❖ Basic operations:

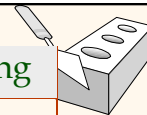
- Selection ( $\sigma$ ) Selects a subset of rows from relation.
- Projection ( $\pi$ ) Deletes unwanted columns from relation.
- Cross-product ( $\times$ ) Allows us to combine two relations.
- Set-difference ( $-$ ) Tuples in reln. 1, but not in reln. 2.
- Union ( $\cup$ ) Tuples in reln. 1 and in reln. 2.

## ❖ Additional operations:

- Intersection, join, division, renaming: Not essential, but (very!) useful. (I.e., don't add expressive power, but...)

## ❖ Since each operation returns a relation, **operations can be composed!** (Algebra is “closed”.)

# Projection



- ❖ Removes attributes that are not in *projection list*.
- ❖ *Schema* of result contains exactly the fields in the projection list, with the same names that they had in the (only) input relation.
- ❖ Relational projection operator has to eliminate *duplicates!* (Why??)
  - Note: real systems typically don't do duplicate elimination unless the user explicitly asks for it. (Q: Why not?)


sname	rating
yuppy	9
lubber	8
guppy	5
rusty	10

$\pi_{sname, rating}(S2)$

age
35.0
55.5

$\pi_{age}(S2)$

## Selection



sid	sname	rating	age
28	yuppy	9	35.0
58	rusty	10	35.0

- ❖ Selects rows that satisfy a *selection condition*.
- ❖ No duplicates in result! (Why?)
- ❖ *Schema* of result identical to schema of (only) input relation.
- ❖ *Result* relation can be the *input* for another relational algebra operation! (This is *operator composition*.)

$$\sigma_{rating > 8}(S_2)$$

sname	rating
yuppy	9
rusty	10

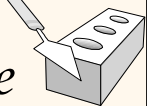
$$\pi_{sname, rating}(\sigma_{rating > 8}(S_2))$$

## Union, Intersection, Set-Difference

- ❖ All of these operations take two input relations, which must be *union-compatible*:
  - Same number of fields.
  - “Corresponding” fields are of the same type.
- ❖ What is the *schema* of result?

sid	sname	rating	age
22	dustin	7	45.0

$$S_1 - S_2$$



sid	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0
44	guppy	5	35.0
28	yuppy	9	35.0

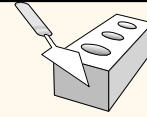
$$S_1 \cup S_2$$

sid	sname	rating	age
31	lubber	8	55.5
58	rusty	10	35.0

$$S_1 \cap S_2$$

Q: Any issues w/ duplicates?

## Cross-Product



- ❖ Each row of S1 is paired with each row of R1.
- ❖ **Result schema** has one field per field of S1 and R1, with field names 'inherited' if possible.
  - **Conflict:** Both S1 and R1 have a field called *sid*.

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	22	101	10/10/96
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	22	101	10/10/96
31	lubber	8	55.5	58	103	11/12/96
58	rusty	10	35.0	22	101	10/10/96
58	rusty	10	35.0	58	103	11/12/96

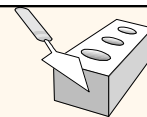
Result relation name

Attribute renaming list

Source expression (anything!)

- **Renaming operator:**  $\rho (C(1 \rightarrow sid1, 5 \rightarrow sid2), S1 \times R1)$

## Renaming



- ❖ **Conflict:** S1 and R1 both had *sid* fields, giving:

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	22	101	10/10/96
...	...	...	...	...	...	...
58	rusty	10	35.0	58	103	11/12/96

- ❖ Several renaming options available:

$\rho (S1R1(1 \rightarrow sid1), S1 \times R1)$  ← Positional renaming

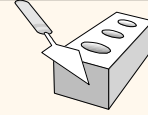
$\rho (TempS1(sid \rightarrow sid1), S1)$  ← Name-based renaming

$TempS1 \times R1$  ← Generalized projection

$(\pi_{sid \rightarrow sid1, sname, rating, age}(S1)) \times R1$



## Joins



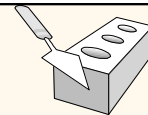
❖ Condition Join:  $R \bowtie_c S = \sigma_c(R \times S)$

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	58	103	11/12/96

$$S1 \bowtie_{sid < sid} R1$$

- ❖ *Result schema* same as that of cross-product.
- ❖ Fewer tuples than cross-product, so might be able to compute more efficiently
- ❖ Sometimes (often!) called a *theta-join*.

## More Joins



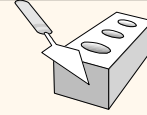
- ❖ Equi-Join: A special case of condition join where the condition  $c$  contains only *equalities*.

sid	sname	rating	age	bid	day
22	dustin	7	45.0	101	10/10/96
58	rusty	10	35.0	103	11/12/96

$$S1 \bowtie_{sid} R1$$

- ❖ *Result schema* similar to cross-product, but only one copy of fields for which equality is specified.
- ❖ Natural Join: An equijoin on *all* commonly named fields.

## Division



- ❖ Not supported as a primitive operator, but useful for expressing queries like:

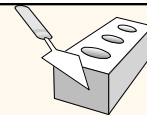
*Find sailors who have reserved all boats.*

- ❖ Let  $A$  have 2 fields,  $x$  and  $y$ ;  $B$  have only field  $y$ :

- $A/B = \{ \langle x \rangle \mid \exists \langle x, y \rangle \in A \ \forall \langle y \rangle \in B \}$
- i.e.,  **$A/B$  contains all  $x$  tuples (sailors) such that for every  $y$  tuple (boat) in  $B$ , there is an  $xy$  tuple in  $A$ .**
- Or: If the set of  $y$  values (boats) associated with an  $x$  value (sailor) in  $A$  contains all  $y$  values in  $B$ , the  $x$  value is in  $A/B$ .

- ❖ In general,  $x$  and  $y$  can be any lists of fields;  $y$  is the list of fields in  $B$ , and  $x \cup y$  is the list of fields of  $A$ .

## Examples of Division $A/B$



sno	pno
s1	p1
s1	p2
s1	p3
s1	p4
s2	p1
s2	p2
s3	p2
s4	p2
s4	p4

$A$

pno
p2

$B1$

sno
s1
s2
s3
s4

$A/B1$

pno
p2
p4

$B2$

sno
s1
s4

$A/B2$

pno
p1
p2
p4

$B3$

sno
s1

$A/B3$

## Expressing $A/B$ Using Basic Operators

(Advanced Topic ☺)

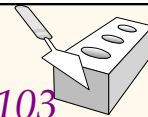


- ❖ Division not an essential op; just useful shorthand.  
(Also true of joins, but joins are so common/important that relational database systems implement joins specially.)
- ❖ *Idea*: For  $A/B$ , compute all  $x$  values that are not “disqualified” by some  $y$  value in  $B$ .
  - $x$  value is *disqualified* if by attaching a  $y$  value from  $B$ , we obtain an  $xy$  tuple that does not appear in  $A$ .

Disqualified  $x$  values:  $\pi_x((\pi_x(A) \times B) - A)$

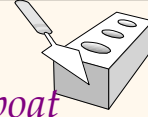
$A/B$ :  $\pi_x(A) - (\text{disqualified } x \text{ values})$

Find names of sailors who've reserved boat #103



- ❖ Solution 1:  $\pi_{sname}((\sigma_{bid=103} \text{Reserves}) \bowtie \text{Sailors})$
- ❖ Solution 2:  $\rho(\text{Temp1}, \sigma_{bid=103} \text{Reserves})$   
 $\rho(\text{Temp2}, \text{Temp1} \bowtie \text{Sailors})$   
 $\pi_{sname}(\text{Temp2})$
- ❖ Solution 3:  $\pi_{sname}(\sigma_{bid=103}(\text{Reserves} \bowtie \text{Sailors}))$

*Find names of sailors who've reserved a red boat*



- ❖ Information about boat color only available in Boats; so need an extra join:

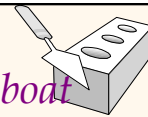
$$\pi_{sname}((\sigma_{color='red'} Boats) \bowtie Reserves \bowtie Sailors)$$

- ❖ A more “efficient” solution:

$$\pi_{sname}(\pi_{sid}((\pi_{bid} \sigma_{color='red'} Boats) \bowtie Res) \bowtie Sailors)$$

*A query optimizer can find this, given the first solution!*

*Find sailors who've reserved a red or a green boat*



- ❖ Can identify all red or green boats, then find sailors who've reserved one of these boats:

$$\rho(Tempboats, (\sigma_{color='red' \vee color='green'} Boats))$$

$$\pi_{sname}(Tempboats \bowtie Reserves \bowtie Sailors)$$

- ❖ Can also define Tempboats using union! (How?)
- ❖ What happens if  $\vee$  is replaced by  $\wedge$  in this query?

Find sailors who've reserved a red and a green boat



- ❖ Previous approach won't work! Must identify sailors who've reserved red boats, sailors who've reserved green boats, then find the intersection (note that *sid* is a key for *Sailors*):

$$\rho (Tempred, \pi_{sid}((\sigma_{color='red'} Boats) \bowtie Reserves))$$

$$\rho (Tempgreen, \pi_{sid}((\sigma_{color='green'} Boats) \bowtie Reserves))$$

$$\pi_{sname}((Tempred \cap Tempgreen) \bowtie Sailors)$$

Find the names of sailors who've reserved all boats



- ❖ Uses division; schemas of the input relations to / must be carefully chosen:

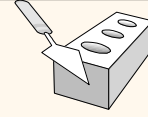
$$\rho (Tempids, (\pi_{sid,bid} Reserves) / (\pi_{bid} Boats))$$

$$\pi_{sname} (Tempids \bowtie Sailors)$$

- ❖ To find sailors who've reserved all 'Interlake' boats:

$$\dots / \pi_{bid}(\sigma_{bname='Interlake'} Boats)$$

## *Relational Algebra Summary*



- ❖ The relational model has (several) rigorously defined query languages that are both simple and powerful.
- ❖ Relational algebra is more operational; very useful as an internal representation for query evaluation plans.
- ❖ Several ways of expressing a given query; a query optimizer should choose the most efficient version. (Take CS122C? 😊)
- ❖ Next up: *Relational Calculus*