

# Relational Algebra

Introduction to Databases

Sina Meraji

Thanks to Ryan Johnson, John Mylopoulos, Arnold Rosenbloom and Renee Miller for material in these slides

## Why the relational model?

- Sounds good: matches how we think about data
- Real reason: *data independence*!
- Earlier models tied to physical data layout
  - Procedural access to data (low-level, explicit access)
  - Relationships stored in data (linked lists, trees, etc.)
  - Change in data layout => application rewrite
- Relational model
  - Declarative access to data (system optimizes for you)
  - Relationships specified by queries (schemas help, too)
  - Develop, maintain apps and data layout separately

*Similar battle today with languages*

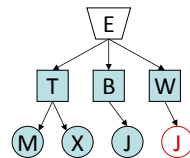
## Comparing data models

### Student job example

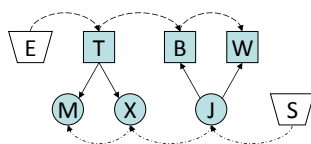
Mary (M) and Xiao (X) both work at Tim Hortons (T)

Jaspreet (J) works at both Bookstore (B) and Wind (W)

### Hierarchical (tree)



### Network (graph)



### Relational (table)

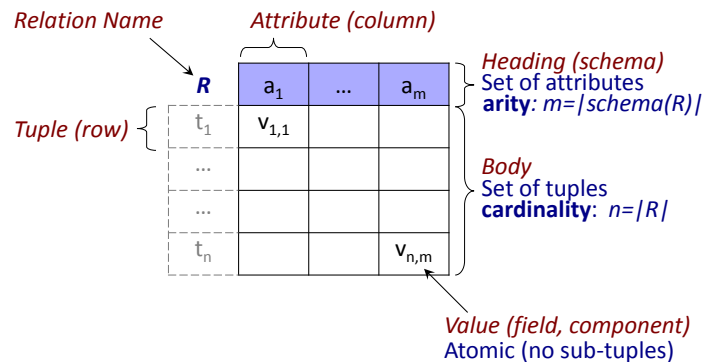
E	S	R
B ...	J ...	M T
T ...	M ...	X T
W ...	X ...	J B
		J W

## What is the relational model?

- Logical representation of data
  - Two-dimensional tables (relations)
- Formal system for manipulating relations
  - Relational algebra
- Result
  - High-level (logical, declarative) description of data
  - Mechanical rules for rewriting/optimizing low-level access
  - Formal methods to reason about soundness

*Relational algebra is the key*

## Relations and tuples



*Set-based: arbitrary row/col ordering*

*Logical: physical layout might be \*very\* different!*

## What is an algebra?

- **Operands (values)**
  - Variables, constants
  - Closed domain
- **Operators**
  - + “Addition”
  - \* “Multiplication”
- **Expressions:**
  - Combine operations with parenthesis (explicit)
  - OR using either precedence (implied)
- **Laws**
  - Identify semantically equivalent expressions
  - Commutativity, associativity, etc.

*Offers formal, sound, mechanical rewriting*

## Example algebra: integer arithmetic

- **Domain: integers**  
... -100, ... -1, 0, 1, ... 100, ...
- **Operators: - , + , \* , ...**  
*'-' is unary negation*
- **Expressions**
  - $((2*a) + ((5*(c + (-d))) + e))$
- **Laws**
  - $a*b = b*a$  *Commutative*
  - $a*(b*c) = (a*b)*c$  *Associative*
  - $a*(b+c) = a*b + a*c$  *Distributive*

*Allows compilers to reason about, optimize*

## Relational algebra

- **Values**
  - Finite relations (cardinality and arity both bounded)
  - Attributes may or may not be typed
- **Operators**
  - Unary:  $\sigma, \pi, \rho$
  - “Additive” (set):  $\cup, \cap, -$
  - “Multiplicative:”  $\times, \bowtie$
  - [details to come]
- **Expressions**
  - Same as arithmetic, but called “queries”
- **Laws**
  - Allow “query rewriting”
  - Basis for query optimization
  - [details to come]

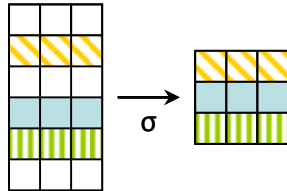
*Expressive power equivalent to 1<sup>st</sup> order logic*



10

## Unary operators: select ( $\sigma$ )

- $\sigma_P(R)$  outputs tuples of R which satisfy P
- same schema as R



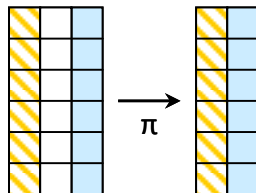
*Removes unwanted rows from relation*



12

## Unary operators: project ( $\pi$ )

- $\pi_Y(R)$  outputs a subset Y of the set of attributes X of relation R



*Removes unwanted columns from relation*



11

## Unary operators: select ( $\sigma$ ) example

**Employees**

Surname	FirstName	Age	Salary
Smith	Mary	25	2000
Black	Lucy	40	3000
Verdi	Nico	36	4500
Smith	Mark	40	3900

**$\sigma_{\text{Age} < 30 \vee \text{Salary} > 4000}$  (Employees)**

Surname	FirstName	Age	Salary
Smith	Mary	25	2000
Verdi	Nico	36	4500



13

## Unary operators: project ( $\pi$ ) example

**Employees**

Surname	FirstName	Department	Head
Smith	Mary	Sales	De Rossi
Black	Lucy	Sales	De Rossi
Verdi	Mary	Personnel	Fox
Smith	Mark	Personnel	Fox

**$\pi_{\text{Surname, FirstName}}$  (Employees)**

Surname	FirstName
Smith	Mary
Black	Lucy
Verdi	Mary
Smith	Mark

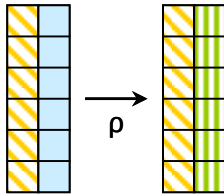
**$\pi_{\text{Department, Head}}$  (Employees)**

Department	Head
Sales	De Rossi
Personnel	Fox



## Unary operators: rename ( $\rho$ )

- $\rho_{S(A,B,C)}(R)$  renames attributes of R to A,B,C and calls the result S
  - $\rho_S(R)$  renames relation R (same attributes)
  - $\rho_{A=X, C=Y}(R)$  (or  $\rho_{A, C \rightarrow X, Y}(R)$ ) renames attributes A and C only

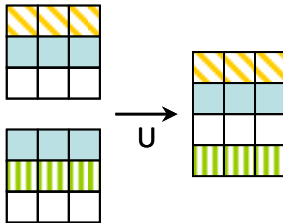


*Modifies schema only - same values*



## Additive operators ( $\cup$ , $\cap$ , $-$ )

- Standard set operators
- Operate on tuples within input relations, but not on schema



## Unary operators: rename ( $\rho$ ) example

**Paternity**

Father	Child
Adam	Cain
Adam	Abel
Abraham	Isaac
Abraham	Ishmael

**$\rho_{\text{Father} \rightarrow \text{Parent}}(\text{Paternity})$**

Parent	Child
Adam	Cain
Adam	Abel
Abraham	Isaac
Abraham	Ishmael



## Additive operators: Union ( $\cup$ )

**Graduates**

Number	Surname	Age
7274	Robinson	37
7432	O'Malley	39
9824	Darkes	38

**Managers**

Number	Surname	Age
9297	O'Malley	56
7432	O'Malley	39
9824	Darkes	38

**$\text{Graduates} \cup \text{Managers}$**

Number	Surname	Age
7274	Robinson	37
7432	O'Malley	39
9824	Darkes	38
9297	O'Malley	56

## Additive operators: Intersection ( $\cap$ )

Graduates

Number	Surname	Age
7274	Robinson	37
7432	O'Malley	39
9824	Darkes	38

Managers

Number	Surname	Age
9297	O'Malley	56
7432	O'Malley	39
9824	Darkes	38

Graduates  $\cap$  Managers

Number	Surname	Age
7432	O'Malley	39
9824	Darkes	38

## Additive operators: Difference (-)

Graduates

Number	Surname	Age
7274	Robinson	37
7432	O'Malley	39
9824	Darkes	38

Managers

Number	Surname	Age
9297	O'Malley	56
7432	O'Malley	39
9824	Darkes	38

Graduates - Managers

Number	Surname	Age
7274	Robinson	37

## A Meaningful but Impossible Union

Paternity

Father	Child
Adam	Cain
Adam	Abel
Abraham	Isaac
Abraham	Ishmael

Maternity

Mother	Child
Eve	Cain
Eve	Seth
Sarah	Isaac
Hagar	Ishmael

Paternity  $\cup$  Maternity ???

- The problem: **Father** and **Mother** are different names, but both represent a parent
- Solution: rename attributes!

## Union with Renaming

Paternity

Father	Child
Adam	Cain
Adam	Abel
Abraham	Isaac
Abraham	Ishmael

Maternity

Mother	Child
Eve	Cain
Eve	Seth
Sarah	Isaac
Hagar	Ishmael

$\rho_{\text{Father} \rightarrow \text{Parent}}(\text{Paternity}) \cup \rho_{\text{Mother} \rightarrow \text{Parent}}(\text{Maternity})$

Parent	Child
Adam	Cain
Adam	Abel
Abraham	Isaac
Abraham	Ishmael
Eve	Cain
Eve	Seth
Sarah	Isaac
Hagar	Ishmael

## Union with Renaming (Many Attributes)

**Employees**

Surname	Branch	Salary
Patterson	Rome	45
Trumble	London	53

**Staff**

Surname	Factory	Wages
Cooke	Chicago	33
Bush	Monza	32

$\rho_{\text{Branch, Salary} \rightarrow \text{Location, Pay}}(\text{Employees}) \cup \rho_{\text{Factory, Wages} \rightarrow \text{Location, Pay}}(\text{Staff})$

Surname	Location	Pay
Patterson	Rome	45
Trumble	London	53
Cooke	Chicago	33
Bush	Monza	32

## Cartesian product (x) example

**Employees**

Employee	Project
Smith	A
Black	A
Black	B

**Projects**

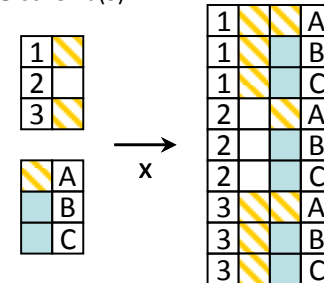
Code	Name
A	Venus
B	Mars

**Employees x Projects**

Employee	Project	Code	Name
Smith	A	A	Venus
Black	A	A	Venus
Black	B	A	Venus
Smith	A	B	Mars
Black	A	B	Mars
Black	B	B	Mars

## Cartesian product (x)

- The outcome of combining every record in R with every record in S
- $T = R \times S$  contains every pairwise combination of R and S tuples
  - schema(T) = schema(R) U schema(S)
  - $|T| = |R| * |S|$

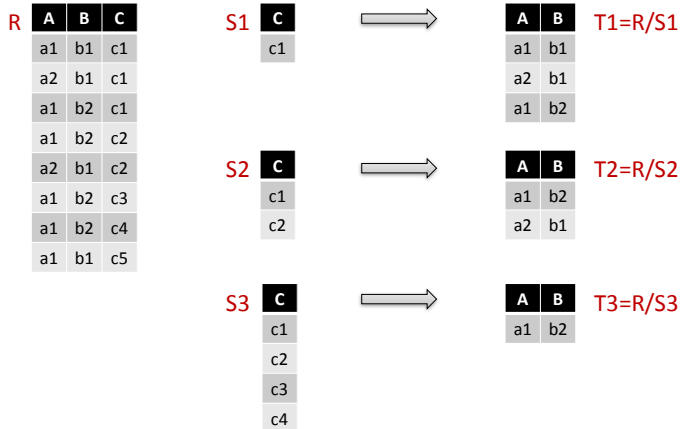


*Input schemas must \*not\* overlap*

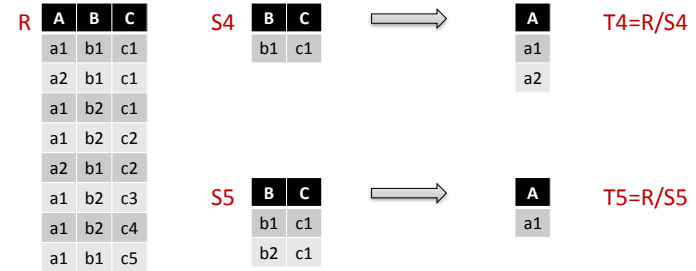
## Division (/)

- Let **R** and **S** be relations with schemas **A1, ..., An**, **B1, ..., Bn** and **B1, ..., Bn** respectively. The result of **R/S** is a relation **T** with
  - Schema **A1, ..., An** (attribute names in **R** but not in **S**)
  - Tuples **t** such that, for every tuple **s** of **S**, the tuple **t || s** (the concatenation of **t** and **s**) is in relation **R**
  - T** contains the largest possible set of tuples s. t.  $S \times T \subseteq R$
- Analogy to integer division:
  - For integers,  $A / B$  is: the largest int  $Q$  s.t.  $Q \times B \leq A$
  - For relations,  $A / B$  is: the largest relation  $Q$  s.t.  $Q \times B \subseteq A$

## Division example



## Division example (cont.)



## Division in RA

- Consider two relations A(x,y), B(y) and suppose we want to specify the query

"Find all x's that are associated through A with all B's"

- This can be expressed as:

$$A/B = \pi_x(A) - \pi_x((\pi_x(A) \times B) - A)$$

- Often useful when the query is about "every" or "all" (but don't just look for these keywords!)
- Doesn't extend the expressiveness of Relational Algebra (convenient to use in many situations)

## Division in RA example

- Assume
  - Take(x,y) - "student x has taken course y",
  - CS(z) - "z is a CS course"
- We want "All students who have taken all CS courses"
  - $\pi_x(\text{Take}) \times \text{CS}$   
(Relation of all <student, CS course> pairs)
  - $(\pi_x(\text{Take}) \times \text{CS}) - \rho_{y \rightarrow z}(\text{Take})$   
(Relation of all <student, CS course> pairs that did NOT occur)
  - $\pi_x((\pi_x(\text{Take}) \times \text{CS}) - \rho_{y \rightarrow z}(\text{Take}))$   
(Relation of all <students> who have NOT taken all CS courses)
  - $\pi_x(\text{Take}) - \pi_x((\pi_x(\text{Take}) \times \text{CS}) - \rho_{y \rightarrow z}(\text{Take}))$   
(Relation of all <students> who have taken all CS courses)

*Work a simple example at home*

## Division Example

Take	CS		R1	R2	R3	R4
x y	z		x z	x z	x	x
S1 C1	C1		S1 C1	S2 C2	S2	S1
S1 C2	C2		S1 C2	S2 C3	S3	S4
S1 C3	C3		S1 C3	S3 C2	S5	
S2 C1			S2 C1	S5 C1		
S3 C1			S2 C2	S5 C3		
S3 C3			S2 C3			
S4 C1			S3 C1			
S4 C2			S3 C2			
S4 C3			S3 C3			
S5 C2			S4 C1			
			S4 C2			
			S4 C3			
			S5 C1			
			S5 C2			
			S5 C3			

$\pi_x(\text{Take}) - \pi_x((\pi_x(\text{Take}) \times \text{CS}) \div \rho_{y \rightarrow z}(\text{Take}))$

## Natural join ( $\bowtie$ )

- T = R  $\bowtie$  S merges tuples from R and S having equal values where their schemas overlap (**join attributes**)
  - T Schema: Union of schemas  $\text{schema}(R) \cap \text{schema}(S) \neq \emptyset$
  - $|T| \leq |R| * |S|$ , usually  $\approx \max(|R|, |S|)$  "join cardinality"
- Special cases
  - No schema overlap:  $\times$
  - Full schema overlap:  $\cap$

1		
2		
3		
4		
5		
6		

A
B
C

 $\bowtie$ 

1	A
3	A
4	B
4	C

Equivalent to  $\pi(\sigma(R \times S))$

## Join

- The most used operator in relational algebra
- Used to **establish connections** among data in different relations, taking advantage of the "value-based" nature of the relational model
- Two main versions of the join:
  - natural join: takes **attribute names** into account
  - theta join: takes **attribute values** into account
- Both join operations denoted by the symbol  $\bowtie$

## Natural join ( $\bowtie$ ) example

Employee	Department
Smith	sales
Black	production
White	production

Department	Head
production	Mori
sales	Brown

$r_1 \bowtie r_2$

Employee	Department	Head
Smith	sales	Brown
Black	production	Mori
White	production	Mori





34

## Properties of Natural Join

- Commutative:

$$R \bowtie S = S \bowtie R$$

- Associative:

$$(R \bowtie S) \bowtie T = R \bowtie (S \bowtie T)$$

- N-ary joins without ambiguity:

$$R_1 \bowtie R_2 \bowtie \dots \bowtie R_n$$



36

## Theta join

- Written as  $T = R \bowtie_{\theta} S$

- Outputs pairwise combinations of tuples which satisfy  $\theta$
- Join cardinality:  $|T| \leq |R| * |S|$

- Most general join

- Arbitrary join predicate (not just equality)

- Equivalent to  $\sigma_{\theta}(R \times S)$

- Schemas must not overlap
- Does not project away any attributes



35

## Example of N-ary Join Operation

$r_1$

Employee	Department
Smith	sales
Black	production
Brown	marketing
White	production

$r_2$

Department	Division
production	A
marketing	B
purchasing	B

$r_3$

Division	Head
A	Mori
B	Brown

$r_1 \bowtie r_2 \bowtie r_3$

Employee	Department	Division	Head
Black	production	A	Mori
Brown	marketing	B	Brown
White	production	A	Mori



37

## Theta join example

**Car**

Car	CarPrice
CarA	20000
CarB	30000
CarC	50000

**Boat**

Boat	BoatPrice
BoatA	10000
BoatB	40000
BoatC	60000

**Car**  $\bowtie_{\text{CarPrice} > \text{BoatPrice}}$  **Boat**

Car	CarPrice	Boat	BoatPrice
CarA	20000	BoatA	10000
CarB	30000	BoatA	10000
CarC	50000	BoatA	10000
CarC	50000	BoatB	40000



38

## Equijoin

- Special case of theta join
- Written as  $R \bowtie_{A=X, B=Y, \dots} S$ 
  - Attribute names in R and S can differ
  - Still compare values for equality
- Like natural join, but using arbitrary attributes
  - Very common due to *foreign keys* in relations
- Equivalent to  $R \bowtie \rho(S)$



39

## Equijoin example

**Employees**

Employee	Project
Smith	A
Black	A
Black	B

**Projects**

Code	Name
A	Venus
B	Mars

**Employees  $\bowtie_{\text{Project=Code}}$  Projects**

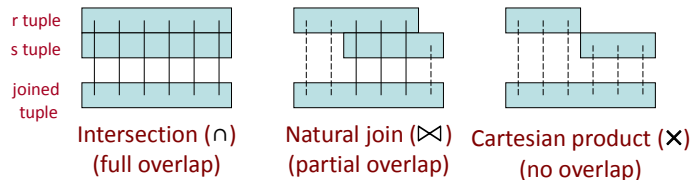
Employee	Project	Code	Name
Smith	A	A	Venus
Black	A	A	Venus
Black	B	B	Mars



40

## Comparison: $\times$ vs. $\cap$ vs. $\bowtie$

- Same general operation
  - Test “overlapping” parts of tuples for equality
  - Combine “matching” pairs (ignore others)
- Differ in degree of schema overlap



*“Generalized intersection”*



41

## Mathematical power vs. efficiency

- Note that  $\times$  expresses both  $\cap$  and  $\bowtie$ 
  - => Mathematically, intersection and joins are unnecessary
- Why bother with them? Two big reasons
- Notation
  - $\pi(\sigma(R \times \rho(S)))$  vs.  $R \cap S$
  - Cartesian product seldom useful *Why not?*
- **Performance**
  - Efficient algorithms compute result directly
  - =>  $|R| * |S|$  rows vs.  $\min(|R|, |S|)$  *Consider  $|R|=|S|=10^6$*

## Summary of Operators

Operation	Name	Symbol	Precedence
choose rows	select	$\sigma$	1
choose columns	project	$\pi$	
rename relation/attribute	rename	$\rho$	
combine tables	natural join	$\bowtie$	2
	theta join	$\bowtie_{\text{condition}}$	
	cartesian product	$\times$	
set operations	intersection	$\cap$	3
	union	$\cup$	
	subtraction	$-$	
assignment	assignment	$:=$	-

## Expressing Integrity Constraints

- Our text (sec 2.5) defines two ways to express an integrity constraint in Relational Algebra. Suppose R and S are expressions in RA. We can write an IC in either of these ways:

$R = \emptyset$  (expresses the fact that R is the empty set.)

$R \subseteq S$  (expresses the fact that R is a subset of S.)

- Equivalent (we don't need the second form) but it's convenient:
  - Saying  $R = \emptyset$  is equivalent to saying  $R \subseteq \emptyset$ .
  - Saying  $R \subseteq S$  is equivalent to saying  $R - S = \emptyset$ .

## Expressing ICs Examples

Course (Dept, CourseNum, Title, Credits)

Section (CRN, Dept, CourseNum, Room, Time, InstructorID)

- Referential integrity constraints

Dept and CourseNum form a foreign key within Section

$\pi_{\text{Dept, CourseNum}}(\text{Section}) \subseteq \pi_{\text{Dept, CourseNum}}(\text{Course})$

- Key constraints

Two tuples which agree on CRN must also agree on Dept, CourseNum, Room, Time, and InstructorID. One of the constraints implied is:

$(\rho_{S1}(\text{Section}) \bowtie_{(S1.CRN=S2.CRN \text{ and } S1.Dept=S2.Dept)} \rho_{S2}(\text{Section})) = \emptyset$

- Domain constraints

E.g., "Course Numbers must be in the range 100-999"

$\sigma_{\text{CourseNum} < 100 \text{ or } \text{CourseNum} > 999}(\text{Course}) = \emptyset$

## Coming next...

- Working Examples
- Tips and Tricks
- You
  - Attend the lectures
  - Read the material and practice
  - Assignment 1



46

## TIPS & TRICKS FOR RELATIONAL ALGEBRA QUERIES

### Tips and Tricks for R.A.

- Ask yourself which relations need to be involved?
  - Ignore the rest
- Every time you combine relations confirm that:
  - (a) attributes that should match will be made to match
  - (b) attributes that will be made to match should match



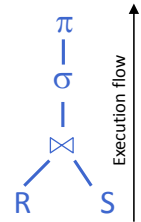
48



47

## Evaluating R.A. Queries

- R.A. is **procedural**
  - an R.A. query itself suggests a procedure for constructing the result (i.e., how to implement the query)
- R.A. suggests a **query execution plan**
  - Two RA expressions might yield the same result but suggest different query execution plans
  - which is best depends on the relation cardinality, defined indices, join ordering, etc.
- DBMSs: **query optimization** takes place
  - Optimizer rewrites queries to be more efficient
  - topic in csc443



### Tips and Tricks for R.A. (cont.)

- Is there an intermediate relation that would help you get the final answer?
  - Draw it out with actual data in it
- Break the answer down by defining intermediate relations using assignment:
  - Use good names for the new relations
  - Name the attributes on the Left-Hand-Side each time, so you don't forget what you have in hand
  - Add a comment that explains exactly what the relation contains



49

## Tips for Specific R.A. Queries

- To show “**max**” (min is analogous):
  - Pair tuples (self-join) and find those that are not the max
  - Then subtract from all to find the max[es]
- To show “**k or more**”:
  - Make all combos of k different tuples that meet the required condition
- To show “**exactly k**”:
  - Find “k or more”
  - Find “(k+1) or more”
  - Then subtract “(k+1) or more” from “k or more”
- To show “**every**” (i.e., division):
  - Make all combos that could have occurred
  - Subtract those that did occur to find those that didn’t always; these are the failures
  - Subtract the failures from all to get the answer



## WORKING EXAMPLES (SET-SEMANTICS)

## A Sample Database

**Employees**

Number	Name	Age	Salary
101	Mary Smith	34	40
103	Mary Bianchi	23	35
104	Luigi Neri	38	61
105	Nico Bini	44	38
210	Marco Celli	49	60
231	Siro Bisi	50	60
252	Nico Bini	44	70
301	Steve Smith	34	70
375	Mary Smith	50	65

**Supervision**

Head	Emp
210	101
210	103
210	104
231	105
301	210
301	231
375	252

In fact, only the **database schema** and **integrity constraints** are needed:

**Employees**(Number, Name, Age, Salary)

**Supervision**(Head, Emp)

$\pi_{\text{Head}}(\text{Supervision}) \subseteq \pi_{\text{Number}}(\text{Employees})$

$\pi_{\text{Emp}}(\text{Supervision}) \subseteq \pi_{\text{Number}}(\text{Employees})$

## Example 1

“Find the numbers, names and ages of employees earning more than 40K”

Employees(Number, Name, Age, Salary)

Supervision(Head, Emp)

$\pi_{\text{Number, Name, Age}}(\sigma_{\text{Salary} > 40} \text{Employees})$

Number	Name	Age
104	Luigi Neri	38
210	Marco Celli	49
231	Siro Bisi	50
252	Nico Bini	44
301	Steve Smith	34
375	Mary Smith	50

## Example 2

“Find the registration numbers of the supervisors of the employees earning more than 40K.”

Employees(Number,Name,Age,Salary)  
Supervision(Head,Emp)

$\pi_{\text{Head}}(\text{Supervision} \bowtie_{\text{Emp=Number}} (\sigma_{\text{Salary}>40} \text{Employees}))$

Head
210
301
375

R(Head,Emp,Number,Name,Age,Salary)

## Example 3

“Find the names and salaries of the supervisors of the employees earning more than 40K.”

Employees(Number,Name,Age,Salary)  
Supervision(Head,Emp)

$\pi_{\text{NameH,SalH}}(\rho_{\text{Number,Name,Salary,Age} \rightarrow \text{NumH,NameH,SalH,AgeH}} \text{Employees})$

$\bowtie_{\text{NumberH=Head}} (\text{Supervision} \bowtie_{\text{Number=Emp}} (\sigma_{\text{Salary}>40} \text{Employees})))$

NameH	SalaryH
Marco Celli	60
Steve Smith	70
Mary Smith	65

R(NumH,NameH,AgeH,SalH,Head,Emp,Number,Name,Age,Salary)

## Example 4

“Find the employees earning more than their respective supervisors; return registration numbers, names and salaries of the employees and their supervisors.”

Employees(Number,Name,Age,Salary)  
Supervision(Head,Emp)

$\pi_{\text{Number,Name,Salary,NumH,NameH,SalH}}(\sigma_{\text{Salary}>\text{SalH}}(\rho_{\text{Number,Name,Salary,Age} \rightarrow \text{NumH,NameH,SalH,AgeH}} \text{Employees} \bowtie_{\text{NumH=Head}} (\text{Supervision} \bowtie_{\text{Emp=Number}} \text{Employees})))$

Number	Name	Salary	NumH	NameH	SalH
104	Luigi Neri	61	210	Marco Celli	60
252	Nico Bini	70	375	Mary Smith	65

R(NumH,NameH,AgeH,SalH,Head,Emp,Number,Name,Age,Salary)

## Example 5

“Find registration numbers and names of supervisors, *all* of whose employees earn more than 40K.”

Employees(Number,Name,Age,Salary)  
Supervision(Head,Emp)

$\pi_{\text{Number,Name}}(\text{Employees} \bowtie_{\text{Number=Head}} (\pi_{\text{Head}}(\text{Supervision}) - \pi_{\text{Head}}(\text{Supervision} \bowtie_{\text{Number=Emp}} (\sigma_{\text{Salary}>40} \text{Employees}))))$

Number	Name
301	Steve Smith
375	Mary Smith



58

## Example 6

"Find registration numbers of supervisors, who supervise *all* employees earning more than 40K."

Employees(Number, Name, Age, Salary)  
Supervision(Head, Emp)

$\pi_{\text{Number}}(\text{Supervision} / \pi_{\text{Number}}(\sigma_{\text{Salary} > 40}(\text{Employees})))$



59

## Example 7

"Find the employees earning **maximum** salary."

Employees(Number, Name, Age, Salary)  
Supervision(Head, Emp)

$A = \pi_{\text{Salary}}(\text{Emp}) - \pi_{\text{Salary}}(\sigma_{\text{Salary1} > \text{Salary}(\text{Emp} \bowtie \rho_{\text{everything}}(\text{Emp}))})$

A provides the maximum salary; the rest is easy ...

$\pi_{\text{Number}}(\text{Emp} \bowtie A)$

- 💡 How to find **minimum** salary?
- 💡 How to find **average** salary?



60

## Example 8

"Find all locations that have **at least two** employees earning more than 40K"

Employees(Number, Name, Loc, Salary)  
Supervision(Head, Emp)

Assume that in the schema we have 'Loc' instead of 'Age' for this example only.

$\sigma_{\text{Salary} > 40}(\text{Emp}) \bowtie \sigma_{\text{Salary1} > 40}(\rho_{\text{Number, Name, Salary} \rightarrow \text{Num1, Name1, Sal1}}(\text{Emp}))$

This results in a relation

$R(\text{Number}, \text{Name}, \text{Salary}, \text{Loc}, \text{Num1}, \text{Name1}, \text{Sal1})$

Now we select tuples where  $\text{Number} \neq \text{Num1}$  and project on Loc

$\pi_{\text{Loc}}(\sigma_{\text{Number} \neq \text{Num1}}(R))$

- 💡 How to find locations that have **at least three** employees ...?
- 💡 How to find locations that have **exactly two** employees ...?



61

## Another Series of Examples

Films(Film#, Title, Director, Year, ProdCost)

Artists(Actor#, Surname, FirstName, Sex, Birthday, Nationality)

Roles(Film#, Actor#, Character)



## Example 1

**Films**(Film#, Title, Director, Year, ProdCost)

**Artists**(Actor#, Surname, FirstName, Sex, Birthday, Nationality)

**Roles**(Film#, Actor#, Character)

“Find the titles of films starring Henry Fonda”

$\pi_{\text{Title}}(\text{Films} \bowtie (\sigma_{(\text{FirstName}=\text{"Henry"}) \wedge (\text{Surname}=\text{"Fonda"})}(\text{Artists} \bowtie \text{Roles})))$



## Example 3

**Films**(Film#, Title, Director, Year, ProdCost)

**Artists**(Actor#, Surname, FirstName, Sex, Birthday, Nationality)

**Roles**(Film#, Actor#, Character)

“Find the actors who have played two characters in the same film; show the title of each such film, first name and surname of the actor and the two characters”

$\pi_{\text{Title}, \text{FirstName}, \text{Surname}, \text{Character1}, \text{Character2}}(\rho_{\text{Film\#}, \text{Actor\#}, \text{Character} \rightarrow \text{Film\#1}, \text{Actor\#1}, \text{Character1}}(\text{Roles}) \bowtie_{(\text{Film\#1}=\text{Film\#}) \wedge (\text{Actor\#1}=\text{Actor\#}) \wedge (\text{Character1} \neq \text{Character2})} \text{Roles} \bowtie \text{Artists} \bowtie \text{Films})$



## Example 2

**Films**(Film#, Title, Director, Year, ProdCost)

**Artists**(Actor#, Surname, FirstName, Sex, Birthday, Nationality)

**Roles**(Film#, Actor#, Character)

“Find the titles of all films in which the director is also an actor”

$\pi_{\text{Title}}(\sigma_{(\text{Director}=\text{Actor\#})}(\text{Films} \bowtie \text{Roles}))$



## Example 4

**Films**(Film#, Title, Director, Year, ProdCost)

**Artists**(Actor#, Surname, FirstName, Sex, Birthday, Nationality)

**Roles**(Film#, Actor#, Character)

“Find the titles of the films in which the actors are all of the same sex”

$\pi_{\text{Title}}(\text{Films}) - \pi_{\text{Title}}(\text{Films} \bowtie \sigma_{\text{sex} \neq \text{sex1} \wedge \text{Film\#1} = \text{Film\#}}((\text{Artists} \bowtie \text{Roles}) \bowtie \rho_{\text{Ac\#, Char, Sur, Fir, Sex, BD, N} \rightarrow \text{Ac\#1, Char1, Sur1, Fir1, Sex1, BD1, N1}}(\text{Artists} \bowtie \text{Roles})))$





66

## RELATION OPERATIONS ON BAGS



67

## Limitations of relational algebra

- Relational algebra is set-based
- Real-life applications need more
  - Expensive (and often unnecessary) to eliminate **duplicates**
  - Important (and often expensive) to **order** output
  - Need a way to apply **scalar expressions** to values
  - What's **\*not\*** there often as important as what is

*Answer: non-set extensions*



68

## Extension: bag semantics

- In practice, relations are bags (multisets)
  - Members are allowed to appear more than once
  - Sometimes people purposefully insert duplicates
  - Projections produce duplicates
- Example: **{1,2,1,1,3}** is a **bag** (still unordered!)
- Most operators still work
  - **Select**, **Rename** remain unchanged
  - **Project** no longer eliminates duplicates
  - **Set operations** need tweaks
  - **Joins** tend to multiply the number of duplicates
- Some laws no longer apply



69

## Bag versions of set operations

- Union
  - Concatenation (except unordered)
  - $\{1, 1, 2, 3\} \cup \{2, 2, 3, 4\} = \{1, 1, 2, 3, 2, 2, 3, 4\}$
- Intersection
  - Take minimum count of each value
  - $\{1, 1, 2, 3\} \cap \{2, 2, 3, 4\} = \{2, 3\}$
- Difference
  - Each occurrence on right can cancel one occurrence on left
  - $\{1, 1, 2, 3\} - \{1, 2, 3, 4\} = \{1\}$
- Union, intersection no longer distribute
 

$\{1\} \cap (\{1\} \cup \{1\})$	vs.	$(\{1\} \cap \{1\}) \cup (\{1\} \cap \{1\})$
$\{1\} \cap \{1, 1\}$	vs.	$\{1\} \cup \{1\}$
$\{1\}$	$\neq$	$\{1, 1\}$

## Bag-projection ( $\pi$ ) and duplicates

Make	Model	Color
Toyota	Prius	Gray
Toyota	Prius	Red
Honda	Accord	Green
Honda	Accord	Red
Honda	Accord	Red
Ford	Echo	Red
Ford	Echo	Gray
Ford	Echo	White

- Consider a relation R modeling cars for sale
- Bag-projection ( $\pi$ ) does **not** eliminate duplicate tuples (as in set-projection)
  - What does  $\pi_{\text{Make}}(R)$  return?
  - How to eliminate duplicates?

## Duplicate elimination ( $\delta$ )

Make	Model	Color
Toyota	Prius	Gray
Toyota	Prius	Red
Honda	Accord	Green
Honda	Accord	Red
Honda	Accord	Red
Ford	Echo	Red
Ford	Echo	Gray
Ford	Echo	White

- Consider a relation R modeling cars for sale
- $\delta$  turns a **bag** into a **set**
- $\delta(\pi_{\text{Make}}(R))$  is a **set**

### Make

Honda
Ford
Toyota

*Duplicates important for summaries ("how many")*

## Summarizing groups of tuples (1)

Student	Year	Dept	Course	Grade
Xiao	2009	CS	A08	B-
Xiao	2009	CS	A48	B
Xiao	2009	CS	A65	B+
Xiao	2009	Math	A23	B
Xiao	2009	Math	A30	B+
Xiao	2009	Math	A37	A
Xiao	2010	CS	B07	B
Xiao	2010	CS	B09	B-
Xiao	2010	CS	B36	B-
Xiao	2010	CS	B58	B
Xiao	2010	Math	B24	A-
Xiao	2010	Math	B41	B
Xiao	2010	Stats	B52	B-
Xiao	2011	CS	C24	B+
Xiao	2011	CS	C43	A-
Xiao	2011	CS	C69	A

All courses Xiao has taken

All courses Xiao took in 2010

All math courses Xiao took in 2010

## Summarizing groups of tuples (2)

Student	Year	Dept	Course	Grade
Xiao	2009	CS	A08	B-
Xiao	2009	CS	A48	B
Xiao	2009	CS	A65	B+
Xiao	2009	Math	A23	B
Xiao	2009	Math	A30	B+
Xiao	2009	Math	A37	A
Xiao	2010	CS	B07	B
Xiao	2010	CS	B09	B-
Xiao	2010	CS	B36	B-
Xiao	2010	CS	B58	B
Xiao	2010	Math	B24	A-
Xiao	2010	Math	B41	B
Xiao	2010	Stats	B52	B-
Xiao	2011	CS	C24	B+
Xiao	2011	CS	C43	A-
Xiao	2011	CS	C69	A

How to summarize this??

Student	Dept	Year	Course	Grade
Xiao	CS	2009	?	?
Xiao	Math	2009	?	?
Xiao	CS	2010	?	?
Xiao	Math	2010	?	?
Xiao	Stats	2010	B52	B-
Xiao	CS	2011	?	?

These columns are easy... equal for every tuple in a group

Show the best grade? Worst grade? Average?

## Summarizing groups of tuples (3)

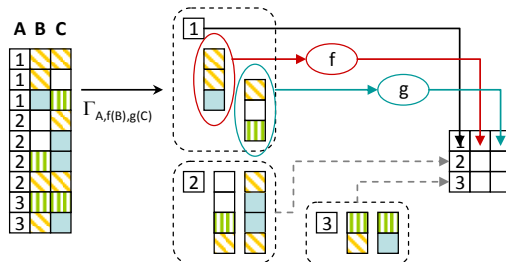
- Description #1: want to output a single tuple which summarizes a set of related tuples
- Description #2: want to “collapse” a set of tuples into a single, “representative” tuple
- Questions
  - How to identify related tuples (set to collapse)?  
=> **Grouping key**: a subset of attributes to test for equality
  - How to collapse a column into a value (summarize it)?  
=> Use an **aggregation function** (sum, count, avg, min, max, ...)

## Grouping ( $\Gamma$ )

- Duplicates useful when computing statistics
  - min, max, sum, count, average, ...
- $\Gamma_{A,B,C,f(X),g(Y),h(Z)}(R)$  computes aggregate values using some attributes as a grouping key
  - Implicit projection (drops unreferenced attributes)
  - A, B, C is the **grouping key**
  - X, Y, Z are **attributes** to aggregate
  - f, g, h are **aggregating** functions to apply
- Aggregating function should be commutative
  - $f(x,y) = f(y,x)$

## Grouping ( $\Gamma$ )

- All tuples having the same key go to same group
  - One output tuple for each unique key
  - Output “group total” for each non-key attribute in group



Example instance:

A is Employee level (1: Employee, 2: Manager, 3: Executive)  
B is Age, so f(B) can represent average Age  
C is Salary, so g(C) can represent average Salary

## Duplicates and grouping

- Consider a relation R modeling cars for sale
- $\Gamma_{\text{Make,count(*)}}(R)$  returns?
  - The number of cars of each make

Make	Model	Color
Toyota	Prius	Gray
Toyota	Prius	Red
Honda	Accord	Green
Honda	Accord	Red
Honda	Civic	Red
Ford	Echo	Red
Ford	Echo	Gray
Ford	Echo	White

Make	Count
Toyota	2
Honda	3
Ford	3

*Duplicates important for summaries (“how many”)*

## Sorting ( $\tau$ )

- $\tau_L(R)$  sorts tuples in R on list of attributes L
  - If L is A1, A2, ..., An tuples sorted first by A1. Ties are broken based on A2;...; Ties that remain after An broken arbitrarily.
  - Default: ascending order; With '-' in front: descending order

- Example:  $\tau_{\text{Count, Make}}(R)$

Make	Count		Make	Count	
Toyota	2		Ford	3	Descending count
Honda	3		Honda	3	
Ford	3	$\tau$	Toyota	2	Alphabetical order when count is equal

- Sorted relations not in the R.A. value domain!
  - =>  $\tau$  must be root operator of query tree\*\*

## The dangling tuple problem

- Consider the following query
  - $\tau_{\text{Total}}(\rho_{\text{Name, Total}}(\Gamma_{\text{Name, sum(Value)}(\text{Emp} \bowtie \text{Sales})))$
  - “List employees and their total sales in descending order”

Emp			Sales		
EID	Name		EID	Value	...
1	Mary	$\bowtie$	1	20	...
2	Xiao		3	10	...
3	Jaspreet		3	15	...

Name	Total
Jaspreet	25
Mary	20
Xiao?	

- Join hides fact that Xiao has no sales!
  - Challenge: rewrite query to include Xiao's zero

What's not there can be very important

## Extending the “inner” join

- All joins so far resemble intersection
  - => Tuples with no match discarded (“dangling”)
- Sometimes desirable to output dangling tuples
  - List all employees and their sales total (even if zero)
- Problem: what to use as missing half of tuple?
  - Introduce special value  $\perp$  (null)
  - Pad dangling tuples as needed to match schema
  - Note:  $\perp$  technically outside R.A. value domain

## Outer join ( $\bowtie$ )

- $T = R \bowtie S$  computes the “outer” join of R (left) and S (right)
  - Like normal join, but all tuples from R and S appear in output
  - Pad (left, right, or all) dangling tuples with  $\perp$  or NULL
    - LEFT — Tuples in inner join padded with tuples in R that have no matching tuples in S.
    - RIGHT — Tuples in inner join padded with tuples in S that have no matching tuples in R.
    - FULL — Tuples in inner join padded with tuples in R that have no matching tuples in S and tuples in S that have no matching tuples in R.
  - Natural, equi-, and theta- variants still apply
  - $|T| \geq \max(|R|, |S|)$

	$\bowtie$	$\bowtie_L$	$\bowtie_R$	$\bowtie_{\text{FULL}}$
1	A			
2	$\perp$			
3	A			
4	B			
4	C			
$\perp$	D			
5	$\perp$			

## Outer join ( $\bowtie$ ) examples

$r_1$	<table><tr><th>Employee</th><th>Department</th></tr><tr><td>Smith</td><td>sales</td></tr><tr><td>Black</td><td>production</td></tr><tr><td>White</td><td>production</td></tr></table>	Employee	Department	Smith	sales	Black	production	White	production	$r_2$	<table><tr><th>Department</th><th>Head</th></tr><tr><td>production</td><td>Mori</td></tr><tr><td>purchasing</td><td>Brown</td></tr></table>	Department	Head	production	Mori	purchasing	Brown	
Employee	Department																	
Smith	sales																	
Black	production																	
White	production																	
Department	Head																	
production	Mori																	
purchasing	Brown																	
$r_1 \bowtie_{\text{LEFT}} r_2$	<table><tr><th>Employee</th><th>Department</th><th>Head</th></tr><tr><td>Smith</td><td>sales</td><td>NULL</td></tr><tr><td>Black</td><td>production</td><td>Mori</td></tr><tr><td>White</td><td>production</td><td>Mori</td></tr></table>	Employee	Department	Head	Smith	sales	NULL	Black	production	Mori	White	production	Mori					
Employee	Department	Head																
Smith	sales	NULL																
Black	production	Mori																
White	production	Mori																
$r_1 \bowtie_{\text{RIGHT}} r_2$	<table><tr><th>Employee</th><th>Department</th><th>Head</th></tr><tr><td>Black</td><td>production</td><td>Mori</td></tr><tr><td>White</td><td>production</td><td>Mori</td></tr><tr><td>NULL</td><td>purchasing</td><td>Brown</td></tr></table>	Employee	Department	Head	Black	production	Mori	White	production	Mori	NULL	purchasing	Brown					
Employee	Department	Head																
Black	production	Mori																
White	production	Mori																
NULL	purchasing	Brown																
$r_1 \bowtie_{\text{FULL}} r_2$	<table><tr><th>Employee</th><th>Department</th><th>Head</th></tr><tr><td>Smith</td><td>Sales</td><td>NULL</td></tr><tr><td>Black</td><td>production</td><td>Mori</td></tr><tr><td>White</td><td>production</td><td>Mori</td></tr><tr><td>NULL</td><td>purchasing</td><td>Brown</td></tr></table>	Employee	Department	Head	Smith	Sales	NULL	Black	production	Mori	White	production	Mori	NULL	purchasing	Brown		
Employee	Department	Head																
Smith	Sales	NULL																
Black	production	Mori																
White	production	Mori																
NULL	purchasing	Brown																

## Extended projection

- $\pi_{x=E}(R)$  computes column x from expression E
  - Arithmetic ( $z=3*x + y$ )
  - String manipulation (substring, capitalization)
  - Some conditional expressions
- Example:

$\tau_{\text{Total}}(\pi_{\text{Name, Total}=\tau_{\text{T:0}}(\rho_{\text{Name, T}}(\Gamma_{\text{Name, sum(Value)}(\text{Emp} \bowtie \text{Sales}))))$

Emp			Sales					
EID	Name		EID	Value	...		Name	Total
1	Mary	$\bowtie$	1	20	...	=	Jaspreet	25
2	Xiao		3	10	...		Mary	20
3	Jaspreet		3	15	...		Xiao	0

## Outer join in action

- Consider the following query
  - $\tau_{\text{Total}}(\rho_{\text{Name, Total}}(\Gamma_{\text{Name, sum(Value)}(\text{Emp} \bowtie \text{Sales})))$
  - “List employees and their total sales in descending order”

Emp			Sales					
EID	Name		EID	Value	...		Name	Total
1	Mary	$\bowtie$	1	20	...	=	Jaspreet	25
2	Xiao		3	10	...		Mary	20
3	Jaspreet		3	15	...		Xiao	⊥

## Coming next...

- SQL (Structured Query Algebra)