

# CSC 365

## Introduction to Database Systems

- Operands in Relational Algebra are relations (or expressions that yield relations)
- A few additional derived operators are not listed: additional join types (Equijoin, Semijoin, Antijoin) and Division

(*derived operators are marked with \**)

Operator Name	Symbol
Selection	$\sigma$
Projection	$\pi$
Rename	$\rho$
Union	$\cup$
Difference	$-$
Intersection *	$\cap$
Cartesian Product	$\times$
Natural Join *	$\bowtie$
Theta Join *	$\bowtie_{\theta}$

- Several types of join operations
  - Natural join
  - Theta join
  - **Equijoin**
  - **Semijoin**
  - **Antijoin**
- Selectively pair tuples from two relations.
- Joins can be expressed in terms of other basic operators we have seen (cartesian product, selection, projection)

A theta join where the  $\theta$  condition is a conjunction of equality comparisons

FACULTY(ID, Name, Dept, Salary)

CHAIR(Dept, FacultyID)

Find salaries of department chairs:

$\pi_{FACULTY.Dept, FACULTY.Salary} (FACULTY \bowtie_{FACULTY.ID = CHAIR.FacultyID} CHAIR)$

A specialized form of natural join. The result includes attributes from only one of the relations (in the case of left semijoin:  $R \ltimes S$ , attributes from R) Given:

EMPLOYEE(EmpID, Name, DeptID)

DEPARTMENT(DeptID, DeptName)

$\text{EMPLOYEE} \ltimes \text{DEPARTMENT}$  has the same schema as EMPLOYEE, lists all employees who are assigned to a department.

$$R \ltimes S = \pi_{R.A1, R.A2, \dots, R.An} (R \bowtie S)$$

The result of  $R \triangleright S$  (sometimes written as  $R \bowtie \overline{S}$ ) includes only those tuples from  $R$  for which there is **no tuple in  $S$**  that is equal on the common attribute names. Antijoin is the complement of Semijoin:

$$R \triangleright S = R - R \bowtie S$$

Consider EMPLOYEE(Name, Dept) and CUSTOMER(Name, Address).  
EMPLOYEE  $\triangleright$  CUSTOMER would list all employees who *don't share* a name with any customers.

- Useful shorthand, sometimes seen when discussing distributed databases, where one key goal is to minimize network transfer
- Note: Semijoin is **not** the same as SQL's `OUTER JOIN`
- Semijoin appears in SQL as `EXISTS` or `IN`
- Antijoin appears in SQL as `NOT EXISTS` or `NOT IN`

Division allows us to compactly represent queries such as:

- Find the names of people who have visited every State Park in California
- Find the names of Netflix customers who have watched every film which features your favorite movie star.



- Five fundamental operators
  - Selection
  - Projection
  - Union
  - Difference
  - Cartesian Product
- Derived operators
  - Intersection
  - Joins: natural, theta, semi, anti
  - Division
- Utility operator
  - Rename

$$R \cap S = R - (R - S) = S - (S - R) = R \bowtie S$$

$$R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$$

$$\sigma_{\theta_1 \wedge \theta_2}(R) = \sigma_{\theta_1}(\sigma_{\theta_2}(R)) = \sigma_{\theta_2}(\sigma_{\theta_1}(R))$$

$$\sigma_{\theta_1 \vee \theta_2}(R) = \sigma_{\theta_1}(R) \cup \sigma_{\theta_2}(R)$$

Selection distributes over the union, intersection, and set difference operations	$\sigma_{\theta} (R \cup S) = \sigma_{\theta} (R) \cup \sigma_{\theta} (S)$ $\sigma_{\theta}(R - S) = \sigma_{\theta}(R) - S = \sigma_{\theta}(R) - \sigma_{\theta}(S)$
Multiple applications of the same select have no effect (selection is idempotent)	$\sigma_{\theta} (\sigma_{\theta} (R)) = \sigma_{\theta} (R)$
Selection is commutative (order of selections does not matter)	$\sigma_{\theta_1} (\sigma_{\theta_2} (R)) = \sigma_{\theta_2} (\sigma_{\theta_1} (R))$
Projection is distributive over set union	$\pi_L (R \cup S) = \pi_L (R) \cup \pi_L (S)$
Projection <i>does not</i> distribute over intersection or difference	<del> <math display="block">\pi_L (R \cap S) = \pi_L (R) \cap \pi_L (S)</math> <math display="block">\pi_L (R - S) = \pi_L (R) - \pi_L (S)</math> </del>

- Laws that highlight equivalent expressions are useful during query optimization (an important task of the DBMS)
- Two equivalent expressions may have *very* different computational costs (when considering I/O & CPU)

Three ways to represent relational algebra expressions:

- One long expression with parentheses
- Tree representation
- Sequence of assignment statements (AKA "linear notation")

We can represent a relational algebra expression as a tree:

- Internal nodes are operators
- Leaf nodes are relations

## AIRPLANE

<u>TailNum</u>	Make	Model	MaxSpeed
C97W	Boeing	797	<i>null</i>
R53Q	Cessna	FG	220
T80H	Airbus	A380	634
G59K	Airbus	A320	450
P88T	Piper	Arrow	180
K30W	Boeing	707	450

## FLIGHT

<u>TailNum</u>	<u>PilotID</u> -----	<u>CopilotID</u> -----	Runway	<u>Date</u>
R53Q	K407	D342	S-2	9/1/17
T80H	K407	null	W-2	9/21/17
C97W	D342	<i>null</i>	W-2	8/9/21
T80H	D342	K407	W-3	9/9/17

## PILOT

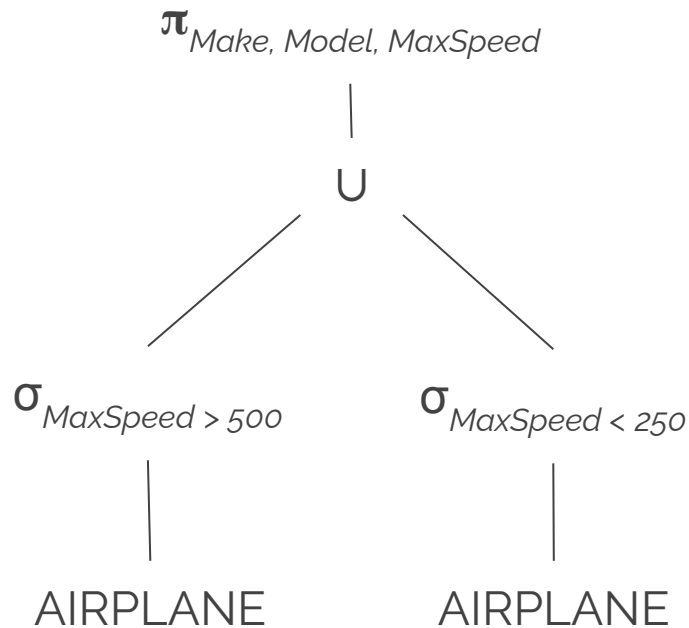
<u>PilotID</u>	Name
D342	Charlie
K407	Juliett
H452	Piper

$\pi_{\text{Make, Model, MaxSpeed}}(\sigma_{\text{MaxSpeed} > 500}(\text{AIRPLANE}))$   
U

$\sigma_{\text{MaxSpeed} < 250}(\text{AIRPLANE})$

Line breaks have no  
significance. This could be  
shown on one long line.

Equivalent expression?





$$\pi_{\text{Make, Model, MaxSpeed}}(\sigma_{\text{MaxSpeed} > 500}(\text{AIRPLANE}) \cup \sigma_{\text{MaxSpeed} < 250}(\text{AIRPLANE}))$$

Linear Notation:

$$R(\text{TailNum, Make, Model, MaxSpeed}) := \sigma_{\text{MaxSpeed} > 500}(\text{AIRPLANE})$$

$$S(\text{TailNum, Make, Model, MaxSpeed}) := \sigma_{\text{MaxSpeed} < 250}(\text{AIRPLANE})$$

$$T(\text{TailNum, Make, Model, MaxSpeed}) := R \cup S$$

$$\text{Answer}(\text{Make, Model, MaxSpeed}) := \pi_{\text{Make, Model, MaxSpeed}}(T)$$

Attribute names may be omitted for brevity.

Given the database schema:

Relation	Sample Tuple
PRODUCT(maker, <u>model</u> , type)	(Dell, 1001, pc)
PC( <u>model</u> , speed, ram, hd, price)	(1001, 3.1, 8192, 512, 2114)
LAPTOP( <u>model</u> , speed, ram, hd, screen, price)	(2001, 1.7, 4096, 256, 13.3, 1699)
PRINTER( <u>model</u> , color, type, price)	(3001, true, laser, 239)

## PRODUCT

<b>maker</b>	<b><u>model</u></b>	<b>type</b>
Dell	1001	pc
Dell	2001	laptop
Orange	1002	pc
HP	1003	pc
HP	2002	laptop

## PC

<b><u>model</u></b>	<b>speed</b>	<b>ram</b>	<b>hd</b>	<b>price</b>
1001	3.1	8	512	1200
1002	2.8	16	512	1300
1003	2.9	8	1024	1500

## LAPTOP

<b><u>model</u></b>	<b>speed</b>	<b>ram</b>	<b>hd</b>	<b>screen</b>	<b>price</b>
2001	1.7	8	512	15	1599
2002	2.0	16	1024	13.3	999

- (1) What PC models have a speed of at least 3.00?
- (2) Find the model number and price of all products (of any type) made by manufacturer *Orange*

(3) Find the manufacturers (makers) that sell laptops, but not PCs.

*Hint: Recall relational algebra's set operators ( $\cup$ ,  $\cap$ ,  $-$ )*

(3) Find the manufacturers (makers) that sell laptops, but not PCs.

*Hint: Recall relational algebra's set operators ( $\cup$ ,  $\cap$ ,  $-$ )*

$$\pi_{maker}(\text{PRODUCT} \bowtie \text{LAPTOP}) - \pi_{maker}(\text{PRODUCT} \bowtie \text{PC})$$

(4) Find hard-disk sizes that occur in two or more PCs.

*Hint: It's entirely acceptable to join a relation to itself (a "self-join")*

(4) Find hard-disk sizes that occur in two or more PCs.

*Hint: It's entirely acceptable to join a relation to itself (a "self-join")*

$$\pi_{hd}(\sigma_{PC1.model \neq PC2.model \wedge PC1.hd = PC2.hd}(\rho_{PC1}(PC) \times \rho_{PC2}(PC)))$$



(5) Find the names of all manufacturers that sell products of all types (pc and laptop in our example; the query should not be specifically limited to these values)

Informally, division is related to the Cartesian product as follows:

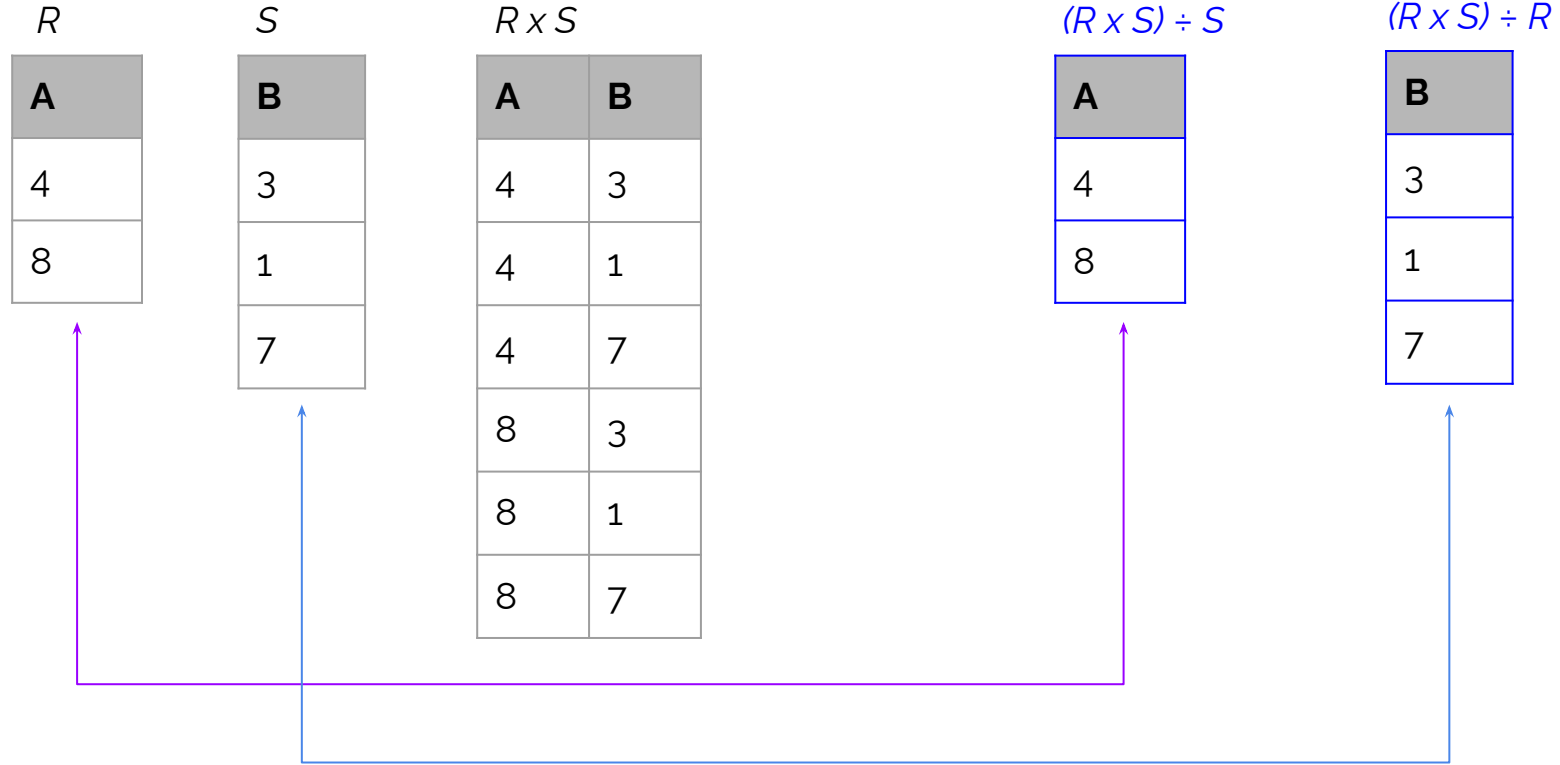
$$C = A \times B$$

$$A = C \div B$$

$$B = C \div A$$

*(This is an oversimplification that does not hold in the case of empty tables.)*

# Relational Algebra - Division Example



*R*

A
4
8

*S*

B
3
1
7

*T*

A	B
4	3
4	1
<del>4</del>	<del>7</del>
8	3
8	1
8	7

*T* ÷ *S*

A
<del>4</del>
8

*T* ÷ *R*

B
3
1
<del>7</del>

In this example,  
the row (4,7) has  
been removed  
from T

Division can be represented as a combination of projection, difference and cross product.

$$R \div S =$$

$$\pi_{\text{Schema}(R) - \text{Schema}(S)}(R) - \pi_{\text{Schema}(R) - \text{Schema}(S)}((\pi_{\text{Schema}(R) - \text{Schema}(S)}(R) \times S) - R)$$

(here,  $\text{Schema}(R)$  and  $\text{Schema}(S)$  represent the attribute sets of relation  $R$  and  $S$ , respectively)

	Relational Model	SQL
Structure	Attributes, Tuples, Relations	Columns, Rows, Tables
Manipulation	Relational Algebra	DML, Queries
Constraints	Primary Key, Referential (Foreign Key), Domain, Relational Algebra expressions to define constraints	CREATE TABLE ALTER TABLE Triggers