

# Relational Algebra and SQL

Slides by:

**Joe Hellerstein**

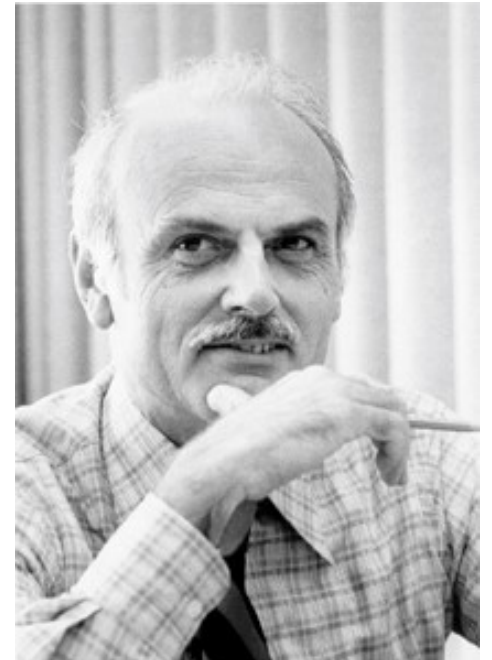
[hellerstein@berkeley.edu](mailto:hellerstein@berkeley.edu)

# A bit of computing history

- Pre-1969: databases were more like data *structures*
  - i.e. hackery

# The Relational Model

- A mathematical abstraction of a database, independent of data structures
  - Indeed, allows you to modify data structures at will without affecting program correctness!
  - So-called *Data Independence*



Edgar F. "Ted" Codd  
(1923 - 2003)  
Turing Award 1981

# Codd's Main Contributions

1. *Relational model*: Mathematical relations with typed attributes, primary keys and consistency constraints, *independent* of physical properties like sort order or indexes. (1969)
2. *A Relational Algebra* of simple operations on relations (1972)
  - In the spirit of abstract algebra (groups, rings, fields, etc)
  - Inspiration for functional libraries like Pandas
3. *A Relational Calculus* of truth expressions over relations (1972)
  - Inspiration for declarative languages like SQL, Datalog, as well as visual languages

# Codd's Theorem (1972)

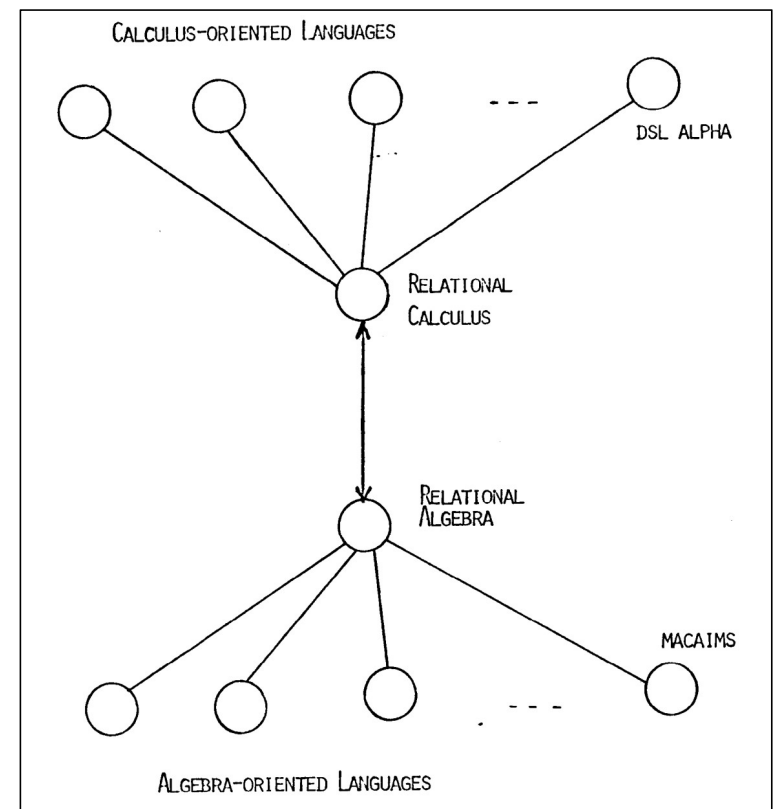
- The relational calculus and the relational algebra have *equivalent expressive power*.

RELATIONAL COMPLETENESS OF DATA BASE SUBLANGUAGES

by

E. F. Codd

IBM Research Laboratory  
San Jose, California



# Historical Perspective

- **1969:** Codd's Relational Model paper
- 1974: IBM System R and Berkeley Ingres research projects begin
- **1979:** Oracle released first commercial SQL system for DEC Vax minicomputer
- 1981: Ted Codd receives Turing Award
- 1983: IBM DB2 released for MVS mainframe
- 1984-87: Teradata, Informix SQL and Sybase released
- 1988: Berkeley Postgres project begins
- 1989: Microsoft SQL Server released (derived from Sybase)
- **1992:** First meaningful SQL standard
- 1995: PostgreSQL released ("Postgres 95"), MySQL released
- 2000: Sqlite released
- 2004: Google MapReduce paper
- 2010: Apache Hive (SQL on Hadoop) released
- **2012:** Pandas library popularized

# Relational Terminology

- *Database*: Set of Relations
- *Relation (Table)*:
  - *Schema (metadata)*
    - A unique name for the relation
    - A list of  $k$  distinct Attribute names, each associated with a type.
    - Optional *constraints* (key constraints)
  - *Instance (data)*
    - Set of  $k$ -tuples satisfying the schema
- *Attribute (Column, Field)*
- *Tuple (Row, Record)*

The schema of a database is the set of schemas of its relations.

# Boat Club Schema

sailors(sid integer, sname text, rating integer, age float)

boats(bid integer, bname text, color text)

reserves(sid integer, bid integer, day date)



# Boat Club

## Example Instances

### Boats

<u>bid</u>	bname	color
101	Interlake	blue
102	Interlake	red
104	Marine	red
103	Clipper	green

Note: primary keys underlined

### R1

<u>sid</u>	<u>bid</u>	<u>day</u>
22	101	10/10/16
58	103	11/12/96

### S1

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

### S2

<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

# Why learn Relational Algebra

- Intuitive for programmers
  - Imperative: apply this, then apply that
  - Set-oriented: no need for for-loops, low-level reasoning
- Basis of functional libraries like Pandas
  - Pandas (over-?) complicates things
  - Nice to have a clean foundation
- Simple optimization rules
  - Will help you think about writing efficient data-centric programs
- Common currency
  - Most data folk know the relational algebra operators

# Relational Algebra Preliminaries

- Algebra of *operators* on *relational instances*

$$\pi_{S.name}(\sigma_{R.bid=100 \wedge S.rating>5}(R \bowtie_{R.sid=S.sid} S))$$

- **Closed**: result is also a relational instance
  - Enables rich composition!
- **Typed**: input schema and operator determines output
  - Why is this important?
- Pure relational algebra has **set semantics**
  - **No duplicate** tuples in a relation instance
  - vs. SQL, which has *multiset* (bag) semantics

# Relational Algebra Operators

Unary Operators: operate on **single** relation instance

- **Projection (  $\pi$  )**: Retains only desired columns (vertical)
- **Selection (  $\sigma$  )**: Selects a subset of rows (horizontal)
- **Renaming (  $\rho$  )**: Rename attributes and relations.

Binary Operators: operate on **pairs** of relation instances

- **Union (  $\cup$  )**: Tuples in  $r1$  or in  $r2$ .
- **Intersection (  $\cap$  )**: Tuples in  $r1$  and in  $r2$ .
- **Set-difference (  $-$  )**: Tuples in  $r1$ , but not in  $r2$ .
- **Cross-product (  $\times$  )**: Allows us to combine two relations.
- **Joins (  $\bowtie_{\theta}$  ,  $\bowtie$  )**: Combine relations that satisfy predicates

# Projection ( $\pi$ ) *Selects a subset of columns (vertical)*

$$\pi_{\text{sname, age}}(S2)$$

List of Attributes

Relational Instance **S2**

<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



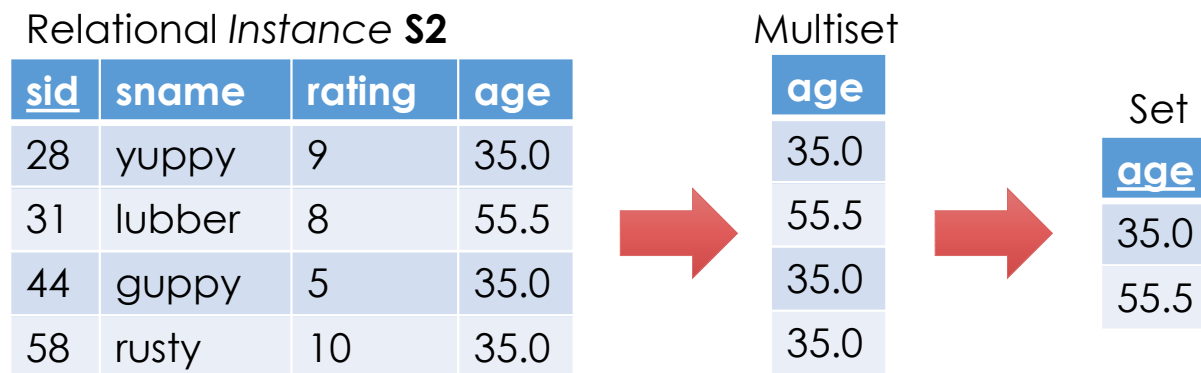
sname	age
yuppy	35.0
lubber	55.5
guppy	35.0
rusty	35.0

- Schema determined by schema of attribute list
  - Names and types correspond to input attributes

# Projection ( $\pi$ )

Selects a subset of columns (vertical)

$$\pi_{\text{age}}(S2)$$



- Set semantics → results in fewer rows
  - SQL systems don't automatically remove duplicates
  - Why?

# Selection( $\sigma$ )


Selects a subset of rows (horizontal)

$$\sigma_{\text{rating} > 8}(S2)$$

Selection Condition (Boolean Expression)

Relational Instance **S2**

<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



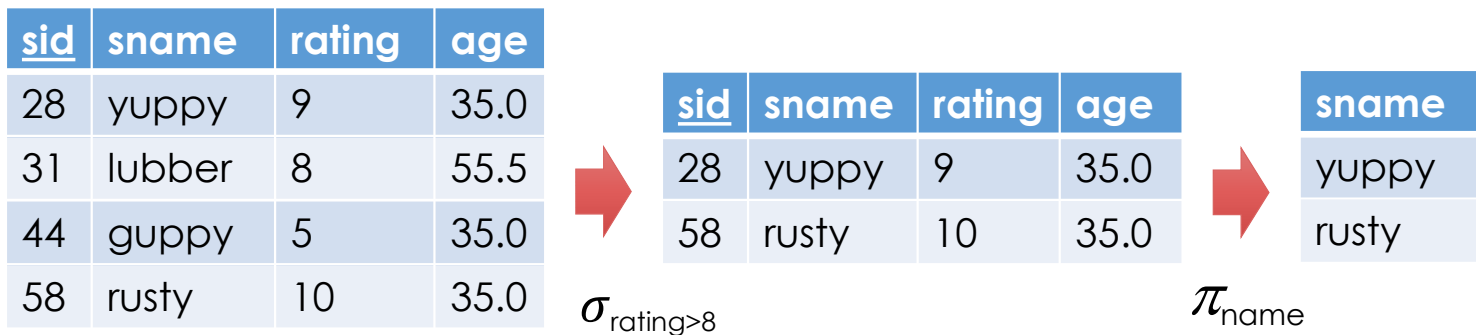
<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
58	rusty	10	35.0

- Output schema same as input
- Duplicate Elimination?

# Composing Select and Project

- Names of sailors with rating > 8

$$\pi_{\text{sname}}(\sigma_{\text{rating} > 8}(S2))$$



- What about:

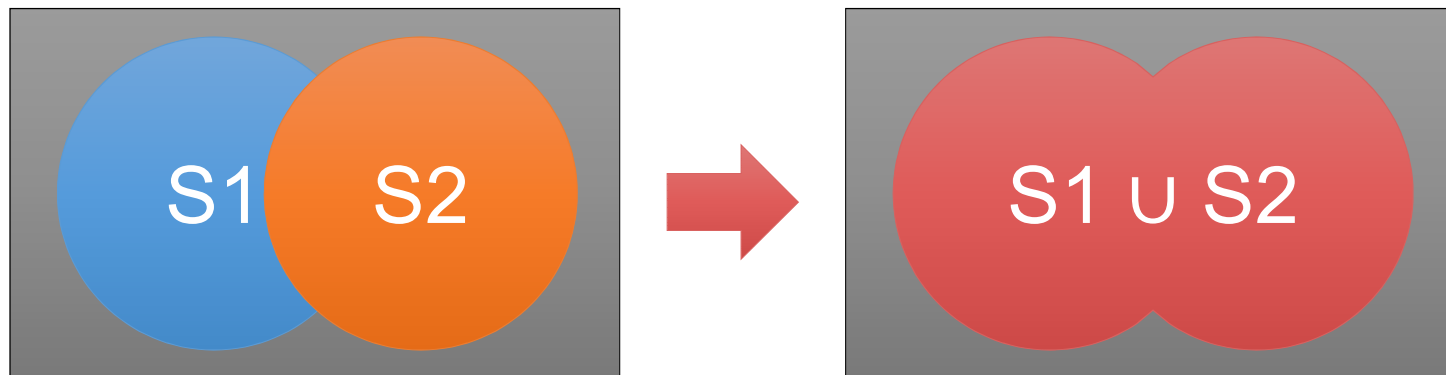
$$\sigma_{\text{rating} > 8}(\pi_{\text{sname}}(S2))$$

Invalid types. Input to  $\sigma_{\text{rating} > 8}$  does not contain *rating*.



# Union ( $\cup$ )

$S1 \cup S2$



Two input relations, must be *compatible*:

- Same number of fields.
- Fields in the same position have same **type**

# Union (U)

Relational Instance **S1**

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

Relational Instance **S2**

<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

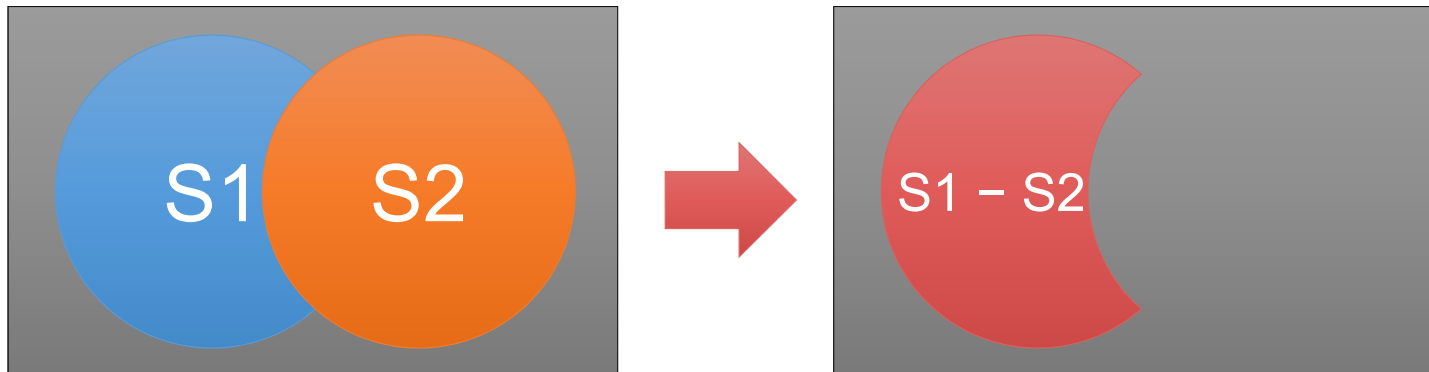
**S1 U S2**

<u>sid</u>	sname	rating	age
22	dustin	7	45
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

Duplicate elimination?

# Set Difference ( - )

$$S1 - S2$$



Same as with union, both input relations must be *compatible*.

# Set Difference ( - )

Relational Instance **S1**

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

Relational Instance **S2**

<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

## **S1 - S2**

<u>sid</u>	sname	rating	age
22	dustin	7	45

Symmetric?

## **S2 - S1**

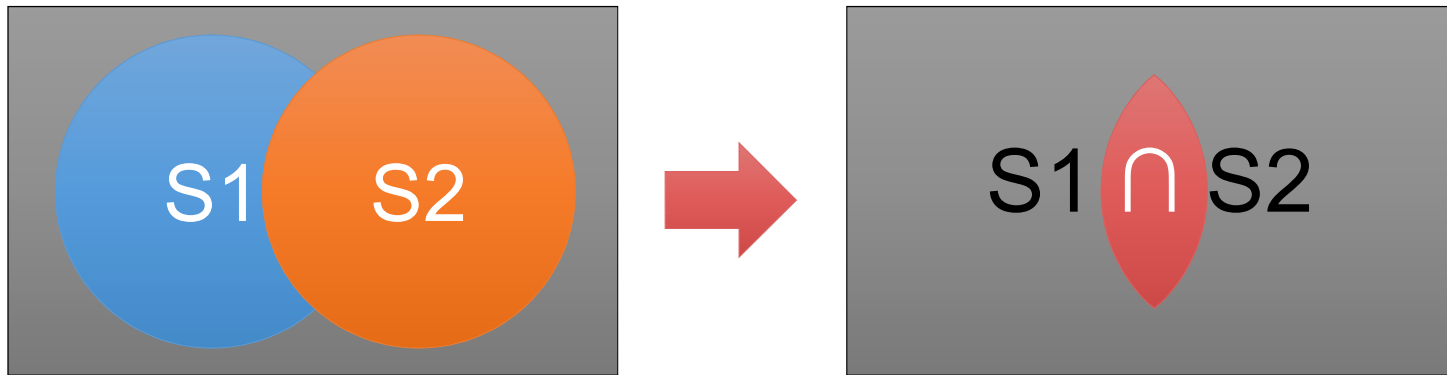
<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
44	guppy	5	35.0

Duplicate elimination?

- Not required

# Intersection ( $\cap$ )

$$s1 \cap s2$$



Same as with union, both input relations must be *compatible*.

# Intersection ( $\cap$ )

Relational Instance **S1**

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

Relational Instance **S2**

<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

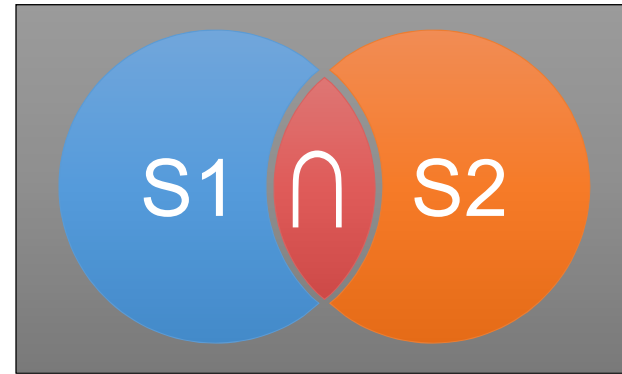
**S1  $\cap$  S2**

<u>sid</u>	sname	rating	age
31	lubber	8	55.5
58	rusty	10	35.0

Is intersection essential?

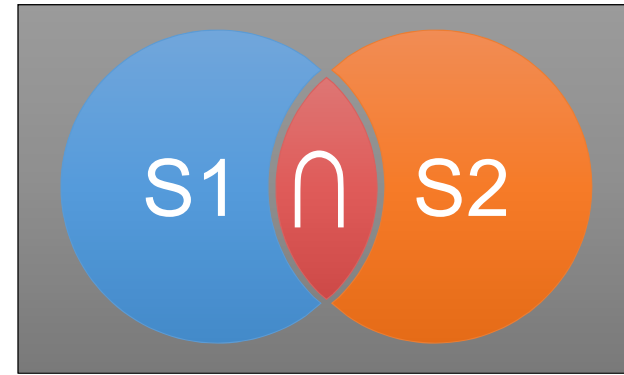
- Implement it with earlier ops. ?

Intersection ( $\cap$ )

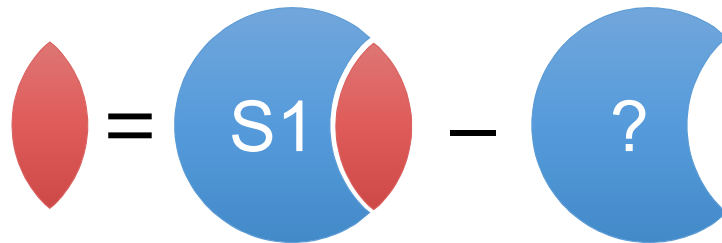


$$S1 \cap S2 = S1 - ?$$

Intersection ( $\cap$ )

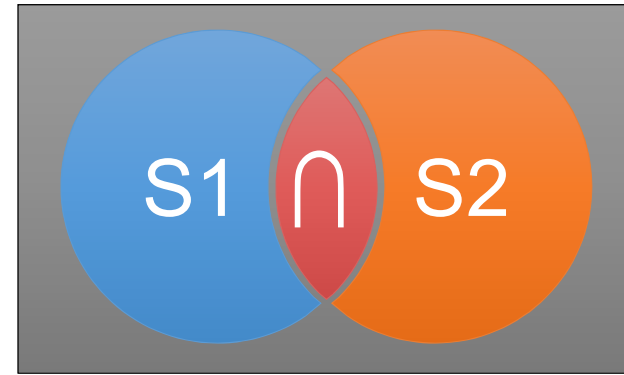


$$S1 \cap S2 = S1 - ?$$

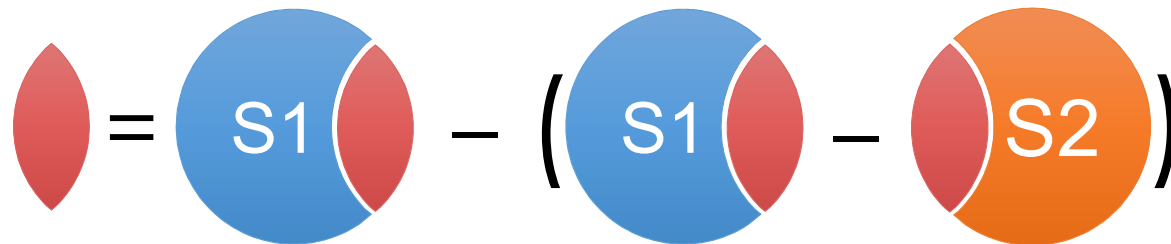




Intersection ( $\cap$ )



$$S1 \cap S2 = S1 - (S1 - S2)$$



# Cross-Product (×)

**R1 × S1:** Each row of **R1** paired with each row of **S1**

**R1:**

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

× **S1:**

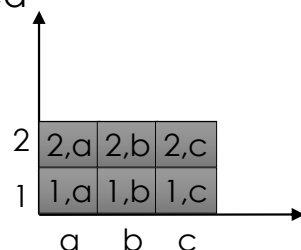
<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

=

**R1 × S1**

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

Sometimes also called  
**Cartesian Product:**



How many rows in the result?

| R1 | \* | R2 |

Schema compatibility?

No requirements.

One field per field in original schemas.

What about duplicate names?

Renaming operator

# Renaming ( $\rho = \text{"rho"}$ )

*Renames relations and their attributes:*

$$\rho(\text{Temp1}(1 \rightarrow \text{sid1}, 4 \rightarrow \text{sid2}), R1 \times S1)$$

Output Relation Name
Renaming List  
position  $\rightarrow$  New Name
Input Relation

**R1 × S1**

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



**Temp1**

sid1	bid	day	sid2	sname	rating	age
22	101	10/10/96	22	dustin	7	45.0
22	101	10/10/96	31	lubber	8	55.5
22	101	10/10/96	58	rusty	10	35.0
58	103	11/12/96	22	dustin	7	45.0
58	103	11/12/96	31	lubber	8	55.5
58	103	11/12/96	58	rusty	10	35.0

- Relational algebra can also be defined *positionally*, without names.
- Difficult to read ...

$$\pi_{f5}(\sigma_{f6 > f8}(S2))$$

# Compound Operator: Join

- Joins are compound operators (like intersection):
  - **Cross product** followed by **selection** and possibly **projection** (for natural join)
- Hierarchy of common kinds:
  - **Theta Join (  $\bowtie_{\theta}$  )**: *join on logical expression  $\theta$* 
    - **Equi-Join**: *theta join with conjunction equalities*
      - **Natural Join (  $\bowtie$  )**: *equi-join on all matching column names*
- Note: we should use a join, not a cross-product, if we can!  
Easier to read, clarifies opportunities for using efficient join algorithms.

# Theta Join ( $\bowtie_{\theta}$ )

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

**Example:** *Pair each sailor with older sailors.*

$$\mathbf{S1} \bowtie_{\text{age} < \text{age}} \mathbf{S1}$$

**S1:**

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

# Theta Join ( $\bowtie_{\theta}$ )

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

**Example:** Pair each sailor with older sailors.

$$S1 \bowtie_{\text{age} < \text{age}} S1$$

$$S1 \times S1$$

**S1:**

sid	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

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

# Theta Join ( $\bowtie_{\theta}$ )

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

**Example:** Pair each sailor with older sailors.

$$\mathbf{S1} \bowtie_{\text{age} < \text{age}} \mathbf{S1}$$

**S1:**

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

S1				S1			
sid	sname	rating	age	sid	sname	rating	age
22	dustin	7	45.0	31	lubber	8	55.5
22	dustin	7	45.0	58	rusty	10	35.0
31	lubber	8	55.5	58	rusty	10	35.0

- Result schema same as that of cross-product.
- Special Case:
  - **Equi-Join:** theta join with conjunction equalities
  - Special special case **Natural Join** ...

# Natural Join ( $\bowtie$ )

Special case of **equi-join** in which equalities are specified for all matching attributes, and duplicate attributes are projected away

$$R \bowtie S = \pi_{\text{unique attr.}} \sigma_{\text{eq. matching attr.}} (R \times S)$$

- Compute  $R \times S$
- **Select** rows where attributes appearing in both relations have equal values
- **Project** onto the set of all unique attributes.



# Natural Join ( $\bowtie$ )

$$\mathbf{R} \bowtie \mathbf{S} = \pi_{\text{unique attr.}} \sigma_{\text{eq. matching attr.}} (\mathbf{R} \times \mathbf{S})$$

Example:

R1:

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

S1:

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

R1  $\bowtie$  S1

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

# Natural Join ( $\bowtie$ )

$$R \bowtie S = \pi_{\text{unique attr.}} \sigma_{\text{eq. matching attr.}} (R \times S)$$

**Example:**

**R1:**

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

**S1:**

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

**R1  $\bowtie$  S1**

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

Commonly used for foreign key joins (as above).

## Exercise:

*Find names of sailors who've reserved boat #103*

➤ Solution 1:

<b>Boats</b> ( <u>bid</u> , bname, color) <b>Sailors</b> ( <u>sid</u> , sname, rating, age) <b>Reserves</b> ( <u>sid</u> , <u>bid</u> , <u>day</u> )
--

$$\pi_{\text{sname}}( \sigma_{\text{bid}=103}(\text{Sailors} \bowtie \text{Reserves}) )$$

➤ Solution 2:

$$\pi_{\text{sname}}(\text{Sailors} \bowtie \sigma_{\text{bid}=103}(\text{Reserves}))$$

## Exercise:

Find names of sailors who've reserved a red boat

➤ Solution 1:

**Boats**(bid, bname, color)  
**Sailors**(sid, sname, rating, age)  
**Reserves**(sid, bid, day)

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

➤ More “efficient” Solution 2:

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

In general many possible equivalent expressions: **algebra**...

# Relational Algebra Rules

## ➤ Selections:

- $\sigma_{c1 \wedge \dots \wedge cn}(R) \equiv \sigma_{c1}(\dots(\sigma_{cn}(R))\dots)$  (cascade)
- $\sigma_{c1}(\sigma_{c2}(R)) \equiv \sigma_{c2}(\sigma_{c1}(R))$  (commute)

## ➤ Projections:

- $\pi_{a1}(R) \equiv \pi_{a1}(\dots(\pi_{a1, \dots, an-1}(R))\dots)$  (cascade)

## ➤ Cartesian Product

- $R \times (S \times T) \equiv (R \times S) \times T$  (associative)
- $R \times S \equiv S \times R$  (commutative)
- *Applies for joins as well but be careful with join predicates ...*

<b>Boats</b> ( <u>bid</u> , bname, color) <b>Sailors</b> ( <u>sid</u> , sname, rating, age) <b>Reserves</b> ( <u>sid</u> , <u>bid</u> , <u>day</u> )
--

# Caution with Join Ordering

- Consider the following:

Boats ⋈<sub>bid</sub> Reserves ⋈<sub>sid</sub> Sailors

- Commute and Associate:

Boats ⋈<sub>bid</sub> (Sailors ⋈<sub>sid</sub> Reserves)

- Incompatible join predicate:

(Boats ⋈<sub>bid</sub> Sailors) ⋈<sub>sid</sub> Reserves

<b>Boats</b> ( <u>bid</u> , bname, color) <b>Sailors</b> ( <u>sid</u> , sname, rating, age) <b>Reserves</b> ( <u>sid</u> , <u>bid</u> , <u>day</u> )
--

# Caution with Join Ordering

- Consider the following:

Boats ⋈<sub>bid</sub> Reserves ⋈<sub>sid</sub> Sailors

- Commute and Associate:

Boats ⋈<sub>bid</sub> (Sailors ⋈<sub>sid</sub> Reserves)

- Incompatible join predicate:

(Boats × Sailors) ⋈<sub>sid, bid</sub> Reserves

# More Relational Algebra Rules

Commuting of selection operators

- $\sigma_c(R \times S) \equiv \sigma_c(R) \times S$  (c only has fields in R)
- $\sigma_c(R \bowtie S) \equiv \sigma_c(R) \bowtie S$  (c only has fields in R)

Commuting of projection operators

- $\pi_a(R \times S) \equiv \pi_{a_1}(R) \times \pi_{a_2}(S)$ 
  - $a_1$  is subset of a that mentions R and  $a_2$  is subset of a that mentions S
  - Similar result holds for joins



# A Standard Extension

- Group By / Aggregation Operator ( $\gamma$ ):

$$\gamma_{age, AVG(rating)}(Sailors)$$

- With selection (HAVING clause):

$$\gamma_{age, AVG(rating), COUNT(*) > 2}(Sailors)$$

## Recall Codd also had a Relational Calculus

- A *declarative* logic language
  - Find all tuples such that the following properties hold ...
  - Says “what” the output should be, not “how” to get it.
- SQL is based on the relational calculus
  - Even though, under the hood, database engines translate to algebra expressions!

# SQL Language

- Two sublanguages:
  - DDL – Data Definition Language
    - Define and modify schema
  - DML – Data Manipulation Language
    - Queries can be written intuitively.
- Relational Database Management System (RDBMS) responsible for efficient evaluation.
  - Choose and run algorithms for declarative queries

# We will learn SQL interactively

- Frontend: psql command line, Jupyter Notebook
- Backend: PostgreSQL