

# Relational Algebra and SQL

Slides by:

**Joe Hellerstein**

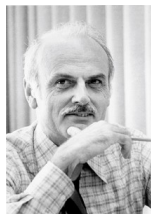
[hellerstein@berkeley.edu](mailto:hellerstein@berkeley.edu)

## 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

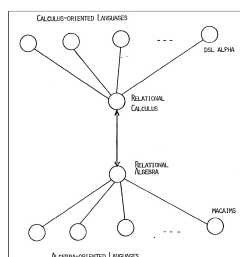
## Codd's Theorem (1972)

- 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
- **1977**: 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

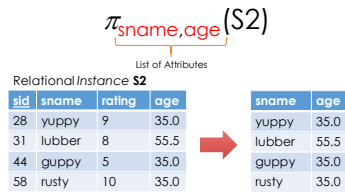
Unary Operators: operate on **single** relation instance

- **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  $r_1$  or in  $r_2$ .
- **Intersection (  $\cap$  )**: Tuples in  $r_1$  and in  $r_2$ .
- **Set-difference (  $-$  )**: Tuples in  $r_1$ , but not in  $r_2$ .
- **Cross-product (  $\times$  )**: Allows us to combine two relations.
- **Joins (  $\bowtie$ ,  $\Join$  )**: Combine relations that satisfy predicates

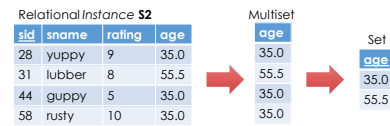
## Projection ( $\pi$ ) Selects a subset of columns (vertical)



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

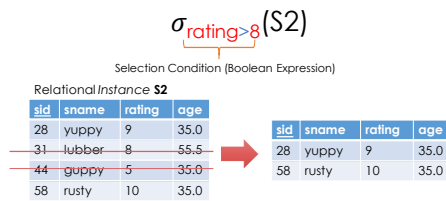
## Projection ( $\pi$ ) Selects a subset of columns (vertical)

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



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

## Selection ( $\sigma$ ) Selects a subset of rows (horizontal)

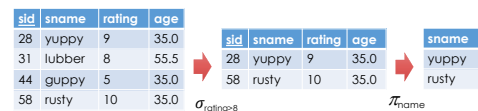


- Output schema same as input
- Duplicate Elimination?

## 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 ( $\cup$ )

**S1  $\cup$  S2**

Relational Instance **S1**

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

Relational Instance **S2**

sid	sname	rating	age
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 ( - )

 $S1 - S2$ Relational Instance  $S1$ 

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

sid	sname	rating	age
22	dustin	7	45

Symmetric?

 $S2 - S1$ Relational Instance  $S2$ 

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

sid	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 ( $\cap$ ) $S1 \cap S2$ Relational Instance  $S1$ 

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

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

Relational Instance  $S2$ 

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

Is intersection essential?

- Implement it with earlier ops. ?

Intersection ( $\cap$ )

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

Intersection ( $\cap$ )

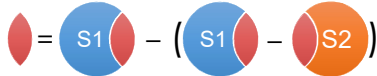
$$S1 \cap S2 = S1 - ?$$



Intersection ( $\cap$ )



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



Cross-Product ( $\times$ )

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

**R1:**

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

**R2:**

sid	bid	day	sid	iname	rating	age
22	101	10/10/96	22	dustin	7	45.0
			31	lubber	8	55.5
			58	rusty	10	35.5
58	103	11/12/96	22	dustin	7	45.0
			31	lubber	8	55.5
			58	rusty	10	35.5

**S1:**

sid	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.5

**R1 × S1**

sid	bid	day	sid	iname	rating	age
22	101	10/10/96	22	dustin	7	45.0
			31	lubber	8	55.5
			58	rusty	10	35.5
58	103	11/12/96	22	dustin	7	45.0
			31	lubber	8	55.5
			58	rusty	10	35.5

How many rows in the result?  $|R1| \times |S1|$

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

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 → New Name

Temp1

**R1 × S1**

sid	bid	day	sid	sname	rating	age
101	101	10/10/16	101	dustin	7	45.0
102	101	10/10/16	102	tubber	8	55.5
103	101	10/10/16	103	dustin	7	45.0
104	101	10/10/16	104	tubber	8	55.5
105	101	10/10/16	105	dustin	7	45.0
106	101	10/10/16	106	tubber	8	55.5
107	101	10/10/16	107	dustin	7	45.0
108	101	10/10/16	108	tubber	8	55.5
109	101	10/10/16	109	dustin	7	45.0
110	101	10/10/16	110	tubber	8	55.5
111	101	10/10/16	111	dustin	7	45.0
112	101	10/10/16	112	tubber	8	55.5
113	101	10/10/16	113	dustin	7	45.0
114	101	10/10/16	114	tubber	8	55.5
115	101	10/10/16	115	dustin	7	45.0
116	101	10/10/16	116	tubber	8	55.5
117	101	10/10/16	117	dustin	7	45.0
118	101	10/10/16	118	tubber	8	55.5
119	101	10/10/16	119	dustin	7	45.0
120	101	10/10/16	120	tubber	8	55.5
121	101	10/10/16	121	dustin	7	45.0
122	101	10/10/16	122	tubber	8	55.5
123	101	10/10/16	123	dustin	7	45.0
124	101	10/10/16	124	tubber	8	55.5
125	101	10/10/16	125	dustin	7	45.0
126	101	10/10/16	126	tubber	8	55.5
127	101	10/10/16	127	dustin	7	45.0
128	101	10/10/16	128	tubber	8	55.5
129	101	10/10/16	129	dustin	7	45.0
130	101	10/10/16	130	tubber	8	55.5
131	101	10/10/16	131	dustin	7	45.0
132	101	10/10/16	132	tubber	8	55.5
133	101	10/10/16	133	dustin	7	45.0
134	101	10/10/16	134	tubber	8	55.5
135	101	10/10/16	135	dustin	7	45.0
136	101	10/10/16	136	tubber	8	55.5
137	101	10/10/16	137	dustin	7	45.0
138	101	10/10/16	138	tubber	8	55.5
139	101	10/10/16	139	dustin	7	45.0
140	101	10/10/16	140	tubber	8	55.5
141	101	10/10/16	141	dustin	7	45.0
142	101	10/10/16	142	tubber	8	55.5
143	101	10/10/16	143	dustin	7	45.0
144	101	10/10/16	144	tubber	8	55.5
145	101	10/10/16	145	dustin	7	45.0
146	101	10/10/16	146	tubber	8	55.5
147	101	10/10/16	147	dustin	7	45.0
148	101	10/10/16	148	tubber	8	55.5
149	101	10/10/16	149	dustin	7	45.0
150	101	10/10/16	150	tubber	8	55.5
151	101	10/10/16	151	dustin	7	45.0
152	101	10/10/16	152	tubber	8	55.5
153	101	10/10/16	153	dustin	7	45.0
154	101	10/10/16	154	tubber	8	55.5
155	101	10/10/16	155	dustin	7	45.0
156	101	10/10/16	156	tubber	8	55.5
157	101	10/10/16	157	dustin	7	45.0
158	101	10/10/16	158	tubber	8	55.5
159	101	10/10/16	159	dustin	7	45.0
160	101	10/10/16	160	tubber	8	55.5
161	101	10/10/16	161	dustin	7	45.

- 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 ( $\bowtie_\theta$ )

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

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

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

**S1:**

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

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

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

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

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

s1:	s1				s1			
	id	name	rating	age	id	name	rating	age
22	dustin	7	45.0		22	dustin	7	45.0
31	lubber	8	55.5		31	lubber	8	55.5
58	rusty	10	35.0		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:				S1				S1			
sid	sname	rating	age	sid	sname	rating	age	sid	sname	rating	age
22	dustin	7	45.0	22	dustin	7	45.0	31	lubber	8	55.5
31	lubber	8	55.5	22	dustin	7	45.0	58	rusty	10	35.0
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 ( $\bowtie$ )

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

Example:

R1

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

S1

sid	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

Natural Join ( $\bowtie$ )

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

Example:

**R1:**

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

**S1:**

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

**R1 ⋈ S1**

Commonly used for foreign key joins (as above).

Commonly used for foreign key joins (as above).

## Exercise:

Find names of sailors who've reserved boat #103

➤ Solution 1:

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

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

➤ 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:**
  - $\sigma_{c_1 \wedge \dots \wedge c_n}(R) \equiv \sigma_{c_1}(\dots(\sigma_{c_n}(R))\dots)$  (cascade)
  - $\sigma_{c_1}(\sigma_{c_2}(R)) \equiv \sigma_{c_2}(\sigma_{c_1}(R))$  (commute)
- **Projections:**
  - $\pi_{a_1}(R) \equiv \pi_{a_1}(\dots(\pi_{a_1}(\dots \pi_{a_{n-1}}(R))\dots))$  (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 ...

## Caution with Join Ordering

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

- Consider the following:

Boats  $\bowtie_{bid}$  Reserves  $\bowtie_{sid}$  Sailors

- Commute and Associate:

Boats  $\bowtie_{bid}$  (Sailors  $\bowtie_{sid}$  Reserves)

- Incompatible join predicate:

(Boats  $\bowtie_{bid}$  Sailors)  $\bowtie_{sid}$  Reserves

## Caution with Join Ordering

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

- Consider the following:

Boats  $\bowtie_{bid}$  Reserves  $\bowtie_{sid}$  Sailors

- Commute and Associate:

Boats  $\bowtie_{bid}$  (Sailors  $\bowtie_{sid}$  Reserves)

- Incompatible join predicate:

(Boats  $\bowtie_{bid}$  Sailors)  $\bowtie_{sid, bid}$  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