

# **CSE 4020/5260**

## **Database Systems**

Instructor: Fitzroy Nembhard, Ph.D.

### **Week 4 & 5**

## **The Relational Model**



# Distribution

---

- All slides included in this class are for the exclusive use of students and instructors associated with Database Systems (CSE 4020/5260) at the Florida Institute of Technology
- Redistribution of the slides is not permitted without the written consent of the author.

# The Relational Model

- Structure of Relational Databases

- Relational Algebra

*Reading:*

*=> Chapter 2*

*=> Chapter 6, sections 1 & 2 (3 is optional).*

# Structure of a Relational Database

- *A relational database consists of a collection of **tables**, each of which is assigned a **unique name and** stores information about a set of entities*
- *A table has columns (attributes), each presents a specific type of information about the table*
  - *$R = (A_1, A_2, \dots, A_n)$  is a relation schema*  
*Example:          instructor (ID, name, dept\_name, salary)*
- *Each row of a table records information about one entity, called also a tuple*
- *A relationship between  $n$  values is represented mathematically by an  $n$ -tuple of values*

# Basic Structure

- Formally, given sets  $D_1, D_2, \dots, D_n$  a relation  $r$  is a subset of  $D_1 \times D_2 \times \dots \times D_n$
- Thus, a relation is a set of tuples  $(a_1, a_2, \dots, a_n)$  where each  $a_i \in D_i$
- Example:

*cust-name*        = {Jones, Smith, Curry, Lindsay}  
*cust-street*     = {Main, North, Park}  
*cust-city*        = {Harrison, Rye, Pittsfield}

$r = \{(Jones, Main, Harrison),$   
       $(Smith, North, Rye),$   
       $(Curry, North, Rye),$   
       $(Lindsay, Park, Pittsfield)\}$

# Relations are Unordered

---

- Since a relation is a *set*, the order of tuples is irrelevant and may be thought of as arbitrary.
- In a real DBMS, tuple order is typically very important and not arbitrary.
- Historically, this was/is a point of contention for the theorists.



# Table vs. Relation

- In a DBMS, a relation is represented or stored as a table.

- The Relation:

{ (10101,Srinivasan,Comp. Sci, 65000),  
(12121,Wu,Finance,90000),  
(15151,Mozart,Music,40000),  
:  
(98345,Kim,Elec. Eng., 80000) }

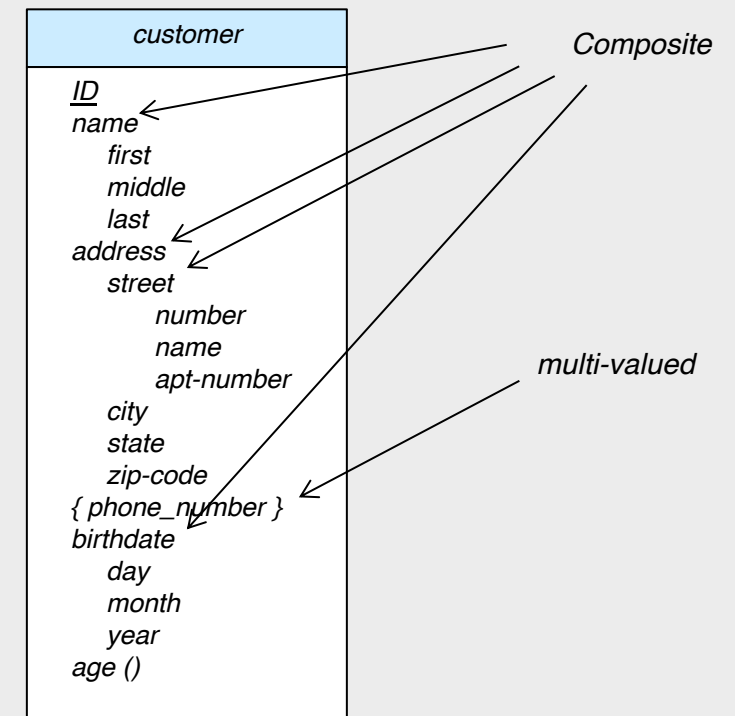
- The Table:

<i>ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>
10101	Srinivasan	Comp. Sci.	65000
12121	Wu	Finance	90000
15151	Mozart	Music	40000
22222	Einstein	Physics	95000
32343	El Said	History	60000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000

# Attribute Types

- Each attribute of a relation has a name.
- The set of allowed values for each attribute is called the domain of the attribute.
- Attribute values are required to be atomic, that is, indivisible.
- This will differ from ER modeling, which will have:
  - Multi-valued attributes
  - Composite attributes

## Recall ER- Modeling





# The Evil Value “Null”

- The special value *null* is an implicit member of every domain.
- Thus, tuples can have a *null* value for some of their attributes.
- A null value can be interpreted in several ways:
  - value is unknown
  - value does not exist
  - value is known and exists, but just hasn't been entered yet
- The null value causes complications in the definition of many operations.
- We shall consider their effect later.

# Relation Schema

- Let  $A_1, A_2, \dots, A_n$  be attributes. Then  $R = (A_1, A_2, \dots, A_n)$  is a relation schema.

*Customer-schema = (customer-name, customer-street, customer-city)*

- Sometimes referred to as a *relational schema* or *relational scheme*.

# Database

- A database consists of multiple relations: (example)

*account* - account information

*depositor* - depositor information, i.e., who deposits into which accounts

*customer* - customer information

- Storing all information as a single relation is possible:

*bank(account-number, balance, customer-name, ..)*

- This results in:

- Repetition of information (e.g. two customers own an account)
- The need for null values (e.g. represent a customer without an account).

# Relational Schemes

- Banking enterprise: (keys underlined)

*customer* (*customer-name*, *customer-street*, *customer-city*)

*branch* (*branch-name*, *branch-city*, *assets*)

*account* (*account-number*, *branch-name*, *balance*)

*loan* (*loan-number*, *branch-name*, *amount*)

*depositor* (*customer-name*, *account-number*)

*borrower* (*customer-name*, *loan-number*)

# Relational Schemes

## ■ University enterprise:

*classroom* (building, room-number, capacity)

*department* (dept-name, building, budget)

*course* (course-id, title, dept-name, credits)

*instructor* (ID, name, depart-name, salary)

*section* (course-id, sec-id, semester, year, building, room-number, time-slot-id)

*teaches* (ID, course-id, sec-id, semester, year)

*student* (ID, name, dept-name, tot-cred)

*takes* (ID, course-id, sec-id, semester, year, grade)

*advisor* (s-ID, i-ID)

*time-slot* (time-slot-id, day, start-time, end-time)

*prereq* (course-id, prereq-id)

# Relational Schemes

## ■ Employee enterprise:

*employee(person-name, street, city)*

*works(person-name, company-name, salary)*

*company(company-name, city)*

*manages(person-name, manager-name)*

# Query Languages

- Language in which user requests information from the database.
- Recall there are two categories of DML languages
  - procedural
  - non-procedural
- Query languages can be categorized as
  - Imperative (*example: Python, C and Java*)
  - Functional (*example: FQL, Geoquery, Kleisli, XMorph*)
  - Declarative (*example: Ruby, R and Haskell*)
- “Pure” languages:
  - Relational Algebra (functional query language, according to the current version of the book)
  - Tuple Relational Calculus (declarative)
  - Domain Relational Calculus (declarative)
- Pure languages form the underlying basis of “real” query languages.



# Imperative Language Code Example

- Find the name of the customer with customer\_id 192-83-7465:

```
try{
    Statement st = connection.createStatement();
    ResultSet rs = st.executeQuery("select customer.customer_name" +
                                   "from customer" +
                                   "where customer.customer_id = '192-83-7465'");

    while(rs.next){
        String s = rs.getString(1);
    }
} catch (SQLException e){}
```

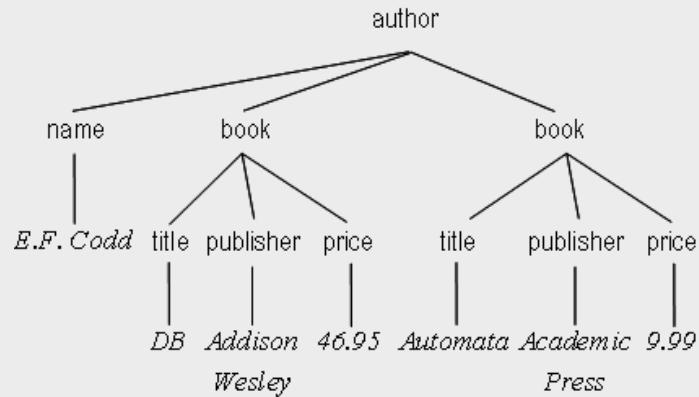
- Resultset declares what data is needed, which are included in the line of the SQL query:

**select** customer.customer-name **from** customer **where** customer.customer-id = '192-83-7465'

- The while loop states the way to retrieve the data.

# Functional Language Code Example

- XMorph is an XML data transformation language and implementation. It is a functional query language (i.e., similar to a query algebra rather than a tuple calculus)

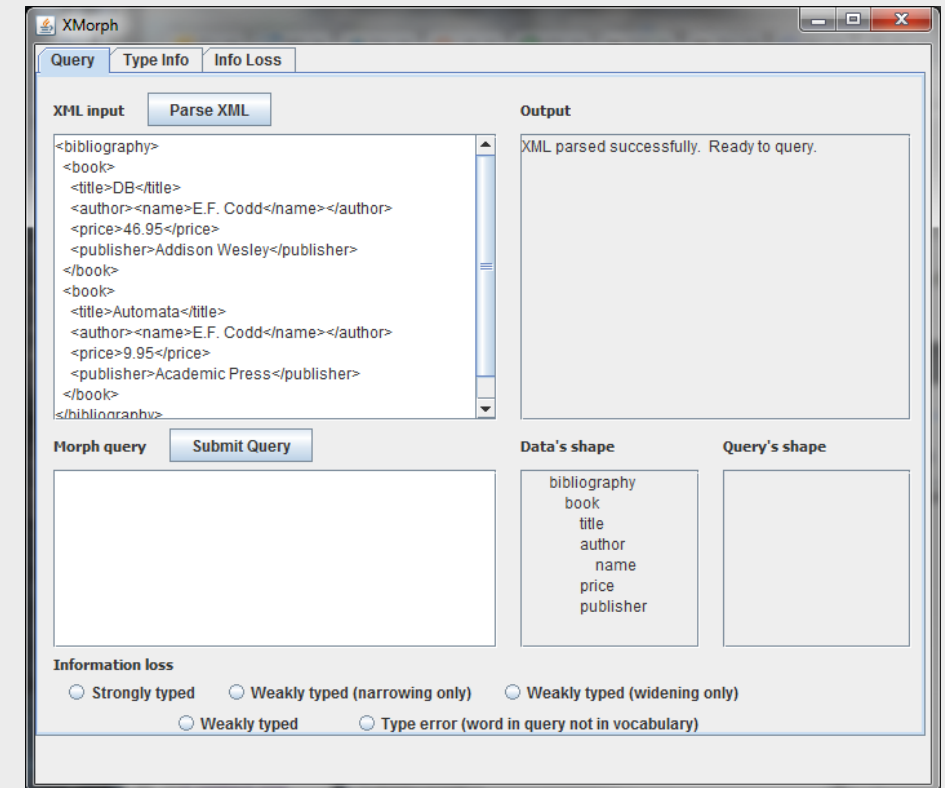


- Query 1: List authors by title

```
MORPH author [ name title ]
```

- Query 2: List titles only for the author named 'E. F. Codd':

```
MORPH author [
  name, WHERE value == 'E.F. Codd'
  title
]
```



Source: <https://cs.usu.edu/people/CurtisDyreson/XMorph>

# Declarative Language Code Example

- Haskell Database Connectivity (HDBC) is a library that provides a common abstraction or interface to different database engines like sqlite, mysql and postgres. The communication to each database engine is handled by the database driver.

```
//Connect to an SQLite Database
import Database.HDBC
import Database.HDBC.Sqlite3

> conn <- connectSqlite3 "zotero.sqlite"
conn :: Connection
>

//A quick query

> :t quickQuery
quickQuery
  :: IConnection conn =>
    conn -> String -> [SqlValue] -> IO [[SqlValue]]
>

//Run the Query

> quickQuery conn "SELECT tagID, name FROM tags WHERE tagID > 80 LIMIT 15" []
```

Source: <https://caiorss.github.io/Functional-Programming/haskell/DatabaseHDBC.html>

# Relational Algebra

- Functional query language (according to the book)

- Six basic operators:

- select ( $\sigma$ )
- project ( $\Pi$ )
- union ( $\cup$ )
- set difference ( $-$ )
- cartesian product ( $\times$ )
- rename ( $\rho$ )

# Relational Algebra

- Each operator takes one or more relations as input and results in a new relation.
- Each operation defines:
  - Requirements or constraints on its parameters.
  - Attributes in the resulting relation, including their types and names.
  - Which tuples will be included in the result.

# Select Operation – Example

## ■ Relation $r$

$A$	$B$	$C$	$D$
$\alpha$	$\alpha$	1	7
$\alpha$	$\beta$	5	7
$\beta$	$\beta$	12	3
$\beta$	$\beta$	23	10

## ■ $\sigma_{A=B \wedge D > 5}(r)$

$A$	$B$	$C$	$D$
$\alpha$	$\alpha$	1	7
$\beta$	$\beta$	23	10

# Tuple Relational Calculus – A Quick Look

- A declarative query language based on mathematical logic, where each query is of the form:

$$\{ t \mid P(t) \}$$

- Read as "the set of all tuples  $t$  such that predicate  $P$  is true for  $t$ "
- $P$  is a formula similar to that of predicate calculus
- $\exists t \in r (Q(t)) \equiv$  "there exists" a tuple  $t$  in relation  $r$  such that  $Q(t)$  is true
- $\forall t \in r (Q(t)) \equiv Q(t)$  is true "for all" tuples  $t$  in relation  $r$

$$\{ t \mid \exists t \in \text{customer} (t[\text{customer-name}] = \text{"Smith"}) \}$$

$$\{ t \mid \forall u \in \text{account} (u[\text{balance}] > 1000 \wedge u[\text{branch-name}] = \text{"Perryridge"}) \}$$

- While we do not cover Tuple Relational Calculus in detail in this class, its syntax is used to define many of the operations done in relational algebra.



# Tuple Relational Calculus – A Quick Example

- Find the *loan-number*, *branch-name*, and *amount* for loans of over \$1200

$$\{t \mid t \in \textit{loan} \wedge t[\textit{amount}] > 1200\}$$

## Schema

*branch* (*branch-name*, *branch-city*, *assets*)

*customer* (*customer-name*, *customer-street*, *customer-city*)

*account* (*account-number*, *branch-name*, *balance*)

*loan* (*loan-number*, *branch-name*, *amount*)

*depositor* (*customer-name*, *account-number*)

*borrower* (*customer-name*, *loan-number*)

# Domain Relational Calculus – A Quick Look

- A declarative query language based on mathematical logic, where each query is of the form:

$$\{ \langle x_1, x_2, x_3, \dots, x_n \rangle \mid P(x_1, x_2, x_3, \dots, x_n) \}$$

where,  $\langle x_1, x_2, x_3, \dots, x_n \rangle$  represents resulting domains variables and

$P(x_1, x_2, x_3, \dots, x_n)$  represents the condition or formula equivalent to the Predicate calculus.

- While we do not cover Domain Relational Calculus in detail in this class, some real query languages are built upon its constructs.

# Domain Relational Calculus – A Quick Example

- Find the *loan-number*, *branch-name*, and *amount* for loans of over \$1200

$$\{ \langle l, b, a \rangle \mid \langle l, b, a \rangle \in \text{loan} \wedge a > 1200 \}$$

## Schema

*branch* (*branch-name*, *branch-city*, *assets*)

*customer* (*customer-name*, *customer-street*, *customer-city*)

*account* (*account-number*, *branch-name*, *balance*)

*loan* (*loan-number*, *branch-name*, *amount*)

*depositor* (*customer-name*, *account-number*)

*borrower* (*customer-name*, *loan-number*)

# Difference Between Tuple Relational Calculus and Domain Relational Calculus

Tuple Relational Calculus (TRC)	Domain Relational Calculus (DRC)
In TRC, the variables represent the tuples from a specified relation.	In DRC, the variables represent the values drawn from a specified domain.
A tuple is a single element of relation. In a database, a tuple is equivalent to a row.	A domain is equivalent to column data type and any constraints on the value of the data.
Filtering variable uses a tuple of a relation.	Filtering is done based on the domain of the attributes.
Notation : $\{T \mid P(T)\}$ or $\{T \mid \text{Condition}(T)\}$	Notation : $\{a_1, a_2, a_3, \dots, a_n \mid P(a_1, a_2, a_3, \dots, a_n)\}$
Example : $\{T \mid \text{EMPLOYEE}(T) \text{ AND } T.\text{DEPT\_ID} = 10\}$	Example : $\{I \mid \text{EMPLOYEE} \langle I \rangle \text{ DEPT\_ID} = 10\}$

# Select Operation

- Notation:

$$\sigma_p(r)$$

where  $p$  is a selection predicate and  $r$  is a relation (or more generally, a relational algebra expression).

- Defined as:

$$\sigma_p(\mathbf{r}) = \{t \mid t \in r \textbf{ and } p(t)\}$$

where  $p$  is a formula in propositional logic consisting of terms connected by:  $\wedge$  (**and**),  $\vee$  (**or**),  $\neg$  (**not**), and where each term can involve the comparison operators:  $=$ ,  $\neq$ ,  $>$ ,  $\geq$ ,  $<$ ,  $\leq$

*\* Note that, in the book's notation, the predicate  $p$  cannot contain a subquery.*

# Select Operation, Cont.

## ■ Example:

$\sigma_{\text{branch-name}=\text{"Perryridge"}}(\text{account})$

$\sigma_{\text{customer-name}=\text{"Smith"} \wedge \text{customer-street} = \text{"main"}}(\text{customer})$

- Logically, one can think of selection as performing a table scan, but technically this may or may not be the case, i.e., an index may be used; that's why relational algebra is most frequently referred to as non-procedural.

### Schema

*branch* (branch-name, branch-city, assets)

*customer* (customer-name, customer-street, customer-city)

*account* (account-number, branch-name, balance)

*loan* (loan-number, branch-name, amount)

*depositor* (customer-name, account-number)

*borrower* (customer-name, loan-number)

# Project Operation – Example

■ Relation  $r$ :

$A$	$B$	$C$
$\alpha$	10	1
$\alpha$	20	1
$\beta$	30	1
$\beta$	40	2

■  $\Pi_{A,C}(r)$

$A$	$C$
$\alpha$	1
$\alpha$	1
$\beta$	1
$\beta$	2

=

$A$	$C$
$\alpha$	1
$\beta$	1
$\beta$	2



# Project Operation

- Notation:

$$\Pi_{A_1, A_2, \dots, A_k} (r)$$

where  $A_1, A_2$  are attribute names and  $r$  is a relation.

- The result is defined as the relation of  $k$  columns obtained by erasing or excluding the columns that are not listed.
- Duplicate rows are removed from the result, since relations are sets.

- Example:

$$\Pi_{\text{account-number, balance}} (\text{account})$$

Note, however, that account is not actually modified.

## Schema

*branch* (branch-name, branch-city, assets)  
*customer* (customer-name, customer-street, customer-city)  
*account* (account-number, branch-name, balance)  
*loan* (loan-number, branch-name, amount)  
*depositor* (customer-name, account-number)  
*borrower* (customer-name, loan-number)

# Project Operation

- The projection operation can also be used to reorder attributes.

$\Pi_{branch-name, balance, account-number} (account)$

As before, however, note that *account* is not actually modified; the order of the attributes is modified only in the result of the expression.

## Schema

*branch* (*branch-name*, *branch-city*, *assets*)  
*customer* (*customer-name*, *customer-street*, *customer-city*)  
*account* (*account-number*, *branch-name*, *balance*)  
*loan* (*loan-number*, *branch-name*, *amount*)  
*depositor* (*customer-name*, *account-number*)  
*borrower* (*customer-name*, *loan-number*)

# Union Operation – Example

■ Relations  $r, s$ :

$A$	$B$
$\alpha$	1
$\alpha$	2
$\beta$	1

$r$

$A$	$B$
$\alpha$	2
$\beta$	3

$s$

$r \cup s$

$A$	$B$
$\alpha$	1
$\alpha$	2
$\beta$	1
$\beta$	3

# Union Operation

■ Notation:  $r \cup s$

■ Defined as:

$$r \cup s = \{t \mid t \in r \text{ or } t \in s\}$$

■ Union can only be taken between *compatible* relations.

- $r$  and  $s$  must have the *same arity* (same number of attributes)
- attribute domains of  $r$  and  $s$  must be compatible (e.g., 2nd attribute of  $r$  deals with “the same type of values” as does the 2nd attribute of  $s$ )

■ Example: find all customers with either an account or a loan

$$\Pi_{customer-name} (depositor) \cup \Pi_{customer-name} (borrower)$$

## Schema

*branch* (branch-name, branch-city, assets)  
*customer* (customer-name, customer-street, customer-city)  
*account* (account-number, branch-name, balance)  
*loan* (loan-number, branch-name, amount)  
*depositor* (customer-name, account-number)  
*borrower* (customer-name, loan-number)

# Union Operation – Example

## ■ Relations *Paternity*, *Maternity*:

<i>Father</i>	<i>Child</i>
<i>Adam</i>	<i>Cain</i>
<i>Adam</i>	<i>Abel</i>
<i>Abraham</i>	<i>Isaac</i>
<i>Abraham</i>	<i>Ishmael</i>

*Paternity*

<i>Mother</i>	<i>Child</i>
<i>Eve</i>	<i>Cain</i>
<i>Eve</i>	<i>Seth</i>
<i>Sarah</i>	<i>Isaac</i>
<i>Hagar</i>	<i>Ishmael</i>

*Maternity*

$\text{Paternity} \cup \text{Maternity}?$

# Set Difference Operation

■ Relations  $r, s$ :

$A$	$B$
$\alpha$	1
$\alpha$	2
$\beta$	1

$r$

$A$	$B$
$\alpha$	2
$\beta$	3

$s$

$r - s$

$A$	$B$
$\alpha$	1
$\beta$	1

# Set Difference Operation, Cont.

- Notation  $r - s$

- Defined as:

$$r - s = \{t \mid t \in r \text{ and } t \notin s\}$$

- Set difference can only be taken between *compatible* relations.

  - $r$  and  $s$  must have the *same arity*

  - attribute domains of  $r$  and  $s$  must be compatible

- Note that there is no requirement that the attribute names be the same.

  - So, what about attribute names in the result?

  - Similarly for union.



# Cartesian-Product Operation

■ Relations  $r, s$ :

$A$	$B$
-----	-----

$\alpha$	1
$\beta$	2

$r$

$C$	$D$	$E$
-----	-----	-----

$\alpha$	10	$a$
$\beta$	10	$a$
$\beta$	20	$b$
$\gamma$	10	$b$

$s$

$r \times s$ :

$A$	$B$	$C$	$D$	$E$
$\alpha$	1	$\alpha$	10	$a$
$\alpha$	1	$\beta$	10	$a$
$\alpha$	1	$\beta$	20	$b$
$\alpha$	1	$\gamma$	10	$b$
$\beta$	2	$\alpha$	10	$a$
$\beta$	2	$\beta$	10	$a$
$\beta$	2	$\beta$	20	$b$
$\beta$	2	$\gamma$	10	$b$

# Cartesian-Product Operation, Cont.

- Notation  $r \times s$

- Defined as:

$$r \times s = \{tq \mid t \in r \textbf{ and } q \in s\}$$

- In some cases, the attributes of  $r$  and  $s$  are disjoint, i.e., that  $R \cap S = \emptyset$ .

- If the attributes of  $r$  and  $s$  are not disjoint:

- Each attribute's name has its originating relations name as a prefix.
- If  $r$  and  $s$  are the same relation, then the rename operation can be used.

# Rename Operation

- The rename operator allows the results of an expression to be renamed.
- The operator appears in two forms:

$\rho_X(E)$  - returns the expression  $E$  under the name  $X$

$\rho_X(A1, A2, \dots, An)(E)$  - returns the expression  $E$  under name  $X$ , with attributes renamed to  $A1, A2, \dots, An$

- Typically used to resolve a name class or ambiguity.

# Rename Operation – Example

## ■ Relations *Paternity*, *Maternity*:

<i>Father</i>	<i>Child</i>
<i>Adam</i>	<i>Cain</i>
<i>Adam</i>	<i>Abel</i>
<i>Abraham</i>	<i>Isaac</i>
<i>Abraham</i>	<i>Ishmael</i>

*Paternity*

<i>Mother</i>	<i>Child</i>
<i>Eve</i>	<i>Cain</i>
<i>Eve</i>	<i>Seth</i>
<i>Sarah</i>	<i>Isaac</i>
<i>Hagar</i>	<i>Ishmael</i>

*Maternity*

<i>Parent</i>	<i>Child</i>
<i>Adam</i>	<i>Cain</i>
<i>Adam</i>	<i>Abel</i>
<i>Abraham</i>	<i>Isaac</i>
<i>Abraham</i>	<i>Ishmael</i>
<i>Eve</i>	<i>Cain</i>
<i>Eve</i>	<i>Seth</i>
<i>Sarah</i>	<i>Isaac</i>
<i>Hagar</i>	<i>Ismael</i>

$Paternity \cup Maternity?$

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

# Composition of Operations

- Expressions can be built using multiple operations

A	B
$\alpha$	1
$\beta$	2

$r$

C	D	E
$\alpha$	10	a
$\beta$	10	a
$\beta$	20	b
$\gamma$	10	b

$s$

$r \times s$

A	B	C	D	E
$\alpha$	1	$\alpha$	10	a
$\alpha$	1	$\beta$	10	a
$\alpha$	1	$\beta$	20	b
$\alpha$	1	$\gamma$	10	b
$\beta$	2	$\alpha$	10	a
$\beta$	2	$\beta$	10	a
$\beta$	2	$\beta$	20	b
$\beta$	2	$\gamma$	10	b

$\sigma_{A=C}(r \times s)$

A	B	C	D	E
$\alpha$	1	$\alpha$	10	a
$\beta$	2	$\beta$	20	a
$\beta$	2	$\beta$	20	b

# Formal (recursive) Definition of a Relational Algebraic Expression

- A basic expression in relational algebra consists of one of the following:
  - A relation in the database
  - A constant relation
  
- Let  $E_1$  and  $E_2$  be relational-algebra expressions. Then the following are also relational-algebra expressions:
  - $E_1 \cup E_2$
  - $E_1 - E_2$
  - $E_1 \times E_2$
  - $\sigma_p(E_1)$ ,  $P$  is a predicate on attributes in  $E_1$
  - $\Pi_s(E_1)$ ,  $S$  is a list consisting of attributes in  $E_1$
  - $\rho_x(E_1)$ ,  $x$  is the new name for the result of  $E_1$

# Banking Example

- Recall the relational schemes from the banking enterprise:

*branch* (*branch-name*, *branch-city*, *assets*)

*customer* (*customer-name*, *customer-street*, *customer-city*)

*account* (*account-number*, *branch-name*, *balance*)

*loan* (*loan-number*, *branch-name*, *amount*)

*depositor* (*customer-name*, *account-number*)

*borrower* (*customer-name*, *loan-number*)

# Example Queries

- Find all loans of over \$1200 (a bit ambiguous).

$\sigma_{amount > 1200} (loan)$

- Find the loan number for each loan with an amount greater than \$1200.

$\Pi_{loan-number} (\sigma_{amount > 1200} (loan))$

## Schema

*branch* (branch-name, branch-city, assets)

*customer* (customer-name, customer-street, customer-city)

*account* (account-number, branch-name, balance)

*loan* (loan-number, branch-name, amount)

*depositor* (customer-name, account-number)

*borrower* (customer-name, loan-number)



# Example Queries

- Find the names of all customers who have a loan, an account, or both.

$$\Pi_{customer-name} (borrower) \cup \Pi_{customer-name} (depositor)$$

- Find the names of all customers who have a loan and an account.

$$\Pi_{customer-name} (borrower) \cap \Pi_{customer-name} (depositor)$$

## Schema

*branch* (branch-name, branch-city, assets)

*customer* (customer-name, customer-street, customer-city)

*account* (account-number, branch-name, balance)

*loan* (loan-number, branch-name, amount)

*depositor* (customer-name, account-number)

*borrower* (customer-name, loan-number)

# Example Queries

- Find the names of all customers who have a loan at the Perryridge branch.

$\Pi_{customer-name} (\sigma_{branch-name="Perryridge"} (\sigma_{borrower.loan-number = loan.loan-number} (borrower \times loan)))$

## Schema

*branch* (branch-name, branch-city, assets)

*customer* (customer-name, customer-street, customer-city)

*account* (account-number, branch-name, balance)

*loan* (loan-number, branch-name, amount)

*depositor* (customer-name, account-number)

*borrower* (customer-name, loan-number)

# Example Queries

- Alternative - Find the names of all customers who have a loan at the Perryridge branch.

$\Pi_{customer-name}(\sigma_{loan.loan-number = borrower.loan-number}(borrower \times \sigma_{branch-name = "Perryridge"}(loan)))$

## Schema

*branch* (branch-name, branch-city, assets)

*customer* (customer-name, customer-street, customer-city)

*account* (account-number, branch-name, balance)

*loan* (loan-number, branch-name, amount)

*depositor* (customer-name, account-number)

*borrower* (customer-name, loan-number)

# Example Queries

- Find the names of all customers who have a loan at the Perryridge branch but no account at any branch of the bank.

$$\Pi_{customer-name} (\sigma_{branch-name = "Perryridge"} (\sigma_{borrower.loan-number = loan.loan-number} (borrower \times loan))) \\ - \Pi_{customer-name}(depositor)$$

- A general query writing strategy – start with something simpler, and then enhance.

## Schema

*branch* (branch-name, branch-city, assets)  
*customer* (customer-name, customer-street, customer-city)  
*account* (account-number, branch-name, balance)  
*loan* (loan-number, branch-name, amount)  
*depositor* (customer-name, account-number)  
*borrower* (customer-name, loan-number)

# Example Queries

## ■ Find the largest account balance:

- Requires comparing each account balance to every other account balance.
- Accomplished by performing a Cartesian product between account and itself.
- Unfortunately, this results in ambiguity of attribute names.
- Resolved by renaming one instance of the *account* relation as *d*.

$$\Pi_{balance}(account) - \Pi_{account.balance}(\sigma_{account.balance < d.balance} (account \times \rho_d(account)))$$

### Schema

*branch* (branch-name, branch-city, assets)  
*customer* (customer-name, customer-street, customer-city)  
*account* (account-number, branch-name, balance)  
*loan* (loan-number, branch-name, amount)  
*depositor* (customer-name, account-number)  
*borrower* (customer-name, loan-number)

# Additional Operations

- The following operations do not add any “power,” or rather, capability to relational algebra queries, but simplify common queries.
  - Set intersection
  - Natural join
  - Theta join
  - Outer join
  - Division
  - Assignment
- All of the above can be defined in terms of the six basic operators.

# Set-Intersection Operation

■ Notation:  $r \cap s$

■ Defined as:

$$r \cap s = \{ t \mid t \in r \textbf{ and } t \in s \}$$

■ Assume:

- $r, s$  have the *same arity*
- attributes of  $r$  and  $s$  are compatible

■ In terms of the 6 basic operators:

$$r \cap s = r - (r - s)$$

# Set-Intersection Operation, Cont.

■ Relation  $r$ ,  $s$ :

A	B
$\alpha$	1
$\alpha$	2
$\beta$	1

$r$

A	B
$\alpha$	2
$\beta$	3

$s$

■  $r \cap s$

A	B
$\alpha$	2



# Natural-Join Operation

■ Notation:  $r \bowtie s$

■ Let  $r$  and  $s$  be relations on schemas  $R$  and  $S$  respectively.

■  $r \bowtie s$  is a relation that:

- Has all attributes in  $R \cup S$
- For each pair of tuples  $t_r$  and  $t_s$  from  $r$  and  $s$ , respectively, if  $t_r$  and  $t_s$  have the same value on all attributes in  $R \cap S$ , add a “joined” tuple  $t$  to the result.

■ Joining two tuples  $t_r$  and  $t_s$  creates a third tuple  $t$  such that:

- $t$  has the same value as  $t_r$  on attributes in  $R$
- $t$  has the same value as  $t_s$  on attributes in  $S$

# Natural-Join Example

- Relational schemes for relations  $r$  and  $s$ , respectively:

$$R = (A, B, C, D)$$

$$S = (E, B, D) \quad \text{-- Note the common attributes, which is typical.}$$

- Resulting schema for  $r \bowtie s$ :

$$(A, B, C, D, E)$$

- In terms of the 6 basic operators  $r \bowtie s$  is defined as:

$$\Pi_{r.A, r.B, r.C, r.D, s.E} (\sigma_{r.B = s.B \wedge r.D = s.D} (r \times s))$$

- More generally, computing the natural join equates to a Cartesian product, followed by a selection, followed by a projection.

# Natural Join Example

■ Relations  $r$ ,  $s$ :

$A$	$B$	$C$	$D$
$\alpha$	1	$\alpha$	a
$\beta$	2	$\gamma$	a
$\gamma$	4	$\beta$	b
$\alpha$	1	$\gamma$	a
$\delta$	2	$\beta$	b

$r$

$B$	$D$	$E$
1	a	$\alpha$
3	a	$\beta$
1	a	$\gamma$
2	b	$\delta$
3	b	$\epsilon$

$s$

■ Contents of  $r \bowtie s$ :

$A$	$B$	$C$	$D$	$E$
$\alpha$	1	$\alpha$	a	$\alpha$
$\alpha$	1	$\alpha$	a	$\gamma$
$\alpha$	1	$\gamma$	a	$\alpha$
$\alpha$	1	$\gamma$	a	$\gamma$
$\delta$	2	$\beta$	b	$\delta$

# Natural Join – Another Example

- Find the names of all customers who have a loan at the Perryridge branch.

Original Expression:

$$\Pi_{customer-name} (\sigma_{branch-name="Perryridge"} (\sigma_{borrower.loan-number = loan.loan-number} (borrower \times loan)))$$

Using the Natural Join Operator:

$$\Pi_{customer-name} (\sigma_{branch-name = "Perryridge"} (borrower \bowtie loan))$$

- Specifying the join explicitly makes it look nicer, plus it helps the query optimizer.

## Schema

*branch* (branch-name, branch-city, assets)

*customer* (customer-name, customer-street, customer-city)

*account* (account-number, branch-name, balance)

*loan* (loan-number, branch-name, amount)

*depositor* (customer-name, account-number)

*borrower* (customer-name, loan-number)

# Natural Join – Another Example

- Find the instructor IDs for those who teach in the Crawford building.

$$\Pi_{ID}(\sigma_{building = "Crawford"}(teaches \bowtie section))$$

- In this case, the natural join is on four attributes  
– *course\_id*, *section\_id*, *semester*, and *year*.

*section*

<i>course_id</i>	<i>sec_id</i>	<i>semester</i>	<i>year</i>	<i>building</i>	<i>room_number</i>	<i>time_slot_id</i>
BIO-101	1	Summer	2017	Painter	514	B
BIO-301	1	Summer	2018	Painter	514	A
CS-101	1	Fall	2017	Packard	101	H
CS-101	1	Spring	2018	Packard	101	F
CS-190	1	Spring	2017	Taylor	3128	E
CS-190	2	Spring	2017	Taylor	3128	A
CS-315	1	Spring	2018	Watson	120	D
CS-319	1	Spring	2018	Watson	100	B
CS-319	2	Spring	2018	Taylor	3128	C
CS-347	1	Fall	2017	Taylor	3128	A
EE-181	1	Spring	2017	Taylor	3128	C
FIN-201	1	Spring	2018	Packard	101	B
HIS-351	1	Spring	2018	Painter	514	C
MU-199	1	Spring	2018	Packard	101	D
PHY-101	1	Fall	2017	Watson	100	A

*teaches*

<i>ID</i>	<i>course_id</i>	<i>sec_id</i>	<i>semester</i>	<i>year</i>
10101	CS-101	1	Fall	2017
10101	CS-315	1	Spring	2018
10101	CS-347	1	Fall	2017
12121	FIN-201	1	Spring	2018
15151	MU-199	1	Spring	2018
22222	PHY-101	1	Fall	2017
32343	HIS-351	1	Spring	2018
45565	CS-101	1	Spring	2018
45565	CS-319	1	Spring	2018
76766	BIO-101	1	Summer	2017
76766	BIO-301	1	Summer	2018
83821	CS-190	1	Spring	2017
83821	CS-190	2	Spring	2017
83821	CS-319	2	Spring	2018
98345	EE-181	1	Spring	2017

# Theta-Join Operation

■ Notation:  $r \bowtie_{\theta} s$

■ Let  $r$  and  $s$  be relations on schemas  $R$  and  $S$ , respectively, and let  $\theta$  be a predicate.

■ Then,  $r \bowtie_{\theta} s$  is a relation that:

- Has all attributes in  $R \cup S$  including duplicate attributes.
- For each pair of tuples  $t_r$  and  $t_s$  from  $r$  and  $s$ , respectively, if  $\theta$  evaluates to true for  $t_r$  and  $t_s$ , then add a “joined” tuple  $t$  to the result.

■ In terms of the 6 basic operators,  $r \bowtie_{\theta} s$  is defined as:

$$\sigma_{\theta}(r \times s)$$

# Theta-Join Example #1

## ■ Example:

$R = (A, B, C, D)$

$S = (E, B, D)$

## ■ Resulting schema:

$(r.A, r.B, r.C, r.D, s.E, s.B, s.D)$

## Theta Join – Example #2

- Consider the following relational schemes:

*Score* = (ID#, Exam#, Grade)

*Exam* = (Exam#, Average)

- Consider the following query:

*“Find the ID#s for those students who scored less than average on some exam.”*

$$\Pi_{\text{Score.ID\#}} (\text{Score} \bowtie_{\text{Score.Exam\# = Exam.Exam\#} \wedge \text{Score.Grade} < \text{Exam.Average}} \text{Exam})$$

- Note the above could also be done with a natural join, followed by a selection.



## Theta Join – Example #3

- Consider the following relational schemes: (Orlando temperatures)

$Temp-Avgs = (\underline{Year}, Avg-Temp)$

$Daily-Temps-2010 = (\underline{Date}, High-Temp)$

- Consider the following query:

*“Find the days during 2010 where the high temperature for the day was higher than the average for some prior year.”*

- Looks ugly, perhaps, but phrasing the query this way does have benefits for query optimization.

$\Pi_{Date} (Daily-Temps-2010 \bowtie Daily-Temps-2010.High-Temp > Temp-Avgs.Avg-Temp \wedge Temp-Avgs.Year < 2010 \ Temp-Avgs)$

# Outer Join

---

- An extension of the join operation that avoids loss of information.
- Computes the join and then adds tuples from one relation that do not match tuples in the other relation.
- Typically introduces *null* values.

# Outer Join – Example

## ■ Relation *loan*:

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>
L-170	Downtown	3000
L-230	Redwood	4000
L-260	Perryridge	1700

## ■ Relation *borrower*:

<i>customer-name</i>	<i>loan-number</i>
Jones	L-170
Smith	L-230
Hayes	L-155

# Outer Join – Example

## Inner Join

*loan* ⋈ *Borrower*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>	<i>customer-name</i>
L-170	Downtown	3000	Jones
L-230	Redwood	4000	Smith

## Left Outer Join

*loan* ⋈<sub>L</sub> *Borrower*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>	<i>customer-name</i>
L-170	Downtown	3000	Jones
L-230	Redwood	4000	Smith
L-260	Perryridge	1700	<i>null</i>

*Loan*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>
L-170	Downtown	3000
L-230	Redwood	4000
L-260	Perryridge	1700

*Borrower*

<i>customer-name</i>	<i>loan-number</i>
Jones	L-170
Smith	L-230
Hayes	L-155

# Outer Join – Example

## ■ Right Outer Join

*loan* ⋈<sub>⊂</sub> *borrower*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>	<i>customer-name</i>
L-170	Downtown	3000	Jones
L-230	Redwood	4000	Smith
L-155	<i>null</i>	<i>null</i>	Hayes

## ■ Full Outer Join

*loan* ⋈<sub>⊃</sub> *borrower*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>	<i>customer-name</i>
L-170	Downtown	3000	Jones
L-230	Redwood	4000	Smith
L-260	Perryridge	1700	<i>null</i>
L-155	<i>null</i>	<i>null</i>	Hayes

### *Loan*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>
L-170	Downtown	3000
L-230	Redwood	4000
L-260	Perryridge	1700

### *Borrower*

<i>customer-name</i>	<i>loan-number</i>
Jones	L-170
Smith	L-230
Hayes	L-155

# Example Left-Outer Join

- Consider the following relational schemes:

*Student = (SS#, Address, Date-of-Birth)*

*Grade-Point-Average = (SS#, GPA)*

- Consider the following query:

*“Create a list of all student SS#s and their GPAs. Be sure to include all students, including first semester freshman, who do not have a GPA.”*

- Solution:

$\Pi_{SS\#,GPA} (Student \bowtie Grade-Point-Average)$

# Outer Join

- In terms of the 6 basic operators (plus natural join ☺), let  $r(R)$  and  $s(S)$  be relations:

$$r \bowtie s = (r - \Pi_R (r \bowtie s)) \times \{(null, null, \dots, null)\} \cup (r \bowtie s)$$

where  $\{(null, null, \dots, null)\}$  is on the schema  $S - R$

# Division Operation

- Notation:  $r \div s$
- Suited to queries that require “universal quantification,” e.g., include the phrase “*for all*.”



# Division Operation

- Let  $r$  and  $s$  be relations on schemas  $R$  and  $S$  respectively where  $S \subseteq R$ .

Assume without loss of generality that the attributes of  $R$  and  $S$  are:

$$R = (A_1, \dots, A_m, B_1, \dots, B_n)$$

$$S = (B_1, \dots, B_n)$$

The  $A_i$  attributes will be referred to as *prefix* attributes, and the  $B_i$  attributes will be referred to as *suffix* attributes.

The result of  $r \div s$  is a relation on schema

$$R - S = (A_1, \dots, A_m)$$

where:

$$r \div s = \{ t \mid t \in \Pi_{R-S}(r) \wedge \forall u \in s (tu \in r) \}$$

# Division – Example #1

Relations  $r, s$ :

$r$

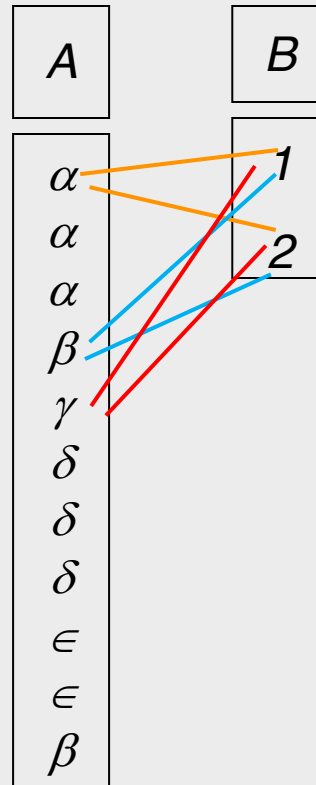
A	B
$\alpha$	1
$\alpha$	2
$\alpha$	3
$\beta$	1
$\gamma$	1
$\delta$	1
$\delta$	3
$\delta$	4
$\epsilon$	6
$\epsilon$	1
$\beta$	2

$s$

B
1
2

$r_{attr} - s_{attr}$

$s$



$r$

A	B
$\alpha$	1
$\alpha$	2
$\alpha$	3
$\beta$	1
$\gamma$	1
$\delta$	1
$\delta$	3
$\delta$	4
$\epsilon$	6
$\epsilon$	1
$\beta$	2

$r \div s$ :

A
$\alpha$
$\beta$

$\gamma, \delta$  and  $\epsilon$  are not returned because when paired with  $s$ , the following are true:

$\{(\gamma, 1), (\gamma, 2)\} \notin r$

$\{(\delta, 1), (\delta, 2)\} \notin r$

$\{(\epsilon, 1), (\epsilon, 2)\} \notin r$

To be in  $r$ , the following must

hold:  $\forall u \in s (tu \in r)$

# Division – Example #2

Relations  $r, s$ :

$r$

$A$	$B$	$C$	$D$	$E$
$\alpha$	a	$\alpha$	a	1
$\alpha$	a	$\gamma$	a	1
$\alpha$	a	$\gamma$	b	1
$\beta$	a	$\gamma$	a	1
$\beta$	a	$\gamma$	b	3
$\gamma$	a	$\gamma$	a	1
$\gamma$	a	$\gamma$	b	1
$\gamma$	a	$\beta$	b	1

$s$

$D$	$E$
a	1
b	1

$r \div s$ :

$A$	$B$	$C$
$\alpha$	a	$\gamma$
$\gamma$	a	$\gamma$

# Division – Example #3

Relations  $r, s$ :

$r$

$A$	$B$	$C$	$D$	$E$
$\alpha$	$a$	$\alpha$	$a$	$1$
$\alpha$	$a$	$\gamma$	$a$	$1$
$\alpha$	$a$	$\gamma$	$b$	$1$
$\beta$	$a$	$\gamma$	$a$	$1$
$\beta$	$a$	$\gamma$	$b$	$3$
$\gamma$	$a$	$\gamma$	$a$	$1$
$\gamma$	$a$	$\gamma$	$b$	$1$
$\gamma$	$a$	$\beta$	$b$	$1$

$s$

$B$	$D$
$a$	$a$
$a$	$b$

$r \div s$ :

$A$	$C$	$E$
$\alpha$	$\gamma$	$1$
$\gamma$	$\gamma$	$1$

## Division Operation (Cont.)

- In terms of the 6 basic operators, let  $r(R)$  and  $s(S)$  be relations, and let  $S \subseteq R$  :

$$r \div s = \Pi_{R-S}(r) - \Pi_{R-S}((\Pi_{R-S}(r) \times s) - \Pi_{R-S,S}(r))$$

To see why:

- $\Pi_{R-S,S}(r)$  simply reorders attributes of  $r$
- $\Pi_{R-S}(\Pi_{R-S}(r) \times s) - \Pi_{R-S,S}(r)$  gives those tuples  $t$  in  $\Pi_{R-S}(r)$  such that for some tuple  $u \in s$ ,  $tu \notin r$ .

- Property:

- Let  $q = r \div s$
- Then  $q$  is the largest relation satisfying  $q \times s \subseteq r$

# Example Queries

- Consider the following query:

*“Find the names of **all customers** who have an account at both the ‘Downtown’ and the ‘Uptown’ branches.”*

## Schema

*branch* (branch-name, branch-city, assets)  
*customer* (customer-name, customer-street, customer-city)  
*account* (account-number, branch-name, balance)  
*loan* (loan-number, branch-name, amount)  
*depositor* (customer-name, account-number)  
*borrower* (customer-name, loan-number)

- Query 1:

$$\Pi_{CN}(\sigma_{BN=\text{“Downtown”}}(\text{depositor} \bowtie \text{account})) \cap \Pi_{CN}(\sigma_{BN=\text{“Uptown”}}(\text{depositor} \bowtie \text{account}))$$

- Query 2:

$$\Pi_{customer-name, branch-name}(\text{depositor} \bowtie \text{account}) \div \rho_{temp(branch-name)}(\{(\text{“Downtown”}), (\text{“Uptown”})\})$$

# Example Queries

- Consider the following (more general) query:

*“Find all customers who have an account at all branches located in the city of Brooklyn.”*

## Schema

*branch* (*branch-name*, *branch-city*, *assets*)  
*customer* (*customer-name*, *customer-street*, *customer-city*)  
*account* (*account-number*, *branch-name*, *balance*)  
*loan* (*loan-number*, *branch-name*, *amount*)  
*depositor* (*customer-name*, *account-number*)  
*borrower* (*customer-name*, *loan-number*)

- How could Query 1 be modified for this scenario?

- How about Query 2?

$$\Pi_{customer-name, branch-name} (depositor \bowtie account) \div \Pi_{branch-name} (\sigma_{branch-city = \text{“Brooklyn”}} (branch))$$

- By the way, what would (should) be the result of the query if there are no Brooklyn branches?

# Assignment Operation

- The assignment operator ( $\leftarrow$ ) provides an easy way to express complex queries.
- Example (for  $r \div s$ ):

$temp1 \leftarrow \Pi_{R-S}(r)$   
 $temp2 \leftarrow \Pi_{R-S}((temp1 \times s) - \Pi_{R-S,S}(r))$   
 $result \leftarrow temp1 - temp2$

\*Do the exercises on the employee/works/company/manages DB!

\*And also the exercises on the university DB!

$$r \div s = \Pi_{R-S}(r) - \Pi_{R-S}((\Pi_{R-S}(r) \times s) - \Pi_{R-S,S}(r))$$



# Extended Relational Algebra Operations

---

- Generalized Projection
- Aggregate Operator

# Generalized Projection

- Extends projection by allowing arithmetic functions in the projection list.

$$\Pi_{F_1, F_2, \dots, F_n}(E)$$

- $E$  is any relational-algebra expression
- Each of  $F_1, F_2, \dots, F_n$  are arithmetic expressions involving constants and attributes in the schema of  $E$ .

**Recall project from the 6 basic operators**

$$\Pi_{A_1, A_2, \dots, A_k}(r)$$

where  $A_1, A_2$  are attribute names and  $r$  is a relation.

# Generalized Projection

- Consider the following relational scheme:

*credit-info* = (customer-name, limit, credit-balance)

- Give a relational algebraic expression for the following query:

*“Determine how much credit is left on each person’s line of credit; Also determine the percentage of their credit line that they have already used.”*

$\Pi_{customer-name, limit - credit-balance, (credit-balance/limit)*100} (credit-info)$

- We can also rename attributes

$\Pi_{customer-name, (limit - credit-balance) \text{ as credit-available}} (credit-info)$

# Aggregate Functions

- An aggregation function takes a collection of values and returns a single value:

<b>avg</b>	- average value
<b>min</b>	- minimum value
<b>max</b>	- maximum value
<b>sum</b>	- sum of values
<b>count</b>	- number of values

- Other aggregate functions are provided by most DBMS vendors.
- Not all aggregate operators are numeric, e.g., some apply to strings.

Hey, Mr. DB Here!  
Check out  
<https://dev.mysql.com/doc/refman/8.0/en/string-functions.html> for  
example string functions.



# The Aggregate Operator

- Aggregation functions are used in the aggregate operator:

$$G_1, G_2, \dots, G_n \mathcal{G} F_1(A_1), F_2(A_2), \dots, F_n(A_n) (E)$$

- $E$  is any relational-algebra expression.
- $G_1, G_2, \dots, G_n$  is a list of attributes on which to group (can be empty).
- Each  $F_i$  is an aggregate function.
- Each  $A_i$  is an attribute name.

# Aggregate Function – Example

■ Relation  $r$ :

$A$	$B$	$C$
$\alpha$	$\alpha$	7
$\alpha$	$\beta$	7
$\beta$	$\beta$	3
$\beta$	$\beta$	10

$g_{sum(c)}(r)$

$sum-C$
27

■ Could also add  $min$ ,  $max$ , and other aggregates to the above expression.

$g_{sum(c), min(c), max(c)}(r)$

$sum-C$	$min-C$	$max-C$
27	3	10

## Grouping – Example

- Grouping is somewhat like sorting, although not identical.
- Relation *account* grouped by *branch-name*:

<i>account-number</i>	<i>branch-name</i>	<i>balance</i>
A-102	Perryridge	400
A-374	Perryridge	900
A-224	Brighton	175
A-161	Brighton	850
A-435	Brighton	400
A-201	Brighton	625
A-217	Redwood	750
A-215	Redwood	750
A-222	Redwood	700

# Aggregate Operation – Example

- Grouping and aggregate functions frequently occur together.
- A list of branch names and the sum of all their account balances:

*branch-name*  $\mathcal{G}$  *sum(balance)* (*account*)

<i>branch-name</i>	<i>balance</i>
Perryridge	1300
Brighton	2050
Redwood	2200



# Aggregate Operation: Grouping on Multiple Attributes

- Consider the following relational scheme:

*History* = (*Student-Name*, *Department*, *Course-Number*, *Grade*)

- Sample data:

<u>Student-Name</u>	<u>Department</u>	<u>Course-Number</u>	<u>Grade</u>
<i>Smith</i>	<i>CSE</i>	<i>1001</i>	<i>90</i>
<i>Jones</i>	<i>MTH</i>	<i>2030</i>	<i>82</i>
<i>Smith</i>	<i>MTH</i>	<i>1002</i>	<i>73</i>
<i>Brown</i>	<i>PSY</i>	<i>4210</i>	<i>86</i>
<i>Jones</i>	<i>CSE</i>	<i>2010</i>	<i>65</i>

:

# Aggregate Operation: Grouping on Multiple Attributes

- Consider the following query:

*“Construct a list of student names and, for each name, list the average course grade for each department in which the student has taken classes.”*

Smith	CSE	87
Smith	MTH	93
Jones	CHM	88
Jones	CSE	75
Brown	PSY	97
	:	

- Recalling the schema:

*History = (Student-Name, Department, Course-Number, Grade)*

- Answer:

*student-name, department  $g_{avg(grade)}$  (History)*

## Aggregate Operation: Grouping on Multiple Attributes

- Adding *count(Course-Number)* would tell how many courses the student had in each department. Similarly, *min* and *max* could be added.

*student-name, department*  $\mathcal{G}$  *avg(grade), count(Course-Number), min(Grade), max(Grade)(History)*

# Aggregate Operation: Grouping on Multiple Attributes

- Would the following two expressions give the same result?

*student-name, department*  $\mathcal{G}$  *avg(grade), count(Course-Number), min(Grade), max(Grade)*(History)

*department, student-name*  $\mathcal{G}$  *avg(grade), count(Course-Number), min(Grade), max(Grade)*(History)

# Aggregate Operation: Naming Attributes

- Note that the aggregated attributes do not have names?

$$g_{sum(c), min(c), max(c)}(r)$$

<i>sum-C</i>	<i>min-C</i>	<i>max-C</i>
27	3	10

# Aggregate Operation: Naming Attributes

- Note that the aggregated attributes do not have names?

$\mathcal{G}_{sum(c), min(c), max(c)}(r)$

?	?	?
27	3	10

- Aggregated attributes can be renamed in the aggregate operator:

*branch-name*  $\mathcal{G}_{sum(balance) \text{ as } sum\text{-}balance}(account)$

# Aggregate Functions and Null Values

- Null values are controversial.
- Various proposals exist in the research literature on whether null values should be allowed and, if so, how they should affect operations.  
(Example: <http://www.dbazine.com/ofinterest/oi-articles/pascal27/>)
- Null values can frequently be eliminated through normalization and decomposition.

# Aggregate Functions and Null Values

- How nulls are treated by relational operators:
  - For duplicate elimination and grouping, null is treated like any other value, i.e., two nulls are assumed to be the same.
  - Aggregate functions (except for *count*) simply ignore null values.
- The above rules are consistent with SQL.
- Note how the second rule can be misleading:
  - Is *avg(grade)* actually a class average?



# Null Values and Expression Evaluation

- Null values also affect how selection predicates are evaluated:
  - The result of any arithmetic expression involving *null* is *null*.
  - Comparisons with *null* returns the special truth value *unknown*.
  - Value of a predicate is treated as *false* if it evaluates to *unknown*.

$$\sigma_{balance*100 > 500} (account)$$

- For more complex predicates, the following three-valued logic is used:

- OR:

(unknown or true)	= true
(unknown or false)	= unknown
(unknown or unknown)	= unknown
- AND:

(true and unknown)	= unknown
(false and unknown)	= false
(unknown and unknown)	= unknown
- NOT:

(not unknown)	= unknown
---------------	-----------

$$\sigma_{(balance*100 > 500) \text{ and } (branch\text{-}name = \text{"Perryridge"})} (account)$$

## Schema

branch (branch-name, branch-city, assets)  
customer (customer-name, customer-street, customer-city)  
account (account-number, branch-name, balance)  
loan (loan-number, branch-name, amount)  
depositor (customer-name, account-number)  
borrower (customer-name, loan-number)

# Null Values and Expression Evaluation, Cont.

- Why doesn't a comparison with *null* simply result in *false*?
- If *false* was used instead of *unknown*, then:

$\text{not } (A < 5)$

would not be equivalent to:

$A \geq 5$

Why would this be a problem?

- How does a comparison with *null* resulting in *unknown* help?

# Modification of the Database

- The database contents can be modified with operations:

- Deletion
- Insertion
- Updating

- These operations can all be expressed using the assignment operator.

- Some can be expressed other ways too.

# Deletion

- A deletion is expressed in relational algebra by:

$$r \leftarrow r - E$$

where  $r$  is a relation and  $E$  is a relational algebra query.

- The deletion of a single tuple is expressed by letting  $E$  be a constant relation containing one tuple.
- Only whole tuples can be deleted, not specific attribute values.

# Deletion Examples

- Forget referential integrity for the moment...

*“Delete all account records with a branch name equal to Perryridge.”*

$account \leftarrow account - \sigma_{branch-name = \text{“Perryridge”}}(account)$

*“Delete all loan records with amount in the range of 0 to 50.”*

$loan \leftarrow loan - \sigma_{amount \geq 0 \text{ and } amount \leq 50}(loan)$

## Schema

*branch* (branch-name, branch-city, assets)

*customer* (customer-name, customer-street, customer-city)

*account* (account-number, branch-name, balance)

*loan* (loan-number, branch-name, amount)

*depositor* (customer-name, account-number)

*borrower* (customer-name, loan-number)

# Deletion Examples

- Now suppose we want to maintain proper referential integrity...

*“Delete all accounts at branches located in Needham”*

- **Version #1:**

$$r_1 \leftarrow \sigma_{\text{branch-city} = \text{“Needham”}} (\text{account} \bowtie \text{branch})$$
$$r_2 \leftarrow \Pi_{\text{account-number}, \text{branch-name}, \text{balance}} (r_1)$$
$$r_3 \leftarrow \Pi_{\text{customer-name}, \text{account-number}} (\text{depositor} \bowtie r_2)$$
$$\text{account} \leftarrow \text{account} - r_2$$
$$\text{depositor} \leftarrow \text{depositor} - r_3$$

## Schema

*branch* (branch-name, branch-city, assets)

*customer* (customer-name, customer-street, customer-city)

*account* (account-number, branch-name, balance)

*loan* (loan-number, branch-name, amount)

*depositor* (customer-name, account-number)

*borrower* (customer-name, loan-number)

# Alternative Versions

## ■ Version #2:

$r_1 \leftarrow \Pi_{\text{branch-name}} (\sigma_{\text{branch-city} = \text{"Needham"}} (\text{branch}))$   
 $r_2 \leftarrow \Pi_{\text{account-number}} (\Pi_{\text{account-number}, \text{branch-name}} (\text{account}) \bowtie r_1)$   
 $\text{account} \leftarrow \text{account} - (\text{account} \bowtie r_2)$   
 $\text{depositor} \leftarrow \text{depositor} - (\text{depositor} \bowtie r_2)$

## ■ Version #3:

$r_1 \leftarrow (\sigma_{\text{branch-city} \neq \text{"Needham"}} (\text{depositor} \bowtie \text{account} \bowtie \text{branch}))$   
 $\text{account} \leftarrow \Pi_{\text{account-number}, \text{branch-name}, \text{balance}} (r_1)$   
 $\text{depositor} \leftarrow \Pi_{\text{customer-name}, \text{account-number}} (r_1)$

### Schema

*branch* (branch-name, branch-city, assets)  
*customer* (customer-name, customer-street, customer-city)  
*account* (account-number, branch-name, balance)  
*loan* (loan-number, branch-name, amount)  
*depositor* (customer-name, account-number)  
*borrower* (customer-name, loan-number)

# Alternative Versions

## ■ Version #4:

$r_1 \leftarrow \text{account} \bowtie \sigma_{\text{branch-city} \triangleleft \text{"Needham"}}(\text{branch})$

$\text{account} \leftarrow \Pi_{\text{account-number}, \text{branch-name}, \text{balance}}(r_1)$

$\text{depositor} \leftarrow \Pi_{\text{customer-name}, \text{account-number}}(\text{depositor} \bowtie r_1)$

## ■ Which version is preferable?

- Note that the last two do not fit the authors' pattern for deletion, i.e., as a set-difference.

### Schema

*branch* (branch-name, branch-city, assets)

*customer* (customer-name, customer-street, customer-city)

*account* (account-number, branch-name, balance)

*loan* (loan-number, branch-name, amount)

*depositor* (customer-name, account-number)

*borrower* (customer-name, loan-number)



# Insertion

- In relational algebra, an insertion is expressed by:

$$r \leftarrow r \cup E$$

where  $r$  is a relation and  $E$  is a relational algebra expression.

- The insertion of a single tuple is expressed by letting  $E$  be a constant relation containing one tuple.

# Insertion Examples

- Insert information in the database specifying that Smith has \$1200 in account A-973 at the Perryridge branch.

$$\begin{aligned} \text{account} &\leftarrow \text{account} \cup \{(A-973, \text{"Perryridge"}, 1200)\} \\ \text{depositor} &\leftarrow \text{depositor} \cup \{(\text{"Smith"}, A-973)\} \end{aligned}$$

- Provide, as a gift, a \$200 savings account for all loan customers at the Perryridge branch. Let the loan number serve as the account number for the new savings account.

$$\begin{aligned} r_1 &\leftarrow (\sigma_{\text{branch-name} = \text{"Perryridge"}}(\text{borrower} \bowtie \text{loan})) \\ \text{account} &\leftarrow \text{account} \cup \Pi_{\text{loan-number}, \text{branch-name}, 200}(r_1) \\ \text{depositor} &\leftarrow \text{depositor} \cup \Pi_{\text{customer-name}, \text{loan-number}}(r_1) \end{aligned}$$

## Schema

*branch* (branch-name, branch-city, assets)  
*customer* (customer-name, customer-street, customer-city)  
*account* (account-number, branch-name, balance)  
*loan* (loan-number, branch-name, amount)  
*depositor* (customer-name, account-number)  
*borrower* (customer-name, loan-number)

# Updating

- Generalized projection is used to change one or more values in a tuple.

$$r \leftarrow \Pi_{F_1, F_2, \dots, F_l}(r)$$

- Each  $F_i$  is either:
  - The  $i$ th attribute of  $r$ , if the  $i$ th attribute is not updated, or,
  - An expression, involving only constants and attributes of  $r$ , which gives a new value for an attribute, when that attribute is to be updated.

# Update Examples

- Make interest payments by increasing all balances by 5 percent.

$account \leftarrow \Pi_{AN, BN, BAL * 1.05} (account)$

where  $AN$ ,  $BN$  and  $BAL$  stand for *account-number*, *branch-name* and *balance*, respectively.

- Pay 6 percent interest to all accounts with balances over \$10,000 and pay 5 percent interest to all others.

$account \leftarrow \Pi_{AN, BN, BAL * 1.06} (\sigma_{BAL > 10000} (account))$   
 $\cup \Pi_{AN, BN, BAL * 1.05} (\sigma_{BAL \leq 10000} (account))$

## Schema

*branch* (*branch-name*, *branch-city*, *assets*)  
*customer* (*customer-name*, *customer-street*, *customer-city*)  
*account* (*account-number*, *branch-name*, *balance*)  
*loan* (*loan-number*, *branch-name*, *amount*)  
*depositor* (*customer-name*, *account-number*)  
*borrower* (*customer-name*, *loan-number*)

# Views

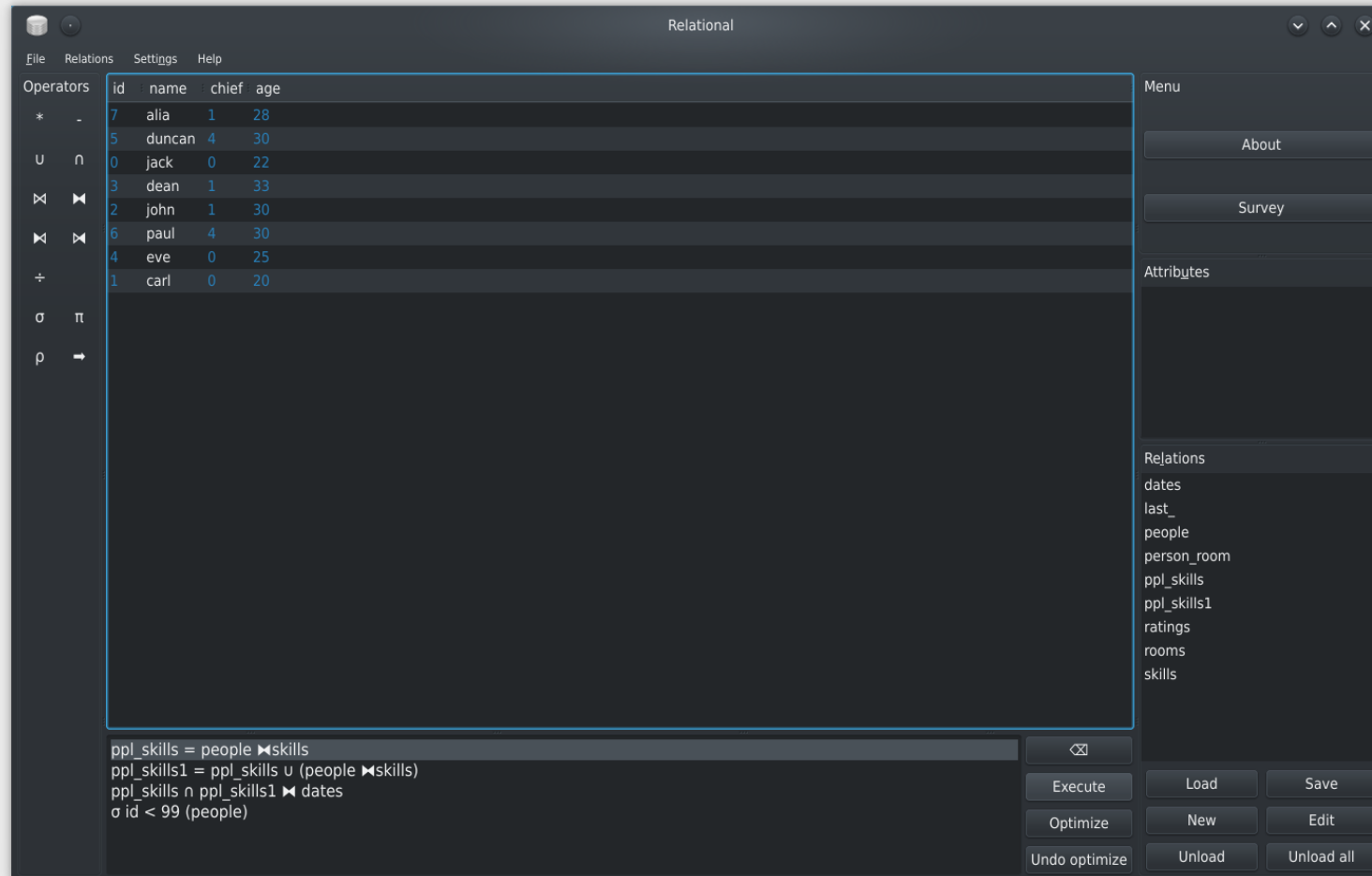
---

- Views are very important, but we will not consider them until chapter 3.

# TOOLS FOR TESTING RELATIONAL ALGEBRA EXPRESSIONS

# The Relational Tool

■ Locate it here: <https://ltworf.github.io/relational/>



The screenshot shows the Relational database tool interface. The main window displays a table with the following data:

id	name	chief	age
7	alia	1	28
5	duncan	4	30
0	jack	0	22
3	dean	1	33
2	john	1	30
6	paul	4	30
4	eve	0	25
1	carl	0	20

On the left, there is a sidebar with operators: \*, -, u, n, join, join, +, sigma, pi, rho, and an arrow. On the right, there is a sidebar with a menu (About, Survey), attributes, and a list of relations: dates, last\_, people, person\_room, ppl\_skills, ppl\_skills1, ratings, rooms, and skills. At the bottom, there is a query editor with the following queries:

```
ppl_skills = people join skills
ppl_skills1 = ppl_skills u (people join skills)
ppl_skills n ppl_skills1 join dates
sigma id < 99 (people)
```

Below the query editor are buttons for Execute, Load, Save, Optimize, New, Edit, Undo optimize, Unload, and Unload all.

# The Relational Tool – Syntax (Cont'd)

## ■ Example Syntax

Symbol	Name	Example	Notes
*	product	$A * B$	
-	difference	$A - B$	
$\cup$	union	$A \cup B$	
$\cap$	intersection	$A \cap B$	
$\div$	division	$A \div B$	
$\bowtie$	join	$A \bowtie B$	
$\ltimes$	left outer join	$A \ltimes B$	All outer joins use a python None value when they have no value to place.
$\rtimes$	right outer join	$A \rtimes B$	
$\Join$	full outer join	$A \Join B$	



# The Relational Tool - Syntax

## ■ Example Syntax

Symbol	Name	Example	Note
$\sigma$	selection	$\sigma \text{ id} == \text{index or rank} > 3 \text{ (A)}$	Expression must be written in python. The variables have the names of the fields in the relation. If the expression contains parenthesis, it <b>must</b> be surrounded by another pair of parenthesis.
$\pi$	projection	$\pi \text{ name, age (A)}$	
$\rho$	rename	$\rho \text{ old\_name} \rightarrow \text{new\_name, age} \rightarrow \text{old (A)}$	

Read more about the syntax here: [https://ltworf.github.io/relational/allowed\\_expressions.html](https://ltworf.github.io/relational/allowed_expressions.html)

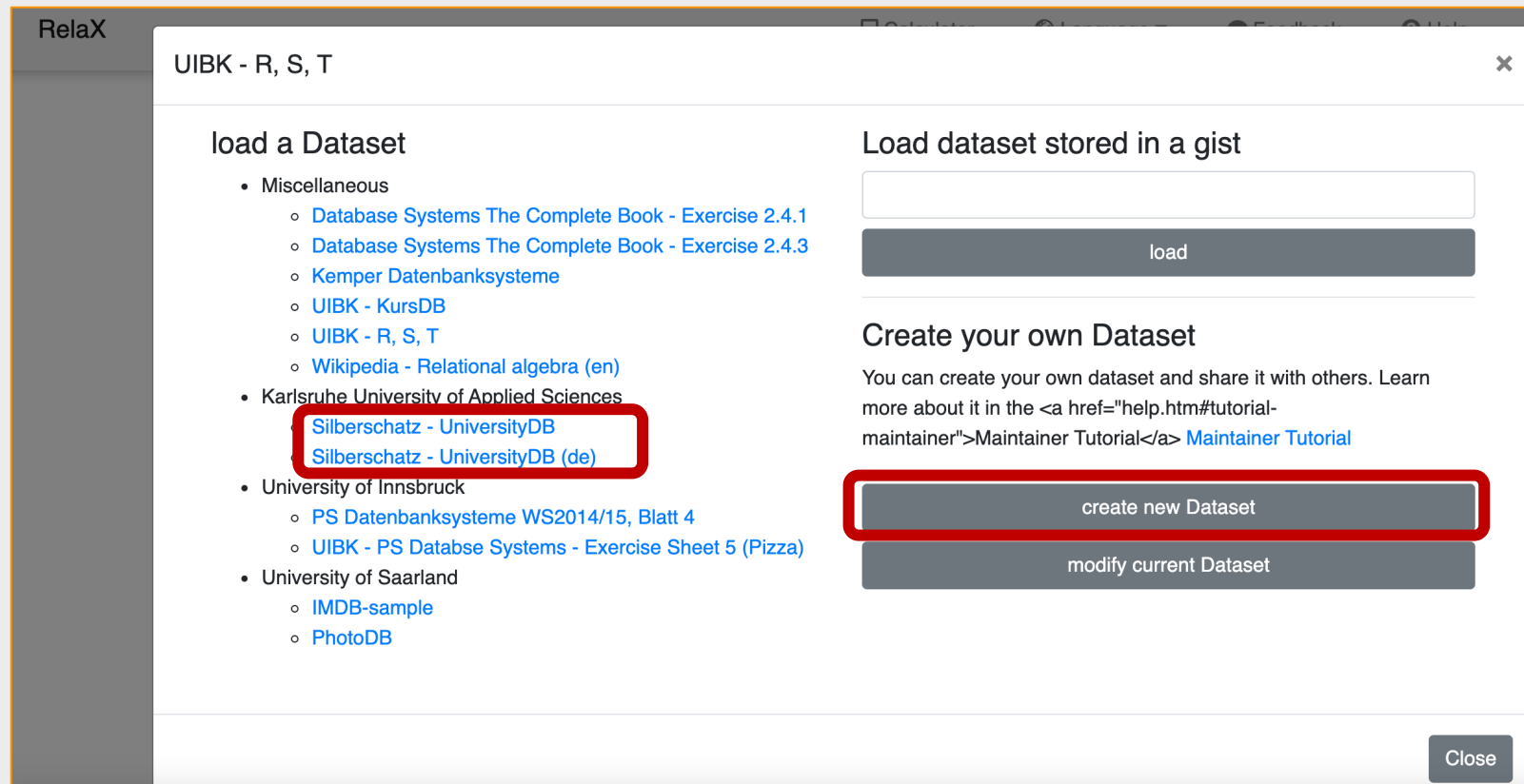
# The Relax Tool

■ Locate it here: <https://dbis-uibk.github.io/relax/landing>

The screenshot displays the Relax tool interface. At the top, the title "RelaX" is on the left, and navigation links for "Calculator", "Language", "Feedback", "Help", and "Imprint" are on the right. The main area is divided into a left sidebar and a central editor. The sidebar, titled "Select DB (UIBK - R, ...)", lists three tables: R (with attributes a: number, b: string, c: string), S (with attributes b: string, d: number), and T (with attributes b: string, d: number). The central editor has tabs for "Relational Algebra", "SQL", and "Group Editor". The "Relational Algebra" tab is active, showing a toolbar with various operators (π, σ, ρ, ←, →, τ, γ, ^, v, ¬, =, ≠, ≥, ≤, ∩, ∪, ÷, -, ×, ⋈, ⋉, ⋊) and a text area with the placeholder "1 | your query goes here ...". Below the text area, a section titled "keyboard shortcuts:" lists: "execute statement: [CTRL]+[RETURN]", "execute selection: [CTRL]+[SHIFT]+[RETURN]", and "autocomplete: [CTRL]+[SPACE]". At the bottom of the editor, there is a blue "execute query" button, a "download" link with a download icon, and a "history" link with a clock icon.

# The RelaX Tool – Choosing a Database (schema)

- Notice the UniversityDB from our textbook in the list.
- You may also add your own relation



# The RelaX Tool – Loading a dataset (relation) - 1

- Click Select DB to select a database or dataset

The screenshot shows the RelaX web application interface. At the top, there is a navigation bar with the title 'RelaX' and several utility links: 'Calculator', 'Language' (with a dropdown arrow), 'Feedback', 'Help' (with a question mark icon), and 'Imprint' (with a person icon). Below the navigation bar, the main content area is divided into three sections. On the left, there is a 'Select DB (UIBK - R,...)' button, which is highlighted with a red rectangular box. Below this button, there are three sections labeled 'R', 'S', and 'T'. Each section lists attributes and their data types: 'R' has 'a number', 'b string', and 'c string'; 'S' has 'b string' and 'd number'; 'T' has 'b string' and 'd number'. In the center, there is a tabbed interface with three tabs: 'Relational Algebra', 'SQL', and 'Group Editor'. The 'Relational Algebra' tab is active, showing a toolbar with various relational algebra symbols (π, σ, ρ, ←, →, τ, γ, ^, v, ¬, =, ≠, ≥, ≤, ∩, ∪, ÷, -, ×, ⋈, ⋉) and a text area containing the placeholder text '1 your query goes here ...'. Below the text area, there is a section titled 'keyboard shortcuts:' with the following information: 'execute statement: [CTRL]+[RETURN]', 'execute selection: [CTRL]+[SHIFT]+[RETURN]', and 'autocomplete: [CTRL]+[SPACE]'. At the bottom of the interface, there is a blue button labeled 'execute query' with a play icon, and two links: 'download' (with a download icon) and 'history' (with a circular arrow icon).

# The RelaX Tool – Loading a dataset (relation) - 2

■ Click “Add new Relation” to add your own dataset

RelaX

Calculator Language Feedback Help Imprint

Select DB (UIBK - R,...)

Relational Algebra SQL Group Editor

**add new relation**

open relation editor

R

- a number
- b string
- c string

S

- b string
- d number

T

- b string
- d number

```
-- this is an example
group: nameOfTheNewGroup
3
4{
a: string, b: number
example, 42
}
7
```

preview

download

# The RelaX Tool – Loading a dataset (relation) - 3



- Paste your data in the cells provided or upload a CSV



RelaX

Relation Editor

works

	Name	employee_na	company_na	salary
Type	string	string	string	number
1	Mac Eger		Google, Inc	75000
2	Bernice Crosbie		Google, Inc	99000
3	Rolanda Pendergast		Howard Johnson's	65000
4	Ambrose Gridley		Comcast	86000
5	Federico Catto		Whole Foods	54789
6	Christinia Harshbarger		Walgreens	47000
7	Marcelo Farrel		First Bank Corporation	75400
8	Caprice Kiker		Google, Inc	120851
9	Pia Balzer		First Bank Corporation	86000
10	Freeman Scranton		Whole Foods	167000
11	Ahmad Locicero		Walgreens	73008
12	Rebecca Sandy		Google, Inc	86000
13	Kate Schutz		Comcast	65320

 Download CSV  Upload CSV

# The RelaX Tool – Loading a dataset (relation) - 4

■ You may also load a relation using the format below

-- this is the works relation

group: works

works = {

employee\_name:string, company\_name:string, salary:number

"Mac Eger", "Google, Inc", 75000

"Bernice Crosbie", "Google, Inc", 99000

"Rolanda Pendergast", "Howard Johnson's", 65000

"}

# The RelaX Tool – Loading a dataset (relation) - 5

- After loading data, click “Use Group in Editor” to execute a query

RelaX

Calculator Language Feedback

preview download

works **use Group in editor**

- works

works.employee_name	works.company_name	works.salary
'Mac Eger'	'Google, Inc'	75000
'Bernice Crosbie'	'Google, Inc'	99000
'Rolanda Pendergast'	'Howard Johnson's'	65000
'Ambrose Gridley'	'Comcast'	86000
'Federico Catto'	'Whole Foods'	54789
'Christinia Harshbarger'	'Walgreens'	47000
'Marcelo Farrel'	'First Bank Corporation'	75400
'Caprice Kiker'	'Google, Inc'	120851
'Pia Balzer'	'First Bank Corporation'	86000
'Freeman Scranton'	'Whole Foods'	167000
'Ahmad Locicero'	'Walgreens'	73008
'Rebecca Sandy'	'Google, Inc'	86000
'Kate Schutz'	'Comcast'	65320
'Ilana Rodriguez'	'Whole Foods'	62000
'Wendell Nodine'	'Tmobile'	56989
'Ayesha Markowitz'	'Walgreens'	78231
'Vincent Sakai'	'First Bank Corporation'	200008
'Leeanne Caffey'	'Small Bank Corporation'	65000
'Jamel Pasha'	'Small Bank Corporation'	72000
'Kamala Harrow'	'Small Bank Corporation'	123000



## The RelaX Tool – Loading a dataset (relation) - 6

## ■ Entering a Query in RelaX

RelaX

Calculator

Language

Feedback

Help

Imprint

Select DB (works)

works

employee\_name string

company\_name string

salary number

Relational Algebra

SQL

Group Editor

$\pi$   $\sigma$   $\rho$   $\leftarrow$   $\rightarrow$   $\tau$   $\gamma$   $\wedge$   $\vee$   $\neg$   $=$   $\neq$   $\geq$   $\leq$   $\cap$   $\cup$   $\div$   $-$   $\times$   $\bowtie$   $\ltimes$   $\ltimes$   $\ltimes$   $\ltimes$   $\ltimes$   $\ltimes$   $\triangleright$   $=$   $--$   $/*$

{}

1  $\pi$  employee\_name (  $\sigma$  company\_name = 'First Bank Corporation' ( works ) )

execute query

download

history

# The RelaX Tool – Loading a dataset (relation) - 7

## Output from RelaX

RelaX

Calculator Language Feedback

execute query download history

```
graph TD; A["Π employee_name  
3 rows"] --- B["σ company_name = 'First Bank Corporation'  
3 rows"]; B --- C["works  
20 rows"]
```

$\Pi_{\text{employee\_name}} ( \sigma_{\text{company\_name} = \text{'First Bank Corporation'}} ( \text{works} ) )$

works.employee_name
'Marcelo Farrel'
'Pia Balzer'
'Vincent Sakai'

< 1 >

# Relational Algebra Practice Problems

Consider the relational database for various employees (EmployeeDB):

- employee (employee-name, street, city)
  - works (employee-name, company-name, salary)
  - company (company-name, city)
  - manages (employee-name, manager-name)
- 
- a) Find the names of all employees who work for First Bank Corporation.
  - b) Find the names and cities of residence of all employees who work for First Bank Corporation.
  - c) Find the names, street address, and cities of residence of all employees who work for First Bank Corporation and earn more than \$10,000 per annum.
  - d) Find the names of all employees in this database who live in the same city as the company for which they work.
  - e) Find the names of all employees who live in the same city and on the same street as do their managers.
  - f) Find the names of all employees in this database who do not work for First Bank Corporation (assume that all employees work for exactly one company.)
  - g) Find the names of all employees who earn more than every employee of Small Bank Corporation.
  - h) Assume the companies may be located in several cities. Find all companies located in every city in which Small Bank Corporation is located.

Hey, Mr. DB Here.  
Wanna learn  
more? Let's  
practice.

