

# Module 2: The Relational Model

## Spring 2022

Cheriton School of Computer Science

CS 348: Intro to Database Management

# Reading Assignments and References

To be read during the Week of May 9–13:

- ▶ Chapter 2 of course textbook.<sup>1</sup> (Material in Sections 2.3 and 2.6 will be covered in later modules.)
- ▶ Section 27.2 of Chapter 27 of course textbook, available online at [db-book.com](http://db-book.com).

## References

1. Abiteboul, Hull and Vianu, *Foundations of Databases*. A book available online at <http://webdam.inria.fr/Alice/>.
2. *Backus-Naur Form*, wiki page.

---

<sup>1</sup>Silberschatz, Korth and Sudarshan, *Database Systems Concepts*, 7<sup>th</sup> edition

# Outline

Unit 1: **Signatures and the Relational Calculus**

Unit 2: Integrity Constraints

Unit 3: Safety and Finiteness

Unit 4: Summary

## A Basic Syntax for Asking Questions and for Answers

To begin with, assume ...

- ▶ Set comprehension syntax for queries:

$$\{\langle \textit{answer} \rangle \mid \langle \textit{condition} \rangle\}.$$

- ▶ Syntax for each  $\langle \textit{answer} \rangle$  is a  $k$ -tuple of *variables*:

$$(x_1, \dots, x_k).$$

- ▶ Answers to a query:

*all  $k$ -tuples  $(c_1, \dots, c_k)$  of constants denoting values for each variable  $x_i$  that satisfy  $\langle \textit{condition} \rangle$ .*

## Asking Questions about Natural Numbers

*What are all pairs of natural numbers that add to 5?*

Question:  $\{(x, y) \mid x + y = 5\}$  or  $\{(x, y) \mid \text{PLUS}(x, y, 5)\}^\dagger$

Answers:  $\{(0, 5), (1, 4), (2, 3), (3, 2), (4, 1), (5, 0)\}$

**Why?** Because  $(0, 5, 5)$ , etc., appear in table PLUS!

*What are all pairs of numbers that add to the same number they subtract to, where  $x + y = x - y$ ?*

Question:  $\{(x, y) \mid \exists z. \text{PLUS}(x, y, z) \wedge \text{PLUS}(z, y, x)\}$

Answers:  $\{(0, 0), (1, 0), \dots\}$  Is  $(5, 5)$  also an answer?

... depends on the content (instance) of table PLUS!

*What is the neutral element of addition?*

Question:  $\{(x) \mid \text{PLUS}(x, x, x)\}$

Answers:  $\{(0)\}$

<sup>†</sup> A relational form for basic conditions.

Table PLUS

A1	A2	R
0	0	0
0	1	1
0	2	2
⋮	⋮	⋮
1	0	1
1	1	2
⋮	⋮	⋮
⋮	⋮	⋮
2	0	2
2	1	3
⋮	⋮	⋮
⋮	⋮	⋮

## Asking Questions about Employees

*Who are all the employees and their departments who work for Bob?*

Question:  $\{(x, y) \mid \text{EMP}(x, y, \text{Bob})\}$

Answers:  $\{(\text{Sue}, \text{CS}), (\text{Bob}, \text{CO})\}$

**Why?** ... because  $(\text{Sue}, \text{CS}, \text{Bob})$ , etc., appear in EMP!

*Who are pairs of employees working for the same boss?*

Q:  $\{(x_1, x_2) \mid \exists y_1, y_2, z. \text{EMP}(x_1, y_1, z) \wedge \text{EMP}(x_2, y_2, z)\}$

A:  $\{(\text{Sue}, \text{Bob}), (\text{Fred}, \text{John}), (\text{Jim}, \text{Eve})\}$  ← Is that all?

*Who are the employees who are their own bosses?*

Q:  $\{(x) \mid \exists y. \text{EMP}(x, y, x)\}$

A:  $\{(\text{Sue}), (\text{Bob})\}$

Table EMP

name	dept	boss
Sue	CS	Bob
Bob	CO	Bob
Fred	PM	Mark
John	PM	Mark
Jim	CS	Fred
Eve	CS	Fred
Sue	PM	Sue

## Relational Databases and the Relational Calculus

Based on *first order predicate logic* (FOL) and *Tarskian semantics*.

Recall example RM database using a common visualization:

AUTHOR	
aid	name
1	Sue
2	John

WROTE	
author	publication
1	1
1	4
2	2
1	2

PUBLICATION	
pubid	title
1	Mathematical Logic
3	Trans. on Databases
2	Principles of DB Systems
4	Query Languages

### Idea

All information is organized in a finite number of relations called *tables*.

Features:

- ▶ simple and clean data model accommodating data independence,
- ▶ declarative DML based on *well-formed formulas* in FOL, and
- ▶ *integrity constraints* also via well-formed formulas.

# Relational Databases

## Components:

**Universe** ▶ a set of **values**  $\mathbf{D}$  (*domain*) with *equality* ( $\approx$ ), and with *constants* for each value.

**Relation** (also called a **table**)

- ▶ **intension**: a relation name (*predicate name*)  $R$ , and *arity*  $k$  of  $R$  (the number of columns), written  $R/k$ , and
- ▶ **extension**: a set of  $k$ -tuples (*interpretation*)  $\mathbf{R} \subseteq \mathbf{D}^k$ .

**Database** ▶ **signature** (metadata): finite set  $\rho$  of predicate names  $R_i$ ; and

- ▶ **instance** (data, *structure*): an extension  $\mathbf{R}_i$  for each  $R_i$ .

## Notation

Signature:  $\rho = (R_1/k_1, \dots, R_n/k_n)$

Instance:  $\mathbf{DB} = (\mathbf{D}, \approx, \mathbf{R}_1, \dots, \mathbf{R}_n)$



## Examples of Relational Databases

- ▶ The integers, with addition *and multiplication*:

(signature)  $\rho = (\text{PLUS}/3, \text{TIMES}/3)$

(data) **DB** = ( $\mathbb{Z}$ ,  $\approx$ , **PLUS**, **TIMES**)

- ▶ The employee database:

(signature)  $\rho = (\text{EMP}/3)$

(data) **DB** = ( $\text{STR}$ ,  $\approx$ , **EMP**)

- ▶ The simple bibliography database:

(signature)  $\rho = (\text{AUTHOR}/2, \text{WROTE}/2, \text{PUBLICATION}/2)$

(data) **DB** = ( $\text{STR} \uplus \mathbb{Z}$ ,  $\approx$ , **AUTHOR**, **WROTE**, **PUBLICATION**)

## Bibliography Relational Database, Version 2

(signature)  $\rho = ($

```
AUTHOR (aid, name),  
WROTE (author, publication),  
PUBLICATION (pubid, title),  
BOOK (pubid, publisher, year),  
JOURNAL-OR-PROCEEDINGS (pubid),  
JOURNAL (pubid, volume, no, year),  
PROCEEDINGS (pubid, year),  
ARTICLE (pubid, appears-in, startpage, endpage)
```

)

Arity is indicated by a sequence of *identifiers*, called **attributes**:

- ▶ Help with understanding semantics; and
- ▶ Used in some DMLs, such as some *relational algebras* and SQL.

## Bibliography Relational Database, Version 2 (cont'd)

(data) **DB** = (STR  $\uplus$  Z,  $\approx$ ,

**AUTHOR** = { (1, Sue), (2, John) },

**WROTE** = { (1, 1), (1, 4), (1, 2), (2, 2) },

**PUBLICATION** = { (1, Mathematical Logic),  
(3, Trans. on Databases),  
(2, Principles of DB Systems),  
(4, Query Languages) },

**BOOK** = { (1, AMS, 1990) },

**JOURNAL-OR-PROCEEDINGS** = { (2), (3) },

**JOURNAL** = { (3, 35, 1, 1990) },

**PROCEEDINGS** = { (2, 1995) },

**ARTICLE** = { (4, 2, 30, 41) }

)

## Simple (Atomic) “Truth”

### Idea

Relationships between values (tuples) that are *present* in an instance are *true*; relationships *absent* are *false*.

In the sample *bibliography* database instance:

- ▶ “John” is the name of an author with id “2” since:  $(2, \text{John}) \in \mathbf{AUTHOR}$ ;
- ▶ “Mathematical Logic” is the title of a publication since:  
 $(1, \text{Mathematical Logic}) \in \mathbf{PUBLICATION}$ ;
- ▶ Moreover, it is a book published by “AMS” in “1990” since:  
 $(1, \text{AMS}, 1990) \in \mathbf{BOOK}$ ;
- ▶ John wrote “Principles of DB Systems” since:  $(2, 2) \in \mathbf{WROTE}$ ;
- ▶ John has **NOT** written “Trans. on Databases” since:  $(2, 3) \notin \mathbf{WROTE}$ ;
- ▶ etc.

# Query Conditions

## Idea

Use **variables** and **valuations** to generalize conditions.

Example:  $\text{AUTHOR}(x, y)$  will be true of any *valuation*  $\{x \mapsto v_1, y \mapsto v_2, \dots\}$  exactly when the 2-tuple of values  $(v_1, v_2)$  occurs in **AUTHOR**.

## Valuation

A valuation is a function  $\theta$  that maps *variable names* to values in the universe:

$$\theta : \{x_1, x_2, \dots\} \rightarrow \mathbf{D}.$$

To denote a modification to  $\theta$  in which variable  $x$  is instead mapped to value  $v$ , one writes:

$$\theta[x \mapsto v].$$

## Query Conditions (cont')

### Idea

Allow more complex conditions to be built from simpler conditions with ...

### Logical connectives:

Conjunction (and):  $\text{AUTHOR}(x, y) \wedge \text{WROTE}(x, z)$

Disjunction (or):  $\text{AUTHOR}(x, y) \vee \text{PUBLICATION}(x, y)$

Negation (not):  $\neg \text{AUTHOR}(x, y)$

### Quantifiers:

Existential (there is...):  $\exists x. \text{author}(x, y)$

### Examples:

- ▶  $\exists z. \text{PLUS}(x, y, z) \wedge \text{PLUS}(z, y, x)$ , or
- ▶  $\exists y_1, y_2, z. \text{EMP}(x_1, y_1, z) \wedge \text{EMP}(x_2, y_2, z)$ .

Summarizing, allow conditions to be *well-formed formulas* (wffs) in the language of FOL.

# Relational Calculus

## Conditions

Given a database signature  $\rho = (R_1/k_1, \dots, R_n/k_n)$ , a set of variable names  $\{x_1, x_2, \dots\}$  and a set of constants  $\{c_1, c_2, \dots\}$ , *conditions* are *formulas* defined by the grammar:

$$\varphi ::= \underbrace{R_i(x_{i,1}, \dots, x_{i,k_i}) \mid x_i = x_j \mid x_i = c_j \mid \varphi_1 \wedge \varphi_2 \mid \exists x_i. \varphi_1 \mid \varphi_1 \vee \varphi_2 \mid \neg \varphi_1}_{\text{conjunctive formulas}}$$
$$\underbrace{\hspace{15em}}_{\text{positive formulas}}$$
$$\underbrace{\hspace{20em}}_{\text{first-order formulas}}$$

A condition is a *sentence* when it has no **free variables**.

The meta-language used to define the grammar is Backus-Naur form. (See wiki for a good overview.)

## Relational Calculus (cont')

### Free Variables

The *free variables* of a formula  $\varphi$ , written  $Fv(\varphi)$ , are defined as follows:

$$Fv(R(x_{i_1}, \dots, x_{i_k})) = \{x_{i_1}, \dots, x_{i_k}\};$$

$$Fv(x_i = x_j) = \{x_i, x_j\};$$

$$Fv(x_i = c_j) = \{x_i\};$$

$$Fv(\varphi \wedge \psi) = Fv(\varphi) \cup Fv(\psi);$$

$$Fv(\exists x_i. \varphi) = Fv(\varphi) - \{x_i\};$$

$$Fv(\varphi \vee \psi) = Fv(\varphi) \cup Fv(\psi); \text{ and}$$

$$Fv(\neg \varphi) = Fv(\varphi).$$



# Semantics for Conditions

## When a Condition is True (Tarski)

The *truth* of a formula  $\varphi$  over a signature  $\rho = (R_1/k_1, \dots, R_n/k_n)$  is defined with respect to

1. a **database instance**  $\mathbf{DB} = (\mathbf{D}, \approx, \mathbf{R}_1, \dots, \mathbf{R}_n)$ , and
2. a **valuation**  $\theta : \{x_1, x_2, \dots\} \rightarrow \mathbf{D}$

as follows:

$\mathbf{DB}, \theta \models R_i(x_{i,1}, \dots, x_{i,k_i})$	if $(\theta(x_{i,1}), \dots, \theta(x_{i,k_i})) \in \mathbf{R}_i$ ;
$\mathbf{DB}, \theta \models x_i = x_j$	if $\theta(x_i) \approx \theta(x_j)$ ;
$\mathbf{DB}, \theta \models x_i = c_j$	if $\theta(x_i) \approx c_j$ ;
$\mathbf{DB}, \theta \models \varphi \wedge \psi$	if $\mathbf{DB}, \theta \models \varphi$ and $\mathbf{DB}, \theta \models \psi$ ;
$\mathbf{DB}, \theta \models \exists x_i. \varphi$	if $\mathbf{DB}, \theta[x_i \mapsto v] \models \varphi$ , for some $v \in \mathbf{D}$ ;
$\mathbf{DB}, \theta \models \varphi \vee \psi$	if $\mathbf{DB}, \theta \models \varphi$ or $\mathbf{DB}, \theta \models \psi$ ; and
$\mathbf{DB}, \theta \models \neg \varphi$	if $\mathbf{DB}, \theta \not\models \varphi$ .

# Equivalences and Syntactic Sugar

## Boolean Equivalences

- ▶  $\neg(\neg\varphi_1) \equiv \varphi_1$
- ▶  $\varphi_1 \vee \varphi_2 \equiv \neg(\neg\varphi_1 \wedge \neg\varphi_2)$
- ▶  $\varphi_1 \rightarrow \varphi_2 \equiv \neg\varphi_1 \vee \varphi_2$
- ▶  $\varphi_1 \leftrightarrow \varphi_2 \equiv (\varphi_1 \rightarrow \varphi_2) \wedge (\varphi_2 \rightarrow \varphi_1)$
- ▶ ...

## First-order Equivalences

- ▶  $\forall x.\varphi \equiv \neg\exists x.\neg\varphi$

## Additional Syntactic Sugar

- ▶  $R(\dots, c, \dots) \equiv \exists x.(R(\dots, x, \dots) \wedge x = c)$ , where  $x$  is *fresh*
- ▶  $\exists x_1, \dots, x_n.\varphi \equiv \exists x_1.\dots.\exists x_n.\varphi$
- ▶  $R(\dots, -, \dots) \equiv \exists x.R(\dots, x, \dots)$ , where  $x$  is *fresh*

## Relational Calculus (cont'd)

### Relational Calculus (RC) Query

A *query* in the relational calculus is a set comprehension of the form

$$\{(x_1, \dots, x_k) \mid \varphi\},$$

where  $\{x_1, \dots, x_k\} = \text{Fv}(\varphi)$  (are the free variables of  $\varphi$ ).

Also:

- ▶ a **conjunctive query** is where  $\varphi$  is a conjunctive formula, and
- ▶ a **positive query** is where  $\varphi$  is a positive formula.

### Query Answers

The *answers* to a query  $\{(x_1, \dots, x_k) \mid \varphi\}$  over **DB** is the **relation**

$$\{(\theta(x_1), \dots, \theta(x_k)) \mid \mathbf{DB}, \theta \models \varphi\}.$$

*Answers to queries*: valuations applied to tuples of variables that make the formula true with respect to a database.

## Example Justification of an Answer to an RC Query

*Who are pairs of employees working for the same boss?*

Q:  $\{(x_1, x_2) \mid \exists y_1, y_2, z. \text{EMP}(x_1, y_1, z) \wedge \text{EMP}(x_2, y_2, z)\}$

A:  $\{(Jim, Eve), \dots\}$

Because:

1. **DB**,  $\theta_1 (= \{x_1 \mapsto Jim, y_1 \mapsto CS, z \mapsto Fred, \dots\}) \models \text{EMP}(x_1, y_1, z)$
2. **DB**,  $\theta_2 (= \{x_2 \mapsto Eve, y_2 \mapsto CS, z \mapsto Fred, \dots\}) \models \text{EMP}(x_2, y_2, z)$
3. **DB**,  $\theta_3 (= \{x_1 \mapsto Jim, y_1 \mapsto CS, x_2 \mapsto Eve, y_2 \mapsto CS, z \mapsto Fred, \dots\})$   
 $\models \text{EMP}(x_1, y_1, z) \wedge \text{EMP}(x_2, y_2, z)$
4. **DB**,  $\theta_4 (= \{x_1 \mapsto Jim, x_2 \mapsto Eve, \dots\})$   
 $\models \exists y_1, y_2, z. \text{EMP}(x_1, y_1, z) \wedge \text{EMP}(x_2, y_2, z)^\dagger$
5.  $(\theta_4(x_1), \theta_4(x_2)) = (Jim, Eve)$

where  $\rho = (\text{EMP}/3)$ , and **DB** = (**STR**,  $\approx$ , **EMP**).

$^\dagger$  Check that  $\{x_1, x_2\} = \text{Fv}(\exists y_1, y_2, z. \text{EMP}(x_1, y_1, z) \wedge \text{EMP}(x_2, y_2, z))$ .

Table EMP

name	dept	boss
Sue	CS	Bob
Bob	CO	Bob
Fred	PM	Mark
John	PM	Mark
Jim	CS	Fred
Eve	CS	Fred
Sue	PM	Sue

## More Examples of RC Queries

Over signature  $\rho = (\text{EMP}(\text{name}, \text{dept}, \text{boss}))$ :

1. *Who are the bosses that manage at least two employees?*

$$\{(b) \mid \exists e_1, e_2. (\exists d_1. \text{EMP}(e_1, d_1, b)) \wedge (\exists d_2. \text{EMP}(e_2, d_2, b)) \wedge \neg(e_1 = e_2)\}$$

(or more simply, with the aid of some syntactic sugar)

$$\{(b) \mid \exists e_1, e_2. \text{EMP}(e_1, -, b) \wedge \text{EMP}(e_2, -, b) \wedge \neg(e_1 = e_2)\}$$

2. *Who are the bosses that do not manage more than two employees?*

$$\{(b) \mid \exists e_1. \text{EMP}(e_1, -, b) \\ \wedge \neg \exists e_2, e_3. \text{EMP}(e_2, -, b) \wedge \text{EMP}(e_3, -, b) \\ \wedge \neg(e_1 = e_2 \vee e_1 = e_3 \vee e_2 = e_3)\}$$

Choose variable names suggestive of what values or (indirectly) entities they refer to, e.g.:

- ▶ “ $e_1$ ” refers indirectly to an employee, and
- ▶ “ $b$ ” refers indirectly to a boss.

## Exercises

1. Over the PLUS-TIMES relational database, with  
signature  $\rho = (\text{PLUS}/3, \text{TIMES}/3)$ , and  
instance  $\mathbf{DB} = (\mathbb{N}, \approx, \mathbf{PLUS}, \mathbf{TIMES})$ :<sup>†</sup>
  - 1.1 *What are all composite numbers?*
  - 1.2 *What are all prime numbers?*
2. Over the bibliography relational database, 2nd version:
  - 2.1 *What are all publication titles?*
  - 2.2 *What are the publication titles that are journals or proceedings?*
  - 2.3 *What are the titles of all books?*
  - 2.4 *What are the publications without authors?*
  - 2.5 *What are all the ordered pairs of coauthor names?*
  - 2.6 *What are all publication titles written by a single author?*

<sup>†</sup> Much harder over the integers  $\mathbb{Z}$ .

# Outline

Unit 1: Signatures and the Relational Calculus

Unit 2: **Integrity Constraints**

Unit 3: Safety and Finiteness

Unit 4: Summary

## Asking Questions about Natural Numbers (revisited)

What is the neutral element of addition?

Question:  $\{(x) \mid \text{PLUS}(x, x, x)\}$

Answers:  $\{(0)\}$

But shouldn't the query really be

$$\{(x) \mid \forall y. \text{PLUS}(x, y, y) \wedge \text{PLUS}(y, x, y)\} \quad (*)$$

### Observation

(\*) is the same as

$$\{(x) \mid \forall y. \text{PLUS}(x, y, y)\} \quad (**)$$

because PLUS is **commutative**! And (\*\*) is the same as

$$\{(x) \mid \text{PLUS}(x, x, x)\}$$

because PLUS is **monotone**!

PLUS should satisfy **integrity constraints** that are the *laws of arithmetic* for natural numbers.

Table PLUS

A1	A2	R
0	0	0
0	1	1
0	2	2
⋮	⋮	⋮
1	0	1
1	1	2
⋮	⋮	⋮
2	0	2
2	1	3
⋮	⋮	⋮



## Integrity Constraints for Addition

Sentences that should always be *true* for any extension of PLUS over the domain of natural numbers:

- ▶ *Addition is commutative:*

$$\begin{aligned} &\forall x, y, z. \text{PLUS}(x, y, z) \rightarrow \text{PLUS}(y, x, z) \\ &\neg \exists x, y, z. \text{PLUS}(x, y, z) \wedge \neg \text{PLUS}(y, x, z) \end{aligned}$$

- ▶ *PLUS is a relational representation of a binary function:*

$$\begin{aligned} &\forall x, y, z_1, z_2. \text{PLUS}(x, y, z_1) \wedge \text{PLUS}(x, y, z_2) \rightarrow z_1 = z_2 \\ &\neg \exists x, y, z_1, z_2. \text{PLUS}(x, y, z_1) \wedge \text{PLUS}(x, y, z_2) \wedge \neg(z_1 = z_2) \end{aligned}$$

- ▶ *Addition is a total function:*

$$\begin{aligned} &\forall x, y. \exists z. \text{PLUS}(x, y, z) \\ &\neg \exists x, y. \neg \exists z. \text{PLUS}(x, y, z) \end{aligned}$$

- ▶ *Addition is monotone in both arguments (harder), etc., etc.*

## Integrity Constraints for Employees

Sentences that should always be *true* for any extension of table

EMP (name, dept, boss) :

- ▶ *Every boss is an employee:*

$$\begin{aligned} \forall e, d_1, b_1. \text{EMP}(e, d_1, b_1) &\rightarrow \exists d_2, b_2. \text{EMP}(b_1, d_2, b_2) \\ \forall b. \text{EMP}(-, -, b) &\rightarrow \text{EMP}(b, -, -)^{\dagger} \end{aligned}$$

- ▶ *Every boss manages a unique department:*

$$\begin{aligned} \forall e_1, e_2, d_1, d_2, b. \text{EMP}(e_1, d_1, b) \wedge \text{EMP}(e_2, d_2, b) &\rightarrow d_1 = d_2 \\ \forall d_1, d_2. (\exists b. \text{EMP}(-, d_1, b) \wedge \text{EMP}(-, d_2, b)) &\rightarrow d_1 = d_2 \end{aligned}$$

- ▶ *No boss has someone else as their boss:*

$$\begin{aligned} \forall e, b_1, b_2. \text{EMP}(e, -, b_1) \wedge \text{EMP}(b_1, -, b_2) &\rightarrow b_1 = b_2 \\ \forall b_1, b_2. (\text{EMP}(-, -, b_1) \wedge \text{EMP}(b_1, -, b_2)) &\rightarrow b_1 = b_2 \end{aligned}$$

<sup>†</sup> Exercise: Show why this is equivalent.

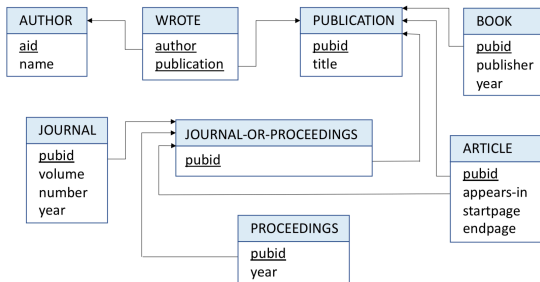
## Integrity Constraints Generally

A relational *signature* captures only the structure of relations.

Valid database instances satisfy additional *integrity constraints* in the form of sentences over the signature.

- ▶ Values of a particular attribute belong to a prescribed *data type*.
- ▶ Values of attributes are unique among tuples in a relation (**keys**).
- ▶ Values appearing in one relation must also appear in another relation (**referential integrity** or **foreign keys**).
- ▶ Values cannot appear simultaneously in certain relations (**disjointness**).
- ▶ Values in a relation must appear in at least one of another set of relations (**coverage**).
- ▶ etc.

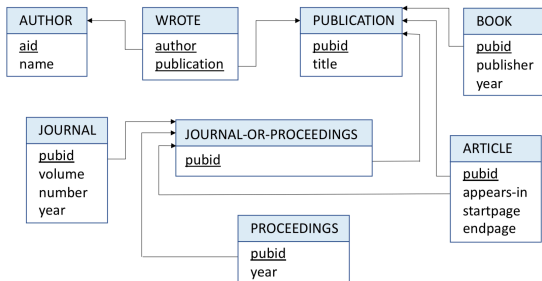
# Bibliography Integrity Constraints



## Typing Constraints / Domain Constraints

- ▶ *Author id's are integers.*
- ▶ *Author names are strings.*
- ▶ *Publication id's are integers.*
- ▶ *Publication titles are strings.*
- ▶ *etc.*

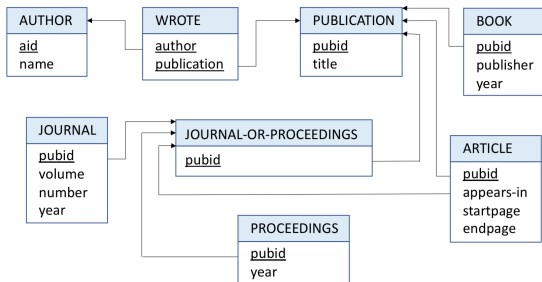
## Bibliography Integrity Constraints (cont'd)



### Uniqueness of Values / Identification (keys)

- ▶ *Author id's are unique and determine author names.*
- ▶ *Publication id's are unique as well.*
- ▶ *Articles can be identified by their publication id.*
- ▶ *Articles can also be identified by the publication id of the collection they have appeared in and their starting page number.*

## Bibliography Integrity Constraints (cont'd)



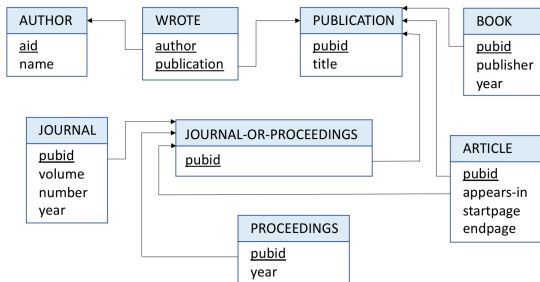
## Referential Integrity / Foreign Keys

- ▶ *Books, journals, proceedings and articles are publications.*
- ▶ *The components of a WROTE tuple must be an author and a publication.*

## Disjointness

- ▶ *Books are different from journals.*
- ▶ *Books are also different from proceedings.*

## Bibliography Integrity Constraints (cont'd)



### Coverage

- ▶ *Every publication is either a book, a journal, a proceedings, or an article.*
- ▶ *Every article appears in a journal or in a proceedings.*

## Views and Integrity Constraints

The extension of a table can be determined by an integrity constraint.

*The extension of table JOURNAL-OR-PROCEEDINGS is the union of the publication id's occurring in table JOURNAL and in table PROCEEDINGS.*

$$\forall p. \text{JOURNAL-OR-PROCEEDINGS}(p) \leftrightarrow (\text{JOURNAL}(p) \vee \text{PROCEEDINGS}(p))$$

### View

Given a signature  $\rho$ , a table  $R$  occurring in  $\rho$  is a *view* when the **relational database schema** contains exactly one integrity constraint of the form:

$$\forall x_1, \dots, x_k. R(x_1, \dots, x_k) \leftrightarrow \varphi,$$

where  $\{x_1, \dots, x_k\} = \text{Fv}(\varphi)$ . Condition  $\varphi$  is called the *view definition* of  $R$ , and  $R$  is said to *depend on* any table mentioned in  $\varphi$ .

*No table occurring in a schema is allowed to depend on itself, either directly or indirectly.*



# Relational Database Schemata and Consistency

## Relational Database Schema

A *relational database schema* is a pair  $\langle \rho, \Sigma \rangle$ , where  $\rho$  is a signature, and where  $\Sigma$  is a finite set of integrity constraints that are sentences over  $\rho$ .

## Relational Databases and Consistency

A *relational database* consists of a relational database schema  $\langle \rho, \Sigma \rangle$  and an instance **DB** of its signature  $\rho$ .

The relational database is *consistent* if and only if, for any integrity constraint  $\varphi \in \Sigma$  and any valuation  $\theta$ :

$$\mathbf{DB}, \theta \models \varphi.$$

# Outline

Unit 1: Signatures and the Relational Calculus

Unit 2: Integrity Constraints

Unit 3: **Safety and Finiteness**

Unit 4: Summary

## Story so far ...

*databases*  $\Leftrightarrow$  *relational structures*

*queries*  $\Leftrightarrow$  *set comprehensions*  
*with conditions as formulas in FOL*<sup>†</sup>

*integrity constraints*  $\Leftrightarrow$  *sentences in FOL*

So are there any remaining issues?

Yes!

Relational databases and RC<sup>‡</sup> queries should also have the following properties:

- ▶ The extension of any relation in a signature should be *finite*; and
- ▶ Queries should be **safe**: their answers should be *finite* when database instances are finite.

<sup>†</sup> *first order predicate logic*

<sup>‡</sup> *relational calculus*

## Unsafe Queries

The set of answers to each of the following queries over the bibliography RDB is not finite:

Case 1  $\{(x, y) \mid x = y\}$

Case 2  $\{(pid, pub, year) \mid \text{BOOK}(pid, pub, year) \vee \text{PROCEEDINGS}(pid, year)\}$

Case 3  $\{(aname) \mid \neg \exists aid. \text{AUTHOR}(aid, aname)\}$

## Domain Independence

An RC query  $\{(x_1, \dots, x_k) \mid \varphi\}$  is *domain independent* when, for any pair of instances  $\mathbf{DB}_1 = (\mathbf{D}_1, \approx, \mathbf{R}_1, \dots, \mathbf{R}_k)$  and  $\mathbf{DB}_2 = (\mathbf{D}_2, \approx, \mathbf{R}_1, \dots, \mathbf{R}_k)$  and any  $\theta$ ,  $\mathbf{DB}_1, \theta \models \varphi$  if and only if  $\mathbf{DB}_2, \theta \models \varphi$ .

## Theorem

Let  $(R_1, \dots, R_k)$  be the signature of a relational database. Answers to domain independent queries contain only values *that occur in the extension  $\mathbf{R}_i$  of any relation  $R_i$* .

**safety**  $\Leftrightarrow$  *domain independence and finite database instances*

# Safety and Query Satisfiability

## Theorem

Satisfiability of RC queries<sup>†</sup> is undecidable;

- ▶ co recursively enumerable in general, and
- ▶ recursively enumerable for finite databases.

<sup>†</sup> *Is there a database for which the answer is non-empty?*

## Proof

Reduction from PCP (see Abiteboul *et. al.* book, p.122-126).

## Theorem

Domain independence of RC queries is undecidable.

## Proof

The query  $\{(x, y) \mid (x = y) \wedge \varphi\}$  is satisfiable if and only if it is not domain independent.

# Range Restricted RC

## Range Restricted Conditions and Queries

Given a database signature  $\rho = (R_1/k_1, \dots, R_n/k_n)$ , a set of variable names  $\{x_1, x_2, \dots\}$  and a set of constants  $\{c_1, c_2, \dots\}$ , *range restricted conditions* are *formulas* defined by the grammar:

$$\begin{array}{l} \varphi ::= R_i(x_{i,1}, \dots, x_{i,k_i}) \\ \quad | \varphi_1 \wedge (x_i = x_j) \quad \text{where } \{x_i, x_j\} \cap \text{Fv}(\varphi_1) \neq \emptyset \text{ (case 1)} \\ \quad | x_i = c_j \\ \quad | \varphi_1 \wedge \varphi_2 \\ \quad | \exists x_i. \varphi_1 \\ \quad | \varphi_1 \vee \varphi_2 \quad \text{where } \text{Fv}(\varphi_1) = \text{Fv}(\varphi_2) \text{ (case 2)} \\ \quad | \varphi_1 \wedge \neg \varphi_2 \quad \text{where } \text{Fv}(\varphi_1) = \text{Fv}(\varphi_2) \text{ (case 3)} \end{array}$$

A *range restricted RC query* has the form  $\{(x_1, \dots, x_n) \mid \varphi\}$ , where  $\{x_1, \dots, x_n\} = \text{Fv}(\varphi)$  and where  $\varphi$  is a range restricted condition.

A query language for the relational model is *relationally complete* if the language is at least as expressive as the range restricted RC.

## Range Restricted RC (cont'd)

### Theorem

Every range restricted RC query is an RC query and is domain independent.

### Proof Outline

Both claims follow by simple inductions on the form of a range restricted condition.  
Exercise: Details.

Do we lose expressiveness by requiring conditions in RC queries to be range restricted?

### Theorem

Every domain independent RC query has an equivalent formulation as a range restricted RC query.

### Proof Outline

1. Restrict every variable in  $\varphi$  to the *active domain*, and
2. express the active domain using a *unary query* over the database instance.

Exercise: Details.

## Computational Properties of Query Answering

- ▶ There is an algorithm for computing the answers to any range restricted RC query.  
⇒ range restricted RC is not *Turing complete*.
- ▶ The **data complexity**, that is, complexity in the size of the database for a *fixed query* is
  - ⇒ in PTIME,
  - ⇒ in LOGSPACE, and
  - ⇒  $AC_0$  (i.e., constant time on polynomially many CPUs in parallel).
- ▶ The **combined complexity**, that is, complexity in size of the query and the database, is
  - ⇒ in PSPACE,
  - (since queries can express NP-hard problems such as SAT).



# Outline

Unit 1: Signatures and the Relational Calculus

Unit 2: Integrity Constraints

Unit 3: Safety and Finiteness

Unit 4: **Summary**

# Query Evaluation versus Query Satisfiability

## Query Evaluation

Given an RC query  $\{(x_1, \dots, x_k) \mid \varphi\}$  and a finite database instance **DB**, find all answers to the query.

## Query Satisfiability

Given an RC query  $\{(x_1, \dots, x_k) \mid \varphi\}$ , determine whether there is a finite database instance **DB** for which the answer is non-empty.

- ▶ Much harder problem, in fact, undecidable.
- ▶ Same as the problem of **query containment** which is fundamental in *query compilation*.
- ▶ Can be solved for fragments of RC.

## Query Subsumption

A query  $\{(x_1, \dots, x_k) \mid \varphi_1\}$  *subsumes* a query  $\{(x_1, \dots, x_k) \mid \varphi_2\}$  *with respect to a relational database schema*  $\langle \rho, \Sigma \rangle$  if, for every instance **DB** of the schema such that  $\mathbf{DB}, \theta \models \psi$  for every  $\psi \in \Sigma$ :

$$\{(\theta(x_1), \dots, \theta(x_k)) \mid \mathbf{DB}, \theta \models \varphi_2\} \subseteq \{(\theta(x_1), \dots, \theta(x_k)) \mid \mathbf{DB}, \theta \models \varphi_1\}$$

- ▶ *Fundamental* in query compilation, e.g., query simplification.
- ▶ Equivalent to determining if the following is satisfiable:

$$\{(x_1, \dots, x_k) \mid \varphi_2 \wedge \neg \varphi_1\}.$$

- ▶ Also equivalent to proving the following in FOL:

$$\left( \bigwedge_{\psi \in \Sigma} \psi \right) \rightarrow (\forall x_1, \dots, x_k. \varphi_2 \rightarrow \varphi_1).$$

- ▶ Again, undecidable in general, but decidable for fragments of RC.

## What queries cannot be expressed in RC?

Recall that range restricted RC is not Turing complete  
⇒ there are computable queries that cannot be expressed.

### Built In Operations

- ▶ ordering, arithmetic, string operations, etc.

### Counting and Aggregation

- ▶ Cardinality of Sets (*parity*)

### Reachability, Connectivity, ...

- ▶ *paths in a graph (binary relation)*

Data model extensions relating to incompleteness and inconsistency:

- ▶ tuples with *unknown* (but existing) values;
- ▶ incomplete relations and *open world assumption*; and
- ▶ conflicting information (e.g., from different data sources).

# The Final Story

<i>databases</i>	$\Leftrightarrow$	<i>relational structures</i>
<i>queries</i>	$\Leftrightarrow$	<i>set comprehensions with conditions as formulas in FOL<sup>†</sup></i>
<i>integrity constraints</i>	$\Leftrightarrow$	<i>sentences in FOL</i>
<i>safety</i>	$\Leftrightarrow$	<i>range restricted RC<sup>‡</sup> and finite database instances</i>

<sup>†</sup> *first order predicate logic*

<sup>‡</sup> *relational calculus*