Updated 2/4/2022

#### Topic 1: Data models and query languages Unit 2: Logic & relational calculus (continued) Lecture 05

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp22)

https://northeastern-datalab.github.io/cs7240/sp22/

2/1/2022

#### Pre-class conversations

- Last class recapitulation
- Project ideas
- Slides from today likely only posted tomorrow night. Sorry!
- today:
  - logic continued (likely next time algebra and the connection)
  - logic becomes super important today again, lack of expertise; thus lots of practice today <sup>(i)</sup>
  - in particular the concept of "undecidability": intuition for why things can quickly get complicated without giving proofs

*Topic 1: Data Models and Query Languages* 

- Lecture 1 (Tue 1/18): Course introduction, SQL, PostgreSQL setup, SQL Activities
- Lecture 2 (Fri 1/21): SQL
- Lecture 3 (Tue 1/25): SQL
- Lecture 4 (Fri 1/28): Logic & Relational Calculus
- Lecture 5 (Tue 2/1): Relational Algebra & Codd's Theorem
- Lecture 6 (Fri 2/4): Relational Algebra & Codd's Theorem
- Lecture 7 (Tue 2/8): Datalog & Recursion,
- Lecture 8 (Fri 2/11): Datalog & Recursion,
- Lecture 9 (Tue 2/15): Alternative Data Models

Pointers to relevant concepts & supplementary material:

- Unit 1. SQL: [SAMS'12], [CS 3200], [Cow'03] Ch3 & Ch5, [Complete'08] Ch6, [Silberschatz+'20] Ch3.8
- Unit 2. Logic & Relational Calculus: First-Order Logic (FOL), relational calculus (RC): [Barland+'08] 4.1.2 & 4.2.1 & 4.4, [Genesereth+] Ch6, [Halpern+'01], [Cow'03] Ch4.3 & 4.4, [Elmasri, Navathe'15] Ch8.6 & Ch8.7, [Silberschatz+'20] Ch27.1 & Ch27.2, [Alice'95] Ch3.1-3.3 & Ch4.2 & Ch4.4 & Ch5.3-5.4, [Barker-Plummer+'11] Ch11
- Unit 3. Relational Algebra & Codd's Theorem: Relational Algebra (RA), Codd's theorem: [Cow'03] Ch4.2,
   [Complete'08] Ch2.4 & Ch5.1-5.2, [Elmasri, Navathe'15] Ch8, [Silberschatz+'20] Ch2.6, [Alice'95] Ch4.4 & Ch5.4
- Unit 4. Datalog & Recursion: Datalog, recursion, Stratified Datalog with negation, Datalog evaluation strategies, Stable Model semantics, Answer Set Programming (ASP): [Complete'08] Ch5.3, [Cow'03] Ch 24, [Koutris'19] L9 & L10, [G., Suciu'10]
- Unit 5. Alternative Data Models: NoSQL: [Hellerstein, Stonebraker'05], [Sadalage, Fowler'12], [Harrison'16]

PRELIMINARY

#### Queries and the connection to logic

- Why logic?
- A crash course in FOL
- Relational Calculus
  - Syntax and Semantics
  - Domain Independence and Safety

#### Entire Story in One Slide

- 1. RC = FOL over DB
- 2. RC can express "bad queries" that depend not only on the DB, but also on the domain from which values are taken (domain dependence)
- 3. We cannot test whether an RC query is "good," but we can use a "good" subset of RC that captures all "good" queries (safety)
- 4. "Good" RC and RA can express the same queries! (equivalence = Codd's theorem)

#### Relational Calculus (RC)

- RC is, essentially, first-order logic (FOL) over the schema relations
  - A query has the form "find all tuples  $(x_1, ..., x_k)$  that satisfy an FOL condition"
- RC is a declarative query language
  - Meaning: a query is not defined by a sequence of operations, but rather by a condition that the result should satisfy

#### Queries and the connection to logic

- Why logic?
- A crash course in FOL
- Relational Calculus
  - Syntax and Semantics
  - Domain Independence and Safety

#### **RC** Query



Person(id, gender, country) Parent(parent, child) Spouse(person1, person2)

assume symmetric relation  $(a,b) \in Spouse \Leftrightarrow (b,a) \in Spouse$ 

# $\left\{ \begin{array}{l} (\mathbf{x},\mathbf{u}) \mid \operatorname{Person}(\mathbf{u}, '\operatorname{female}', '\operatorname{Canada}') \wedge & (a,b) \in \operatorname{Spouse} \Leftrightarrow (b,a) \in (a,b) \in \operatorname{Spouse} \Leftrightarrow (b,a) \in (b,a) \in (a,b) \in \operatorname{Spouse} \Leftrightarrow (b,a) \in (a,b) \in \operatorname{Spouse} \otimes (b,a) \in (a,b) \in (a,$



which relatives does ? ?

#### **RC** Symbols

- Constant values: a, b, c, female, Canada, ...
  - Values that may appear in table cells (optionally with quotation marks)
- Variables: x, y, z, ...
  - Range over the values that may appear in table cells
- Relation symbols: R, S, T, Person, Parent, ...
  - Each with a specified arity
  - Will be fixed by the relational schema at hand
  - No attribute names, only attribute positions (= unnamed perspective)!
- Unlike general FOL, no function symbols!

RC Formulas (atomic and non-atomic)

- Atomic formulas:
  - $R(t_1,...,t_k)$ 
    - R is a k-ary relation, Each t<sub>i</sub> is a variable or a constant
    - Semantically it states that  $(t_1, ..., t_k)$  is a tuple in R

x op u

- x is a variable, u is a variable/constant, op is one of >, <, =, ≠</li>
- Simply binary predicates, predefined interpretation
- Formula:
  - Atomic formula
  - If  $\phi$  and  $\psi$  are formulas then these are formulas:  $\phi \land \psi \quad \phi \lor \psi \quad \phi \rightarrow \psi \quad \phi \rightarrow \psi \quad \neg \phi \quad \exists x \phi \quad \forall x \phi$

Person(x, 'female', 'Canada')

x=y, y≠w, z>5, z='female'

#### Free Variables

- Intuitively: free variable are not bound to quantifiers
- Formally:
  - A free variable of an atomic formula is a variable that occurs in the atomic formula
  - A free variable of  $\phi \land \psi, \phi \lor \psi, \phi \rightarrow \psi$  is a free variable of either  $\phi$  or  $\psi$
  - A free variable of  $\neg \phi$  is a free variable of  $\phi$
  - A free variable of  $\exists x \phi$  and  $\forall x \phi$  is a free variable y of  $\phi$  such that  $y \neq x$
- We write  $\phi(x_1,...,x_k)$  to state that  $x_1,...,x_k$  are the free variables of formula  $\phi$  (in some order)

#### Back to our earlier example



#### This is a formula!

#### Person(u,)'female', 'Canada') $\land$ $\exists z, y [Parent(z, y) \land Parent(y, x) \land$ $\exists w [Parent(z, w) \land y \neq w \land (u = w \lor Spouse(u, w))] ]$



## what are the free ?

#### Back to our earlier example



#### Person(u, 'female', 'Canada') $\land$ $\exists z, y [Parent(z, y) \land Parent(y, x) \land$ $\exists w [Parent(z, w) \land y \neq w \land (u = w \lor Spouse(u, w))] ]$



Notation:

 $\varphi(x,u)$  / CanadianAunt(u,x)

RC query



# $\left\{ \begin{array}{l} (\mathbf{x}, \mathbf{u}) \mid \text{Person}(\mathbf{u}, \text{'female'}, \text{'Canada'}) \land \\ \exists z, y \left[ \text{Parent}(z, y) \land \text{Parent}(y, x) \land \\ \exists w \left[ \text{Parent}(z, w) \land y \neq w \land (\mathbf{u} = w \lor \text{Spouse}(\mathbf{u}, w)) \right] \right\}$



{ 
$$(x_1,...,x_k) | \phi(x_1,...,x_k)$$
 }

 $\varphi(x,u)$  / CanadianAunt(u,x)

Relation Calculus Query

some condition on the variables  $COND(x_1,...,x_k)$ 

• An RC query is an expression of the form

 $\{(x_1, \dots, x_k) \mid \varphi(x_1, \dots, x_k)\}$ where  $\varphi(x_1, \dots, x_k)$  is an RC formula

 An RC query is *over* a relational schema S if all the relation symbols belong to S (with matching arities)

#### DRC vs. TRC

- There are two common variants of RC:
  - DRC: Domain Relational Calculus (what we have seen so far)
  - TRC: Tuple Relational Calculus
- DRC applies vanilla FO: terms interpreted as attribute values, relations have arity but no attribute names (= unnamed perspective)
- TRC is more "database friendly": terms interpreted as tuples with named attributes
- There are easy conversions between the two formalisms

 $\begin{array}{l} \begin{array}{l} \text{Our Example in TRC} \\ \text{optionally "t(nephew, aunt)"} \\ \\ \begin{array}{l} \textbf{t} & \exists a \in \text{Person [a.gender = 'female' $\ensuremath{\Lambda}$ a.country = 'Canada'] $\ensuremath{\Lambda}$ \\ \\ \exists p,q,w \in \text{Parent [p.child $\ext{t.nephew}$ $\ensuremath{\Lambda}$ q.child $= p.parent $\ensuremath{\Lambda}$ \\ \\ \text{w.parent = $q.parent $\ensuremath{\Lambda}$ w.child $\neq$ q.child $\ensuremath{\Lambda}$ a.id $= \ensuremath{\textbf{t.aunt}}$ \\ \end{array} \end{array}$ 

(w.child =  $\mathbf{a}$ .id V  $\exists \mathbf{s} [\mathbf{s} \in \text{Spouse } \land \mathbf{s}$ .person1 = w.child  $\land \mathbf{s}$ .person2 =  $\mathbf{a}$ .id])]}



tuple variables like in SQL instead of domain variables: {+ | COND(+)}

often used short forms: $\forall x \in \mathbb{R}[\phi]$ same as $\forall x[x \in \mathbb{R} \Rightarrow \phi]$  $\exists x \in \mathbb{R}[\phi]$ same as $\exists x[x \in \mathbb{R} \land \phi]$ 

83

Based on material by Benny Kimelfeld and Oded Shmueli for 236363 Database Management Systems, Technion, 2018. However, notice I prefer and follow here the notation of [Ramakrishnan, Gehrke' 03] and [Elmasri, Navathe'15] of using a.country = 'Canada', instead of the alternative notation a[country]='Canada' used by [Silberschatz, Korth, Sudarshan 2010] Wolfgang Gatterbauer. Principles of scalable data management: <a href="https://northeastern-datalab.github.io/cs7240/">https://northeastern-datalab.github.io/cs7240/</a>

#### Different TRC notations

Find persons that frequent <u>only</u> bars that serve <u>only</u> drinks they like.

(Find persons who like all drinks that are served in all the bars they visit.)

(Find persons for which there does not exist a bar they frequent that serves a drink they do not like.)

 $\{q(person) \mid \exists f \in Frequents \ [f.person=q.person \land \neg(\exists f2 \in Frequents \ [f2.person=f.person \land mq Preferred notation \neg(\exists I \in Likes, \exists s \in Serves \ [I.drink=s.drink \land f2.bar=s.bar \land f2.person=I.person])]\}$ 

 $\{F. person \mid F \in Frequents.(\nexists F2 \in Frequents.(F2. person=F. person \land (\nexists L \in Likes, \nexists S \in Serves.(L.drink=S.drink \land F2. bar=S.bar \land F2. person=L. person))) \}$ 

- {t: Person |  $\exists f \in Frequents [t(Person)=f(Person) \land \neg \exists f 2 \in Frequents [F2(person)=F(person) \land [Deutsch 2019] \neg (\exists I \in Likes \exists s \in Serves) [I(Drink)=s(Drink) \land f2(Bar)=s(Bar) \land f2(Person)=I(Person)]]}$
- {f.Person | Frequents(f) AND (NOT(∃f2)(Frequents(f2) AND f2.person=f.person ∧ [Elmasri 2015] (NOT(∃I)(∃s)(Likes(I) AND Serves(s) AND I.drink=s.drink AND f2.bar=s.bar AND f2.person=I.person)))}
- $\{\mu^{(1)} \mid (\exists \rho^{(2)}) (\text{Frequents}(\rho) \land \rho[1] = \mu[1] \land \neg((\exists \lambda^{(2)}) (\text{Frequents}(\lambda) \land \lambda[1] = \rho[1] \land \text{[Ullman 1988]} \land ((\exists \nu^{(2)}) (\exists \theta^{(2)}) (\text{Likes}(\nu) \land \text{Serves}(\theta) \land \nu(2) = \theta(2) \land \lambda(2) = \theta(1) \land \lambda(1) = \nu(1)))) \}$
- {P|  $\exists F \in Frequents$  (F.person=P.person  $\land \neg \exists F2 \in Frequents(F2.person=F.person \land \neg (\exists L \in Likes \exists S \in Serves (L.drink=S.drink \land F2.bar=S.bar \land F2.person=L.person)))}$

[Ramakrishnan 2003]

331



#### Queries and the connection to logic

- Why logic?
- A crash course in FOL
- Relational Calculus
  - Syntax and Semantics
  - Domain Independence and Safety

#### Bringing in the Domain

- Let **S** be a schema, D a database over **S**, and Q an RC query over **S**
- Then D gives an interpretation for the underlying FOL
  - Predicates  $\rightarrow$  relations; constants copied; no functions



#### Bringing in the Domain

- Let **S** be a schema, D a database over **S**, and Q an RC query over **S**
- Then D gives an interpretation for the underlying FOL
  - Predicates → relations; constants copied; no functions
  - Not yet! We need to answer first: What is the <u>domain</u>?
- The active domain ADom (of D and Q) is the set of all the values that occur in either D or Q
- The query Q is evaluated over D with respect to a domain Dom that contains the active domain (Dom ⊇ ADom)
- Denote by Q<sup>Dom</sup>(D) the result of evaluating Q over D relative to the domain
   Dom

#### Domain Independence

- Let  ${\boldsymbol{S}}$  be a schema, and let Q be an RC query over  ${\boldsymbol{S}}$
- We say that Q is domain independent if for every database D over  $\boldsymbol{S}$  and  $\ldots$



#### Domain Independence

- Let  ${\boldsymbol{S}}$  be a schema, and let Q be an RC query over  ${\boldsymbol{S}}$
- We say that Q is domain independent if for every database D over S and every two domains Dom<sub>1</sub> and Dom<sub>2</sub> that contain the active domain, we have:

$$Q^{\text{Dom}}(D) = Q^{\text{Dom}}(D) = Q^{\text{ADom}}(D)$$

#### Bad News...

- We would like be able to tell whether a given RC query is domain independent, and then reject "bad queries"
- Alas, this problem is undecidable!
  - That is, there is no algorithm that takes as input an RC query and returns true iff the query is domain independent

Based on material by Benny Kimelfeld and Oded Shmueli for 236363 Database Management Systems, Technion, 2018. First observed in "Di Paola. The Recursive Unsolvability of the Decision Problem for the Class of Definite Formulas, JACM 1969. <u>https://doi.org/10.1145/321510.321524</u>" Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

#### Good News

- Domain-independent RC has an "effective syntax", that is:
  - A syntactic restriction of RC in which every query is domain independent
  - Restricted queries are said to be safe
- Safety can be tested automatically (and efficiently)
  - Most importantly, for every domain independent RC query there exists an equivalent safe RC query!

#### Safety

- We do not formally define the safe syntax in this course
- Details on the safe syntax can be found in Ch 5.4 of [Alice'95]: Foundations of Databases by Abiteboul, Hull and Vianu
  - Example:
    - In  $\exists x \phi$ , the variable x should be guarded by  $\phi$
    - Every variable  $x_i$  is guarded by  $R(x_1,...,x_k)$
    - In  $\phi \land (x=y)$ , the variable x is guarded if and only if either x or y is guarded by  $\phi$
    - ... and so on



Based on material by Benny Kimelfeld and Oded Shmueli for 236363 Database Management Systems, Technion, 2018.

An accessible overview of issues involving safety is: Topor, Safety and Domain Independence, Encyclopedia of Database Systems. <u>https://doi.org/10.1007/978-0-387-39940-9\_1255</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

ADom =  $\{1, 2, 3, 'female', 'Canada'\}$ Dom = ADom U  $\{'elefant', 'car', 'lemon', \pi, ...\}$ 

 $\{(\mathbf{x}) \mid \neg \operatorname{Person}(\mathbf{x}, '\operatorname{female'}, '\operatorname{Canada'})\}$ 

 $\{(x,y) | \exists z [Spouse(x,z) \land y=z] \}$ 

 $\{ (x,y) | \exists z [Spouse(x,z) \land y \neq z] \}$ 

Based on material by Benny Kimelfeld and Oded Shmueli for 236363 Database Management Systems, Technion, 2018. Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2)

103



x could be also 'Canada' or 'female' or ...

Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2) Example fixes: ...  $\Lambda \exists y, z. \operatorname{Person}(x) y, z$  (...  $\Lambda \exists y. \operatorname{Person}(x, y, 'Canada')$ 

Not DI

 $\{ (\mathbf{x}, \mathbf{y}) | \exists z [Spouse(\mathbf{x}, z) \land \mathbf{y} = z] \}$ 

 $\{(x) \mid \neg Person(x, 'female', 'Canada')\}$ 

 $\{ (\mathbf{x}, \mathbf{y}) | \exists z [Spouse(\mathbf{x}, z) \land \mathbf{y} \neq z] \}$ 

 $\{ (x) \mid \neg Person(x, 'female', 'Canada') \}$ 

Example fixes: ...  $\Lambda \exists y, z. Person(x, y, z)$ 

x could be also 'Canada' or 'female' or ...

some as  $\{(x,y) \mid \text{Spouse}(x,y)\} = \text{Spouse}(x,y)$ 



...  $\Lambda$  **Jy**. Person(x,y, 'Canada')

Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2)

Not DI

DI

 $\{ (x,y) | \exists z [Spouse(x,z) \land y \neq z] \}$ 

 $\left\{ (\mathfrak{Y}, \mathfrak{Y}) \mid \exists z \ [Spouse(\mathfrak{Y}, \mathfrak{Y}) \land \mathfrak{Y}=\mathfrak{Y}] \right\}$ 



Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2)

Not DI

DI

Example fixes: ...  $\land \exists y, z. \operatorname{Person}(x, y, z)$ ...  $\land \exists y. \operatorname{Person}(x, y, 'Canada')$  $\{(x) \mid \neg \operatorname{Person}(x, 'female', 'Canada') \}$ x could be also 'Canada' or 'female' or ...

 $\left\{ (x,y) | \exists z [Spouse(x,z) \land y=z] \right\}$ same as  $\{(x,y) | Spouse(x,y)\}$ 

 $\left\{ (\mathbf{x}, \mathbf{y}) | \exists z [Spouse(\mathbf{x}, z) \land \mathbf{y} \neq z] \right\}$ 



Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2)

Not DI

DI

Example fixes: ...  $\Lambda \exists y, z. \operatorname{Person}(x, y, z)$ ...  $\Lambda \exists y. \operatorname{Person}(x, y, 'Canada')$  $\{ (x) \mid \neg \operatorname{Person}(x, 'female', 'Canada') \}$ x could be also 'Canada' or 'female' or ...

 $\left\{ (x,y) | \exists z [Spouse(x,z) \land y=z] \right\}$ same as  $\{(x,y) | Spouse(x,y)\}$ 

 $\left\{ (\mathbf{x}, \mathbf{y}) | \exists z [Spouse(\mathbf{x}, z) \land \mathbf{y} \neq z] \right\}$ 

Based on material by Benny Kimelfeld and Oded Shmueli for 236363 Database Management Systems, Technion, 2018. Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

Not DI



Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2)

#### $\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \exists y [\neg \operatorname{Likes}(x, y)] \}$

#### $\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \forall y [\neg \operatorname{Likes}(x, y)] \}$

#### $\left\{ (\mathbf{x}) | \exists z, w \operatorname{Person}(\mathbf{x}, z, w) \land \forall y [\neg \operatorname{Likes}(\mathbf{x}, y)] \land \exists y [\neg \operatorname{Likes}(\mathbf{x}, y)] \right\}$



Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2)

Person('Alice', 'female', 'Canada') Person('Beate', 'female', 'Canada') Person('Cecile', 'female', 'Canada')

Likes('Alice', 'Beate') Likes('Alice', 'Cecile') Likes('Alice', 'Alice')

 $\mathbf{ADom} = \mathbf{?}$ 

#### $\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \exists y [\neg \operatorname{Likes}(x, y)] \}$

$$\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \forall y [\neg \operatorname{Likes}(x, y)] \}$$

$$\left\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \forall y [\neg \operatorname{Likes}(x, y)] \land \exists y [\neg \operatorname{Likes}(x, y)] \right\}$$

Person('Alice', 'female', 'Canada') Likes('Alice', 'Beate')
Person('Beate', 'female', 'Canada') Likes('Alice', 'Cecile')
Person('Cecile', 'female', 'Canada') Likes('Alice', 'Alice')

**ADom** = {'Alice', 'Beate', 'Cecile', 'female', 'Canada')

#### $\{ (\mathbf{x}) | \exists z, w \operatorname{Person}(\mathbf{x}, z, w) \land \exists y [\neg \operatorname{Likes}(\mathbf{x}, y)] \}$

$$\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \forall y [\neg \operatorname{Likes}(x, y)] \}$$

$$\left\{ (\mathbf{x}) | \exists z, w \operatorname{Person}(\mathbf{x}, z, w) \land \forall y [\neg \operatorname{Likes}(\mathbf{x}, y)] \land \exists y [\neg \operatorname{Likes}(\mathbf{x}, y)] \right\}$$



Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2)

Person('Alice', 'Alice', 'Alice') Person('Beate', 'Beate', 'Beate') Person('Cecile', 'Beate', 'Beate')

D

Likes('Alice', 'Beate') Likes('Alice', 'Cecile') Likes('Alice', 'Alice') Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2)

ADom = {'Alice', 'Beate', 'Cecile')
Dom = {'Alice', 'Beate', 'Cecile', 'Dora')

 $\{ (\mathbf{x}) | \exists z, w \operatorname{Person}(\mathbf{x}, z, w) \land \exists y [\neg \operatorname{Likes}(\mathbf{x}, y)] \}$ 

 $\left\{ (\mathbf{x}) | \exists z, w \operatorname{Person}(\mathbf{x}, z, w) \land \forall y [\neg \operatorname{Likes}(\mathbf{x}, y)] \right\}$ 

 $\left\{ (\mathbf{x}) | \exists z, w \operatorname{Person}(\mathbf{x}, z, w) \land \forall y [\neg \operatorname{Likes}(\mathbf{x}, y)] \land \exists y [\neg \operatorname{Likes}(\mathbf{x}, y)] \right\}$ 



Person('Alice', 'Alice', 'Alice')

Person('Beate', 'Beate', 'Beate')

D



Likes('Alice', 'Beate')

Likes('Alice', 'Cecile')

Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2)

Person('Cecile', 'Beate', 'Beate') Likes('Alice', 'Alice') **ADom** = {'Alice', 'Beate', 'Cecile') = {'Alice', 'Beate', 'Cecile', 'Dora') Dom Person(y,\_,\_) Example fix: ...  $\Lambda \exists u, v [Person(y, u, v)]$  $\{ (\mathbf{x}) | \exists z, w \operatorname{Person}(\mathbf{x}, z, w) \land \exists y [\neg \operatorname{Likes}(\mathbf{x}, y)] \}$ Not DI Alice is in the output if  $Dom \supset ADom$  (e.g., Dora is in Dom)  $\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \forall y [\neg \operatorname{Likes}(x, y)] \}$  $\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \forall y [\neg \operatorname{Likes}(x, y)] \land \exists y [\neg \operatorname{Likes}(x, y)] \}$ 

Likes('Alice', 'Beate')

Likes('Alice', 'Cecile')

Likes('Alice', 'Alice')

Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2)

Person('Alice', 'Alice', 'Alice') Person('Beate', 'Beate', 'Beate') Person('Cecile', 'Beate', 'Beate')

D

**ADom** = {'Alice', 'Beate', 'Cecile') = {'Alice', 'Beate', 'Cecile', 'Dora') Dom Person(y,\_,\_) Example fix: ...  $\Lambda \exists u, v [Person(y, u, v)]$  $\{ (\mathbf{x}) | \exists z, w \operatorname{Person}(\mathbf{x}, z, w) \land \exists y [\neg \operatorname{Likes}(\mathbf{x}, y)] \}$ Not DI Alice is in the output if  $Dom \supset ADom$  (e.g., Dora is in Dom)  $\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \forall y [\neg \operatorname{Likes}(x, y)] \}$ D x never occurs in Likes(x, ): Beate, Cecile  $\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \forall y [\neg \operatorname{Likes}(x, y)] \land \exists y [\neg \operatorname{Likes}(x, y)] \}$ 

Which One is Domain Independent? Person(id, gender, country) Likes(person1, person2) Spouse(person1, person2) Likes('Alice', 'Beate') Person('Alice', 'Alice', 'Alice') Person('Beate', 'Beate', 'Beate') Likes('Alice', 'Cecile') Person('Cecile', 'Beate', 'Beate') Likes('Alice', 'Alice') **ADom** = {'Alice', 'Beate', 'Cecile') = {'Alice', 'Beate', 'Cecile', 'Dora') Dom Person(y,\_,\_) Example fix: ...  $\Lambda \exists u, v [Person(y, u, v)]$  $\{ (\mathbf{x}) | \exists z, w \operatorname{Person}(\mathbf{x}, z, w) \land \exists y [\neg \operatorname{Likes}(\mathbf{x}, y)] \}$ Not DI Alice is in the output if  $Dom \supset ADom$  (e.g., Dora is in Dom)  $\{ (x) | \exists z, w \operatorname{Person}(x, z, w) \land \forall y [\neg \operatorname{Likes}(x, y)] \}$ x never occurs in Likes(x, ): Beate, Cecile  $\{(x) | \exists z, w \operatorname{Person}(x, z, w) \land \forall y [\neg \operatorname{Likes}(x, y)] \land \exists y [\neg \operatorname{Likes}(x, y)] \}$ D implication (absorption) if Dom  $\neq \emptyset$ , which is necessary for there to be Person(x,\_,)

What is the meaning of following unsafe expressions?

- $\left\{ \mathbf{x} \mid \exists \mathbf{y}. \mathbf{R}(\mathbf{x}) \right\}$  $\left\{ \mathbf{x} \mid \mathbf{x} \ge 10 \right\}$
- $\left\{ \mathbf{x} \mid \forall \mathbf{y} \, \mathbf{R}(\mathbf{x}, \mathbf{y}) \right\}$

## What is the meaning of following unsafe expressions? $\begin{cases} x \mid \exists y. R(x) \end{cases}$ logically equivalent to $\{ x \mid R(x) \} = R(x) \\ \{ x \mid x \ge 10 \}$ ? $\{ x \mid \forall y R(x,y) \}$ ?

# What is the meaning of following unsafe expressions? $\{x \mid \exists y. R(x)\}$ logically equivalent to $\{x \mid \mathcal{R}(x)\} = \mathcal{R}(x)$ $\{x \mid x \ge 10\}$ What if $\mathcal{D}om=\mathbb{N}$ ? $\mathcal{D}I: \{x \mid A(x) \land x \ge 10\}$ $\{x \mid \forall y R(x,y)\}$ ?

What is the meaning of following unsafe expressions?  $\left\{ \mathbf{x} \mid \exists \mathbf{y}. \mathbf{R}(\mathbf{x}) \right\}$ logically equivalent to  $\{x \mid R(x)\} = R(x)$  $\left\{ \mathbf{x} \mid \mathbf{x} \ge 10 \right\}$  $\mathcal{DI:} \left\{ x \mid A(x) \land x \ge 10 \right\}$ What if Dom=N? $\mathcal{DI} : \left\{ x \mid \forall y \left[ A(y) \rightarrow R(x,y) \right] \right\}$  $\left\{ \mathbf{x} \mid \forall \mathbf{y} \, \mathbf{R}(\mathbf{x}, \mathbf{y}) \right\}$  $\mathcal{D}$ :  $\mathbb{R}(a',a')$  $ADom = \{a'\}$ Dom={'a','Chile'} <u>, ?</u>

What is the meaning of following unsafe expressions?  $\left\{ \mathbf{x} \mid \exists \mathbf{y}. \mathbf{R}(\mathbf{x}) \right\}$ logically equivalent to  $\{x \mid R(x)\} = R(x)$  $\left\{ \mathbf{x} \mid \mathbf{x} \ge 10 \right\}$  $\mathcal{DI:} \left\{ x \mid A(x) \land x \ge 10 \right\}$ what if Dom=N? $\left\{ \mathbf{x} \mid \forall \mathbf{y} \, \mathbf{R}(\mathbf{x}, \mathbf{y}) \right\}$  $\mathcal{DI} : \left\{ x \mid \forall y \left[ A(y) \rightarrow R(x,y) \right] \right\}$  $\mathcal{D}$ :  $\mathbb{R}(a',a')$  $ADom = \{a'\}$ what if relation A is empty? Dom={'a','Chile'}

What is the meaning of following unsafe expressions?  $\left\{ \mathbf{x} \mid \exists \mathbf{y}. \mathbf{R}(\mathbf{x}) \right\}$ logically equivalent to  $\{x \mid R(x)\} = R(x)$  $\left\{ \mathbf{x} \mid \mathbf{x} \ge 10 \right\}$  $\mathcal{DI:} \left\{ \mathbf{x} \mid \mathbf{A}(\mathbf{x}) \land \mathbf{x} \ge 10 \right\}$ what if Dom=N? $\mathcal{DI} : \left\{ x \mid \forall y \left[ A(y) \rightarrow R(x,y) \right] \right\}$  $\left\{ \mathbf{x} \mid \forall \mathbf{y} \, \mathbf{R}(\mathbf{x}, \mathbf{y}) \right\}$  $\mathcal{D}$ :  $\mathbb{R}(a',a')$  $ADom = \{ a \}$ what if relation A is empty? Dom={'a','Chile'} 1. always true for  $A=\emptyset$  $\left\{ x \mid \forall y \left( \neg A(y) \lor R(x,y) \right] \right\}$ 

What is the meaning of following unsafe expressions?  $\left\{ \mathbf{x} \mid \exists \mathbf{y}. \mathbf{R}(\mathbf{x}) \right\}$ logically equivalent to  $\{x \mid \mathbb{R}(x)\} = \mathbb{R}(x)$  $\left\{ \mathbf{x} \mid \mathbf{x} \ge 10 \right\}$  $\mathcal{DI:} \left\{ \mathbf{x} \mid \mathbf{A}(\mathbf{x}) \land \mathbf{x} \ge 10 \right\}$ what if Dom=N?not DI:  $\{ x \mid \forall y [A(y) \rightarrow R(x,y)] \}$  $\left\{ \mathbf{x} \mid \forall \mathbf{y} \, \mathbf{R}(\mathbf{x}, \mathbf{y}) \right\}$  $\mathcal{D}$ :  $\mathbb{R}(a',a')$ ADom={'a'} what if relation A is empty? Dom={'a','Chile'} Neutral element for  $\forall$  is TRUE 1. always true for  $A=\emptyset$  $\sum : 0 + x = x$  $\left\{ x \mid \forall y \left( \neg A(y) \lor R(x,y) \right] \right\}$  $\prod$ : 1 · x = x V: FALSE V x = x $\exists$  :  $x_1 \lor x_2 \lor \dots \lor FALSE$ 2. alternative way  $\forall$ :  $x_1 \land x_2 \land \dots \land TRUE$  $\Lambda$ : TRUE  $\Lambda x = x$ to see that MIN: MIN( $\infty$ , x) = x

Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/



Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/

What is the meaning of following unsafe expressions?  $\left\{ \mathbf{x} \mid \exists \mathbf{y}. \mathbf{R}(\mathbf{x}) \right\}$ logically equivalent to  $\{x \mid R(x)\} = R(x)$  $\left\{ \mathbf{x} \mid \mathbf{x} \ge 10 \right\}$  $\mathcal{DI}: \left\{ \mathbf{x} \mid \mathbf{A}(\mathbf{x}) \land \mathbf{x} \ge 10 \right\}$ what if Dom=N? Not DI:  $\{ x \mid \forall y [A(y) \rightarrow R(x,y)] \}$  $\{ \mathbf{x} \mid \forall \mathbf{y} \mathbf{R}(\mathbf{x},\mathbf{y}) \}$ DI:

What is the meaning of following unsafe expressions?  $\left\{ \mathbf{x} \mid \exists \mathbf{y}. \mathbf{R}(\mathbf{x}) \right\}$ logically equivalent to  $\{x \mid R(x)\} = R(x)$  $\left\{ \mathbf{x} \mid \mathbf{x} \ge 10 \right\}$  $\mathcal{DI:} \left\{ \mathbf{x} \mid \mathbf{A}(\mathbf{x}) \land \mathbf{x} \ge 10 \right\}$ what if Dom=N? Not DI:  $\{ x \mid \forall y [A(y) \rightarrow R(x,y)] \}$  $\{ \mathbf{x} \mid \forall \mathbf{y} \mathbf{R}(\mathbf{x},\mathbf{y}) \}$  $\begin{array}{c} \text{or } \mathbb{R}(x, \_) \quad \exists z [\mathbb{R}(x, z) \land ...] \\ \mathbb{DI}: \quad \left\{ \begin{array}{c} x \mid A(x) \land \forall y [A(y) \rightarrow \mathbb{R}(x, y)] \right\} \end{array} \right.$  $\left\{ \mathbf{x} \mid \mathbf{A}(\mathbf{x}) \land \nexists \mathbf{y} \left[ \mathbf{A}(\mathbf{y}) \land \neg \mathbf{R}(\mathbf{x}, \mathbf{y}) \right] \right\}$ 





A encodes the directed edges of a graph ("arcs")

Based on an example by Dan Suciu from CSE 554, 2011. Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



What do these queries return ?

$$\left\{ x \mid \exists y. A(x,y) \right\}$$

Nodes that have at least one child:





A encodes the directed edges of a graph ("arcs")

Based on an example by Dan Suciu from CSE 554, 2011. Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/

 $\left\{ \begin{array}{l} x \mid \exists y, z, u. [A(x, y) \land A(y, z) \land A(z, u)] \right\} \\ \end{array} \\ \left\{ \begin{array}{l} (x, y) \mid \forall z. [A(x, z) \rightarrow A(y, z)] \right\} \\ \end{array} \right\}$ 



What do these queries return ?

$$\left\{ \mathbf{x} \mid \exists \mathbf{y}. A(\mathbf{x}, \mathbf{y}) \right\}$$

Nodes that have at least one child:  $\{1,2,3\}$ 

 $\left\{ \mathbf{x} \mid \exists \mathbf{y}, \mathbf{z}, \mathbf{u}. [A(\mathbf{x}, \mathbf{y}) \land A(\mathbf{y}, \mathbf{z}) \land A(\mathbf{z}, \mathbf{u})] \right\}$ 

 $\left\{ (\mathbf{x}, \mathbf{y}) \mid \forall \mathbf{z} . [\mathbf{A}(\mathbf{x}, \mathbf{z}) \to \mathbf{A}(\mathbf{y}, \mathbf{z})] \right\}$ 



A encodes the directed edges of a graph ("arcs")

Based on an example by Dan Suciu from CSE 554, 2011. Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/



What do these queries return ?

$$\left\{ \mathbf{x} \mid \exists \mathbf{y}. \mathbf{A}(\mathbf{x},\mathbf{y}) \right\}$$

Nodes that have at least one child:  $\{1,2,3\}$ 



A encodes the directed edges of a graph ("arcs")

 $\left\{ \begin{array}{l} x \mid \exists y, z, u. [A(x, y) \land A(y, z) \land A(z, u)] \right\}$ Nodes that have a great-grand-child:

$$\left\{ (\mathbf{x},\mathbf{y}) \mid \forall \mathbf{z}.[\mathbf{A}(\mathbf{x},\mathbf{z}) \to \mathbf{A}(\mathbf{y},\mathbf{z})] \right\}$$



What do these queries return ?

$$\left\{ \mathbf{x} \mid \exists \mathbf{y}. \mathbf{A}(\mathbf{x},\mathbf{y}) \right\}$$

Nodes that have at least one child:  $\{1,2,3\}$ 



A encodes the directed edges of a graph ("arcs")

$$\left\{ \mathbf{x} \mid \exists \mathbf{y}, \mathbf{z}, \mathbf{u}. [A(\mathbf{x}, \mathbf{y}) \land A(\mathbf{y}, \mathbf{z}) \land A(\mathbf{z}, \mathbf{u})] \right\}$$

Nodes that have a great-grand-child: {1,2}  $\{ (x,y) \mid \forall z.[A(x,z) \rightarrow A(y,z)] \}$   $\begin{array}{l} \forall z \in A(x,z) \rightarrow A(y,z) = 0 \\ \forall z \in A(y$ 



What do these queries return ?

$$\left\{ \mathbf{x} \mid \exists \mathbf{y}. \mathbf{A}(\mathbf{x}, \mathbf{y}) \right\}$$

Nodes that have at least one child:  $\{1,2,3\}$ 

 $\begin{array}{c|cccc}
A: & 1 & 2 \\
& 2 & 1 \\
& 2 & 3 \\
& 1 & 4 \\
& 3 & 4 \\
\end{array}$ 

A encodes the directed edges of a graph ("arcs")

$$\left\{ \mathbf{x} \mid \exists \mathbf{y}, \mathbf{z}, \mathbf{u}. [A(\mathbf{x}, \mathbf{y}) \land A(\mathbf{y}, \mathbf{z}) \land A(\mathbf{z}, \mathbf{u})] \right\}$$

Nodes that have a great-grand-child: {1,2}

 $\begin{cases} (x,y) \mid \forall z.[A(x,z) \rightarrow A(y,z)] \\ \text{Every child of x is a child of y.} \end{cases} \quad \begin{array}{l} \text{Which of the} \\ \text{following tuples} \\ \text{fulfill the condition?} \\ (1,1) \quad (4,4) \quad (1,3) \quad (3,1) \quad (4,1) \end{cases}$ 



What do these queries return ?

$$\left\{ \mathbf{x} \mid \exists \mathbf{y}. A(\mathbf{x}, \mathbf{y}) \right\}$$

Nodes that have at least one child:  $\{1,2,3\}$ 



4

A encodes the directed edges of a graph ("arcs")

$$\left\{ \mathbf{x} \mid \exists \mathbf{y}, \mathbf{z}, \mathbf{u}. [A(\mathbf{x}, \mathbf{y}) \land A(\mathbf{y}, \mathbf{z}) \land A(\mathbf{z}, \mathbf{u})] \right\}$$

Nodes that have a great-grand-child: {1,2}

 $\begin{cases} (x,y) \mid \forall z.[A(x,z) \rightarrow A(y,z)] \\ \forall vz.[A(x,z) \rightarrow A(y,z)] \\ Every child of x is a child of y. fulfill the condition? \\ (1,1) (4,4) (1,3) (3,1) (4,1) \\ (1,1),(2,2),(3,1),(3,3),(4,1),(4,2),(4,3),(4,4) \\ \end{cases}$ 

Likes(person, drink) Frequents(person, bar) Serves(bar, drink)



What does the following query compute?

$$\{x \mid \forall y. [Frequents(x,y) \rightarrow \exists z. [Serves(y,z) \land Likes(x,z)]\}$$

Likes(person, drink) Frequents(person, bar) Serves(bar, drink)



What does the following query compute?

$$\{x \mid \forall y. [Frequents(x,y) \rightarrow \exists z. [Serves(y,z) \land Likes(x,z)]\}$$

Find drinkers that frequent <u>only</u> bars that serve <u>some</u> drink they like.

Is this query domain independent? If not, how to fix?

Likes(person, drink) Frequents(person, bar) Serves(bar, drink)



#### What does the following query compute?

 $\exists w.[Frequents(x,w) \land ... \land Are those two options to \\\exists w.[Likes(x,w) \land ... \land make it safe identical \\ x \mid \forall y.[Frequents(x,y) \rightarrow \exists z.[Serves(y,z) \land Likes(x,z)] \end{cases}$ 

Find drinkers that frequent <u>only</u> bars that serve <u>some</u> drink they like.

Careful! This query is not domain independent.

Likes(person, drink) Frequents(person, bar) Serves(bar, drink)



#### What does the following query compute?

 $\exists w.[Frequents(x,w) \land ... \\ \exists w.[Likes(x,w) \land ... \\ likes a drink but does not frequent a bar be returned? \\ x | \forall y.[Frequents(x,y) \rightarrow \exists z.[Serves(y,z) \land Likes(x,z)] \}$ 

Find drinkers that frequent <u>only</u> bars that serve <u>some</u> drink they like.

Challenge: write this query without the  $\forall$  quantifier! And then in SQL

#### The person/bar/drinks example

Challenge: write these in SQL. Solutions at: <u>https://demo.queryvis.com</u> Likes(person, drink) Frequents(person, bar) Serves(bar, drink)



Find persons that frequent some bar that serves some drink they like.

Find persons that frequent only bars that serve some drink they like.

 $\left\{ x \mid \forall y. [Frequents(x,y) \rightarrow \exists z. [Serves(y,z) \land Likes(x,z)] \right\}$ 

Find persons that frequent some bar that serves only drinks they like.

Find persons that frequent <u>only</u> bars that serve <u>only</u> drinks they like.

(= Find persons who like all drinks that are served in all the bars they visit.)

(= Find persons for which there does not exist a bar they frequent that serves a drink they do not like.)

