Updated 2/23/2022

# Topic 2: Complexity of Query Evaluation Unit 1: Conjunctive Queries Lecture 11

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp22)

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

2/22/2022

Updated 2/23/2022

# Topic 2: Complexity of Query Evaluation Unit 1: Conjunctive Queries Lecture 11

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp22)

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

2/22/2022

Let L be a database query language.

• The Query Evaluation Problem:

• The Query Equivalence Problem:

#### • The Query Containment Problem:

Let L be a database query language.

- The Query Evaluation Problem:
  - "Given a query q in L and a database instance D, evaluate q(D)"
  - That's the main problem in query processing.
- The Query Equivalence Problem:

• The Query Containment Problem:

#### Let L be a database query language.

- The Query Evaluation Problem:
  - "Given a query q in L and a database instance D, evaluate q(D)"
  - That's the main problem in query processing.
- The Query Equivalence Problem:
  - "Given two queries q and q' in L, is it the case that  $q \equiv q'$ ?"
    - i.e., is it the case that, for every database instance D, we have that q(D) = q'(D)?
  - This problem underlies query optimization: transform a given query to an equivalent more efficient one.
- The Query Containment Problem:

#### Let L be a database query language.

- The Query Evaluation Problem:
  - "Given a query q in L and a database instance D, evaluate q(D)"
  - That's the main problem in query processing.
- The Query Equivalence Problem:
  - "Given two queries q and q' in L, is it the case that  $q \equiv q'$ ?"
    - i.e., is it the case that, for every database instance D, we have that q(D) = q'(D)?
  - This problem underlies query optimization: transform a given query to an equivalent more efficient one.
- The Query Containment Problem:
  - "Given two queries q and q' in L, is it the case that  $q \subseteq q'$ ?"

Based on Phokion Kolaitis' "Logic and Databases" series at Simons Institute, 2016. <u>https://simons.berkeley.edu/talks/logic-and-databases</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



 $\subseteq a'(\mathcal{D})$ 

# Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
  - CQ equivalence and containment
  - Graph homomorphisms
  - Homomorphism beyond graphs
  - CQ containment
  - CQ minimization
- T2-2: Equivalence Beyond CQs
  - Union of CQs, and inequalities
  - Union of CQs equivalence under bag semantics
  - Nested queries
  - Tree pattern queries

Why bother about Query Containment

 The Query Containment Problem and Query Equivalence Problem are closely related to each other:

-  $q \equiv q'$  if and only if ? -  $q \subseteq q'$  if and only if ?





Why bother about Query Containment

- The Query Containment Problem and Query Equivalence Problem are closely related to each other:
  - $q \equiv q'$  if and only if
    - $q \subseteq q'$  and  $q \supseteq q'$



Why bother about Query Containment

- The Query Containment Problem and Query Equivalence Problem are closely related to each other:
  - $q \equiv q'$  if and only if •  $q \subseteq q'$  and  $q \supseteq q'$
  - $q \subseteq q'$  if and only if
    - $q \equiv (q \cap q')$



• Thm: The Query Equivalence Problem for relational calculus (RC) queries is...



• Thm: The Query Equivalence Problem for relational calculus (RC) queries is...

... undecidable 😕

A decision problem is <u>undecidable</u> if it is impossible to construct an algorithm that always leads to a correct yes-or-no answer.

- Proof: using <u>Trakhtenbrot's Theorem</u> (1949):
  - The <u>Finite Validity Problem</u> (problem of validity in FOL on the class of all finite models) is undecidable.
     a formula is <u>valid</u> if it comes out as true (or "satisfied") under all admissible assignments of meaning to that formula within the intended semantics for the logical language

what problem do we have to reduce to what other problem

Tip:  $A \leq B$ : reduction from A to B. Means: B could be used to solve A. But A is hard ...



• Thm: The Query Equivalence Problem for relational calculus (RC) queries is...

A decision problem is undecidable if it is impossible to construct an algorithm that always leads to a correct yes-or-no answer.

Proof: using Trakhtenprot's Theorem (1949):

... undecidable 😕

- The Finite Validity Problem problem of validity in FOL on the class of all finite models) is -a formula is valid if it comes out as true (or "satisfied") under all admissible assignments undecidable. of meaning to that formula within the intended semantics for the logical language Finite Validity Problem ≤ Query Equivalence Problem

Tip:  $A \preccurlyeq B$ : reduction from A to B. Means: B could be used to solve A. But A is hard ...

Corollary: The Query Containment Problem for RC is undecidable.



how



• Thm: The Query Equivalence Problem for relational calculus (RC) queries is...

... undecidable 😕

A decision problem is <u>undecidable</u> if it is impossible to construct an algorithm that always leads to a correct yes-or-no answer.

- Proof: using Trakhtenbrot's Theorem (1949):
  - The <u>Finite Validity Problem</u> (problem of validity in FOL on the class of all finite models) is a formula is valid if it comes out as true (or "satisfied") under all admissible assignments undecidable. of meaning to that formula within the intended semantics for the logical language
     Finite Validity Problem ≤ Query Equivalence Problem
  - - Take a fixed finitely valid RC sentence  $\psi$ , and assume you can solve the query equivalence problem. Then for every RC sentence  $\varphi$ , we could solve validity: Tip: A  $\leq$  B: reduction from A to B.  $\varphi$  is finitely valid  $\Leftrightarrow \varphi \equiv \psi$ . Means: B could be used to solve A. But A is hard ...
- Corollary: The Query Containment Problem for RC is undecidable.



• Thm: The Query Equivalence Problem for relational calculus (RC) queries is...

... undecidable 😕

A decision problem is undecidable if it is impossible to construct an algorithm that always leads to a correct yes-or-no answer.

- Proof: using Trakhtenbrot's Theorem (1949):
  - The <u>Finite Validity Problem</u> (problem of validity in FOL on the class of all finite models) is a formula is valid if it comes out as true (or "satisfied") under all admissible assignments undecidable. of meaning to that formula within the intended semantics for the logical language
     Finite Validity Problem ≤ Query Equivalence Problem
  - - Take a fixed finitely valid RC sentence  $\psi$ , and assume you can solve the query equivalence problem. Then for every RC sentence  $\varphi$ , we could solve validity: Tip: A  $\leq$  B: reduction from A to B.  $\varphi$  is finitely valid  $\Leftrightarrow \varphi \equiv \psi$ . Means: B could be used to solve A. But A is hard ...
- Corollary: The Query Containment Problem for RC is undecidable.

- Proof: Query Equivalence  $\leq$  Query Containment, since  $q \equiv q' \Leftrightarrow (q \subseteq q' \text{ and } q' \supseteq q)$ 

Updated 2/25/2022

# Topic 2: Complexity of Query Evaluation Unit 1: Conjunctive Queries (continued) Lecture 12

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp22)

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

2/25/2022

# Complexity of the Query Evaluation Problem

- The Query Evaluation Problem for Relational Calculus:
  - Given a RC formula  $\varphi$  and a database instance D, find  $\varphi^{adom}(D)$ .
- Theorem: The Query Evaluation Problem for Relational Calculus is ... ... PSPACE-complete.
  - PSPACE: decision problems, can be solved using an amount of memory that is polynomial in the input length (~ in polynomial amount of space).
  - PSPACE-complete: PSPACE + every other PSPACE problem can be transformed to it in polynomial time (PSPACE-hard)
- Proof: We need to show both
  - This problem is in PSPACE.
  - This problem is PSPACE-hard. (We only focus on this task)

Complexity of the Query Evaluation Problem

- Theorem: The Query Evaluation Problem for Relational Calculus is PSPACE-hard.
- Reduction uses QBF Quantified Boolean Formulas
  - Given QBF  $\forall x_1 \exists x_2 \dots \forall x_k \psi$ ,
  - is it true or false (notice every variable is <u>quantified = bound</u> at beginning of <u>sentence</u>, there are no free variables)
- Proof
  - Show that QBF  $\leq p$  Query Evaluation for Relational Calculus

# Complexity of the Query Evaluation Problem

Proof: Show that QBF  $\leq p$  Query Evaluation for Relational Calculus

- Given QBF  $\forall x_1 \exists x_2 \dots \forall x_k \psi$ ,
- Let V and P be two unary relation symbols
- Obtain  $\psi^*$  from  $\psi$  by replacing  $x_i$  by  $P(x_i)$ , and  $\neg x_i$  by  $\neg P(x_i)$
- Let D be the database instance with V = {0,1}, P={1}.
- Then the following statements are equivalent:
  - $\forall x_1 \exists x_2 \dots \forall x_k \psi$  is true
  - $\forall x_1 (V(x_1) \rightarrow \exists x_2 (V(x_2) \land (... \forall x_k (V(x_k) \rightarrow \psi^*))...) \text{ is true on } D.$

# Sublanguages of Relational Calculus

 Question: Are there interesting sublanguages of relational calculus for which the Query Containment Problem and the Query Evaluation Problem are "easier" than the full relational calculus?

- Answer:
  - Yes, the language of Conjunctive Queries (CQs) is such a sublanguage.
  - Moreover, conjunctive queries are the most frequently asked queries against relational databases.

# Conjunctive Queries (CQs)

- Definition:
  - A CQ is a query expressible by a RC formula in prenex normal form built from atomic formulas  $R(y_1,...,y_n)$ , and ∧ and ∃ only.

{  $(x_1,...,x_k): \exists z_1 ... \exists z_m \phi(x_1,...,x_k, z_1,...,z_k)$  },

- where  $\phi(x_1, ..., x_k, z_1, ..., z_k)$  is a conjunction of atomic formulas of the form  $R(y_1, ..., y_m)$ .
- <u>Prenex formula</u>: prefix (quantifiers & bound variables), then quantifier-free part
- Equivalently, a CQ is a query expressible by a RA expression of the form
  - $\pi_X(\sigma_{\Theta}(R_1 \times ... \times R_n))$ , where
  - Θ is a conjunction of equality atomic formulas (equijoin).
- Equivalently, a CQ is a query expressible by an SQL expression of the form
  - SELECT <list of attributes>
    - FROM <list of relation names>

# Conjunctive Queries (CQs)

- Definition:
  - A CQ is a query expressible by a RC formula in prenex normal form built from atomic formulas  $R(y_1,...,y_n)$ , and ∧ and ∃ only.

{  $(x_1,...,x_k): \exists z_1 ... \exists z_m \phi(x_1,...,x_k, z_1,...,z_k)$  },

- where  $\phi(x_1, ..., x_k, z_1, ..., z_k)$  is a conjunction of atomic formulas of the form  $R(y_1, ..., y_m)$ .
- Equivalently, a CQ can be written as a logic-programming rule:

 $Q(x_1,...,x_k) := R_1(u_1), ..., R_n(u_n)$ , where

- Each variable x<sub>i</sub> occurs in the right-hand side of the rule.
- Each **u**<sub>i</sub> is a tuple of variables (not necessarily distinct)
- The variables occurring in the right-hand side (the body), but not in the left-hand side (the head) of the rule are existentially quantified (but the quantifiers are not displayed).





• Return paths of Length 2: (binary output)

RC: ? RA: ? Datalog: ?

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

Examples of Conjunctive Queries

• Return paths of Length 2: (binary output)

RC:  $\{(x, y) \mid \exists z [E(x, z) \land E(z, y)]\}$ 

RA:

Datalog:



HOMOM.

150 M.









- Return paths of Length 2: (binary output)
  - RC:  $\{(x, y) \mid \exists z [E(x, z) \land E(z, y)]\}$
  - RA:  $\pi_{1,4}(\sigma_{\$2=\$3}(E \times E))$  unnamed perspective Datalog: ?

*E*(from, to)



• Return paths of Length 2: (binary output)

RC: 
$$\{(x, y) \mid \exists z [E(x, z) \land E(z, y)]\}$$

RA:  $\pi_{1,4}(\sigma_{\$2=\$3}(E \times E))$  unnamed perspective

Datalog: Q(x,y) := E(x,z), E(z,y)

• Cycle of Length 3: (Boolean query)

Datalog:

RC:





- Return paths of Length 2: (binary output)
  - RC:  $\{(x, y) \mid \exists z [E(x, z) \land E(z, y)]\}$
  - RA:  $\pi_{1,4}(\sigma_{\$2=\$3}(E \times E))$  unnamed perspective

Datalog: Q(x,y) := E(x,z), E(z,y)

• Cycle of Length 3: (Boolean query) RC:  $\exists z \exists z \exists z [E(x, y) \land E(y, z) \land E(z, x)]$ Datalog: **?** 





• Return paths of Length 2: (binary output)

RC: 
$$\{(x, y) \mid \exists z [E(x, z) \land E(z, y)]\}$$

RA:  $\pi_{1,4}(\sigma_{\$2=\$3}(E \times E))$  unnamed perspective

Datalog: Q(x,y) := E(x,z), E(z,y)

• Cycle of Length 3: (Boolean query)

RC: 
$$\exists z \exists z \exists z [E(x,y) \land E(y,z) \land E(z,x)]$$

Datalog: Q := E(x, y), E(y, z), E(z, x)

# **Conjunctive Queries**

- Every natural join is a conjunctive query with ... ... no existentially quantified variables
- Example: Given P(A,B,C), R(B,C,D)
  - P  $\bowtie$  R = {(x,y,z,w): P(x,y,z)  $\land$  R(y,z,w)}
  - q(x,y,z,w) := P(x,y,z), R(y,z,w)

(no variables are existentially quantified)

- SELECT P.A, P.B, P.C, R.D
  FROM P, R
  WHERE P.B = R.B AND P.C = R.C
- Conjunctive queries are also known as SPJ-queries (SELECT-PROJECT-JOIN queries)

## Conjunctive Query Evaluation and Containment

- Definition: Two fundamental problems about CQs
  - Conjunctive Query Evaluation (CQE):
    - Given a conjunctive query q and an instance D, find q(D).
  - Conjunctive Query Containment (CQC):
    - Given two k-ary conjunctive queries q<sub>1</sub> and q<sub>2</sub>, is it true that q<sub>1</sub> ⊆ q<sub>2</sub>?
       (i.e., for every instance D, we have that q1(D) ⊆ q2(D))
    - Given two Boolean conjunctive queries q1 and q2, is it true that q<sub>1</sub> ⊨ q<sub>2</sub>? (that is, for all D, if D ⊨ q<sub>1</sub>, then D ⊨ q<sub>2</sub>)?
- Notice that CQC is logical implication.
- Later today: connection to homomorphisms

# Vardi's Taxonomy of the Query Evaluation Problem

M.Y Vardi, "The Complexity of Relational Query Languages", 1982

- Definition: Let L be a database query language.
  - The combined complexity of L is the decision problem:
    - given an L-sentence and a database instance D, is  $\varphi$  true on D?
    - In symbols, does  $D \models \varphi$  (does D satisfy  $\varphi$ )?



- given a database instance D, does  $D \models \phi$ ?
- The query complexity of L is the family of the following decision problems P<sub>D</sub>, where D is a database instance:
  - given an L-sentence  $\varphi$ , does D  $\models \varphi$ ?

Vardi's Taxonomy of the Query Evaluation Problem

Vardi's "empirical" discovery:

- For most query languages L:
  - The data complexity of L is of lower complexity than both the combined complexity of L and the query complexity of L.
  - The query complexity of L can be as hard as the combined complexity of L.

#### Taxonomy of the Query Evaluation Problem for Relational Calculus

#### **Complexity Classes**



The Query Evaluation Problem for Relational Calculus

Problem	Complexity
Combined	PSPACE-complete
Complexity	
Query Complexity	• In PSPACE
	<ul> <li>can be PSPACE-</li> </ul>
	complete
Data Complexity	In LOGSPACE
	$\wedge$

#### Summary

- Relational Algebra and Relational Calculus have "essentially" the same expressive power.
- The Query Equivalence Problem for Relational Calculus is undecidable.
  - Therefore also the Query Containment Problem
- The Query Evaluation Problem for Relational Calculus:
  - Data Complexity is in LOGSPACE
  - Combined Complexity is PSPACE-complete
  - Query Complexity is PSPACE-complete.

# Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
  - CQ equivalence and containment
  - Graph homomorphisms
  - Homomorphism beyond graphs
  - CQ containment
  - CQ minimization
- T2-2: Equivalence Beyond CQs
  - Union of CQs, and inequalities
  - Union of CQs equivalence under bag semantics
  - Nested queries
  - Tree pattern queries

## Injective, Surjective, and Bijective functions $f: X \rightarrow Y$


### Injective, Surjective, and Bijective functions $f: X \rightarrow Y$



Sources: https://en.wikipedia.org/wiki/Bijection\_injection\_and\_surjection\_atalab.github.io/cs7240/

#### Injective, Surjective, and Bijective functions $f: X \to Y$



Sources: https://en.wikipedia.org/wiki/Bijection, injection, and surjection Wolfgang Catterburg, Windples of sediable data management: https://hortheastern-datalab.github.io/cs7240/

#### Injective, Surjective, and Bijective functions $f: X \rightarrow Y$

Function



Injective function

> Bijective function

maps each <u>argument</u> (element from its <u>domain</u>) to exactly one <u>image</u> (element in its <u>codomain</u>)  $\forall x \in X, \exists ! y \in Y[y = f(x)]$ }

("one-to-one"): each element of the codomain is mapped to by <u>at most one</u> element of the domain (i.e. distinct elements of the domain map to distinct elements in the codomain)

$$\forall x, x' \in X. [f(x) = f(x') \Rightarrow x = x'] \forall x, x' \in X. [x \neq x' \Rightarrow f(x) \neq f(x')]$$

("onto"): each element of the codomain is mapped to by <u>at least one</u> element of the domain (i.e. the image and the codomain of the function are equal)  $\forall y \in Y, x \in X[y = f(x)]$ 

### Injective, Surjective, and Bijective functions $f: X \rightarrow Y$



Sources: https://en.wikipedia.org/wiki/Bijection\_injection\_and\_surjection\_atalab.github.io/cs7240/





Wolfgang Gatterbauer. P6i7240esPoincipleblefcoadablendgtamentalatope//inottpes//inottpes//inortbleastepredatabab./distrubio/cs7240/





≻

≻€

►



not a mapping (or function)!

injective function (or one-to-one): maps distinct elements of its domain to <u>distinct elements of its codomain</u>

Wolfgang Gatterbauer. Ebinables Poincipleblef doed a bendgemean and the period of the asterbard at the binable of the



not a mapping (or function)!

injective function (or one-to-one): maps distinct elements of its domain to <u>distinct elements of its codomain</u>

surjective (or onto): every element y in the codomain Y of f has at least one element x in the domain that maps to it



not a mapping (or function)!

injective function (or one-to-one): maps distinct elements of its domain to <u>distinct elements of its codomain</u>

surjective (or onto): every element y in the codomain Y of f has at least one element x in the domain that maps to it

injective & surjective



not a mapping (or function)!

injective function (or one-to-one): maps distinct elements of its domain to <u>distinct elements of its codomain</u>

surjective (or onto): every element y in the codomain Y of f has at least one element x in the domain that maps to it

injective & surjective

neighter



not a mapping (or function)!

injective function (or one-to-one): maps distinct elements of its domain to <u>distinct elements of its codomain</u>

surjective (or onto): every element y in the codomain Y of f has at least one element x in the domain that maps to it

injective & surjective

neighter

not even a mapping!



Wolfgang Gatterbauer. PSin240esPoinciplablefcoadablandgtamentalattope//inbitpe//inbi

# Bijection, Injection, and Surjection



Sources: <a href="https://www.intechopen.com/books/protein-interactions/relating-protein-structure-and-function-through-a-bijection-and-its-implications-on-protein-structur">https://www.intechopen.com/books/protein-interactions/relating-protein-structure-and-function-through-a-bijection-and-its-implications-on-protein-structur,</a>

#### Bijection, Injection, and Surjection



A function not injective not surjective

An injective function not surjective

A surjective function not injective

A bijective function injective + surjective

Not a function

### We make a detour to Graph matching

• Finding a correspondence between the nodes and the edges of two graphs that satisfies some (more or less stringent) constraints



- A graph homomorphism *h* from graph  $G(V_G, E_G)$  to  $H(V_H, E_H)$ , is a mapping from  $V_G$  to  $V_H$  such that  $\{x, y\} \in E_G$  implies  $\{h(x), h(y)\} \in E_H$ 
  - "edge-preserving": if two nodes in G are linked by an edge, then they are mapped to two nodes in H that are also linked







- A graph homomorphism *h* from graph  $G(V_G, E_G)$  to  $H(V_H, E_H)$ , is a mapping from  $V_G$  to  $V_H$  such that  $\{x, y\} \in E_G$  implies  $\{h(x), h(y)\} \in E_H$ 
  - "edge-preserving": if two nodes in G are linked by an edge, then they are mapped to two nodes in H that are also linked



Wolfgang Gatterbauer, CS 7240: Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/



- A graph homomorphism *h* from graph  $G(V_G, E_G)$  to  $H(V_H, E_H)$ , is a mapping from  $V_G$  to  $V_H$  such that  $\{x, y\} \in E_G$  implies  $\{h(x), h(y)\} \in E_H$ 
  - "edge-preserving": if two nodes in G are linked by an edge, then they are mapped to two nodes in H that are also linked





- A graph homomorphism *h* from graph  $G(V_G, E_G)$  to  $H(V_H, E_H)$ , is a mapping from  $V_G$  to  $V_H$  such that  $\{x, y\} \in E_G$  implies  $\{h(x), h(y)\} \in E_H$ 
  - "edge-preserving": if two nodes in G are linked by an edge, then they are mapped to two nodes in H that are also linked





- Graphs  $G(V_G, E_G)$  and  $H(V_H, E_H)$  are isomorphic iff there is an invertible h from  $V_G$  to  $V_H$  s.t.  $\{x, y\} \in E_G$  iff  $\{h(u), h(v)\} \in E_H$ 
  - We need to find a one-to-one correspondence



Wolfgange Gatterbauer. CS 7240: Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/



- Graphs  $G(V_G, E_G)$  and  $H(V_H, E_H)$  are isomorphic iff there is an invertible h from  $V_G$  to  $V_H$  s.t.  $\{x, y\} \in E_G$  iff  $\{f(u), f(v)\} \in E_H$ 
  - We need to find a one-to-one correspondence



Wolfgange Gatterbauer, CS 7240: Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



- Graphs  $G(V_G, E_G)$  and  $H(V_H, E_H)$  are isomorphic iff there is an invertible h from  $V_G$  to  $V_H$  s.t.  $\{x, y\} \in E_G$  iff  $\{f(u), f(v)\} \in E_H$ 
  - We need to find a one-to-one correspondence





- Graphs  $G(V_G, E_G)$  and  $H(V_H, E_H)$  are isomorphic iff there is an invertible h from  $V_G$  to  $V_H$  s.t.  $\{x, y\} \in E_G$  iff  $\{f(u), f(v)\} \in E_H$ 
  - We need to find a one-to-one correspondence





Is there an isomorphismYes: $h: \{(1,a), (2,b), (3,d), (4,c), (5,e)\}$ from G to H?bijection = surjective and injective mapping

Wolfgang Gatterbauer, CS 7240: Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/

Another perspective on absorption and subsumption

#### Absorption (or the challenge with self-joins)



f is true if there is an edge

 $f = \exists x, y. \ x \land y \land (x, y) \in E$ 



#### Absorption (or the challenge with self-joins)



f is true if there is an edge

 $f = \exists x, y. \ x \land y \land (x, y) \in E$ 



 $f = ab \lor ac$ 



#### Absorption (or the challenge with self-joins)



f is true if there is an edge

 $f = \exists x, y. \ x \land y \land (x, y) \in E$ 



 $f = ab \lor ac$ 



#### Absorption

#### Absorption

 $(A \lor B) \land A = A$  $(A \min B) \max A = A$  $A \le B \Leftrightarrow A \min B = A$  $(A \lor B) \land (A \Rightarrow B) = A$  $(A \min B) = A, \text{ if } A \le B$  $\sim A \land B = A$ 

Two binary operations, V and A, are said to be connected by the absorption law if:  $a \lor (a \land b) = a \land (a \lor b) = a$ .

A set equipped with two **commutative**, **associative** and **idempotent** binary operations V ("join") and  $\Lambda$  ("meet") that are connected by the absorption law is called a **lattice**. Examples of lattices include Boolean algebras, the set of sets with union and intersection operators, and ordered sets with min and max operations. <u>https://en.wikipedia.org/wiki/Absorption\_law</u>

Updated 3/2/2021

# Topic 2: Complexity of Query Evaluation Unit 1: Conjunctive Queries Lecture 12

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp21)

https://northeastern-datalab.github.io/cs7240/sp21/ 3/2/2021

# Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
  - CQ equivalence and containment
  - Graph homomorphisms
  - Homomorphism beyond graphs
  - CQ containment
  - CQ minimization
- T2-2: Equivalence Beyond CQs
  - Union of CQs, and inequalities
  - Union of CQs equivalence under bag semantics
  - Nested queries
  - Tree pattern queries

# Graph Homomorphism beyond graphs

**Definition :** Let G and H be graphs. A homomorphism of G to H is a function  $f: V(G) \rightarrow V(H)$  such that

5

 $(x,y) \in E(G) \Rightarrow (f(x),f(y)) \in E(H).$ 

We sometimes write  $G \rightarrow H$  (G  $\rightarrow$  H) if there is a homomorphism (no homomorphism) of G to H

Definition of a homomorphism naturally extends to:

- digraphs (directed graphs)
- edge-colored graphs
- relational systems
- constraint satisfaction problems (CSPs)

















Can this assignment be extended to a homomorphism?

No, this assignment requires a loop on vertex 1 (in H)

Fxample bytRickBrewaterperaphiamemerphiametutorialtreestern-datalab.github.io/cs7240/





Fxampale bytRinkaBrewaterperaphiamemerahiamethtamethtamethtamethtametheastern-datalab.github.io/cs7240/



**Definition:** Let G and H be graphs. A homom. of G to H is a function f:  $V(G) \rightarrow V(H)$  s.t. that

 $(x,y) \in E(G) \Rightarrow (f(x),f(y)) \in E(H).$ 



Example bytRick Brewsterp Graph damomer phism thtorial tragge for the astern-datalab.github.io/cs7240/



**Definition:** Let G and H be graphs. A homom. of G to H is a function f:  $V(G) \rightarrow V(H)$  s.t. that

 $(x,y) \in E(G) \Rightarrow (f(x),f(y)) \in E(H).$ 




### An example



The quotient set of the partition (set of equivalence classes of the partition) is a subgraph of H.





### Some observations

### When does $G \rightarrow K_3$ hold? ( $K_3 = 3$ -clique = triangle)



## Some observations

When does  $G \rightarrow K_3$  hold? ( $K_3 = 3$ -clique = triangle)

iff G is 3-colorable

When does  $G \rightarrow K_n$  hold? ( $K_n = n$ -clique)



Some observations

When does  $G \rightarrow K_3$  hold? ( $K_3 = 3$ -clique = triangle)

iff G is 3-colorable

When does  $G \rightarrow K_n$  hold? ( $K_n = n$ -clique) iff G is n-colorable

Thus homomorphisms generalize colorings: Notation:  $G \rightarrow H$  is an H-coloring of G.

What is the complexity of testing for the existence of a homomorphism?



94



When does  $G \rightarrow K_3$  hold? ( $K_3 = 3$ -clique = triangle) iff G is 3-colorable

- When does  $G \rightarrow K_n$  hold? ( $K_n = n$ -clique) iff G is n-colorable
- Thus homomorphisms generalize colorings: Notation:  $G \rightarrow H$  is an H-coloring of G.

What is the complexity of testing for the existence of a homomorphism?

### Some observations





### The complexity of H-coloring

- Let H be a fixed graph. H-coloring Instance: A graph G.
- Question: Does G admit an H-coloring?

Theorem [Hell, Nesetril'90]: If H is bipartite or contains a loop, then H-colouring is polynomial time solvable; otherwise, H is NP-complete.





### Repeated variable names



In sentences with multiple quantifiers, <u>distinct variables do not need</u> <u>to range over distinct objects</u>! (cp. homomorphism vs. isomorphism)

which of the following formulas imply each other?

 $\forall x. \forall y. E(x,y)$ 

 $\forall x. E(x,x)$ 

 $\exists x. \exists y. E(x,y)$ 

 $\exists x. E(x,x)$ 

### Repeated variable names

In sentences with multiple quantifiers, <u>distinct variables do not need</u> <u>to range over distinct objects</u>! (cp. homomorphism vs. isomorphism)





### Repeated variable names



In sentences with multiple quantifiers, <u>distinct variables do not need</u> <u>to range over distinct objects</u>! (cp. homomorphism vs. isomorphism)





# A more abstract (general) view on homomorphisms

Homomorphisms on Binary Structures

- Definition (Binary algebraic structure): A binary algebraic structure is a set together with a binary operation on it. This is denoted by an ordered pair (S,\*) in which S is a set and \* is a binary operation on S.
- Definition (homomorphism of binary structures): Let (S,\*) and (S',∘) be binary structures. A homomorphism from (S,\*) to (S',∘) is a map h: S → S' that satisfies, for all x, y in S:

 $h(x \star y) = h(x) \circ h(y)$ 

• We can denote it by  $h: (S, \star) \longrightarrow (S', \circ)$ .

• Let  $f(x) = e^x$ . Is f a homomorphism b/w two binary structures?



- Let  $f(x) = e^x$ . Is f a homomorphism b/w two binary structures?
  - Yes, from the real numbers with addition ( $\mathbb{R}$ ,+) to  $f(x+y) = f(x) \cdot f(y)$
  - the positive real numbers with multiplication  $(\mathbb{R}^+, \cdot)$   $f:(\mathbb{R}, +) \longrightarrow (\mathbb{R}^+, \cdot)$
  - It is even an isomorphism!

The exponential map exp :  $\mathbb{R} \to \mathbb{R}^+$  defined by  $\exp(x) = e^x$ , where *e* is the base of the natural logarithm, is an isomorphism from  $(\mathbb{R}, +)$  to  $(\mathbb{R}^+, \times)$ . Exp is a bijection since it has an inverse function (namely  $\log_e$ ) and exp preserves the group operations since  $e^{x+y} = e^x e^y$ . In this example both the elements and the operations are different yet the two groups are isomorphic, that is, as groups they have identical structures.

• Let  $g(x) = e^{ix}$ . Is g also a homomorphism?

### ?

- Let  $f(x) = e^x$ . Is f a homomorphism b/w two binary structures?
  - Yes, from the real numbers with addition ( $\mathbb{R}$ ,+) to  $f(x+y) = f(x) \cdot f(y)$
  - the positive real numbers with multiplication  $(\mathbb{R}^+, \cdot)$   $f:(\mathbb{R}, +) \to (\mathbb{R}^+, \cdot)$
  - It is even an isomorphism!

The exponential map exp :  $\mathbb{R} \to \mathbb{R}^+$  defined by  $\exp(x) = e^x$ , where *e* is the base of the natural logarithm, is an isomorphism from  $(\mathbb{R}, +)$  to  $(\mathbb{R}^+, \times)$ . Exp is a bijection since it has an inverse function (namely  $\log_e$ ) and exp preserves the group operations since  $e^{x+y} = e^x e^y$ . In this example both the elements and the operations are different yet the two groups are isomorphic, that is, as groups they have identical structures.

- Let  $g(x) = e^{ix}$ . Is g also a homomorphism?
  - Yes, from the real numbers with addition ( $\mathbb{R}$ ,+) to
  - the unit circle in the complex plane with rotation



$$G = \mathbb{R} \text{ under } +$$

$$H = \{ z \in \mathbb{C} : |z| = 1 \}$$

$$= \text{Group under } \times$$

$$Hint:$$

$$Hint:$$

$$Every \ z \in \mathbb{C} \text{ with } |z| = 1$$

$$can be written as \ z = e^{i\theta}.$$

$$f: G \to H$$

$$x \mapsto e^{ix}$$

$$Show \ f(x + y) = f(x) \times f(y)$$

$$e^{i(x+y)} = e^{ix} \times e^{iy}$$

$$e^{ix+iy} = e^{ix} \times e^{iy}$$

$$e^{ix} \times e^{iy} = e^{ix} \times e^{iy}$$

$$f(0) = f(2\pi) = 1, \ f(2\pi\pi) = 1$$

$$f \text{ is not } 1-1$$



WolfgengoBatterbaverEPlen's perof scalable to a calendary agement Attps://horhteastervalable to a calendary and the comparent of the calendary and the calen

### Isomorphism

- **Definition**: A homomorphism of binary structures is called an isomorphism iff the corresponding map of sets is:
  - one-to-one (injective) and
  - onto (surjective).



### Some homomorphisms



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

### Pointers to related work

- Kolaitis. Logic and Databases. Logical Structures in Computation Boot Camp, Berkeley 2016. <u>https://simons.berkeley.edu/talks/logic-and-databases</u>
- Abiteboul, Hull, Vianu. Foundations of Databases. Addison Wesley, 1995. <u>http://webdam.inria.fr/Alice/</u>, Ch 2.1: Theoretical background, Ch 6.2: Conjunctive queries & homomorphisms & query containment, Ch 6.3: Undecidability of equivalence for calculus.
- Kolaitis, Vardi. Conjunctive-Query Containment and Constraint Satisfaction. JCSS 2000. <u>https://doi.org/10.1006/jcss.2000.1713</u>
- Vardi. The Complexity of Relational Query Languages. STOC 1982. <u>https://doi.org/10.1145/800070.802186</u>
- Vardi. Constraint satisfaction and database theory: a tutorial. PODS 2000. <u>https://doi.org/10.1145/335168.335209</u>

### Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
  - CQ equivalence and containment
  - Graph homomorphisms
  - Homomorphism beyond graphs
  - CQ containment
  - CQ minimization
- T2-2: Equivalence Beyond CQs
  - Union of CQs, and inequalities
  - Union of CQs equivalence under bag semantics
  - Nested queries
  - Tree pattern queries

### Query Equivalence

Two queries  $q_1$ ,  $q_2$  are equivalent, denoted  $q_1 \equiv q_2$ , if for every database instance D, we have  $q_1(D) = q_2(D)$ .

Query  $q_1$  is contained in query  $q_2$ , denoted  $q_1 \subseteq q_2$ , if for every database instance D, we have  $q_1(D) \subseteq q_2(D)$ 

#### Corollary

$$q_1 \equiv q_2$$
 is equivalent to  $(q_1 \subseteq q_2 \text{ and } q_1 \supseteq q_2)$ 

If queries are Boolean, then query containment = logical implication:  $q_1 \Leftrightarrow q_2$  is equivalent to

### Query Equivalence

Two queries  $q_1, q_2$  are equivalent, denoted  $q_1 \equiv q_2$ , if for every database instance D, we have  $q_1(D) = q_2(D)$ .

Query  $q_1$  is contained in query  $q_2$ , denoted  $q_1 \subseteq q_2$ for every database instance D, we have  $q_1(D) \subseteq q_2($ 

#### Corollary

<u>Corollary</u>  $q_1 \equiv q_2$  is equivalent to  $(q_1 \subseteq q_2 \text{ and } q_1 \supseteq q_2)$ 

If queries are Boolean, then query containment = logical implication:  $q_1 \Leftrightarrow q_2$  is equivalent to  $q_1 \Rightarrow q_2$  and  $q_1 \leftarrow q_2$ )



A homomorphism *h* from Boolean  $q_2$  to  $q_1$  is a function  $h: \operatorname{var}(q_2) \to \operatorname{var}(q_1) \cup \operatorname{const}(q_1)$  such that: for every atom  $R(x_1, x_2, ...)$  in  $q_2$ , there is an atom  $R(h(x_1), h(x_2), ...)$  in  $q_1$ 

#### **Example**

$q_1(x) := R(x,y), R(y,y), R(y,z)$	
q <sub>2</sub> (s) :- R(s,u), R(u,w), R(s,v), R(v,w), R(u,v)	
	$v$ $w$ $q_2(x)$
$q_1(x)$ Wolfgang Gatterbauer. Principles of sole data management: <u>https://northeastern-datalab.github.io/cs7240/</u>	u s 114



A homomorphism *h* from Boolean  $q_2$  to  $q_1$  is a function *h*: var $(q_2) \rightarrow$  var $(q_1) \cup$  const $(q_1)$  such that: for every atom  $R(x_1, x_2, ...)$  in  $q_2$ , there is an atom  $R(h(x_1), h(x_2), ...)$  in  $q_1$ 

#### **Example**

$q_1(x) := R(x,y), R(y,y), R($	r,z) s,v), R(v,w), R(u,v)	
V Z	Also: h <sub>2→1</sub> ': {s,u,v,w} →{y}; recall [Hell, Nesetril'90]	$v$ $w$ $q_2(x)$
$q_1(x)$ Wolfgang Gatterbauer. Principles of scale data mana	$ \rightarrow_{1}: \{(S,X), (U,Y), (V,Y), (W,Z)\} $ gement: <u>https://northeastern-datalab.github.io/cs7240/</u>	u s 115



A homomorphism *h* from Boolean  $q_2$  to  $q_1$  is a function  $h: var(q_2) \rightarrow var(q_1) \cup const(q_1)$  such that: for every atom  $R(x_1, x_2, ...)$  in  $q_2$ , there is an atom  $R(h(x_1), h(x_2), ...)$  in  $q_1$ 

#### Example

 $q_1(x) := R(x,y), R(y,y), R(y,z)$  $q_2(s) := R(s,u), R(u,w), R(s,v), R(v,w), R(u,v)$ 







A homomorphism h from Boolean  $q_2$  to  $q_1$  is a function h:  $var(q_2) \rightarrow var(q_1) \cup const(q_1)$  such that: for every atom  $R(x_1, x_2, ...)$  in  $q_2$ , there is an atom  $R(h(x_1), h(x_2), ...)$  in  $q_1$ 

#### Example

Wolfgang Gatterbauer. Principles of

 $q_1(x) := R(x,y), R(y,y), R(y,z)$  $q_{2}(s) := R(s,u), R(u,w), R(s,v), R(v,w), R(u,v)$ What about:  $h_{1 \to 2}$ : {(x,s),(y,v),(z,w)} ? Ζ  $q_1(x)$  $h_{2 \rightarrow 1}$ : {(s,x),(u,y),(v,y),(w,z)} x

 $q_2(x)$ S 117



A homomorphism *h* from Boolean  $q_2$  to  $q_1$  is a function *h*: var $(q_2) \rightarrow$  var $(q_1) \cup$  const $(q_1)$  such that: for every atom  $R(x_1, x_2, ...)$  in  $q_2$ , there is an atom  $R(h(x_1), h(x_2), ...)$  in  $q_1$ 

#### Example

 $\begin{aligned} & q_1(x) := R(x,y), \, R(y,y), \, R(y,z) \\ & q_2(s) := R(s,u), \, R(u,w), \, R(s,v), \, R(v,w), \, R(u,v) \ , \ R(v,v) \end{aligned}$ 







A homomorphism *h* from Boolean  $q_2$  to  $q_1$  is a function *h*: var $(q_2) \rightarrow$  var $(q_1) \cup$  const $(q_1)$  such that:

for every atom  $R(x_1, x_2, ...)$  in  $q_2$ , there is an atom  $R(h(x_1), h(x_2), ...)$  in  $q_1$ 



$$\exists x. E(x,x) ? \exists x. \exists y. E(x,y)$$

#### Example





A homomorphism *h* from Boolean  $q_2$  to  $q_1$  is a function *h*: var $(q_2) \rightarrow$  var $(q_1) \cup$  const $(q_1)$  such that:

for every atom  $R(x_1, x_2, ...)$  in  $q_2$ , there is an atom  $R(h(x_1), h(x_2), ...)$  in  $q_1$ 



### Canonical database



Definition (Canonical database)

Given a conjunctive query q, the canonical database  $D_c[q]$  is the database instance where each atom in q becomes a fact in the instance.

Example  $q_1(x) := R(x,y), R(y,y), R(y,z)$  $D_c[q] = ?$ 

### Canonical database

Definition (Canonical database)

Given a conjunctive query q, the canonical database  $D_c[q]$  is the database instance where each atom in q becomes a fact in the instance.

Example  $q_1(x) := R(x,y), R(y,y), R(y,z)$   $D_c[q] = \{R('x', 'y'), R('y', 'y'), R('y', 'z')\}$  $\equiv \{R(a,b), R(b,b), R(b,c)\}$ 

Just treat each variable as different constant 😳









We will only look at 2)  $\Rightarrow$  1)

If there is a homomorphism *h* from  $q_2$  to  $q_1$ , then  $q_1 \subseteq q_2$ 

1. Given  $h=h_{2\rightarrow 1}$ , we will show that for any D:  $q_1(D) \Rightarrow q_2(D)$ 

2. For  $q_1(D)$  to hold, there is a valuation v s.t.  $v(q_1) \in D$ 

3. We will show that the composition  $g = v \circ h$  is a valuation for  $q_2$ 



- If there is a homomorphism *h* from  $q_2$  to  $q_1$ , then  $q_1 \subseteq q_2$
- 1. Given  $h=h_{2\rightarrow 1}$ , we will show that for any D:  $q_1(D) \Rightarrow q_2(D)$
- 2. For  $q_1(D)$  to hold, there is a valuation v s.t.  $v(q_1) \in D$
- 3. We will show that the composition  $g = v \circ h$  is a valuation for  $q_2$ 
  - 3a. By definition of h, for every  $R(x_1, x_2, ...)$  in  $q_2$ ,  $R(h(x_1), h(x_2), ...)$  in  $q_1$
  - 3b. By definition of v, for every  $R(x_1, x_2, ...)$  in  $q_2$ ,  $R(v(h(x_1)), v(h(x_2)), ...)$  in D

- If there is a homomorphism *h* from  $q_2$  to  $q_1$ , then  $q_1 \subseteq q_2$
- 1. Given  $h=h_{2\rightarrow 1}$ , we will show that for any D:  $q_1(D) \Rightarrow q_2(D)$
- 2. For  $q_1(D)$  to hold, there is a valuation v s.t.  $v(q_1) \in D$
- 3. We will show that the composition  $g = v \circ h$  is a valuation for  $q_2$ 
  - 3a. By definition of h, for every  $R(x_1, x_2, ...)$  in  $q_2$ ,  $R(h(x_1), h(x_2), ...)$  in  $q_1$ 3b. By definition of v, for every  $R(x_1, x_2, ...)$  in  $q_2$ ,  $R(v(h(x_1)), v(h(x_2)), ...)$  in D

 $g=v \circ h$ 

g(x) = v(h(x))

#### Example

 $q_1() := R(x,y), R(y,y), R(y,z)$  $q_2() := R(s,u), R(u,w), R(s,v), R(v,w), R(u,v)$ 



- If there is a homomorphism *h* from  $q_2$  to  $q_1$ , then  $q_1 \subseteq q_2$
- 1. Given  $h=h_{2\rightarrow 1}$ , we will show that for any D:  $q_1(D) \Rightarrow q_2(D)$
- 2. For  $q_1(D)$  to hold, there is a valuation v s.t.  $v(q_1) \in D$
- 3. We will show that the composition  $g = v \circ h$  is a valuation for  $q_2$ 
  - 3a. By definition of *h*, for every  $R(x_1, x_2, ...)$  in  $q_2$ ,  $R(h(x_1), h(x_2), ...)$  in  $q_1$ 3b. By definition of *v*, for every  $R(x_1, x_2, ...)$  in  $q_2$ ,  $R(v(h(x_1)), v(h(x_2)), ...)$  in *D*

 $g=v \circ h$ 

g(x) = v(h(x))

#### **Example**


# [Chandra and Merlin 1977]

- If there is a homomorphism *h* from  $q_2$  to  $q_1$ , then  $q_1 \subseteq q_2$
- 1. Given  $h=h_{2\rightarrow 1}$ , we will show that for any D:  $q_1(D) \Rightarrow q_2(D)$
- 2. For  $q_1(D)$  to hold, there is a valuation v s.t.  $v(q_1) \in D$
- 3. We will show that the composition  $g = v \circ h$  is a valuation for  $q_2$ 
  - 3a. By definition of h, for every  $R(x_1, x_2, ...)$  in  $q_2$ ,  $R(h(x_1), h(x_2), ...)$  in  $q_1$ 3b. By definition of v, for every  $R(x_1, x_2, ...)$  in  $q_2$ ,  $R(v(h(x_1)), v(h(x_2)), ...)$  in D

 $g=v \circ h$ 

g(x) = v(h(x))

#### Example



Combined complexity of CQC and CQE

#### Corollary:

The following problems are NP-complete:

- 1) Given two (Boolean) conjunctive queries Q and Q', is  $Q \subseteq Q'$ ?
- 2) Given a Boolean conjunctive query Q and an instance D, does  $D \models Q$ ?

siz a

#### Proof:

(a) Membership in NP follows from the Homom. Theorem:  $Q \subseteq Q'$  if and only if there is a homomorphism h:  $Q' \rightarrow Q$ 

(b) NP-hardness follows from 3-Colorability:

G is 3-colorable if and only if  $Q^{K_3} \subseteq Q^{G_{.}}$ 

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

# The Complexity of Database Query Languages

	Relational	CQs
	Calculus	
Query Eval.:	In LOGSPACE	In LOGSPACE
Data Complexity	(hence, in P)	(hence, in P)
Query Eval.:	PSPACE-	NP-complete
Combined Compl.	complete	
Query Equivalence	Undecidable	NP-complete
& Containment		

### Containment of conjunctive queries

 $Q = \{ x \mid \exists x', x'', y ( R(xy) \land R(x'y) \land R(x'', 1)) \}$ 





### Another query

 $\exists x', x'', y (R(xy) \land R(x'y) \land R(x'', 1) \land x' = x'' \land R(z, z))$  $Q' = \exists x', y, z (R(xy) \land R(x'y) \land R(x', 1) \land R(z, z))$ 



Q'(I)

Question

### Definition: $Q' \subseteq Q$ if for all $I, Q'(I) \subseteq Q(I)$ $Q' \equiv Q$ if $Q' \subseteq Q$ and $Q \subseteq Q'$

#### Problem: given Q', Q, test whether $Q' \subseteq Q$

Central issue for query optimization

### If there is a homomorphism from Q to Q', Q' $\subseteq$ Q



# If $Q' \subseteq Q$ , there is a homomorphism from Q to Q'



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

2/23/22

# $Q' \subseteq Q$ iff there is a homomorphism from Q to Q'

#### The problem is NP-complete

# Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
  - CQ equivalence and containment
  - Graph homomorphisms
  - Homomorphism beyond graphs
  - CQ containment
  - CQ minimization
- T2-2: Equivalence Beyond CQs
  - Union of CQs, and inequalities
  - Union of CQs equivalence under bag semantics
  - Nested queries
  - Tree pattern queries

# Islands of Tractability of CQ Evaluation

- Major Research Program: Identify <u>tractable cases</u> of the combined complexity of conjunctive query evaluation.
- Note: Over the years, this program has been pursued by two different research communities:
  - The Database Theory community
  - The Constraint Satisfaction community
- Explanation:

**Constraint Satisfaction Problem** 

≡ (Feder-Vardi, 1993)

Homomorphism Problem

≡ (Chandra-Merlin, 1977)

**Conjunctive Query Evaluation** 

# Beyond Conjunctive Queries

• What can we say about query languages of intermediate expressive power between conjunctive queries and the full relational calculus?

 Conjunctive queries form the sublanguage of relational algebra obtained by using only cartesian product, projection, and selection with equality conditions.

• The next step would be to consider relational algebra expressions that also involve union.

# Beyond Conjunctive Queries

- Definition:
  - A union of conjunctive queries (UCQ) is a query expressible by an expression of the form  $q_1 \cup q_2 \cup ... \cup q_m$ , where each  $q_i$  is a conjunctive query.
  - A monotone query is a query expressible by a relational algebra expression which uses only union, cartesian product, projection, and selection with equality condition.
     Fact:
- Fact:
  - Every union of conjunctive queries is a monotone query.
  - Every monotone query is equivalent to a union of conjunctive queries, but
    - the union may have exponentially many disjuncts.
- (normal form for monotone queries).
  - Monotone queries are precisely the queries expressible by relational calculus expressions using A, V, and E only.

Unions of CQs and Monotone Queries Union of Conjunctive Queries (UCQ)



Given edge relation *E*(*A*,*B*), find paths of length 1 or 2



(unnamed RA)



Union of Conjunctive Queries (UCQ)

Given edge relation *E*(*A*,*B*), find paths of length 1 or 2

RA 
$$E \cup \pi_{1,4}(\sigma_{2=3}(E \times E))$$
  
RC ?

(unnamed RA)

#### Union of Conjunctive Queries (UCQ)

Given edge relation *E*(*A*,*B*), find paths of length 1 or 2

RA 
$$E \bigcup \pi_{1,4}(\sigma_{2=3}(E \times E))$$
 (unnamed RA)  
RC  $E(x_1, x_2) \lor \exists z [E(z, x_2) \land E(z, x_2)]$ 

#### Union of Conjunctive Queries (UCQ)

Given edge relation *E*(*A*,*B*), find paths of length 1 or 2

RA  $E \cup \pi_{1,4}(\sigma_{2=3}(E \times E))$  (unnamed RA) RC  $E(x_1, x_2) \lor \exists z [E(z, x_2) \land E(z, x_2)]$ 

Monotone Query

Assume schema R(A,B), S(A,B), T(B,C), V(B,C)

Is following query monotone  $?(R \cup S) \bowtie (T \cup V)$ 



### Union of Conjunctive Queries (UCQ)

Given edge relation *E*(*A*,*B*), find paths of length 1 or 2

RA 
$$E \cup \pi_{1,4}(\sigma_{2=3}(E \times E))$$
 (unnamed RA)

RC 
$$E(x_1, x_2) \lor \exists z [E(z, x_2) \land E(z, x_2)]$$

Monotone Query

Assume schema R(A,B), S(A,B), T(B,C), V(B,C)

Is following query monotone?  $(R \cup S) \bowtie (T \cup V)$ 

#### Equal to a UCQ?



### Union of Conjunctive Queries (UCQ)

Given edge relation *E*(*A*,*B*), find paths of length 1 or 2

RA 
$$E \cup \pi_{1,4}(\sigma_{2=3}(E \times E))$$
 (unnamed RA)

$$\mathsf{RC} \quad E(x_1, x_2) \lor \exists z [E(z, x_2) \land E(z, x_2)]$$

#### Monotone Query

Assume schema R(A,B), S(A,B), T(B,C), V(B,C) Is following query monotone?  $(R \cup S) \bowtie (T \cup V)) \lor (\top \bowtie \chi)$ Equal to a UCQ?  $(R \bowtie T) \cup (R \bowtie V) \cup (S \bowtie T) \cup (S \bowtie V)$ 



# The Containment Problem for Unions of CQs

THEOREM [Sagiv and Yannakakis 1981] Let  $q_1 \cup q_2 \cup \cdots \cup q_m$  and  $q'_1 \cup q'_2 \cup \cdots \cup q'_n$  be two UCQs. Then the following are equivalent:

1) 
$$q_1 \cup q_2 \cup \cdots \cup q_m \subseteq q'_1 \cup q'_2 \cup \cdots \cup q'_n$$

2) For every  $i \le m$ , there is  $j \le n$  such that  $q_i \subseteq q'_i$ 

#### Proof: Use the Homomorphism Theorem 1. $\Rightarrow$ 2. Since $D_{C}[q_{i}] \models q_{i}$ , we have that $D_{C}[q_{i}] \models q_{1} \cup q_{2} \cup ... \cup q_{m}$ hence $D_{C}[q_{i}] \models q'_{1} \cup q'_{2} \cup ... \cup q'_{n}$ , hence there is some $j \le n$ such that $D_{C}[q_{i}] \models q'_{j}$ , hence (by the Homomorphism Theorem) $q_{i} \subseteq q'_{j}$ .

#### 2. $\Rightarrow$ 1. This direction is obvious.

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

### The Complexity of Database Query Languages

	Relational	CQs	UCQs
	Calculus		
Query Eval.:	In LOGSPACE	In LOGSPACE	In LOGSPACE
Data Complexity	(hence <i>,</i> in P)	(hence, in P)	(hence, in P)
Query Eval.:	PSPACE-	NP-complete	NP-complete
Combined Compl.	complete		
Query Equivalence	Undecidable	NP-complete	NP-complete
& Containment			

### Monotone Queries

- Even though monotone queries have the same expressive power as unions of conjunctive queries, the containment problem for monotone queries has higher complexity than the containment problem for unions of conjunctive queries (syntax/complexity tradeoff)
- Theorem: Sagiv and Yannakakis 1982
  The containment problem for monotone queries is Π<sub>2</sub><sup>p</sup>complete.
- Note: The prototypical  $\Pi_2^p$ -complete problem is  $\forall \exists$ SAT, i.e., the restriction of QBF to formulas of the form

 $\forall x_1 \dots \forall x_m \exists y_1 \dots \exists y_n \varphi.$ 

# The Complexity of Database Query Languages

	Relational	CQs	UCQs	Monotone queries
	Calculus			
Query Eval.:	In LOGSPACE	In LOGSPACE	In Logspace	In LOGSPACE
Data Complexity	(hence <i>,</i> in P)	(hence, in P)	(hence, in P)	(hence, in P)
Query Eval.:	PSPACE-	NP-complete	NP-complete	NP-complete
Combined Compl.	complete			
Query Equivalence	Undecidable	NP-complete	NP-complete	Π <sub>2</sub> <sup>p</sup> -complete
& Containment				

# Conjunctive Queries with Inequalities

- Definition: Conjunctive queries with inequalities form the sublanguage of relational algebra obtained by using only cartesian product, projection, and selection with equality and inequality (≠, <, ≤) conditions.</li>
- Example: Q(x,y):- E(x,z), E(z,w), E(w,y),  $z \neq w$ , z < y.
- Theorem: (Klug 1988, van der Meyden 1992)
  - The query containment problem for conjunctive queries with inequalities is  $\Pi_2^{p}$ -complete.
  - The query evaluation problem for conjunctive queries with inequalities in NP-complete.

# The Complexity of Database Query Languages

	Relational	CQs	UCQs	Monotone queries /
	Calculus			CQs with inequalities
Query Eval.:	In LOGSPACE	In LOGSPACE	In LOGSPACE	In LOGSPACE
Data Complexity	(hence <i>,</i> in P)	(hence, in P)	(hence <i>,</i> in P)	(hence, in P)
Query Eval.:	PSPACE-	NP-complete	NP-complete	NP-complete
Combined Compl.	complete			
Query Equivalence	Undecidable	NP-complete	NP-complete	Π <sub>2</sub> <sup>p</sup> -complete
& Containment				

# Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
  - CQ equivalence and containment
  - Graph homomorphisms
  - Homomorphism beyond graphs
  - CQ containment
  - CQ minimization
- T2-2: Equivalence Beyond CQs
  - Union of CQs, and inequalities
  - Union of CQs equivalence under bag semantics
  - Nested queries
  - Tree pattern queries

Following slides are from Phokion Kolaitis's talk on "Logic and databases" at "Logical structures in Computation Boot Camp", Berkeley 2016: https://simons.berkeley.edu/talks/logic-and-databases

#### **Logic and Databases**

Phokion G. Kolaitis

UC Santa Cruz & IBM Research – Almaden

Lecture 4 – Part 1





Credit: Phokion Kolaitis: <u>https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01</u>

#### **Thematic Roadmap**

- ✓ Logic and Database Query Languages
  - Relational Algebra and Relational Calculus
  - Conjunctive queries and their variants
  - Datalog
- ✓ Query Evaluation, Query Containment, Query Equivalence
  - Decidability and Complexity
- ✓ Other Aspects of Conjunctive Query Evaluation
- Alternative Semantics of Queries
  - Bag Databases: Semantics and Conjunctive Query Containment
  - Probabilistic Databases: Semantics and Dichotomy Theorems for Conjunctive Query Evaluation
  - Inconsistent Databases: Semantics and Dichotomy Theorems

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

#### **Alternative Semantics**

- So far, we have examined logic and databases under classical semantics:
  - The database relations are sets.
  - Tarskian semantics are used to interpret queries definable be first-order formulas.
- Over the years, several different alternative semantics of queries have been investigated. We will discuss three such scenarios:
  - The database relations can be bags (multisets).
  - The databases may be probabilistic.
  - The databases may be inconsistent.

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

#### Sets vs. Multisets

Relation EMPLOYEE(name, dept, salary)

• Relational Algebra Expression:

 $\pi_{\text{salary}} \left( \sigma_{\text{dept} = \text{CS}} \left( \text{EMPLOYEE} \right) \right)$ 

• SQL query:

SELECT salary FROM EMPLOYEE WHERE dpt = 'CS'

- SQL returns a bag (multiset) of numbers in which a number may appear several times, provided different faculty had the same salary.
- SQL does not eliminate duplicates, in general, because:
  - Duplicates are important for aggregate queries (e.g., average)
  - Duplicate elimination takes nlogn time.

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

#### **Relational Algebra Under Bag Semantics**

Operation	Multiplicity	•	R <sub>1</sub>	<u>A B</u>
$\begin{array}{l} \text{Union} \\ \text{R}_1 \cup \ \text{R}_2 \end{array}$	m <sub>1</sub> + m <sub>2</sub>			1 2 1 2 2 3
$\frac{\text{Intersection}}{R_1 \cap R_2}$	min(m <sub>1</sub> , m <sub>2</sub> )	•	R <sub>2</sub>	<u>BC</u> 24 25
Product	$m_1 \times m_2$			20
$R_1 \times R_2$		•	$(R_1 \Join R_2)$	<u>ABC</u>
Projection and Selection	Duplicates are not eliminated			1 2 4 1 2 5 1 2 5

5

Source: Phokion Kolaitis: <u>https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01</u>

**Conjunctive Queries Under Bag Semantics** 

Chaudhuri & Vardi – 1993 Optimization of *Real* Conjunctive Queries

- Called for a re-examination of conjunctive-query optimization under bag semantics.
- In particular, they initiated the study of the containment problem for conjunctive queries under bag semantics.
- This problem has turned out to be *much more challenging* than originally perceived.

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

#### PROBLEMS

Problems worthy of attack prove their worth by hitting back.

in: Grooks by Piet Hein (1905-1996)

Source: Phokion Kolaitis: <u>https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01</u>

#### **Query Containment Under Set Semantics**

Class of Queries	Complexity of Query Containment
Conjunctive Queries	NP-complete Chandra & Merlin – 1977
Unions of Conjunctive Queries	NP-complete Sagiv & Yannakakis - 1980
Conjunctive Queries with $\neq$ , $\leq$ , $\geq$	Π <sub>2</sub> <sup>p</sup> -complete Klug 1988, van der Meyden -1992
First-Order (SQL) queries	Undecidable Trakhtenbrot - 1949

Source: Phokion Kolaitis: <u>https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01</u>

Bag Semantics vs. Set Semantics

• For bags  $R_1, R_2$ :  $R_1 \subseteq_{BAG} R_2$  if  $m(a, R_1) \le m(a, R_2)$ , for every tuple  $a \in \mathbb{R}$ 

- Q<sup>BAG</sup>(D) : Result of evaluating Q on (bag) database D.
- $Q_1 \subseteq_{BAG} Q_2$  if for every (bag) database D, we have that  $Q_1^{BAG}(D) \subseteq_{BAG} Q_2^{BAG}(D)$ .

#### Fact:

- $Q_1 \subseteq_{BAG} Q_2$  implies  $Q_1 \subseteq Q_2$ .
- The converse does not always hold.



Source: Phokion Kolaitis: <u>https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01</u>

#### **Bag Semantics vs. Set Semantics**

**Fact:**  $Q_1 \subseteq Q_2$  does not imply that  $Q_1 \subseteq_{BAG} Q_2$ .

#### **Example:**

- Q<sub>1</sub>(x) :- P(x), T(x)
  Q<sub>2</sub>(x) :- P(x)
- $Q_1 \subseteq Q_2$  (obvious from the definitions)
- $Q_1 \not\subseteq_{BAG} Q_2$
- Consider the (bag) instance  $D = \{P(a), T(a), T(a)\}$ . Then:
  - $Q_1(D) = \{a,a\}$
  - $Q_2(D) = \{a\}, \text{ so } Q_1(D) \notin Q_2(D).$

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01
Query Containment under Bag Semantics

- Chaudhuri & Vardi 1993 stated that: Under bag semantics, the containment problem for conjunctive queries is Π<sub>2</sub><sup>p</sup>-hard.
- Problem:
  - What is the exact complexity of the containment problem for conjunctive queries under bag semantics?
  - Is this problem decidable?

Source: Phokion Kolaitis: <u>https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01</u>

#### **Query Containment Under Bag Semantics**

- 23 years have passed since the containment problem for conjunctive queries under bag semantics was raised.
- Several attacks to solve this problem have failed.
- At least two technically flawed PhD theses on this problem have been produced.
- Chaudhuri and Vardi have withdrawn the claimed  $\Pi_2^p$ -hardness of this problem; no one has provided a proof.

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

#### **Query Containment Under Bag Semantics**

• The containment problem for conjunctive queries under bag semantics remains **open** to date.

- However, progress has been made towards the containment problem under bag semantics for the two main extensions of conjunctive queries:
  - Unions of conjunctive queries
  - Conjunctive queries with  $\neq$

Source: Phokion Kolaitis: <u>https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01</u>

#### **Unions of Conjunctive Queries**

**Theorem** (loannidis & Ramakrishnan – 1995): Under bag semantics, the containment problem for unions of conjunctive queries is **undecidable**.

Hint of Proof:

Reduction from Hilbert's 10<sup>th</sup> Problem.

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

Hilbert's 10<sup>th</sup> Problem

 Hilbert's 10<sup>th</sup> Problem – 1900 (10<sup>th</sup> in Hilbert's list of 23 problems)

15

Given a Diophantine equation with any number of unknown quantities and with rational integral numerical coefficients: To devise a process according to which it can be determined in a finite number of operations whether the equation is solvable in rational integers.

In effect, Hilbert's 10<sup>th</sup> Problem is: Find an algorithm for the following problem: Given a polynomial  $P(x_1,...,x_n)$  with integer coefficients, does it have an all-integer solution?



#### Hilbert's 10<sup>th</sup> Problem

- Hilbert's 10<sup>th</sup> Problem 1900
  - (10<sup>th</sup> in Hilbert's list of 23 problems)



Find an algorithm for the following problem:

Given a polynomial  $P(x_1,...,x_n)$  with integer coefficients, does it have an all-integer solution?

• Y. Matiyasevich – 1971

(building on M. Davis, H. Putnam, and J. Robinson)

Hilbert's 10<sup>th</sup> Problem is undecidable, hence no such algorithm exists.

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

16

#### Hilbert's 10<sup>th</sup> Problem

- Fact: The following variant of Hilbert's 10<sup>th</sup> Problem is undecidable:
  - Given two polynomials p<sub>1</sub>(x<sub>1</sub>,...x<sub>n</sub>) and p<sub>2</sub>(x<sub>1</sub>,...x<sub>n</sub>) with positive integer coefficients and no constant terms, is it true that p<sub>1</sub> ≤ p<sub>2</sub>?
    In other words, is it true that p<sub>1</sub>(a<sub>1</sub>,...,a<sub>n</sub>) ≤ p<sub>2</sub>(a<sub>1</sub>,...a<sub>n</sub>), for all positive integers a<sub>1</sub>,...,a<sub>n</sub>?
- Thus, there is no algorithm for deciding questions like:

$$- \text{ Is } 3x_1^4x_2x_3 + 2x_2x_3 \le x_1^6 + 5x_2x_3^?$$

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

17

#### Unions of Conjunctive Queries

Theorem (loannidis & Ramakrishnan – 1995): Under bag semantics, the containment problem for unions of conjunctive queries is **undecidable**.

#### Hint of Proof:

- Reduction from the previous variant of Hilbert's 10<sup>th</sup> Problem:
  - Use joins of unary relations to encode monomials (products of variables).
  - Use unions to encode sums of monomials.

18

#### Unions of Conjunctive Queries

Example: Consider the polynomial  $3x_1^4x_2x_3 + 2x_2x_3$ 

- The monomial  $x_1 4 x_2 x_3$  is encoded by the conjunctive query  $P_1(w), P_1(w), P_1(w), P_1(w), P_2(w), P_3(w).$
- The monomial  $x_2x_3$  is encoded by the conjunctive query  $P_{2}(w), P_{3}(w).$
- The polynomial \$x\_1^4x\_2x\_3\$ + 2x\_2x\_3\$ is encoded by the union having:
  three copies of \$P\_1(w), P\_1(w), P\_1(w), P\_1(w), P\_2(w), P\_3(w)\$
  - and
  - two copies of  $P_2(w), P_3(w)$ .

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

#### **Complexity of Query Containment**

<b>Class of Queries</b>	Complexity –	Complexity –
	Set Semantics	<b>Bag Semantics</b>
Conjunctive queries	NP-complete CM – 1977	
Unions of conj. queries	NP-complete SY - 1980	Undecidable IR - 1995
Conj. queries with $\neq$ , $\leq$ , $\geq$	П <sub>2</sub> <sup>p</sup> -complete vdM - 1992	
First-order (SQL) queries	Undecidable Trakhtenbrot - 1949	Undecidable

Source: Phokion Kolaitis: <u>https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01</u>

#### Conjunctive Queries with ≠

Theorem (Jayram, K ..., Vee – 2006): Under bag semantics, the containment problem for conjunctive queries with  $\neq$  is **undecidable**.

In fact, this problem is **undecidable** even if

- the queries use only a single relation of arity 2;
- the number of inequalities in the queries is at most some fixed (albeit huge) constant.

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

#### **Complexity of Query Containment**

Class of Queries	Complexity –	Complexity –
	Set Semantics	<b>Bag Semantics</b>
Conjunctive queries	NP-complete CM – 1977	Open
Unions of conj. queries	NP-complete SY - 1980	Undecidable IR - 1995
Conj. queries with $\neq$ , $\leq$ , $\geq$	Π <sub>2</sub> <sup>p</sup> -complete vdM - 1992	Undecidable JKV - 2006
First-order (SQL) queries	Undecidable Trakhtenbrot - 1949	Undecidable

Source: Phokion Kolaitis: <u>https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01</u>

#### Subsequent Developments

- Some progress has been made towards identifying special classes of conjunctive queries for which the containment problem under bag semantics is decidable.
  - Afrati, Damigos, Gergatsoulis 2010
    - Projection-free conjunctive queries.
  - Kopparty and Rossman 2011
    - A large class of boolean conjunctive queries on graphs.

Source: Phokion Kolaitis: https://simons.berkeley.edu/talks/phokion-kolaitis-2016-09-01

209

NOT ELISTS A (NOT EX

### Pointers to related work

- Kolaitis. Logic and Databases. Logical Structures in Computation Boot Camp, Berkeley 2016. <u>https://simons.berkeley.edu/talks/logic-and-databases</u>
- Abiteboul, Hull, Vianu. Foundations of Databases. Addison Wesley, 1995. <u>http://webdam.inria.fr/Alice/</u>, Ch 2.1: Theoretical background, Ch 6.2: Conjunctive queries & homomorphisms & query containment, Ch 6.3: Undecidability of equivalence for calculus.
- Chandra, Merlin. Optimal implementation of conjunctive queries in relational data bases. STOC 1977. <u>https://doi.org/10.1145/800105.803397</u>

Updated 3/5/2021

# Topic 2: Complexity of Query Evaluation Unit 1: Conjunctive Queries Lecture 13

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp21)

https://northeastern-datalab.github.io/cs7240/sp21/ 3/5/2021

# Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
  - CQ equivalence and containment
  - Graph homomorphisms
  - Homomorphism beyond graphs
  - CQ containment
  - CQ minimization
- T2-2: Equivalence Beyond CQs
  - Union of CQs, and inequalities
  - Union of CQs equivalence under bag semantics
  - Nested queries
  - Tree pattern queries

## Exercise: Find the Homomorphisms





Order of subgoals in the query does not matter (thus ~sets)

q<sub>2</sub>: {E(x,y),E(y,z),E(z,x)}

 $q_3: \{E(x,y), E(y,x)\}$ 

what is the containment relation between these queries ?

 $q_4: \{E(x,y), E(y,x), E(y,y)\}$   $q_5: \{E(x,x)\}$ 

### Exercise: Find the Homomorphisms





 $q_1: \{E(x,y), E(y,z), E(z,w)\}$ 

 $q_4$ : {E(x,y),E(y,x),E(y,y)}  $q_5: \{E(x,x)\}$ Х 🔶



Example bytendreas Piecisles of scalable data management: https://northeastern-datalab.github.io/cs7240/



### Exercise: Find the Homomorphisms







### Exercise: Find the Homomorphisms







 $q_1(x,y) := R(x,u), R(v,u), R(v,y)$   $var(q_1) = \{x, u, v, y\}$ 

 $q_2(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y)$  var $(q_2) = \{x, u, v, w, t, y\}$ 

Are these queries equivalent ?



 $q_{1}(x,y) := R(x,u), R(v,u), R(v,y) \qquad var(q_{1}) = \{x, u, v, y\}$   $q_{2}(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y) \qquad var(q_{2}) = \{x, u, v, w, t, y\}$ 

# $q_1 \rightarrow q_2$ Thus ?

which query contains the other?



 $q_{1}(x,y) := R(x,u), R(v,u), R(v,y) \qquad var(q_{1}) = \{x, u, v, y\}$   $q_{2}(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y) \qquad var(q_{2}) = \{x, u, v, w, t, y\}$ 





 $q_1(x,y) := R(x,u), R(v,u), R(v,y) \qquad var(q_1) = \{x, u, v, y\}$ 

 $q_2(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y)$   $var(q_2) = \{x, u, v, w, t, y\}$ 

Is there any homomorphism  $q_2 \rightarrow q_1$ and then  $q_1 \subseteq q_2$ ?



 $q_{1}(x,y) := R(x,u), R(y,u), R(y,v) \qquad \text{var}(q_{1}) = \{x, u, v, y\}$   $q_{2}(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y) \qquad \text{var}(q_{2}) = \{x, u, v, w, t, y\}$ 



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



# Minimizing Conjunctive Queries

- Goal: minimize the number of joins in a query
- Definition: A conjunctive query Q is minimal if...

?

# Minimizing Conjunctive Queries

- Goal: minimize the number of joins in a query
- Definition: A conjunctive query Q is minimal if there is no conjunctive query Q' such that:
  - 1. Q ≡ Q′
  - 2. Q' has fewer atoms than Q
- The task of CQ minimization is, given a conjunctive query Q, to compute a minimal one that is equivalent to Q

### Minimization by Deletion

**Theorem:** Consider a conjunctive query  $Q_1(x_1,...,x_k) := body_1$ . If  $Q_1$  is equivalent to a conjunctive query  $Q_2(y_1,...,y_k) := body_2$ where  $|body_2| < |body_1|$ , then  $Q_1$  is equivalent to a query  $Q_3(x_1,...,x_k) := body_3$  such that  $body_3 \subseteq body_1$ 

> The above theorem says that to minimize a conjunctive query  $Q_1(\mathbf{x})$  :- body we simply need to remove some atoms from body

Can be shown by exploiting the homomorphism theorem...

Conjunctive query minimization algorithm

Q :-E(x,y), E(y,z) Q':-E(x,y)  $Minimize(Q(\mathbf{x}) := body)$ Repeat { We know  $Q' \rightarrow Q$ Choose an atom  $\alpha \in body$ Thus:  $\mathbf{Q} \subset \mathbf{Q}'$ Remove  $\alpha$  from Q; let Q' be the new query  $\checkmark$ If there is a homomorphism from Q to Q',  $\checkmark$ •  $Q' \subset Q$   $Q \rightarrow Q'$ then body := body  $\setminus \{\alpha\}$ Until no atom can be removed}

#### Notice: the order in which we inspect subgoals doesn't matter

### Minimization Procedure: Example



a,b,c,d are constants

Q(x): R(x,y), R(x,b'), R(a',b'), R(u,c'), R(u,v), S(a',c',d')R(x,y), R(x,b'), R(a',b'), R(u,c'), R(u,v), S(a',c',d')

#### Is this query minimal?

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

### Minimization Procedure: Example



a,b,c,d are constants



92 TROE 1, FACSO

3

Is this query minimal

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

### Minimization Procedure: Example



a,b,c,d are constants



#### Is this query minimal?
### Minimization Procedure: Example



a,b,c,d are constants



Is this query minimal

### Minimization Procedure: Example



a,b,c,d are constants

Q(x) :- R(x,y), R(x,'b'), R('a','b'), R(u,'c'), R(u,v), S('a','c','d')  $\{y \rightarrow b\}$ R(x,'b'), R('a','b'), R(u,'c'), R(u,v), S('a','c','d') Q(x) :- $\{ V \longrightarrow C \}$ S('a','c','d') Q(x) :-R(x,'b'), R('a','b'), R(u,'c'), Minimal query  $\{x \rightarrow a\}$ R('a', b'), R(u,'c') S('a','c','d') Q('a') :=

Mapping  $x \rightarrow a$  is not valid since x is a distinguished variable!

# Uniqueness of Minimal Queries

**Natural question:** does the order in which we remove atoms from the body of the input conjunctive query matter?

**Theorem:** Consider a conjunctive query Q. Let  $Q_1$  and  $Q_2$  be minimal conjunctive queries such that  $Q_1 \equiv Q$  and  $Q_2 \equiv Q$ . Then,  $Q_1$  and  $Q_2$  are isomorphic (i.e., they are the same up to variable renaming)



CHURCH - ROSSER

Therefore, given a conjunctive query Q, the result of Minimization(Q) is unique (up to variable renaming) and is called the core of Q



Is this query still minimal?

### AND E1.university = 'Northeastern' AND E2.university = 'Northeastern'

AND E3.university = 'Northeastern' AND E4.university = 'Northeastern'

Employee(<u>name</u>, university, manager)



### Query Minimization for Views

#### NEU employees managed by NEU emp.:

CREATE VIEW NeuMentors AS SELECT DISTINCT E1.name,E1.manager FROM Employee E1, Employee E2 WHERE E1.manager = E2.name AND E1.university = 'Northeastern' AND E2.university= 'Northeastern'

#### ←This query / view is minimal F1 → F2



name	university	manager
Alice	Northeastern	Bob
Bob	Northeastern	Cecile
Cecile	Northeastern	

### NEU emp. managed by NEU emp. managed by NEU emp .:

SELECT DISTINCT N1.name FROM NeuMentors N1, NeuMentors N2 WHERE N1.manager = N2.name

#### ←This query is minimal



#### View expansion (when you run a SQL query on a view)

SELECT DISTINCT E1.name
FROM Employee E1, Employee E2, Employee E3, Employee E4
WHERE E1.manager = E2.name AND E1.manager = E3.name AND E3.manager = E4.name
AND E1.university = 'Northeastern' AND E2.university = 'Northeastern'
AND E3.university = 'Northeastern' AND E4.university = 'Northeastern'

#### E2 is redundant!



Employee(name, university, manager)



# Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
  - CQ equivalence and containment
  - Graph homomorphisms
  - Homomorphism beyond graphs
  - CQ containment
  - CQ minimization
- T2-2: Equivalence Beyond CQs
  - Union of CQs, and inequalities
  - Union of CQs equivalence under bag semantics
  - Nested queries
  - Tree pattern queries

# Equivalence of nested queries

- Query equivalence is one of the foundational questions in database theory (and practice?)
  - touches on logics and decidability
  - what modifications allow tractability
- Lots of work (and open questions) on query equivalence
  - But not so much on nested queries!
- Related to QueryViz project (<u>http://queryviz.com</u>) and two foundational questions on visual formalism:
  - 1. When can visual formalism *unambiguously* express logical statements?
  - 2. When can equivalent logical statements be transformed to each other by a sequence of visual transformations? (*Query equivalence*)

# Diagrammatic reasoning systems and their expressiveness

(2)

R





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

## Diagrammatic reasoning systems and their expressiveness



THE OGICAL

**Diagrams** are widely used in reasoning about problems in physics, mathematics, and logic, but have traditionally been considered to be only heuristic tools and not valid elements of mathematical proofs. This book challenges this prejudice against visualization in the history of logic and mathematics and provides a formal foundation for work on natural reasoning in a visual mode.

The author presents Venn diagrams as a formal system of representation equipped with its own syntax and semantics and specifies rules of transformation that make this system sound and complete. The system is then extended to the equivalent of a first-order monadic language. The soundness of these diagrammatic systems refutes the contention that graphical representation is misleading in reasoning. The validity of the transformation rules ensures that the correct application of the rules will not lead to fallacies. The book concludes with a discussion of some fundamental differences between graphical systems and linguistic systems.

This groundbreaking work will have important influence on research in logic, philosophy, and knowledge representation. objects. Conjunctive information is more naturally represented by diagrams than by linguistic formulæ. For example, a single Venn diagram can

Still, not all relations can be viewed as membership or inclusion. Shin has been careful throughout her book to restrict herself to monadic systems. Relations per se (polyadic predicates) are not considered. And while it may be true that the formation of a system (such as Venn-II) that is provably both sound and complete would help mitigate the prejudice

perception. In her discussion of perception she shows that disjunctive information is not representable in *any* system. In doing so she relies on

The logical status of diagrams, Sun-Joo Shin, Cambridge university press 1994. <u>https://doi.org/10.1017/CBO9780511574696</u>

# QueryViz

- Motivation: Can we create an automatic system that:
  - unambiguously visualizes the logical intent of a SQL query (thus no two different queries lead to an "identical" visualization; with "identical" to be formalized correctly)
  - for some important subset of nested queries
  - with visual diagrams that allow us to reason about SQL design patterns
- Related:
  - Lot's of interest on conjunctive queries equivalence. Now: For what fragment of nested queries is equivalence decidable (under set semantics)?
- Suggestion:
  - nested queries, with inequalities, without any disjunctions

Wolfgang GStraictersupperset a of a conjunctive squeries ub.io/cs7240/

what is the intend of this query?

```
SELECT
                                L1.drinker
                     FROM
                                Likes L1
                                NOT EXISTS
                     WHERE
                                (SELECT *
                                FROM
                                          Likes L2
                                          L1.drinker <> L2.drinker
                                WHERE
                                          NOT EXISTS
                                AND
                                           (SELECT *
                                          FROM
                                                     Likes L3
                                          WHERE
                                                     L3.drinker = L2.drinker
                                          AND
                                                     NOT EXISTS
                                                      (SELECT *
                                                     FROM
                                                                Likes L4
                                                     WHERE
                                                                L4.drinker = L1.drinker
                                                     AND
                                                                L4.beer = L3.beer)
                                AND
                                          NOT EXISTS
                                           (SELECT
                                          FROM
                                                     Likes L5
                                          WHERE
                                                     L5. drinker = L1. drinker
                                          AND
                                                     NOT EXISTS
                                                      (SELECT
                                                     FROM
                                                                Likes L6
                                                     WHERE
                                                                L6.drinker = L2.drinker
Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>eer= L5.beer)))
```

2019/10/21

Likes		
drinker		
beer		



Unique set query: "Find drinkers that like a unique set of beers."



2019/10/21



Lewentidig & Query Vis: Logic based Diagrams help baren Understand Complicated & Queries Faster, SIGMOD 2020. https://doi.org/10.1145/3318464.3389767 268

"Return any drinker, s.t. there does not exist any other drinker, s.t. there does not exist any beer liked by that other drinker that is not also liked by the returned drinker and there does not exist any beer liked by the returned drinker that is not also liked by the same other drinker."

Let x be a drinker, and S(x) be the set of liked beers by drinker x. Find any drinker x, s.t. there does not exist another drinker x ', x for which:  $S(x') \subseteq S(x)$  and  $S(x') \supseteq S(x)$  Unique set query: "Find drinkers that like a unique set of beers."

0

2

3

 $\{ L1.d \mid \exists L1 \in Likes \land$ **∄L2** ∈ Likes [L2.d <> L1.d ∧  $\nexists L3 \in Likes [L3.d = L1.d \land$  $\nexists L4 \in Likes [L4.d = L2.d \land L4.b = L3.b]] \land$  $\nexists L5 \in Likes [L5.d = L2.d \land$  $\nexists L6 \in Likes [L6.d = L1.d \land L6.b = L5.b]]]$ 

Notice how the logic tree portrays the nesting hierarchy shown in the FOL (TRC) representation of the SQL query.

Each node in the LT represents the root of a scope in the FOL representation. The predicates in each node are the predicates in the root of the scope of a given node (thus the predicates which do not use any additionally quantified variables).



### Atomic predicate classification



#### Allowed:

local op value (local selection pred.) local op local (local join pred.) local op ancestor (connecting join pred.) Not allowed:

ancestor op value (foreign selection pred.) ancestor op ancestor (foreign join pred.)



type

### Focus: one single nesting level

- We first restrict ourselves to
  - equi-joins (no inequalities like T.A < T.B)</li>
  - paths (no siblings = every node can have only one nested child)
  - one single nesting level
  - Boolean queries
  - no foreign predicates
  - only binary relations (thus can be represented as graphs)
  - only one single relation R
  - (and as before only conjunctions)
- Given two such queries, what is a generalization of the homomorphism procedure that works for that fragment?

# Simplifying notation

#### What will become handy, is a short convenient notation for queries

```
SELECT TRUE
        R R1. R R2. R R3
FROM
WHERE R1.B = R2.A
        R2.B = R3.A
AND
NOT EXISTS
     (SELECT *
     FROM
              R R4, R R5, R R6
     WHERE R4.B = R5.A
     AND
             R5.B = R6.A
             R4.A = R1.A
     AND
              R6.A = R2.B)
     AND
```

 $q_0 := R(x,y), R(y,z), R(z,w)$ 

 $q_1(s,t):= R(s,u), R(u,v), R(v,t), s=x, t=y$ 

 $q_0 \coloneqq R(x,y), R(y,z), R(z,w), \neg q_1(x,z)$ 

∃ R1, R2, R3 ∈ R (R1.B=R2.A ∧ R2.B=R3.A ∧ ∄ R4, R5, R6 ∈ R (R4.B=R5.A ∧ R5.B=R6.A ∧ R4.A=R1.A ∧ R6.A = R2.B)



# Simplifying notation

#### What will become handy, is a short convenient notation for queries

```
SELECT TRUE
        R R1. R R2. R R3
FROM
WHERE R1.B = R2.A
AND
        R2.B = R3.A
NOT EXISTS
     (SELECT *
     FROM
             R R4, R R5, R R6
     WHERE R4.B = R5.A
     AND
             R5.B = R6.A
             R4.A = R1.A
     AND
              R6.A = R2.B)
     AND
```

$$q_0 := R(x,y), R(y,z), R(z,w)$$

 $\neg q_1 := R(x,u), R(u,v), R(v,y)$ 

∃ R1, R2, R3 ∈ R (R1.B=R2.A ∧ R2.B=R3.A ∧ ∄ R4, R5, R6 ∈ R (R4.B=R5.A ∧ R5.B=R6.A ∧ R4.A=R1.A ∧ R6.A = R2.B)



# Simplifying notation

#### What will become handy, is a short convenient notation for queries

```
SELECT TRUE
FROM
        R R1, R R2, R R3
WHERE
        R1.B = R2.A
AND
        R2.B = R3.A
NOT EXISTS
     (SELECT *
     FROM
              R R4, R R5, R R6
     WHERE
             R4.B = R5.A
     AND
              R5.B = R6.A
              R4.A = R1.A
     AND
     AND
              R6.A = R2.B)
```



$$q_0 := R(x,y), R(y,z), R(z,w)$$

$$\neg q_1 := R(x,u), R(u,v), R(v,y)$$

Cartesian product: R'(x,y,z,w) = R(x,y), R(y,z), R(z,w)? can be expressed in guarded fragment of FOL (with negation)? But single join already not guarded

See Barany, Cate, Segoufin, "Guarded negatation", JACM 2015

#### guardedness

















### Question

- Find two such nested queries (somehow leveraging the example below) that are equivalent (based on some simple reasoning)
- What is then the \*structured\* procedure to prove equivalence?

### Example

```
\begin{array}{c} q_{1}(x) \coloneqq R(x,y), R(y,y), R(y,z) \\ q_{2}(s) \coloneqq R(s,u), R(u,w), R(s,v), R(u,w), R(u,v) , R(v,v) \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & &
```

# Undecidability 🟵

- Unfortunately, the following problem is already undecidable
  - Consider the class of nested queries with maximal nesting level 2, no disjunctions, our safety restrictions from earlier, set semantics, arbitrary number of siblings
  - Deciding whether any given query is finitely satisfiable is undecidable.
- This follows non-trivially from from following Arxiv paper:
  - "Undecidability of satisfiability in the algebra of finite binary relations with union, composition, and difference" by Tony Tan, Jan Van den Bussche, Xiaowang Zhang, Corr 1406.0349.
     <u>https://arxiv.org/abs/1406.0349</u>



See "Undecidability of satisfiability in the algebra of finite binary relations with union, composition, and difference" by Tan, Van den Bussche, Zhang. https://arxiv.org/abs/1406.0349





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

### Pointers to related work

- Kolaitis. Logic and Databases. Logical Structures in Computation Boot Camp, Berkeley 2016. <u>https://simons.berkeley.edu/talks/logic-and-databases</u>
- Abiteboul, Hull, Vianu. Foundations of Databases. Addison Wesley, 1995. <u>http://webdam.inria.fr/Alice/</u>, Ch 2.1: Theoretical background, Ch 6.2: Conjunctive queries & homomorphisms & query containment, Ch 6.3: Undecidability of equivalence for calculus.
- Chandra, Merlin. Optimal implementation of conjunctive queries in relational data bases. STOC 1977. <u>https://doi.org/10.1145/800105.803397</u>
- Gatterbauer. *Databases will visualize queries too*. PVLDB 2011. <u>http://www.vldb.org/pvldb/vol4/p1498-gatterbauer.pdf</u>

# Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
  - CQ equivalence and containment
  - Graph homomorphisms
  - Homomorphism beyond graphs
  - CQ containment
  - CQ minimization
- T2-2: Equivalence Beyond CQs
  - Union of CQs, and inequalities
  - Union of CQs equivalence under bag semantics
  - Nested queries
  - Tree pattern queries

### Tree pattern queries



Does the query on the left find a match on in the data on the right (i.e. is there a homomorphism from left to right)?

Notice that "a", "b", "c" are labels (not node ids), thus like constants in a query

### Tree pattern queries





294





Profession

name: artist

Country

name: United States

295



Figure 2: A tree pattern finding the artists who were born in the United States. The query returns the person names and the cities where they were born. (Fully circled nodes are return nodes.)

296
## Optimizing tree patterns



How are those two tree patterns related to each other?

ExarolfgaingnGaïterbaioleg Prencipter of foralable/idatampen agement Shttps://hortheastern-datalab/github.10/cos7240/arys. SIGMOD record 2017. https://doi.org/10.1145/3093754.3093759

297



	TREE PATTERN MINIMIZATION
Given:	A tree pattern $p$ and $k \in \mathbb{N}$
Question:	Is there a tree pattern $q$ , equivalent to
	p, such that its size is at most $k$ ?

# Minimality =? Nonredundancy

#### **1.4 History of the Problem**

Although the patterns we consider here have been widely studied [14, 24, 36, 15, 22, 1, 9, 4, 32], their minimization problem remained elusive for a long time. The most important previous work for their minimization was done by Kimelfeld and Sagiv [22] and by Flesca, Furfaro, and Masciari [14, 15].

The key challenge was <u>understanding the relationship be</u>tween *minimality* (M) and *nonredundancy* (NR). Here, a tree pattern is minimal if it has the smallest number of nodes among all equivalent tree patterns. It is nonredundant if none of its leaves (or branches<sup>2</sup>) can be deleted while remaining equivalent. The question was if minimality and nonredundancy are the same ([22, Section 7] and [15, p. 35]):

 $M \stackrel{?}{=} NR$  Problem:

Is a tree pattern minimal if and only if it is nonredundant? Notice that a part of the  $M \stackrel{?}{=} NR$  problem is easy to see: a minimal pattern is trivially also nonredundant (that is,  $M \subseteq NR$ ). The opposite direction is much less clear.

If the problem would have a positive answer, it would mean that the simple algorithmic idea summarised in Algorithm 1 correctly minimizes tree patterns. Therefore, the  $M \stackrel{?}{=} NR$  problem is a natural question about the design of minimization algorithms for tree patterns.

Algorithm 1 Computing a nonredundant subpattern	
<b>Input:</b> A tree pattern $p$ <b>Output:</b> A nonredundant tree pattern $q$ , equivalent to $p$	
while a leaf of $p$ can be removed (remaining equivalent to $p$ ) d	
Remove the leaf	
end while return the resulting pattern	

The  $M \stackrel{?}{=} NR$  problem is also a question about complexity. The main source of complexity of the nonredundancy algorithm lies in testing equivalence between a pattern p and a pattern p', which is generally coNP-complete [24]. If  $M \stackrel{?}{=}$ NR has a positive answer, then TREE PATTERN MINIMIZA-TION would also be coNP-complete.

In fact, the problem was claimed to be coNP-complete in 2003 [14, Theorem 2], but the status of the minimizationand the  $M \stackrel{?}{=} NR$  problems were re-opened by Kimelfeld and Sagiv [22], who found errors in the proofs. Flesca et al.'s journal paper then proved that M = NR for a limited class of tree patterns, namely those where every wildcard node has at most one child [15]. Nevertheless, for tree patterns,

- (a) the status of the  $M \stackrel{?}{=} NR$  problem and
- (b) the complexity of the minimization problem remained open.

Czerwinski, Martens, Niewerth, Parys [PODS 2016]

- (a) There exists a tree pattern that is nonredundant but not minimal. Therefore,  $M \neq NR$ .
- (b) TREE PATTERN MINIMIZATION is  $\Sigma_2^P$ -complete. This implies that even the main idea in Algorithm 1 cannot work unless  $\operatorname{coNP} = \Sigma_2^P$ .

## Tree pattern containment



### Tree pattern containment



but ⊉!



Figure 7: A non-redundant tree pattern p (right) and an equivalent tree pattern q that is smaller (left)  $q \subseteq p$  from previous argument. To be shown:  $q \supseteq p$ To be shown  $q \supseteq p$ : (idea: whenever p matches, then also q) Idea:  $a=\star$  can be matched in 3 ways in a graph



(a) How q can be matched if  $p_1$  can be matched



(b) How q can be matched if  $p_2$  can be matched



(c) How q can be matched if  $p_3$  can be matched