

The homomorphism lattice of finite structures, unique characterization, and exact learnability

A scenic view of a canal in a city, likely Amsterdam. The canal is filled with water and has a boat docked on the left. The right bank is paved with bricks and has a brick building with a modern facade. There are trees and a clear blue sky in the background.

Balder ten Cate

A|C seminar
Oct 6, 2021

About me

Career path:

ILLC → Ivl → INRIA / ENS de Cachan → UC Santa Cruz → LogicBlox → Google → ILLC

About me

Career path:

ILLC → Ivl → INRIA / ENS de Cachan → UC Santa Cruz → LogicBlox → Google → ILLC

Past research topics:

- Extended modal languages; well-behaved fragments of first-order logic and fixpoint logic
- XML query languages; automata and logics for finite trees
- Data management and data interoperability (schema mappings, access constraints)
- Knowledge representation and reasoning (ontology-based data access)
- Declarative languages for data analysis and machine learning
- And, conversely, techniques for learning declarative representations

European reintegration fellowship LLAMA (Logic and Learning - an Algebraic and Model-theoretic Approach)

Outline

1. Conjunctive queries (CQs)
2. Main technical content: Homomorphism lattice of finite structures
3. Applications of the above to some problems regarding learning CQs (and other declarative specification languages)

Main reference:

- Balder ten Cate, Victor Dalmau (2020). **Conjunctive Queries: Unique Characterizations and Exact Learnability**. arxiv.org/abs/2008.06824

Conjunctive Queries

A **conjunctive query (CQ)** is a positive-existential-conjunctive FO formula (that is, CQs form the fragment of FO that uses only \exists and \wedge)

A **union of conjunctive queries (UCQ)** is a disjunction of CQs (that is, UCQs form the fragment of FO that uses only \exists , \wedge , and \vee --- suitably normalized).

Conjunctive Queries

A **conjunctive query (CQ)** is a positive-existential-conjunctive FO formula (that is, CQs form the fragment of FO that uses only \exists and \wedge)

A **union of conjunctive queries (UCQ)** is a disjunction of CQs (that is, UCQs form the fragment of FO that uses only \exists , \wedge , and \vee --- suitably normalized).

Why are CQs (and UCQs) interesting?

- CQs are of fundamental importance in databases as they capture the core SELECT-FROM-WHERE construct in SQL.
- They arise naturally in numerous other areas of computer science as well.

Brief digression

Recall that

- **Modal logic** can be viewed as a fragment of FO, and
- In fact, it is the **bisimulation-invariant fragment** of FO
 - (“Van Benthem Characterization Theorem”)

Similarly,

- **UCQ** is a fragment of FO logic, and
- In fact, it is the **homomorphism-preserved fragment** of FO
 - (“Homomorphism Preservation Theorem”)

Brief digression

Recall that

- **Modal logic** can be viewed as a fragment of FO, and
- In fact, it is the **bisimulation-invariant fragment** of FO
 - (“Van Benthem Characterization Theorem”)

Similarly,

- **UCQ** is a fragment of FO logic, and
- In fact, it is the **homomorphism-preserved fragment** of FO
 - (“Homomorphism Preservation Theorem”)

A **homomorphism** is a map $h : \text{dom}(A) \rightarrow \text{dom}(B)$ that preserves structure (more precise definition will come later in this talk).

Brief digression (ct'd)

Many classic results of model theory “*fail in the finite*” (restricted to finite structures)

- Examples: Compactness, Craig interpolation, Los-Tarski Theorem, ...
- Failure-in-the-finite appears to be the rule rather than the exception.

What is an example of a model-theoretic result that *does* hold in the finite?

Brief digression (ct'd)

Many classic results of model theory “*fail in the finite*” (restricted to finite structures)

- Examples: Compactness, Craig interpolation, Łoś–Tarski Theorem, ...
- Failure-in-the-finite appears to be the rule rather than the exception.

What is an example of a model-theoretic result that *does* hold in the finite?

1. Van Benthem Characterization Theorem (as shown by Rosen, 1997)
2. Homomorphism Preservation Theorem (as shown by Rossman, 2008)

World of modal logic

World of database theory

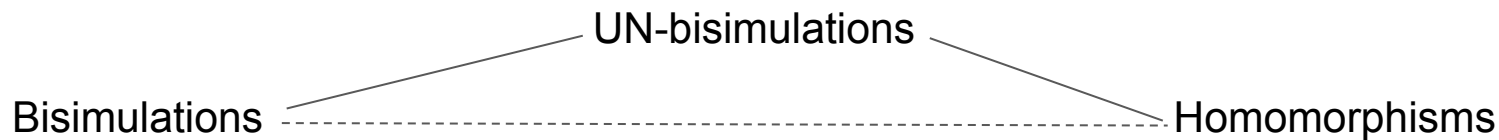
Modal logic ----- UCQs

Modal mu-calculus ----- Monadic Datalog

Bisimulations ----- Homomorphisms

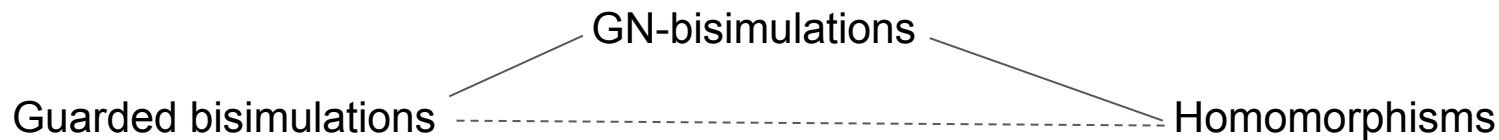
World of modal logic

World of database theory



World of modal logic

World of database theory



End of the Brief digression

(Ask me later if you are interested in this.)

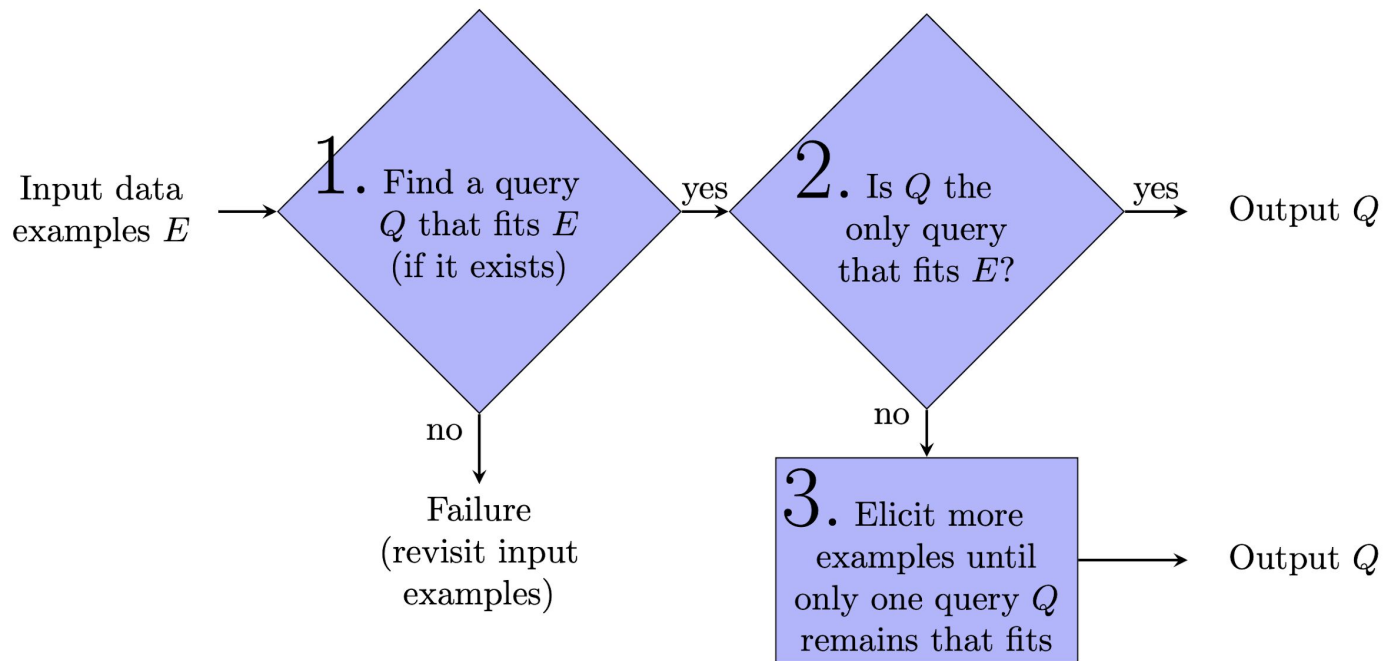
Example-Driven Query Discovery

Scenario: *we are trying to construct a database query from data examples.*

Example-Driven Query Discovery

Scenario: *we are trying to construct a database query from data examples.*

- The **database schema** S and the **arity** k of the query are known.
- A **data example** is a triple: (I, \mathbf{a}, s)
 - I is a database instance (for schema S)
 - \mathbf{a} is a k -tuple
 - s in $\{+, -\}$ indicates if this is a positive or a negative example.
- A **query** Q ...
 - ... *fits* data example $(I, \mathbf{a}, +)$ if $\mathbf{a} \in Q(I)$
 - ... *fits* data example $(I, \mathbf{a}, -)$ if $\mathbf{a} \notin Q(I)$



Remainder of this talk

1. Some aspects of the Homomorphism Lattice of Finite Structures
2. Applications to Example-Driven Query Discovery for CQs.
3. (Time permitting:) Schema Mappings and Description Logics

Structures and Homomorphisms

Fix a finite **schema** \mathbf{S} consisting of relation symbols and constants symbols.

By a **structure** we will mean a finite relational structure over \mathbf{S} .

A **homomorphism** from A to B is a function $h: \text{dom}(A) \rightarrow \text{dom}(B)$ such that

1. For every fact of A , its image (under h) is a fact of B .
2. $h(c^A) = c^B$ for each constant symbol c

The Homomorphism Lattice

The *Homomorphism Lattice* is the partial order $(\text{Str}, \leq_{\text{hom}})$

- Str is the set of all structures
- $A \leq_{\text{hom}} B$ if there is a homomorphism from A to B .

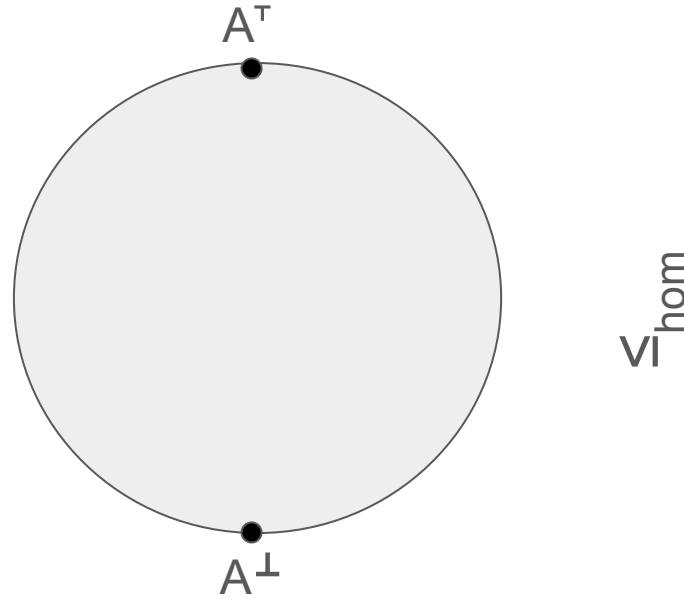
The Homomorphism Lattice

The *Homomorphism Lattice* is the partial order $(\text{Str}, \leq_{\text{hom}})$

- Str is the set of all structures
- $A \leq_{\text{hom}} B$ if there is a homomorphism from A to B .

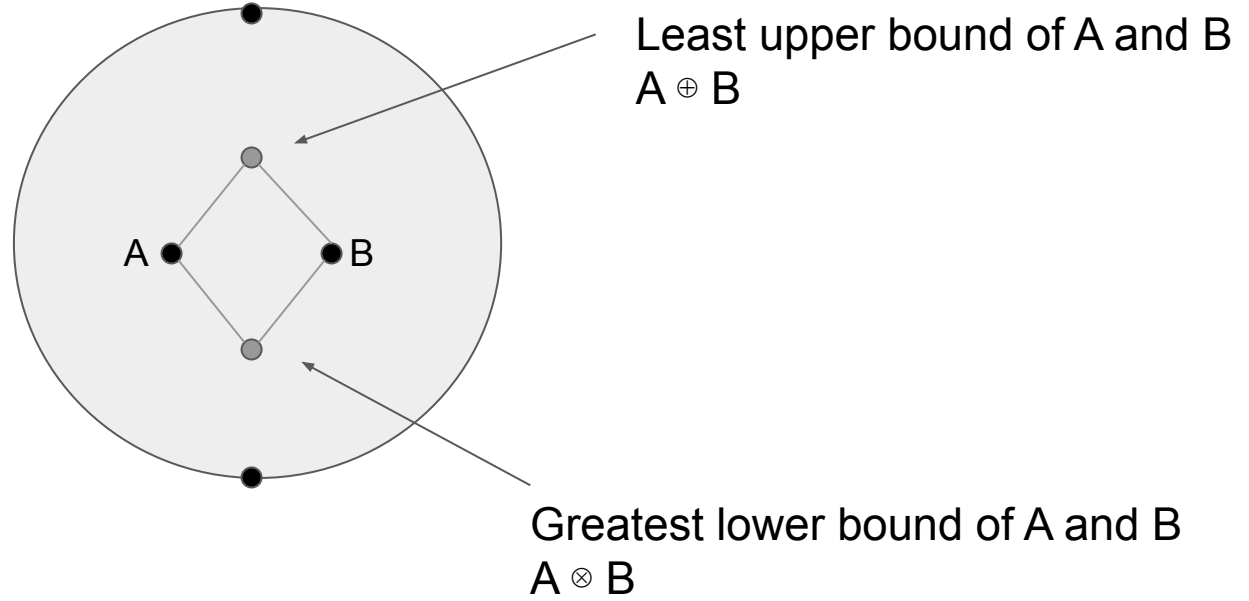
Two structures are *homomorphically equivalent* if $A \leq_{\text{hom}} B$ and $B \leq_{\text{hom}} A$.

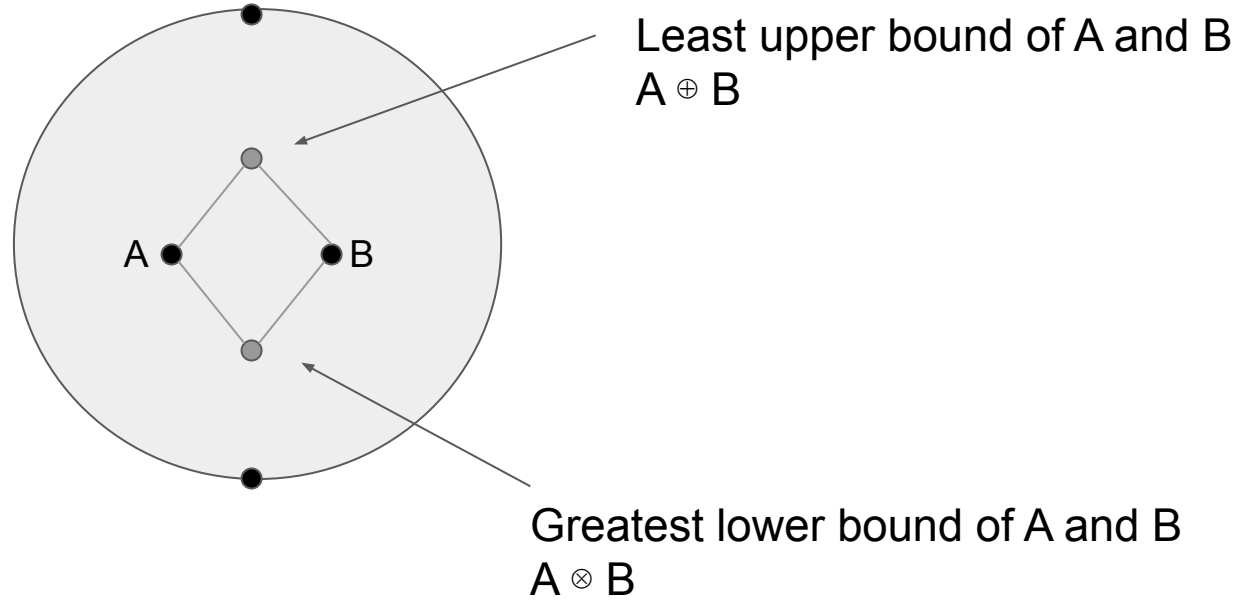
We will not distinguish between homomorphically equivalent structures.



A^\perp : structure with no facts; each constant symbol denotes a different element.

A^T : single-element structure (all constants denote the same element); all facts.





$A \otimes B$ = direct product of A and B

$A \oplus B \sim$ disjoint union of A and B

(in the case with constants the precise construction is a bit more tedious)

The Homomorphism Lattice

Algebra (Str, \oplus , \otimes , \top , \perp , ...)

The Homomorphism Lattice

Algebra (Str, \oplus , \otimes , \top , \perp , ...)

One more relevant operation:

for every A and B there is a structure A^B satisfying:

$$C \leq_{\text{hom}} A^B \quad \text{iff} \quad (C \otimes B) \leq_{\text{hom}} A$$

Small digression

Product Homomorphism Problem (PHP):

Given A_1, \dots, A_n and B , decide whether $(A_1 \otimes \dots \otimes A_n) \leq_{\text{hom}} B$

Theorem (Willard 2010; tC and Dalmau 2015) :

PHP is NExpTime-complete (even for a fixed schema).

Density

Let $\mathbf{S} = \{ R \}$ with R a binary relation. Let $A_1 = \{ R(a,b) \}$ and let $A_2 = \{ R(a,b) , R(b,a) \}$.

Fact 1: $A^\perp <_{\text{hom}} A_1 <_{\text{hom}} A_2$

Density

Let $\mathbf{S} = \{ R \}$ with R a binary relation. Let $A_1 = \{ R(a,b) \}$ and let $A_2 = \{ R(a,b) , R(b,a) \}$.

Fact 1: $A^\perp <_{\text{hom}} A_1 <_{\text{hom}} A_2$

Fact 2: There is *no* B such that $A^\perp < B < A_1$.

Fact 3: For every $B <_{\text{hom}} A_2$ there is a structure B' such that $B <_{\text{hom}} B' <_{\text{hom}} A_2$.

Follows from an extension (due to Nešetřil and Rödl 1989) of Erdős (1959)'s celebrated theorem on the existence of graphs of high girth and chromatic number

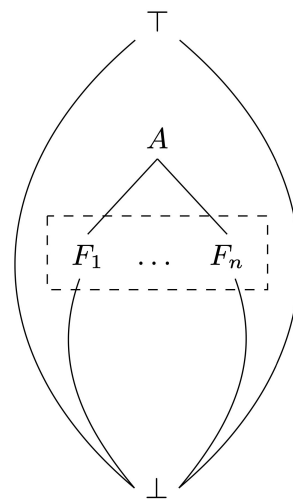
Frontiers

Definition; a *frontier* for A is a finite set of structures $\{F_1, \dots, F_n\}$, such that

1. each $F_i <_{\text{hom}} A$, and
2. whenever $B <_{\text{hom}} A$ then $B \leq_{\text{hom}} F_i$ for some F_i

The frontier separates A from the structures strictly below A .

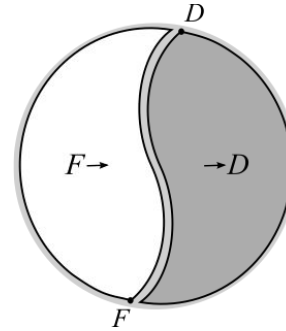
Not every structure has a frontier. A_2 has no frontier.



Dualities

Def 1: A pair of structure (F, D) is a *duality pair* if:

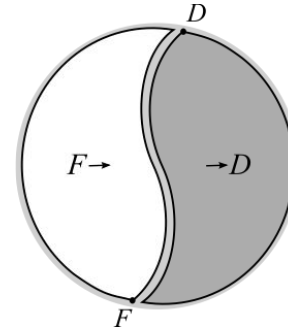
$$\{A \mid F \leq_{\text{hom}} A\} = \{A \mid A \not\leq_{\text{hom}} D\}$$



Dualities

Def 1: A pair of structure (F, D) is a *duality pair* if:

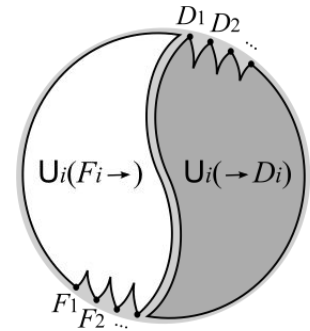
$$\{A \mid F \leq_{\text{hom}} A\} = \{A \mid A \not\leq_{\text{hom}} D\}$$



Def 2: A pair of finite sets $(\{F_1, \dots, F_n\}, \{D_1, \dots, D_m\})$

is a *generalized homomorphism duality* if:

$$\{A \mid F_i \leq_{\text{hom}} A \text{ for some } i \leq n\} = \{A \mid A \not\leq_{\text{hom}} D_i \text{ for all } i \leq m\}$$



Example 1 (duality pair)

Let

- P_{k+1} be the directed path with $k+1$ elements.
- T_k be the linear order with k elements.

Gallai-Hasse-Roy-Vitaver Theorem (~1965) for directed graphs:

(P_{k+1}, T_k) is a duality pair,

i.e., for every directed graph G , it holds that $P_{k+1} \rightarrow G$ if and only if $G \not\rightarrow T_k$

Example 2 (infinitary homomorphism duality)

A graph is 2-colorable if and only if it does not have a cycle of odd length.

$G \rightarrow \mathbf{K}_2$ if and only if $\mathbf{C}_{2n+1} \not\rightarrow G$ for all n

where \mathbf{K}_2 is the 2-element clique and \mathbf{C}_n is the a cycle of length n .

Frontiers and Dualities

Theorem (from Neseštril and Tardiff 2000, cf. also tC and Dalmau 2020):

- If $\{B_1, \dots, B_n\}$ is a **frontier** for A , then

$(\{A\}, \{A^{B_1}, \dots, A^{B_n}\})$ is a **generalized homomorphism duality**

- If $(\{A\}, \{B_1, \dots, B_n\})$ is a **generalized homomorphism duality**,

then $\{(A \otimes B_1), \dots, (A \otimes B_n)\}$ is a **frontier** for A

(*Statement assumes no constants but naturally extends to the case with constants)

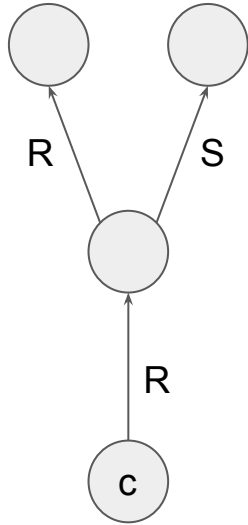
C-acyclicity

A structure is *c-acyclic* if it does not have cycles except for cycles that involve an element named by a constant symbol.

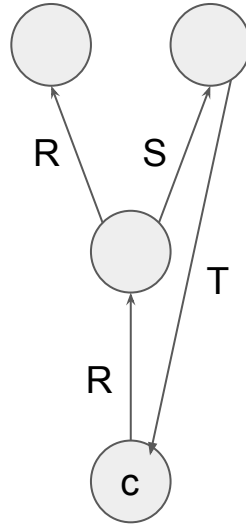
More precise definition:

- The *incidence graph* of a structure A is the bi-partite multi-graph where
 - The nodes of the graph are the elements and facts of A
 - There is an edge between an element and a fact if the element occurs in the fact.
 - If an element occurs multiple times in a fact, each occurrence generates an edge.
- A structure is *c-acyclic* if every cycle in the incidence graph goes through at least one element named by a constant symbol.

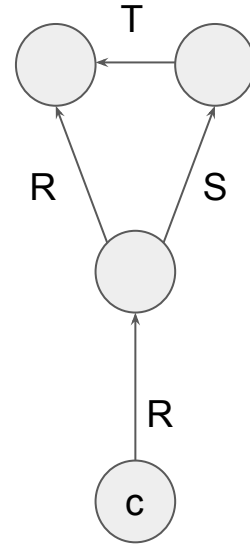
Examples



acyclic



c-acyclic



cyclic

Frontiers: existence and how to construct them

Thm 1: For all structures A , the following are equivalent:

1. A **has a frontier**
2. A is homomorphically equivalent to a **c-acyclic** structure.

Thm 2: For c-acyclic structures, a frontier **can be computed in polynomial time**.

Thm 3: Testing whether a given set of structures $\{ F_1, \dots, F_n \}$ is a frontier for a structure A is **NP-complete**.

Foniok, Nešetřil and Tardiff (2008), tC and Dalmau (2020)

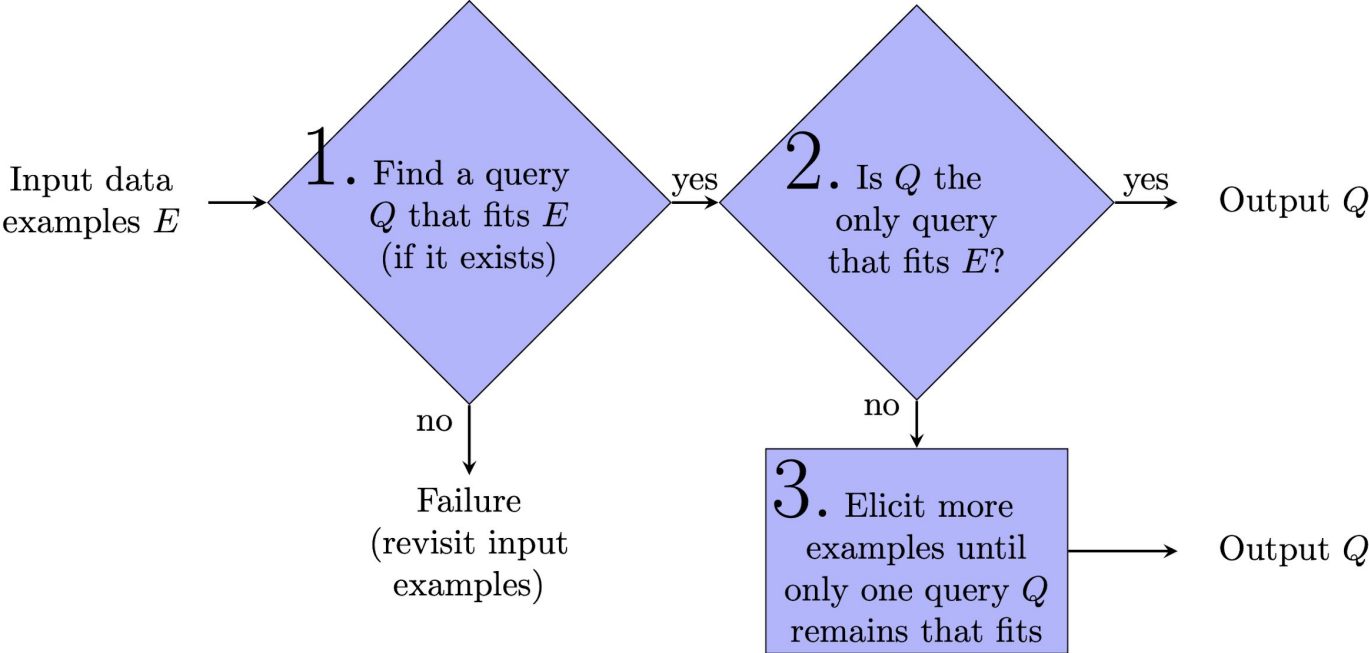
Some further results on frontiers

The class of c -acyclic structures is not “*frontier-closed*” (although every c -acyclic structure has a frontier, it does not necessarily consist of c -acyclic structures).

Thm (tC and Dalmau 2020). The class of **acyclic connected structures (with a single constant)** is polynomial-time frontier-closed.

Thm (from Nešetřil and Ossona de Mendez 2008; cf. tC and Dalmau 2020):
Every class of structures that has bounded expansion admits relativized frontiers.

Example-Driven Query Discovery



Conjunctive Queries

Let's restrict attention to **conjunctive queries (CQs)**.

- Every **k-ary CQs** over a relational schema **S** can be equivalently viewed as a **structure over schema $S \cup \{c_1, \dots, c_k\}$** .
- Every **example** corresponds to a structure over **$S \cup \{c_1, \dots, c_k\}$** as well.
- Q fits a $(I, \mathbf{a}, +)$ iff there is a homomorphism from $q(\mathbf{x})$ to (I, \mathbf{a})
- Q fits a $(I, \mathbf{a}, -)$ iff there is no homomorphism from $q(\mathbf{x})$ to (I, \mathbf{a})

This sets the stage for us to apply results about the homomorphism lattice.

1. The Fitting Problem

Fitting problem:

Input: a finite set E of data examples.

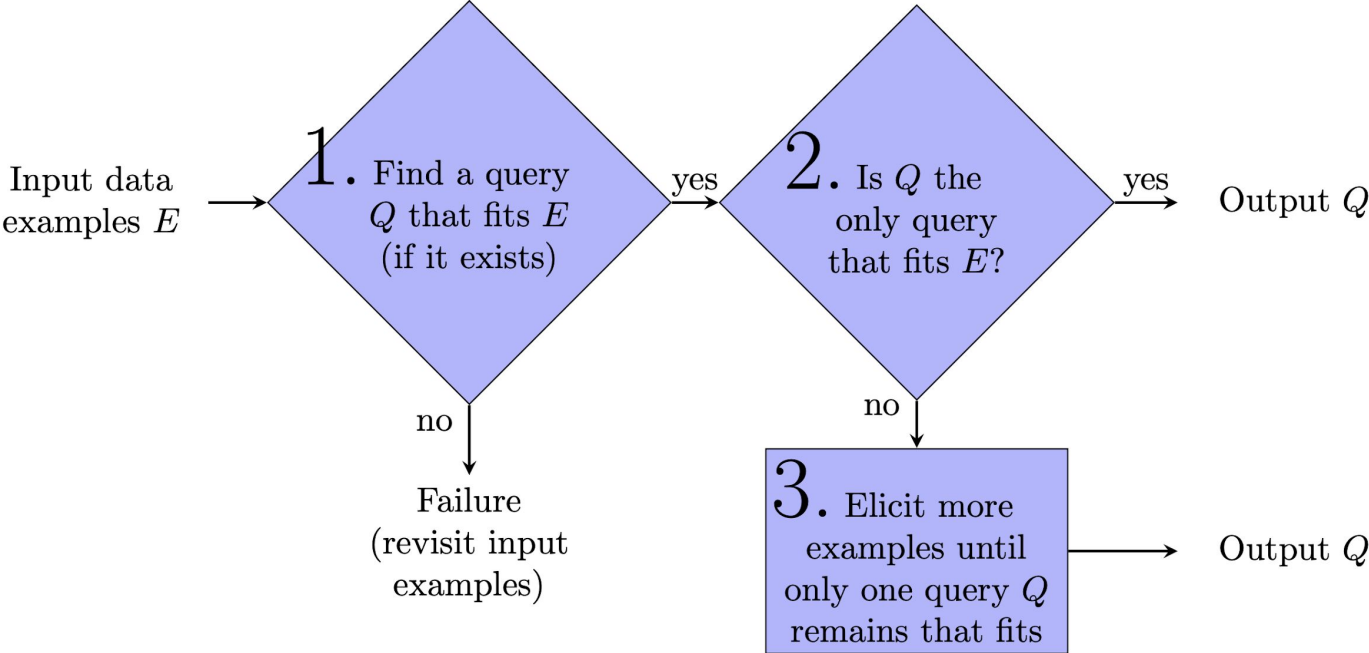
Decide whether there is a CQ that fits all data examples in E .

Lemma: let Q^* be the **direct product of the positive examples in E** . The following are equivalent:

1. There exists a CQ that fits the data examples in E .
2. Q^* fits the data examples in E .
3. Q^* does not have a homomorphism to any negative example in E .
4. Q^* is the **most-specific fitting CQ** for the data examples in E .

Theorem: the fitting problem is coNExpTime-complete. (by reduction from PHP)

Example-Driven Query Discovery



2. The Uniqueness Problem

Definition: E *uniquely characterizes* Q if Q is the only CQ (up to logical equivalence) that fits the examples in E.

Unique characterization problem:

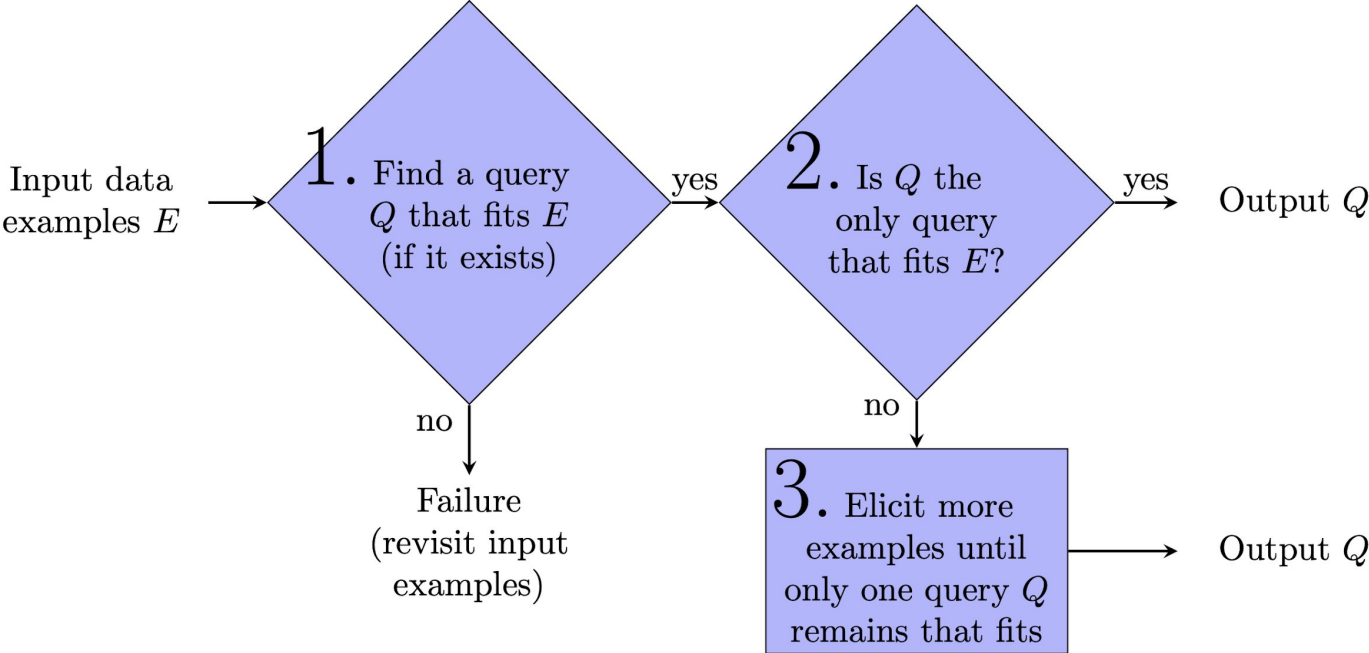
Input: a finite set of data examples E and a most-specific-fitting CQ Q for E.
Decide if E uniquely characterizes Q.

Lemma: The following are equivalent:

1. E uniquely characterizes Q
2. The negative examples in E form a frontier for Q.

Theorem: the unique characterization problem is NP-complete.

Example-Driven Query Discovery



3. Eliciting Further Examples

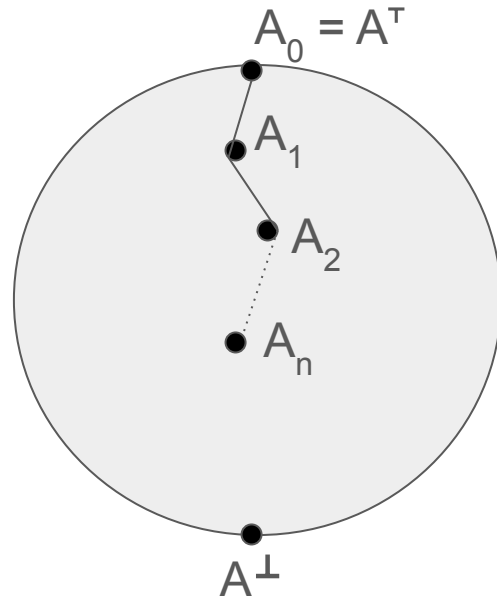
Will finitely many examples even ever be enough to uniquely characterize the target CQ?

Theorem 1: for all CQs Q , the following are equivalent

1. Q is uniquely characterized by a finite set of data examples
2. Q has a frontier.
3. Q is logically equivalent to a c-acyclic CQ.

Theorem 2: the class of c-acyclic CQs is efficiently exactly learnable with membership queries.

Theorem 2 means that, if the “target CQ” is c-acyclic, then there is a PTIME algorithm that, after asking for the label of one or more examples, terminates and identifies it correctly (modulo logical equivalence).



VI^{hom}

Unions of Conjunctive Queries

The situation for UCQs is a little different:

- **Unique characterizations** for UCQs ~ **generalized homomorphism dualities**
- A UCQ is **uniquely characterized by a finite set of examples** iff it is logically equivalent to a **c-acyclic** UCQ.
- For c-acyclic UCQs, a uniquely characterizing set of examples can be effectively constructed but (provably) **not in polynomial time**.
- The class of c-acyclic UCQs is *not* efficiently exactly learnable with membership queries (but is efficiently exactly learnable with membership and equivalence queries).

Schema Mappings

A schema mapping $M=(\mathbf{S},\mathbf{T},\Sigma)$ is a high-level declarative specifications of the relationships between two database schemas.

Two of the most well-studied schema mapping specification languages are **LAV** (“Local-as-View”) and **GAV** (“Global-as-View”) schema mappings.

Schema Mappings

A schema mapping $M=(\mathbf{S},\mathbf{T},\Sigma)$ is a high-level declarative specifications of the relationships between two database schemas.

Two of the most well-studied schema mapping specification languages are [LAV](#) (“Local-as-View”) and [GAV](#) (“Global-as-View”) schema mappings.

Alexe, tC, Kolaitis, and Tan (2011) studied the question when a schema mapping be uniquely characterized by a finite set of data examples.

tC, Dalmau, and Kolaitis (2013) studied efficient exact learnability for GAV schema mappings.

Schema Mappings

Corr. Fix a source schema S and a target schema T .

1. A LAV schema mapping $M=(S,T,\Sigma)$ is **uniquely characterizable** by a finite set of positive and negative schema-mapping examples if and only if M is logically equivalent to a **c-acyclic** LAV schema mapping.
2. If M is c-acyclic, then a uniquely characterizing set of positive and negative schema-mapping examples **can be constructed in PTIME**.
3. The class of c-acyclic LAV schema mappings over S, T is **efficiently exactly learnable with membership queries**.

Question: what happens in the presence of integrity constraints?

Description logics

Description logics (DLs) are formal specification languages used to represent domain knowledge.

The DL concept language *ELI* can be characterized as a notational variant of acyclic connected unary CQs (over a schema with unary and binary relations only).

For example, the *ELI* concept

$\text{logician} \sqcap \exists \text{HasParent} . \exists \text{HasParent} . \text{logician}$

corresponds to $q(x) = \exists y, z . (\text{logician}(x) \wedge \text{HasParent}(x,y) \wedge \text{HasParent}(y,z) \wedge \text{logician}(z))$

Description logics

Corr.

1. Every ELI concept expression is uniquely characterizable by a finite collection of positive and negative QA examples, which is PTIME-computable.
2. Furthermore, the class of ELI concept expressions is efficiently exactly learnable with membership queries.

Question: does this hold also in the presence of an ontology (background theory)?

One recent result

Assume a schema \mathbf{S} consisting of unary and binary relations only.

A *DL-Lite ontology* is a FO theory consisting of

- Inclusion constraints $\alpha(x) \rightarrow \beta(x)$ and/or
- Disjointness constraints of the form $\alpha(x) \rightarrow \neg\beta(x)$, $R(x,y) \rightarrow \neg S(x,y)$, and/or $R(x,y) \rightarrow \neg S(y,x)$

where $\alpha(x)$, $\beta(x)$ are unary projections of relations, i.e., of the form $P(x)$, $\exists y.R(x,y)$ or $\exists y.R(y,x)$

Thm (Funk, Jung and Lutz, DL 2021): Acyclic, connected, unary CQs over \mathbf{S} (in other words, ELI description logic concepts) are efficiently exactly learnable in the presence of DL-Lite ontologies.

Outlook (LLAMA)

- Unique characterizations and exact learnability for modal formulas
- Unique characterizations and exact learnability in the presence of a background theory

Collaboration with Raoul Koudijs, ... (room for more collaborations!)

Thank you!