

CS294-248 Special Topics in Database Theory

Unit 1: Logic and Queries

Dan Suciu

University of Washington

Welcome!

This course is intended for graduate students interested in getting deeper into data management technologies: understanding the underlying theory.

I am a professor at the University of Washington, attending the SIMONS institute [Logic and Algorithms in Database Theory and AI](#), and the recipient of the Theory of Computing Chancellor's Professorship at UC Berkeley.

This course is a one-time offering.

About

- What this course is about: logic, complexity, algorithms, all related to data management. There will be proofs in class.
- What is course is not: a course on data science, data management, or database systems.¹

¹Recommended: CS286: Graduate DB Systems in Spring 2024. (This is a natural graduate-level follow-on to the undergrad CS186 class.)

General Info

- Lectures: Tue/Thu 11-12:20;
<https://berkeley-cs294-248.github.io/>
- Workshops at the Simons Institute: weeks of 9/25, 10/16, 11/13.
- Theory homework assignments – first one is already published.
- Final report and presentation: the week of 12/4.

Tentative Course Outline

Tue	Thu	Unit	Topic	Lecturer
8/29	8/31	U1	Logic and Queries.	
9/5	9/7	U2	Basic Query Evaluation.	
9/12	9/14	U3	Incremental View Maintenance	Dan Olteanu
9/19	9/21	U4	AGM Bound WCOJ	Hung Ngo
9/25-9/29: WS 1: Fine-grained Complexity, Logic, Query Eval				
10/3	10/5	U5	Database Constraints.	
10/10	10/12	U6	Probabilistic databases	
10/16-10/20: WS 2: Probabilistic Circuits and Logic				
10/24	10/26	U7	Semirings, K-Relations.	Val Tannen
10/31		U8	FAQ	Hung Ngo
11/7	11/9	U9	Datalog, Chase.	
11/13-11/17: WS 3: Logic and Algebra for Query Evaluation				
...	...		TBD	

Final Report

- Task: pick a theory problem/result and explain it to a wide audience.
- Write a short report. Suggested length: 3-5 page.
- Give a short presentation (10'), in class, probably on Tuesday 12/5. Details TBD.

Recommended Readings



The “Alice Book” [Abiteboul et al., 1995]

Libkin’s *Finite Model Theory* [Libkin, 2004]

A much shorter tutorial in PODS [Libkin, 2009].

New upcoming book on Database Theory [Arenas et al., 2022].

Basic Definitions

Structures

A **vocabulary** σ is a set of relation symbols R_1, \dots, R_k and function symbols f_1, \dots, f_m , each with a fixed arity.

A **structure** is $\mathbf{D} = (D, R_1^D, \dots, R_k^D, f_1^D, \dots, f_m^D)$, where $R_i^D \subseteq (D)^{\text{arity}(R_i)}$ and $f_j^D : (D)^{\text{arity}(f_j)} \rightarrow D$.

D = the *domain* or the *universe*; always assumed $\neq \emptyset$.

$v \in D$ is called an *element* or a *value* or a *point*.

\mathbf{D} called a *structure* or a *model* or *database*.

Structures

A **vocabulary** σ is a set of relation symbols R_1, \dots, R_k and function symbols f_1, \dots, f_m , each with a fixed arity.

A **structure** is $\mathbf{D} = (D, R_1^D, \dots, R_k^D, f_1^D, \dots, f_m^D)$, where $R_i^D \subseteq (D)^{\text{arity}(R_i)}$ and $f_j^D : (D)^{\text{arity}(f_j)} \rightarrow D$.

D = the *domain* or the *universe*; always assumed $\neq \emptyset$.

$v \in D$ is called an *element* or a *value* or a *point*.

\mathbf{D} called a *structure* or a *model* or *database*.

Structures

A **vocabulary** σ is a set of relation symbols R_1, \dots, R_k and function symbols f_1, \dots, f_m , each with a fixed arity.

A **structure** is $\mathbf{D} = (D, R_1^D, \dots, R_k^D, f_1^D, \dots, f_m^D)$, where $R_i^D \subseteq (D)^{\text{arity}(R_i)}$ and $f_j^D : (D)^{\text{arity}(f_j)} \rightarrow D$.

D = the *domain* or the *universe*; always assumed $\neq \emptyset$.

$v \in D$ is called an *element* or a *value* or a *point*.

\mathbf{D} called a *structure* or a *model* or *database*.

Examples

A **graph** is $G = (V, E)$, $E \subseteq V \times V$.

A **field** is $\mathbb{F} = (F, 0, 1, +, \cdot)$ where

- F is a set.
- 0 and 1 are constants (i.e. functions $F^0 \rightarrow F$).
- $+$ and \cdot are functions $F^2 \rightarrow F$.

An **ordered set** is $\mathbf{S} = (S, \leq)$ where $\leq \subseteq S \times S$.

A **database** is $\mathbf{D} = (\text{Domain}, \text{Customer}, \text{Order}, \text{Product})$.

Discussion

- We don't really need functions, since $f : D^k \rightarrow D$ is represented by its graph $\subseteq D^{k+1}$, but we keep them when convenient.
- If f is a 0-ary function $D^0 \rightarrow D$, then it is a constant D , and we denote it c rather than f .
- D can be a finite or an infinite structure.

First Order Logic

Fix a vocabulary σ and a set of variables x_1, x_2, \dots

Terms:

- Every constant c and every variable x is a term.
- If t_1, \dots, t_k are terms then $f(t_1, \dots, t_k)$ is a term.

Formulas:

- \mathbf{F} is a formula (means *false*).
- If t_1, \dots, t_k are terms, then $t_1 = t_2$ and $R(t_1, \dots, t_k)$ are formulas.
- If φ, ψ are formulas, then so are $\varphi \rightarrow \psi$ and $\forall x(\varphi)$.

Discussion

We were very frugal! We used only **F**, \rightarrow , \forall .

In practice we use several derived operations:

- $\neg\varphi$ is a shorthand for $\varphi \rightarrow \mathbf{F}$.
- $\varphi \vee \psi$ is a shorthand for $(\neg\varphi) \rightarrow \psi$.
- $\varphi \wedge \psi$ is a shorthand for $\neg(\varphi \vee \psi)$.
- $\exists x(\varphi)$ is a shorthand for $\neg(\forall x(\neg\varphi))$.

F often denoted: false or \perp or 0.

$=$ is not always part of the language

Formulas and Sentences

We say that $\forall x(\varphi)$ *binds* x in φ . Every occurrence of x in φ is *bound*. Otherwise, it is *free*.

A **sentence** is a formula φ without free variables.

E.g. formula $\varphi(x, z) = \exists y(E(x, y) \wedge E(y, z))$. Free variables: x, z .

E.g. sentence $\varphi = \exists x \forall z \exists y(E(x, y) \wedge E(y, z))$.

Truth

Let φ be a sentence, and \mathbf{D} a structure

Definition

We say that φ is **true** in \mathbf{D} , written $\boxed{\mathbf{D} \models \varphi}$, if:

- φ is $c = c'$ and c, c' are the same constant.
- φ is $R(c_1, \dots, c_n)$ and $(c_1, \dots, c_n) \in R^{\mathbf{D}}$.
- φ is $\psi_1 \rightarrow \psi_2$ and $\mathbf{D} \not\models \psi_1$, or $\mathbf{D} \models \psi_1$ and $\mathbf{D} \models \psi_2$.
- φ is $\forall y(\psi)$, and, for all $b \in D$, $\mathbf{D} \models \psi[b/y]$.

This definition is boring but important!

Special Case: Propositional Logic

A **nullary relation**, $A()$, is the same as a propositional variable:

- In any structure \mathbf{D} , $A^{\mathbf{D}}$ can be either \emptyset or $\{()\}$.
- If $A^{\mathbf{D}} = \{()\}$ then we say that $A^{\mathbf{D}}$ is **true**.
- If $A^{\mathbf{D}} = \emptyset$ then we say that $A^{\mathbf{D}}$ is **false**.

Sentences over nullary predicates are the same as propositional formulas:

$$A() \wedge (B() \vee \neg A())$$

$$A \wedge (B \vee \neg A)$$

What do these sentences say about *D*?

$$\exists x \exists y \exists z (x \neq y) \wedge (x \neq z) \wedge (y \neq z)$$

$$\exists x \exists y \forall z (z = x) \vee (z = y)$$

What do these sentences say about D ?

$$\exists x \exists y \exists z (x \neq y) \wedge (x \neq z) \wedge (y \neq z)$$

“There are at least three elements”, i.e. $|D| \geq 3$

$$\exists x \exists y \forall z (z = x) \vee (z = y)$$

What do these sentences say about D ?

$$\exists x \exists y \exists z (x \neq y) \wedge (x \neq z) \wedge (y \neq z)$$

“There are at least three elements”, i.e. $|D| \geq 3$

$$\exists x \exists y \forall z (z = x) \vee (z = y)$$

“There are at most two elements”, i.e. $|D| \leq 2$

What do these sentences say about *D*?

$$\forall x \exists y E(x, y) \vee E(y, x)$$

$$\forall x \forall y \exists z E(x, z) \wedge E(z, y)$$

$$\begin{aligned} &\exists x \exists y \exists z (\forall u (u = x) \vee (u = y) \vee (u = z)) \\ &\quad \wedge \neg E(x, x) \wedge E(x, y) \wedge \neg E(x, z) \\ &\quad \wedge \neg E(y, z) \wedge \neg E(y, y) \wedge E(y, z) \\ &\quad \wedge E(z, x) \wedge \neg E(z, y) \wedge \neg E(z, z) \end{aligned}$$

What do these sentences say about D ?

$$\forall x \exists y E(x, y) \vee E(y, x)$$

“There are no isolated nodes”

$$\forall x \forall y \exists z E(x, z) \wedge E(z, y)$$

$$\begin{aligned} &\exists x \exists y \exists z (\forall u (u = x) \vee (u = y) \vee (u = z)) \\ &\quad \wedge \neg E(x, x) \wedge E(x, y) \wedge \neg E(x, z) \\ &\quad \wedge \neg E(y, z) \wedge \neg E(y, y) \wedge E(y, z) \\ &\quad \wedge E(z, x) \wedge \neg E(z, y) \wedge \neg E(z, z) \end{aligned}$$

What do these sentences say about *D*?

$$\forall x \exists y E(x, y) \vee E(y, x)$$

“There are no isolated nodes”

$$\forall x \forall y \exists z E(x, z) \wedge E(z, y)$$

“Every two nodes are connected by a path of length 2”

$$\begin{aligned} &\exists x \exists y \exists z (\forall u (u = x) \vee (u = y) \vee (u = z)) \\ &\quad \wedge \neg E(x, x) \wedge E(x, y) \wedge \neg E(x, z) \\ &\quad \wedge \neg E(y, z) \wedge \neg E(y, y) \wedge E(y, z) \\ &\quad \wedge E(z, x) \wedge \neg E(z, y) \wedge \neg E(z, z) \end{aligned}$$

What do these sentences say about D ?

$$\forall x \exists y E(x, y) \vee E(y, x)$$

“There are no isolated nodes”

$$\forall x \forall y \exists z E(x, z) \wedge E(z, y)$$

“Every two nodes are connected by a path of length 2”

$$\begin{aligned} \exists x \exists y \exists z (&\forall u (u = x) \vee (u = y) \vee (u = z)) \\ &\wedge \neg E(x, x) \wedge E(x, y) \wedge \neg E(x, z) \\ &\wedge \neg E(y, z) \wedge \neg E(y, y) \wedge E(y, z) \\ &\wedge E(z, x) \wedge \neg E(z, y) \wedge \neg E(z, z) \end{aligned}$$

It completely determines the graph: $D = \{a, b, c\}$ and $a \rightarrow b \rightarrow c \rightarrow a$.

Classical Model Theory

Fix a sentence φ , and a set of sentences Σ (may be infinite).

- **Satisfiability:** Σ is satisfiable if $\exists \mathbf{D}$ such that $\mathbf{D} \models \Sigma$. $\text{SAT}(\Sigma)$.
- **Implication:** Σ implies φ if $\forall \mathbf{D}$, $\mathbf{D} \models \Sigma$ implies $\mathbf{D} \models \varphi$. $\Sigma \models \varphi$.
- **Validity:** φ is valid if $\forall \mathbf{D}$, $\mathbf{D} \models \varphi$. We write $\models \varphi$ or $\text{VAL}(\varphi)$.

$$\neg \text{SAT}(\varphi) \text{ iff } \text{VAL}(\neg \varphi)$$

Completeness, Undecidability

Gödel's Completeness Thm: $\Sigma \models \varphi$ iff there exists a finite proof $\Sigma \vdash \varphi$.

Church's Undecidability Thm: VAL is undecidable. Hence, so is SAT.

We will not discuss what a “proof” $\Sigma \vdash \varphi$ means.

- There exists an algorithm that enumerates all valid sentences:

$$\text{VAL} = \{\varphi_0, \varphi_1, \varphi_2, \dots\}$$

- There exists an algorithm that enumerates all unsatisfiable sentences:

$$\text{UNSAT} = \{\varphi_0, \varphi_1, \varphi_2, \dots\}$$

We say that VAL is **recursively enumerable**, r.e., and SAT is **co-r.e.**

Completeness, Undecidability

Gödel's Completeness Thm: $\Sigma \models \varphi$ iff there exists a finite proof $\Sigma \vdash \varphi$.

Church's Undecidability Thm: VAL is undecidable. Hence, so is SAT.

We will not discuss what a “proof” $\Sigma \vdash \varphi$ means.

- There exists an algorithm that enumerates all valid sentences:

$$\text{VAL} = \{\varphi_0, \varphi_1, \varphi_2, \dots\}$$

- There exists an algorithm that enumerates all unsatisfiable sentences:

$$\text{UNSAT} = \{\varphi_0, \varphi_1, \varphi_2, \dots\}$$

We say that VAL is **recursively enumerable**, r.e., and SAT is **co-r.e.**

Finite Model Theory, Databases, Verification

All previous problems, where the models are restricted to be finite:

- Finite satisfiability: $\text{SAT}_{\text{fin}}(\Sigma)$.
- Finite implication: $\Sigma \models_{\text{fin}} \varphi$.
- Finite validity: $\models_{\text{fin}} \varphi$ or $\text{VAL}_{\text{fin}}(\varphi)$.

New problems that make sense only in the finite:

- **Model checking**: Given φ , D , determine whether $D \models \varphi$.
- **Query evaluation**: Given $\varphi(\mathbf{x})$, D , compute $\{\mathbf{a} \mid D \models \varphi[\mathbf{a}/\mathbf{x}]\}$.

Examples

- $F \models \varphi$ for any sentence φ **why?**.

Examples

- $\mathbf{F} \models \varphi$ for any sentence φ **why?**.
- $\Sigma \models \mathbf{F}$ iff Σ is unsatisfiable **why?**.

Examples

- $\mathbf{F} \models \varphi$ for any sentence φ **why?**.
- $\Sigma \models \mathbf{F}$ iff Σ is unsatisfiable **why?**.
- Is $\Phi \stackrel{\text{def}}{=}} (\forall x \exists y E(x, y)) \wedge (\forall x_1 \forall x_2 \forall y (E(x_1, y) \wedge E(x_2, y) \Rightarrow x_1 = x_2))$ satisfiable?

Examples

- $\mathbf{F} \models \varphi$ for any sentence φ **why?**.
- $\Sigma \models \mathbf{F}$ iff Σ is unsatisfiable **why?**.
- Is $\Phi \stackrel{\text{def}}{=}} (\forall x \exists y E(x, y)) \wedge (\forall x_1 \forall x_2 \forall y (E(x_1, y) \wedge E(x_2, y) \Rightarrow x_1 = x_2))$ satisfiable?

Every node has an outgoing edge, and at most one incoming edge

Examples

- $\mathbf{F} \models \varphi$ for any sentence φ **why?**.
- $\Sigma \models \mathbf{F}$ iff Σ is unsatisfiable **why?**.
- Is $\Phi \stackrel{\text{def}}{=} (\forall x \exists y E(x, y)) \wedge (\forall x_1 \forall x_2 \forall y (E(x_1, y) \wedge E(x_2, y) \Rightarrow x_1 = x_2))$ satisfiable?

Every node has an outgoing edge, and at most one incoming edge

Yes: E.g. $E = \{(1, 2), (2, 3), (3, 1)\}$.

Examples

- $\mathbf{F} \models \varphi$ for any sentence φ **why?**.
- $\Sigma \models \mathbf{F}$ iff Σ is unsatisfiable **why?**.
- Is $\Phi \stackrel{\text{def}}{=} (\forall x \exists y E(x, y)) \wedge (\forall x_1 \forall x_2 \forall y (E(x_1, y) \wedge E(x_2, y) \Rightarrow x_1 = x_2))$ satisfiable?
 Every node has an outgoing edge, and at most one incoming edge
Yes: E.g. $E = \{(1, 2), (2, 3), (3, 1)\}$.
- Is $\Phi \wedge (\exists y \forall x \neg E(x, y))$ (where Φ is defined above) satisfiable?

Examples

- $\mathbf{F} \models \varphi$ for any sentence φ **why?**.
- $\Sigma \models \mathbf{F}$ iff Σ is unsatisfiable **why?**.
- Is $\Phi \stackrel{\text{def}}{=} (\forall x \exists y E(x, y)) \wedge (\forall x_1 \forall x_2 \forall y (E(x_1, y) \wedge E(x_2, y) \Rightarrow x_1 = x_2))$ satisfiable?
 Every node has an outgoing edge, and at most one incoming edge
Yes: E.g. $E = \{(1, 2), (2, 3), (3, 1)\}$.
- Is $\Phi \wedge (\exists y \forall x \neg E(x, y))$ (where Φ is defined above) satisfiable?
 ... and there exists a node with no incoming edge

Examples

- $\mathbf{F} \models \varphi$ for any sentence φ **why?**.
- $\Sigma \models \mathbf{F}$ iff Σ is unsatisfiable **why?**.
- Is $\Phi \stackrel{\text{def}}{=} (\forall x \exists y E(x, y)) \wedge (\forall x_1 \forall x_2 \forall y (E(x_1, y) \wedge E(x_2, y) \Rightarrow x_1 = x_2))$ satisfiable?
 Every node has an outgoing edge, and at most one incoming edge
Yes: E.g. $E = \{(1, 2), (2, 3), (3, 1)\}$.
- Is $\Phi \wedge (\exists y \forall x \neg E(x, y))$ (where Φ is defined above) satisfiable?
 ... and there exists a node with no incoming edge
Yes: $E = \{(0, 1), (1, 2), (2, 3), \dots\}$

Examples

- $\mathbf{F} \models \varphi$ for any sentence φ **why?**.
- $\Sigma \models \mathbf{F}$ iff Σ is unsatisfiable **why?**.
- Is $\Phi \stackrel{\text{def}}{=}} (\forall x \exists y E(x, y)) \wedge (\forall x_1 \forall x_2 \forall y (E(x_1, y) \wedge E(x_2, y) \Rightarrow x_1 = x_2))$ satisfiable?
Every node has an outgoing edge, and at most one incoming edge
Yes: E.g. $E = \{(1, 2), (2, 3), (3, 1)\}$.
- Is $\Phi \wedge (\exists y \forall x \neg E(x, y))$ (where Φ is defined above) satisfiable?
... and there exists a node with no incoming edge
Yes: $E = \{(0, 1), (1, 2), (2, 3), \dots\}$
But **Not satisfiable in the finite.**
“Axioms of infinity” [Börger et al., 1997]

$$\boxed{\text{SAT}_{\text{fin}}(\varphi) \Rightarrow \text{SAT}(\varphi)}$$

Finite v.s. Classical Model Theory

In relational databases we are interested in **Finite Model Theory**.

VAL_{fin} , SAT_{fin} differ from VAL , SAT .

Could VAL_{fin} , SAT_{fin} be decidable?

Finite v.s. Classical Model Theory

In relational databases we are interested in **Finite Model Theory**.

VAL_{fin} , SAT_{fin} differ from VAL , SAT .

Could VAL_{fin} , SAT_{fin} be decidable?

There is hope:

- In classical model theory VAL is r.e., SAT is co-r.e.
- In finite model theory SAT_{fin} is r.e. **why?**.

Trakhtenbrot's Undecidability Theorem

Theorem (Trakhtenbrot)

If the vocabulary includes at least one relation of arity ≥ 2 , then SAT_{fin} is undecidable. (We will prove it later.)

Therefore static analysis of arbitrary FO formulas is undecidable; same as for Turing-complete programming languages. This justifies studying fragments of FO, where static analysis is possible.

We will prove Trakhtenbrot's theorem later.

The condition [at least one relation of arity \$\geq 2\$](#) is necessary. See HW 1.

Summary

Classical Model Theory:

- Concerned with satisfiability, validity, provability.
- Major, fundamental results: Gödel's completeness; Church undecidability; the Compactness Theorem; Löwenheim–Skolem; Gödel's incompleteness.

Finite Model Theory:

- Concerned with similar questions, plus evaluation.
- Major, fundamental results: Trakhtenbrot's undecidability; Fagin's 0/1-law; Fagin's $SO=NP$ theorem.

Relational Model

Origins

In 1970-1971 Tedd Codd proposed that databases should be modeled as finite structures, and queries represented by formulas.

A decade of debates followed, where the relational data model had to compete against the established CODASYL model.

This story is now the founding legend, par of the folklore of our community. A great reading is *What Goes Around Comes Around* in [Bailis et al., 2015].

Relational Databases

Fix the schema (vocabulary): R_1, R_2, \dots

A **relational database instance** is a finite structure $\mathbf{D} = (D, R_1^D, R_2^D, \dots)$

We often omit the domain and write $\mathbf{D} = (R_1^D, R_2^D, \dots)$.

The **active domain**, $\text{ADom}(\mathbf{D})$, is the set of constants that occur in R_1^D, R_2^D, \dots

A **query**, $Q(\mathbf{x})$, is an FO formula with free variables \mathbf{x} . We write (with some overloading) $Q(\mathbf{D})$ for the result of Q on a database \mathbf{D} .

The Drinkers-Beer-Bar Example

Introduced by [Ullman, 1980].

Frequents(Drinker,Bar)

Serves(Bar,Beer)

Likes(Drinker,Beer)

The Drinkers-Beer-Bar Example

Introduced by [Ullman, 1980].

Frequents(Drinker, Bar)
Serves(Bar, Beer)
Likes(Drinker, Beer)

Drinkers who frequent some bar who serve some beer that they like:

$$Q(x) = \exists y \exists z (\text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z))$$

The Drinkers-Beer-Bar Example

Introduced by [Ullman, 1980].

Frequents(Drinker, Bar)
Serves(Bar, Beer)
Likes(Drinker, Beer)

Drinkers who frequent some bar who serve some beer that they like:

$$Q(x) = \exists y \exists z (\text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z))$$

Drinkers who frequent only bars who serve only beers that they like:

$$Q(x) = \forall y (\text{Frequents}(x, y) \Rightarrow \forall z (\text{Serves}(y, z) \Rightarrow \text{Likes}(x, z)))$$

The Drinkers-Beer-Bar Example

Introduced by [Ullman, 1980].

Frequents(Drinker, Bar)
Serves(Bar, Beer)
Likes(Drinker, Beer)

Drinkers who frequent some bar who serve some beer that they like:

$$Q(x) = \exists y \exists z (\text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z))$$

Drinkers who frequent only bars who serve only beers that they like:

$$Q(x) = \forall y (\text{Frequents}(x, y) \Rightarrow \forall z (\text{Serves}(y, z) \Rightarrow \text{Likes}(x, z)))$$

The last query is incorrect!

Domain Independence

The last query is actually incorrect! Let's look at a simpler query:

$$Q(x) = \forall y (R(x, y) \Rightarrow S(x, y))$$

Domain Independence

The last query is actually incorrect! Let's look at a simpler query:

$$Q(x) = \forall y(R(x, y) \Rightarrow S(x, y))$$

Recall the “boring” definition: if $c \in D$, $c \notin \Pi_x(R^D)$ then $Q(c)$ is true.
 Q returns values c that are not in the active domain; **domain dependent**.

Domain Independence

The last query is actually incorrect! Let's look at a simpler query:

$$Q(x) = \forall y(R(x, y) \Rightarrow S(x, y))$$

Recall the “boring” definition: if $c \in D$, $c \notin \Pi_x(R^D)$ then $Q(c)$ is true. Q returns values c that are not in the active domain; **domain dependent**.

Definition

An FO formula (query) is **domain independent** if it does not depend on the domain D of the structure (D, R_1^D, R_2^D, \dots) .

Domain Independence

The last query is actually incorrect! Let's look at a simpler query:

$$Q(x) = \forall y(R(x, y) \Rightarrow S(x, y))$$

Recall the “boring” definition: if $c \in D$, $c \notin \Pi_x(R^D)$ then $Q(c)$ is true.
 Q returns values c that are not in the active domain; **domain dependent**.

Definition

An FO formula (query) is **domain independent** if it does not depend on the domain D of the structure (D, R_1^D, R_2^D, \dots) .

Are these queries independent?

$$Q(x, y) = R(x) \vee S(y)$$

Domain Independence

The last query is actually incorrect! Let's look at a simpler query:

$$Q(x) = \forall y(R(x, y) \Rightarrow S(x, y))$$

Recall the “boring” definition: if $c \in D$, $c \notin \Pi_x(R^D)$ then $Q(c)$ is true.
 Q returns values c that are not in the active domain; **domain dependent**.

Definition

An FO formula (query) is **domain independent** if it does not depend on the domain D of the structure (D, R_1^D, R_2^D, \dots) .

Are these queries independent?

$$Q(x, y) = R(x) \vee S(y)$$

Domain dependent.

Domain Independence

The last query is actually incorrect! Let's look at a simpler query:

$$Q(x) = \forall y(R(x, y) \Rightarrow S(x, y))$$

Recall the “boring” definition: if $c \in D$, $c \notin \Pi_x(R^D)$ then $Q(c)$ is true.
 Q returns values c that are not in the active domain; **domain dependent**.

Definition

An FO formula (query) is **domain independent** if it does not depend on the domain D of the structure (D, R_1^D, R_2^D, \dots) .

Are these queries independent?

$$Q(x, y) = R(x) \vee S(y)$$

Domain dependent.

$$Q(x) = R(x) \wedge \exists y(\neg S(x, y))$$

Domain Independence

The last query is actually incorrect! Let's look at a simpler query:

$$Q(x) = \forall y(R(x, y) \Rightarrow S(x, y))$$

Recall the “boring” definition: if $c \in D$, $c \notin \Pi_x(R^D)$ then $Q(c)$ is true.
 Q returns values c that are not in the active domain; **domain dependent**.

Definition

An FO formula (query) is **domain independent** if it does not depend on the domain D of the structure (D, R_1^D, R_2^D, \dots) .

Are these queries independent?

$$Q(x, y) = R(x) \vee S(y)$$

Domain dependent.

$$Q(x) = R(x) \wedge \exists y(\neg S(x, y))$$

Domain dependent.

Domain Independence

Given a formula φ , how do we check whether it is domain independent?

Domain Independence

Given a formula φ , how do we check whether it is domain independent?

Theorem

Checking domain independence is undecidable.

Domain Independence

Given a formula φ , how do we check whether it is domain independent?

Theorem

Checking domain independence is undecidable.

Proof Assuming an algorithm for checking domain independence, we solve SAT_{fin} , which contradicts Trakhtenbrot's theorem:

- Fix some domain-dependent query, say $\varphi = \forall x R(x)$.
- Given an FO sentence Φ , construct a new sentence $\psi \stackrel{\text{def}}{=} \Phi \wedge \varphi$.
- Then ψ is domain independent iff Φ is unsatisfiable.

Domain Independence

Given a formula φ , how do we check whether it is domain independent?

Theorem

Checking domain independence is undecidable.

Proof Assuming an algorithm for checking domain independence, we solve SAT_{fin} , which contradicts Trakhtenbrot's theorem:

- Fix some domain-dependent query, say $\varphi = \forall x R(x)$.
- Given an FO sentence Φ , construct a new sentence $\psi \stackrel{\text{def}}{=} \Phi \wedge \varphi$.
- Then ψ is domain independent iff Φ is unsatisfiable.

Syntactic restriction: Q is **range-restricted** if each var is restricted to (a subset of) ADom :

$$Q(x) = \exists u (R(x, u) \vee S(x, u)) \wedge (\forall y (R(x, y) \Rightarrow S(x, y)))$$

Relational Algebra – Quick Review

Five operators:

- Selection σ
- Projection Π
- Join \bowtie
- Union \cup
- Difference $-$

Relational Algebra – Quick Review

Five operators:

- Selection σ
- Projection Π
- Join \bowtie
- Union \cup
- Difference $-$

Who likes Leffe?

$$Q_1(x) = \text{Likes}(x, \text{'Leffe'})$$

Relational Algebra – Quick Review

Five operators:

- Selection σ
- Projection Π
- Join \bowtie
- Union \cup
- Difference $-$

Who likes Leffe?

$$Q_1(x) = \text{Likes}(x, \text{'Leffe'})$$

$$\begin{array}{c} \Pi_x \\ | \\ \sigma_{y=\text{'Leffe'}} \\ | \\ \text{Likes}(x,y) \end{array}$$

Relational Algebra – Quick Review

Five operators:

- Selection σ
- Projection Π
- Join \bowtie
- Union \cup
- Difference $-$

$$Q_2(x, y, z) = \text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z)$$

Who likes Leffe?

$$Q_1(x) = \text{Likes}(x, \text{'Leffe'})$$

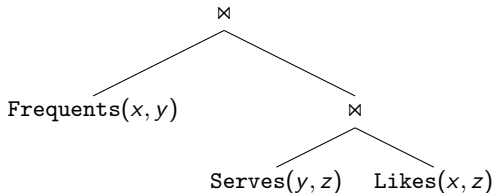
$$\begin{array}{c} \Pi_x \\ | \\ \sigma_{y=\text{'Leffe'}} \\ | \\ \text{Likes}(x, y) \end{array}$$

Relational Algebra – Quick Review

Five operators:

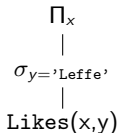
- Selection σ
- Projection Π
- Join \bowtie
- Union \cup
- Difference $-$

$$Q_2(x, y, z) = \text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z)$$



Who likes Leffe?

$$Q_1(x) = \text{Likes}(x, \text{'Leffe'})$$

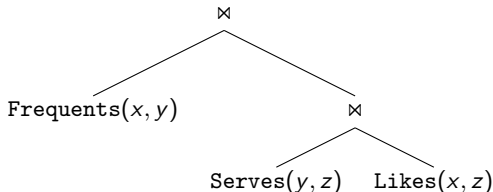


Relational Algebra – Quick Review

Five operators:

- Selection σ
- Projection Π
- Join \bowtie
- Union \cup
- Difference $-$

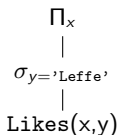
$$Q_2(x, y, z) = \text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z)$$



$$Q_3(x) = A(x) \wedge \forall y (B(y) \Rightarrow C(x, y))$$

Who likes Leffe?

$$Q_1(x) = \text{Likes}(x, \text{'Leffe'})$$

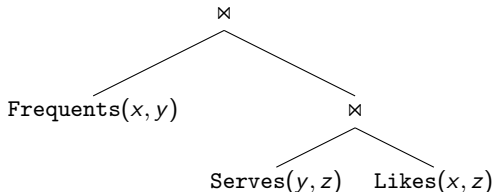


Relational Algebra – Quick Review

Five operators:

- Selection σ
- Projection Π
- Join \bowtie
- Union \cup
- Difference $-$

$$Q_2(x, y, z) = \text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z)$$

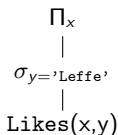


$$Q_3(x) = A(x) \wedge \forall y (B(y) \Rightarrow C(x, y))$$

$$Q_3(x) = A(x) \wedge \neg \exists y (B(y) \wedge \neg C(x, y))$$

Who likes Leffe?

$$Q_1(x) = \text{Likes}(x, \text{'Leffe'})$$

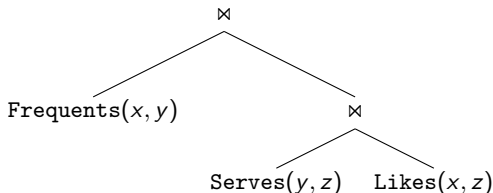


Relational Algebra – Quick Review

Five operators:

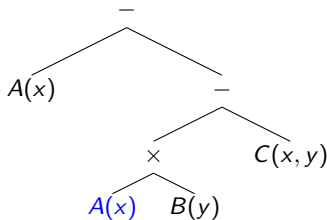
- Selection σ
- Projection Π
- Join \bowtie
- Union \cup
- Difference $-$

$$Q_2(x, y, z) = \text{Frequents}(x, y) \wedge \text{Serves}(y, z) \wedge \text{Likes}(x, z)$$



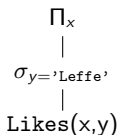
$$Q_3(x) = A(x) \wedge \forall y (B(y) \Rightarrow C(x, y))$$

$$Q_3(x) = A(x) \wedge \neg \exists y (B(y) \wedge \neg C(x, y))$$



Who likes Leffe?

$$Q_1(x) = \text{Likes}(x, \text{'Leffe'})$$



Easier with an **anti-semijoin** (look it up).

FO and RA are Equivalent

Theorem

Domain-independent FO and Relational Algebra express the same class of queries.

Proof: exercise.

FO and RA are Equivalent

Theorem

Domain-independent FO and Relational Algebra express the same class of queries.

Proof: exercise.

Physical independence principle: separation of **What** from **How**.

- Users write **what** they want, in a declarative language (FO).
- System decides **how** to compute the query most efficiently (RA plan).

Summary

- Relational data model is founded on finite model theory.
- **Physical Data Independence** is perhaps the deepest reason why it is still successful 50 years later: separate the **What** from the **How**.
- **What** is in FO. But too abstract for the real world (e.g. domain independence!), hence SQL and its history.
- **Why** is RA. But too limited for the real world, hence extended with aggregates, group-by, dependent joins, anti-semijoins, etc, etc.
- FO used in databases beyond query expressions: for constraints, optimization rules, verification.

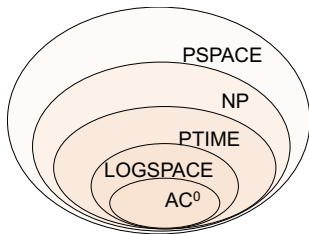
The Query Evaluation Problem

Complexity

A Turing-complete language can express any computable problem.

But FO is restricted. What is the complexity of the problems it can express?

First, we are interested in the **complexity class**.
Later we will study **efficient algorithms**.



The Query Evaluation Problem

Given a query Q and a database instance D , compute $Q(D)$.
This is the bread-and-butter of database engines.

The Query Evaluation Problem

Given a query Q and a database instance \mathbf{D} , compute $Q(\mathbf{D})$.

This is the bread-and-butter of database engines.

Definition (Complexity of Query Evaluation [Vardi, 1982])

Three ways to define the complexity:

- **Data Complexity**. Fix the query Q , complexity is $f(|\mathbf{D}|)$.
- **Query Complexity**. Fix the database \mathbf{D} , complexity is $f(|Q|)$.
A.k.a. **expression complexity**.
- **Combined Complexity**, $f(|Q|, |\mathbf{D}|)$.

The Query Evaluation Problem

Given a query Q and a database instance \mathbf{D} , compute $Q(\mathbf{D})$.

This is the bread-and-butter of database engines.

Definition (Complexity of Query Evaluation [Vardi, 1982])

Three ways to define the complexity:

- **Data Complexity**. Fix the query Q , complexity is $f(|\mathbf{D}|)$.
- **Query Complexity**. Fix the database \mathbf{D} , complexity is $f(|Q|)$.
A.k.a. **expression complexity**.
- **Combined Complexity**, $f(|Q|, |\mathbf{D}|)$.

Which is most important in practice?

Data Complexity of FO is in AC^0

Theorem

The Data Complexity of FO is in AC^0

(Stronger: it is in **uniform AC^0** , but we will ignore this.)

Recall that AC^0 is at the bottom of the hierarchy:

$$AC^0 \subseteq LOGSPACE \subseteq \dots \subseteq PTIME$$

Before we prove the theorem let's prove something simpler:
 The Data Complexity of FO is in PTIME.

Data Complexity of FO is in PTIME: Proof

How do we evaluate this? $Q = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

Data Complexity of FO is in PTIME: Proof

How do we evaluate this? $Q = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

```

some_x = false;
for x = 1,N do:
  if A(x) then:
    all_y = true
    for y = 1,N do:
      if not (B(y) => C(x,y))
        then: all_y = false;
    if all_y then: some_x = true;
return some_x
    
```

Data Complexity of FO is in PTIME: Proof

How do we evaluate this?

$$Q = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$$

```

some_x = false;
for x = 1,N do:
  if A(x) then:
    all_y = true
    for y = 1,N do:
      if not (B(y) => C(x,y))
        then: all_y = false;
    if all_y then: some_x = true;
return some_x
    
```

- Generalizes to any sentence φ .
- Runtime $O(N^k)$, where:
 $N = |\text{ADom}|$
 $k = |\text{Vars}(\varphi)|$
- In PTIME (and in LOGSPACE),
 for fixed φ .

Data Complexity of FO is in PTIME: Proof

How do we evaluate this? $Q = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

```

some_x = false;
for x = 1,N do:
  if A(x) then:
    all_y = true
    for y = 1,N do:
      if not (B(y) => C(x,y))
        then: all_y = false;
    if all_y then: some_x = true;
return some_x
    
```

- Generalizes to any sentence φ .
- Runtime $O(N^k)$, where:
 $N = |\text{ADom}|$
 $k = |\text{Vars}(\varphi)|$
- In PTIME (and in LOGSPACE),
 for fixed φ .

Many texts state that the data complexity is in LOGSPACE, or in PTIME.
 The correct complexity is AC^0 . **Let's prove it**

Definition of AC^0

Definition

A problem is in AC^0 if $\forall N$, there exists a circuit of polynomial size and constant depth, consisting NOT gates and unbounded fan-in AND and OR gates, that computes the problem on inputs of size N encoded using N bits.

Definition of AC^0

Definition

A problem is in AC^0 if $\forall N$, there exists a circuit of polynomial size and constant depth, consisting NOT gates and unbounded fan-in AND and OR gates, that computes the problem on inputs of size N encoded using N bits.

E.g. Given an undirected graph with n nodes, check if it has a triangle.

Definition of AC^0

Definition

A problem is in AC^0 if $\forall N$, there exists a circuit of polynomial size and constant depth, consisting NOT gates and unbounded fan-in AND and OR gates, that computes the problem on inputs of size N encoded using N bits.

E.g. Given an undirected graph with n nodes, check if it has a triangle.
When $n = 4$ there are $N = 6$ possible edges, E_{ij} for $1 \leq i < j \leq 4$

Definition of AC^0

Definition

A problem is in AC^0 if $\forall N$, there exists a circuit of polynomial size and constant depth, consisting NOT gates and unbounded fan-in AND and OR gates, that computes the problem on inputs of size N encoded using N bits.

E.g. Given an undirected graph with n nodes, check if it has a triangle.
When $n = 4$ there are $N = 6$ possible edges, E_{ij} for $1 \leq i < j \leq 4$

$$(E_{12} \wedge E_{23} \wedge E_{13}) \vee (E_{12} \wedge E_{24} \wedge E_{14}) \vee (E_{23} \wedge E_{34} \wedge E_{24}) \vee (E_{34} \wedge E_{14} \wedge E_{13})$$

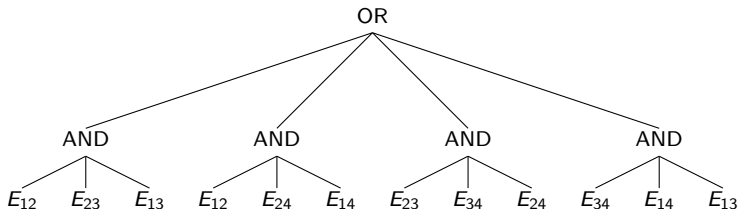
Definition of AC^0

Definition

A problem is in AC^0 if $\forall N$, there exists a circuit of polynomial size and constant depth, consisting NOT gates and unbounded fan-in AND and OR gates, that computes the problem on inputs of size N encoded using N bits.

E.g. Given an undirected graph with n nodes, check if it has a triangle.
When $n = 4$ there are $N = 6$ possible edges, E_{ij} for $1 \leq i < j \leq 4$

$$(E_{12} \wedge E_{23} \wedge E_{13}) \vee (E_{12} \wedge E_{24} \wedge E_{14}) \vee (E_{23} \wedge E_{34} \wedge E_{24}) \vee (E_{34} \wedge E_{14} \wedge E_{13})$$



Data Complexity FO is in AC^0 : Proof

Fix a Boolean query Q in FO. Encode the input database \mathbf{D} using bits:

- Let $N = |\text{ADom}(\mathbf{D})|$.

Data Complexity FO is in AC^0 : Proof

Fix a Boolean query Q in FO. Encode the input database \mathbf{D} using bits:

- Let $N = |\text{ADom}(\mathbf{D})|$.
- Encode a relation of arity k using N^k bits.
E.g. $R(X, Y)$ encoded using Boolean matrix R_{ij} .

Data Complexity FO is in AC^0 : Proof

Fix a Boolean query Q in FO. Encode the input database \mathbf{D} using bits:

- Let $N = |\text{ADom}(\mathbf{D})|$.
- Encode a relation of arity k using N^k bits.
E.g. $R(X, Y)$ encoded using Boolean matrix R_{ij} .
- Expand Q into a Boolean formula $\Phi_{Q, \mathbf{D}}$, called the **lineage** of Q on \mathbf{D} .

Data Complexity FO is in AC^0 : Proof

Fix a Boolean query Q in FO. Encode the input database D using bits:

- Let $N = |\text{ADom}(D)|$.
- Encode a relation of arity k using N^k bits.
E.g. $R(X, Y)$ encoded using Boolean matrix R_{ij} .
- Expand Q into a Boolean formula $\Phi_{Q,D}$, called the **lineage** of Q on D .
- Represent $\Phi_{Q,D}$ using an AC^0 circuit.

Data Complexity FO is in AC^0 : Proof

Fix a Boolean query Q in FO. Encode the input database D using bits:

- Let $N = |\text{ADom}(D)|$.
- Encode a relation of arity k using N^k bits.
E.g. $R(X, Y)$ encoded using Boolean matrix R_{ij} .
- Expand Q into a Boolean formula $\Phi_{Q,D}$, called the **lineage** of Q on D .
- Represent $\Phi_{Q,D}$ using an AC^0 circuit.

Example: $Q = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

Data Complexity FO is in AC^0 : Proof

Fix a Boolean query Q in FO. Encode the input database \mathbf{D} using bits:

- Let $N = |\text{ADom}(\mathbf{D})|$.
- Encode a relation of arity k using N^k bits.
E.g. $R(X, Y)$ encoded using Boolean matrix R_{ij} .
- Expand Q into a Boolean formula $\Phi_{Q,\mathbf{D}}$, called the **lineage** of Q on \mathbf{D} .
- Represent $\Phi_{Q,\mathbf{D}}$ using an AC^0 circuit.

Example: $Q = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

Its **lineage** is: $\Phi_{Q,\mathbf{D}} = \bigvee_{i=1,N} \left(A_i \wedge \bigwedge_{j=1,N} (\neg B_j \vee C_{ij}) \right)$

Data Complexity FO is in AC^0 : Proof

Fix a Boolean query Q in FO. Encode the input database \mathbf{D} using bits:

- Let $N = |\text{ADom}(\mathbf{D})|$.
- Encode a relation of arity k using N^k bits.
E.g. $R(X, Y)$ encoded using Boolean matrix R_{ij} .
- Expand Q into a Boolean formula $\Phi_{Q,\mathbf{D}}$, called the **lineage** of Q on \mathbf{D} .
- Represent $\Phi_{Q,\mathbf{D}}$ using an AC^0 circuit.

Example: $Q = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

Its **lineage** is: $\Phi_{Q,\mathbf{D}} = \bigvee_{i=1,N} \left(A_i \wedge \bigwedge_{j=1,N} (\neg B_j \vee C_{ij}) \right)$

In class: construct a circuit of depth 5 and size $O(N^2)$.

Summary

- Data complexity is in AC^0 ; this implies LOGSPACE, PTIME.
- Expression complexity, combined complexity: PSPACE complete
We will discuss this later.
- AC^0 is the class of highly parallelizable problems.
- “SQL is embarrassingly parallel”

Restricted Query Languages

Motivation

- FO is too rich for powerful optimizations: Trakhtenbrot's theorem is a fundamental limit.
- For fragments of FO static analysis is possible, and they still capture the most important queries in practice.
- Assuming FO consists of $\exists, \forall, \wedge, \vee, \neg, =$, we will obtain fragments by restricting the connectives.

Conjunctive Queries

Definition

A **Conjunctive Query (CQ)** is an expression of the form:

$$Q(\mathbf{x}_0) = \exists \mathbf{y} \bigwedge_i R_i(\mathbf{x}_i)$$

Conjunctive Queries

Definition

A **Conjunctive Query (CQ)** is an expression of the form:

$$Q(\mathbf{x}_0) = \exists \mathbf{y} \bigwedge_i R_i(\mathbf{x}_i)$$

- E.g. $Q(x, y) = \exists z (E(x, z) \wedge E(z, y))$.
- Equivalently: FO formula restricted to $=, \wedge, \exists$ **What fragment of RA?**

Conjunctive Queries

Definition

A **Conjunctive Query (CQ)** is an expression of the form:

$$Q(\mathbf{x}_0) = \exists \mathbf{y} \bigwedge_i R_i(\mathbf{x}_i)$$

- E.g. $Q(x, y) = \exists z (E(x, z) \wedge E(z, y))$.
- Equivalently: FO formula restricted to $=, \wedge, \exists$ **What fragment of RA?**
- CQ has the same expressive power as RA restricted to σ, Π, \bowtie .
- These correspond to SELECT-FROM-WHERE queries in SQL (but we have to be careful what we allow in each clause).

Unions of Conjunctive Queries

Definition

A **Union of Conjunctive Queries (UCQ)** is a formula of the form:

$$Q(\mathbf{x}) = \bigvee_i Q_i(\mathbf{x})$$

where all Q_i 's are CQs, and have the same sets of free variables.

Unions of Conjunctive Queries

Definition

A **Union of Conjunctive Queries (UCQ)** is a formula of the form:

$$Q(\mathbf{x}) = \bigvee_i Q_i(\mathbf{x})$$

where all Q_i 's are CQs, and have the same sets of free variables.

- E.g. $Q(x, y) = E(x, y) \vee \exists z(E(x, z) \wedge E(z, y))$.
- Equivalently, UCQs are FO formulas restricted to $=, \wedge, \exists, \vee$.
- UCQ has the same expressive power as RA restricted to $\sigma, \Pi, \bowtie, \cup$.

Monotone Queries

Given two databases D, D' over the same schema, we write $D \subseteq D'$ if $R_i^D \subseteq R_i^{D'}$ for every relation R_i in the schema.

Definition

A query Q is **monotone** if $D \subseteq D'$ implies $Q(D) \subseteq Q(D')$.

Example: $\exists x, y, z (E(x, y) \wedge E(y, z))$

Non-example: $\exists x \forall y (V(x) \wedge \forall y (V(y) \Rightarrow E(x, y)))$

Monotone Queries

Given two databases D, D' over the same schema, we write $D \subseteq D'$ if $R_i^D \subseteq R_i^{D'}$ for every relation R_i in the schema.

Definition

A query Q is **monotone** if $D \subseteq D'$ implies $Q(D) \subseteq Q(D')$.

Example: $\exists x, y, z (E(x, y) \wedge E(y, z))$ **Non-example:** $\exists x \forall y (V(y) \Rightarrow E(x, y))$

All UCQ queries are monotone. **Exercise**

The only non-monotone operators are:

- negation \neg in FO.
- difference $-$ in RA.

Other Ways to Restrict the Query Language (1/2)

Adding $\neq, <, \leq$ to CQ, UCQ:

- By default they are not allowed in CQ, UCQ.
- If we want them, we write e.g. CQ^{\neq} or UCQ^{\leq} .
- $Q(x, y) = \exists u \exists v (E(x, u) \wedge E(u, v) \wedge E(v, y) \wedge x < u < v < y)$.

Other Ways to Restrict the Query Language (1/2)

Adding $\neq, <, \leq$ to CQ, UCQ:

- By default they are not allowed in CQ, UCQ.
- If we want them, we write e.g. CQ^{\neq} or UCQ^{\leq} .
- $Q(x, y) = \exists u \exists v (E(x, u) \wedge E(u, v) \wedge E(v, y) \wedge x < u < v < y)$.
Is this query monotone?

Other Ways to Restrict the Query Language (1/2)

Adding $\neq, <, \leq$ to CQ, UCQ:

- By default they are not allowed in CQ, UCQ.
- If we want them, we write e.g. CQ^{\neq} or UCQ^{\leq} .
- $Q(x, y) = \exists u \exists v (E(x, u) \wedge E(u, v) \wedge E(v, y) \wedge x < u < v < y)$.
Is this query monotone?

YES!

Other Ways to Restrict the Query Language (2/2)

Restricting the number of variables in FO:

- FO^k : restricted to using only k variables.
- E.g. check in FO^2 if there a path of length ≥ 5 :

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

Other Ways to Restrict the Query Language (2/2)

Restricting the number of variables in FO:

- FO^k : restricted to using only k variables.
- E.g. check in FO^2 if there a path of length ≥ 5 :

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

Theorem ([Grädel et al., 1997])

If $\varphi \in FO^2$ has any model (possibly infinite), then it has a model of size at most exponential in $|\varphi|$. Thus, $SAT_{fin}(FO^2)$ is decidable.

Other Ways to Restrict the Query Language (2/2)

Restricting the number of variables in FO:

- FO^k : restricted to using only k variables.
- E.g. check in FO^2 if there a path of length ≥ 5 :

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

Theorem ([Grädel et al., 1997])

If $\varphi \in FO^2$ has any model (possibly infinite), then it has a model of size at most exponential in $|\varphi|$. Thus, $SAT_{fin}(FO^2)$ is decidable.

Suggested research: what are the implications for a query optimizer?

Other Ways to Restrict the Query Language (2/2)

Restricting the number of variables in FO:

- FO^k : restricted to using only k variables.
- E.g. check in FO^2 if there a path of length ≥ 5 :

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

Theorem ([Grädel et al., 1997])

If $\varphi \in FO^2$ has any model (possibly infinite), then it has a model of size at most exponential in $|\varphi|$. Thus, $SAT_{fin}(FO^2)$ is decidable.

Suggested research: what are the implications for a query optimizer?

What about FO^3 ?

Other Ways to Restrict the Query Language (2/2)

Restricting the number of variables in FO:

- FO^k : restricted to using only k variables.
- E.g. check in FO^2 if there a path of length ≥ 5 :

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

Theorem ([Grädel et al., 1997])

If $\varphi \in FO^2$ has any model (possibly infinite), then it has a model of size at most exponential in $|\varphi|$. Thus, $SAT_{fin}(FO^2)$ is decidable.

Suggested research: what are the implications for a query optimizer?

What about FO^3 ?

To watch how many variables we need to prove Trakhtenbrot's theorem

Conjunctive Queries

Are the most important and most studied fragment. Terminology:

- **Boolean query**: no head vars: $Q() = \exists x \exists y \exists z (E(x, y) \wedge E(y, z)).$
- **Full query**: no existential vars: $Q(x, y, z) = E(x, y) \wedge E(y, z).$
- **Without selfjoins**: every relation name occurs at most once.

$$Q(x) = \exists y \exists z (R(x, y) \wedge S(y, z) \wedge T(z, x)).$$

- We often omit the existential quantifiers, and write for example:

$$Q(x) = R(x, y) \wedge S(y, z) \wedge T(z, x).$$

Summary

- Most of our discussion will be focused on CQ's.
- UCQs come almost for free, or with very little additional effort.
- Let's re-examine query evaluation when the query is restricted to a CQ.

Query Evaluation for CQ

Motivation

We already know that the data complexity is in AC^0 .

What is the expression complexity? The combined complexity?

Will answer both, and also discuss the expression/combined complexity for FO (which we left out).

Importantly: we will define query evaluation for CQ in terms of
Homomorphisms

Equivalent Concepts

- A Conjunctive Query:

$$R(x, y, z) \wedge S(x, u) \wedge S(y, v) \wedge S(z, w) \wedge R(u, v, w)$$

- A database instance:

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

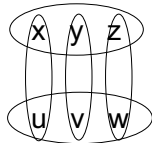
A	B	C
x	y	z
u	v	w

$$S(D, E) =$$

D	E
x	u
y	v
z	w

- A labeled hypergraph, $G = (V, E)$, where

$V = \{x, y, z, u, v, w\}$, $E = \{\{x, y, z\}, \{u, v, w\}, \{x, u\}, \{y, v\}, \{z, w\}\}$
(hyperedges are labeled with R, S respectively).



Equivalent Concepts

- A Conjunctive Query:

$$R(x, y, z) \wedge S(x, u) \wedge S(y, v) \wedge S(z, w) \wedge R(u, v, w)$$

- A database instance:

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

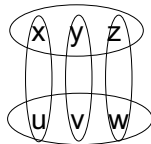
A	B	C
x	y	z
u	v	w

$$S(D, E) =$$

D	E
x	u
y	v
z	w

- A labeled hypergraph, $G = (V, E)$, where

$V = \{x, y, z, u, v, w\}$, $E = \{\{x, y, z\}, \{u, v, w\}, \{x, u\}, \{y, v\}, \{z, w\}\}$
(hyperedges are labeled with R, S respectively).



We will often switch back-and-forth between these equivalent notions

Homomorphisms

$$Q(\mathbf{x}_0) = R_1(\mathbf{x}_1) \wedge \cdots \wedge R_m(\mathbf{x}_m), \quad Q'(\mathbf{y}_0) = S_1(\mathbf{y}_1) \wedge \cdots \wedge S_n(\mathbf{y}_n).$$

Definition

A **homomorphism** $h : Q' \rightarrow Q$ is a function

$h : \text{Const}(Q') \cup \text{Vars}(Q') \rightarrow \text{Const}(Q) \cup \text{Vars}(Q)$ s.t.:

- $\forall c \in \text{Const}(Q'), h(c) = c.$
- $S_j(\mathbf{y}_j) \in \text{Atoms}(Q'), \exists R_i(\mathbf{x}_i) \in \text{Atoms}(Q)$ such that $R_i = S_j$ (they are the same relation name) and $h(\mathbf{y}_j) = \mathbf{x}_i.$
- h maps head vars to head vars: $h(\mathbf{y}_0) = \mathbf{x}_0.$

Homomorphisms

$$Q(\mathbf{x}_0) = R_1(\mathbf{x}_1) \wedge \cdots \wedge R_m(\mathbf{x}_m), \quad Q'(\mathbf{y}_0) = S_1(\mathbf{y}_1) \wedge \cdots \wedge S_n(\mathbf{y}_n).$$

Definition

A **homomorphism** $h : Q' \rightarrow Q$ is a function

$h : \text{Const}(Q') \cup \text{Vars}(Q') \rightarrow \text{Const}(Q) \cup \text{Vars}(Q)$ s.t.:

- $\forall c \in \text{Const}(Q'), h(c) = c.$
- $S_j(\mathbf{y}_j) \in \text{Atoms}(Q'), \exists R_i(\mathbf{x}_i) \in \text{Atoms}(Q)$ such that $R_i = S_j$ (they are the same relation name) and $h(\mathbf{y}_j) = \mathbf{x}_i.$
- h maps head vars to head vars: $h(\mathbf{y}_0) = \mathbf{x}_0.$

Graph homomorphism $h : G' \rightarrow G$ is $h : V \rightarrow V'$ s.t. $\forall e \in E', h(e) \in E.$

Query Evaluation for CQ and Homomorphisms

Computing $Q(\mathbf{D})$ consists of finding all homomorphisms $h : Q \rightarrow D$ and returning $h(\text{Head}(Q))$.

$$Q(x) = R(x) \wedge S(x, y) \wedge T(y, 'a')$$

$R =$	<table><tr><th>x</th></tr><tr><td>1</td></tr><tr><td>2</td></tr></table>	x	1	2	$S =$	<table><tr><th>x</th><th>y</th></tr><tr><td>1</td><td>10</td></tr><tr><td>1</td><td>20</td></tr><tr><td>2</td><td>20</td></tr></table>	x	y	1	10	1	20	2	20	$T =$	<table><tr><th>y</th><th>z</th></tr><tr><td>10</td><td><i>a</i></td></tr><tr><td>10</td><td><i>b</i></td></tr><tr><td>20</td><td><i>a</i></td></tr></table>	y	z	10	<i>a</i>	10	<i>b</i>	20	<i>a</i>
	x																							
	1																							
2																								
x	y																							
1	10																							
1	20																							
2	20																							
y	z																							
10	<i>a</i>																							
10	<i>b</i>																							
20	<i>a</i>																							

We list all homomorphisms:

	$x(= \text{Head}(Q))$	y	a
$h =$	1	10	a
	1	20	a
	2	20	a

Final answer after duplicate elimination: $Q(\mathbf{D}) = \{1, 2\}$.

The Combined Complexity for UCQ is in NP

Theorem

The combined complexity for UCQ is in NP.

Proof: Fix a UCQ $Q = Q_1 \vee Q_2 \vee \dots$ and a database D .

To check $D \models Q$:

- “guess” a CQ Q_i , and
- “guess” a homomorphism $h : Q_i \rightarrow D$

The Expression Complexity for CQ is NP-hard

Theorem

There exists a database D for which the expression complexity of CQ queries is NP complete.

Thus, the expression complexity is also NP-complete.

The Expression Complexity for CQ is NP-hard

Theorem

There exists a database D for which the expression complexity of CQ queries is NP complete.

Thus, the expression complexity is also NP-complete.

Proof Many proofs are possible (will explain shortly why). We will use reduction from 3SAT, because we will reuse it a few times.

The Expression Complexity for CQ is NP-hard

Theorem

There exists a database \mathbf{D} for which the expression complexity of CQ queries is NP complete.

Thus, the expression complexity is also NP-complete.

Proof Many proofs are possible (will explain shortly why). We will use reduction from 3SAT, because we will reuse it a few times.

Given a 3CNF formula Φ we construct Q_Φ, \mathbf{D} such that:

The Expression Complexity for CQ is NP-hard

Theorem

There exists a database D for which the expression complexity of CQ queries is NP complete.

Thus, the expression complexity is also NP-complete.

Proof Many proofs are possible (will explain shortly why). We will use reduction from 3SAT, because we will reuse it a few times.

Given a 3CNF formula Φ we construct Q_Φ, D such that:

$\Phi \text{ is satisfiable iff } \exists h : Q_\Phi \rightarrow D.$

Notice that D is independent of Φ .

Details next.

Reduction from 3SAT to CQ Evaluation

Given a 3CNF formula Φ we construct Q_Φ, \mathbf{D} such that:

Φ is satisfiable iff $\exists h : Q_\Phi \rightarrow \mathbf{D}$.

Q_Φ has one atom for each clause C in Φ :

- If $C = (X_i \vee X_j \vee X_k)$ then Q contains $A(x_i, x_j, x_k)$.

Reduction from 3SAT to CQ Evaluation

Given a 3CNF formula Φ we construct Q_Φ, \mathbf{D} such that:

Φ is satisfiable iff $\exists h : Q_\Phi \rightarrow \mathbf{D}$.

Q_Φ has one atom for each clause C in Φ :

- If $C = (X_i \vee X_j \vee X_k)$ then Q contains $A(x_i, x_j, x_k)$.
- If $C = (X_i \vee X_j \vee \neg X_k)$ then Q contains $B(x_i, x_j, x_k)$.

Reduction from 3SAT to CQ Evaluation

Given a 3CNF formula Φ we construct Q_Φ, \mathbf{D} such that:

Φ is satisfiable iff $\exists h : Q_\Phi \rightarrow \mathbf{D}$.

Q_Φ has one atom for each clause C in Φ :

- If $C = (X_i \vee X_j \vee X_k)$ then Q contains $A(x_i, x_j, x_k)$.
- If $C = (X_i \vee X_j \vee \neg X_k)$ then Q contains $B(x_i, x_j, x_k)$.
- If $C = (X_i \vee \neg X_j \vee \neg X_k)$ then Q contains $C(x_i, x_j, x_k)$.
- If $C = (\neg X_i \vee \neg X_j \vee \neg X_k)$ then Q contains $D(x_i, x_j, x_k)$.

Reduction from 3SAT to CQ Evaluation

Given a 3CNF formula Φ we construct Q_Φ, \mathbf{D} such that:

Φ is satisfiable iff $\exists h : Q_\Phi \rightarrow \mathbf{D}$.

Q_Φ has one atom for each clause C in Φ :

- If $C = (X_i \vee X_j \vee X_k)$ then Q contains $A(x_i, x_j, x_k)$.
- If $C = (X_i \vee X_j \vee \neg X_k)$ then Q contains $B(x_i, x_j, x_k)$.
- If $C = (X_i \vee \neg X_j \vee \neg X_k)$ then Q contains $C(x_i, x_j, x_k)$.
- If $C = (\neg X_i \vee \neg X_j \vee \neg X_k)$ then Q contains $D(x_i, x_j, x_k)$.

\mathbf{D} has 4 tables with 7 tuples each **which tuple is missing?**

$$A = \begin{array}{|c|c|c|} \hline 0 & 0 & 1 \\ \hline \vdots & & \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$$B = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline \vdots & & \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$C = \dots$

$D = \dots$

Reduction from 3SAT to CQ Evaluation

Given a 3CNF formula Φ we construct Q_Φ, \mathbf{D} such that:

Φ is satisfiable iff $\exists h : Q_\Phi \rightarrow \mathbf{D}$.

Q_Φ has one atom for each clause C in Φ :

- If $C = (X_i \vee X_j \vee X_k)$ then Q contains $A(x_i, x_j, x_k)$.
- If $C = (X_i \vee X_j \vee \neg X_k)$ then Q contains $B(x_i, x_j, x_k)$.
- If $C = (X_i \vee \neg X_j \vee \neg X_k)$ then Q contains $C(x_i, x_j, x_k)$.
- If $C = (\neg X_i \vee \neg X_j \vee \neg X_k)$ then Q contains $D(x_i, x_j, x_k)$.

\mathbf{D} has 4 tables with 7 tuples each **which tuple is missing?**

$$A = \begin{array}{|c|c|c|} \hline 0 & 0 & 1 \\ \hline \vdots & & \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$$B = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline \vdots & & \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$C = \dots$

$D = \dots$

In class: Φ is satisfiable iff $\exists h : Q \rightarrow \mathbf{D}$.

Combined Complexity for FO

Recall that the combined complexity of FO is in PSPACE.

Theorem

*There exists a database **D** for which the expression complexity of FO queries is PSPACE complete.*

Thus, the combined complexity is also PSPACE-complete.

Combined Complexity for FO

Recall that the combined complexity of FO is in PSPACE.

Theorem

There exists a database \mathbf{D} for which the expression complexity of FO queries is PSPACE complete.

Thus, the combined complexity is also PSPACE-complete.

Proof: Reduction from the [Quantified Boolean Formula Satisfiability](#):

$$Q_1 X_1 \ Q_2 X_2 \ \cdots \ Q_n X_n \ \Phi$$

where Φ is 3CNF.

Use the same Q_Φ , \mathbf{D} before, but add appropriate quantifiers to Q_Φ :

$$Q X_1 \ Q X_2 \ \cdots \ Q X_n \ Q_\Phi(x_1, \dots, x_n)$$

Discussion: CQ and CSP

The generalized Constraint Satisfaction Problem is:

Definition ([Kolaitis and Vardi, 1998])

Given two classes of finite structures \mathcal{A}, \mathcal{B} , the $CSP(\mathcal{A}, \mathcal{B})$ problem is:
Given $A \in \mathcal{A}, B \in \mathcal{B}$, is there a homomorphism $h : A \rightarrow B$?

Discussion: CQ and CSP

The generalized Constraint Satisfaction Problem is:

Definition ([Kolaitis and Vardi, 1998])

Given two classes of finite structures \mathcal{A}, \mathcal{B} , the $CSP(\mathcal{A}, \mathcal{B})$ problem is:
Given $A \in \mathcal{A}, B \in \mathcal{B}$, is there a homomorphism $h : A \rightarrow B$?

Standard CSP restricts the right-hand side, $CSP(-, B)$.

What is B for 3SAT? For 3-colorability? For Hamiltonian path?

Discussion: CQ and CSP

The generalized Constraint Satisfaction Problem is:

Definition ([Kolaitis and Vardi, 1998])

Given two classes of finite structures \mathcal{A}, \mathcal{B} , the $CSP(\mathcal{A}, \mathcal{B})$ problem is:
Given $A \in \mathcal{A}, B \in \mathcal{B}$, is there a homomorphism $h : A \rightarrow B$?

Standard CSP restricts the right-hand side, $CSP(-, B)$.

What is B for 3SAT? For 3-colorability? For Hamiltonian path?

Query evaluation restricts the left-hand side, $CSP(Q, -)$

“Query evaluation is CSP from the other side.”

Summary

- Evaluating $Q(\mathbf{D})$ consists of finding homomorphisms $h : Q \rightarrow \mathbf{D}$.
- This problem is in NP, in fact it is the very definition of NP.
- If Q is fixed, then the problem is in PTIME in $|\mathbf{D}|$. [Data complexity](#)
- If Q is part of the input (i.e. can be huge) then NP-complete.
[Expression complexity](#)



Abiteboul, S., Hull, R., and Vianu, V. (1995).

Foundations of Databases.

Addison-Wesley.



Arenas, M., Barceló, P., Libkin, L., Martens, W., and Pieris, A. (2022).

Database Theory.

Open source at <https://github.com/pdm-book/community>.

See also

<https://www.theoinf.uni-bayreuth.de/en/news/2021/Book-on-Database-Theory-in-the-works/index.html>.



Bailis, P., Hellerstein, J. M., and Stonebraker, M., editors (2015).

Readings in Database Systems, 5th Edition.



Börger, E., Grädel, E., and Gurevich, Y. (1997).

The Classical Decision Problem.

Perspectives in Mathematical Logic. Springer.



Grädel, E., Kolaitis, P. G., and Vardi, M. Y. (1997).

On the decision problem for two-variable first-order logic.

Bull. Symb. Log., 3(1):53–69.



Kolaitis, P. G. and Vardi, M. Y. (1998).

Conjunctive-query containment and constraint satisfaction.

In Mendelzon, A. O. and Paredaens, J., editors, *Proceedings of the Seventeenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, June 1-3, 1998, Seattle, Washington, USA*, pages 205–213. ACM Press.



Libkin, L. (2004).

Elements of Finite Model Theory.

Texts in Theoretical Computer Science. An EATCS Series. Springer.



Libkin, L. (2009).

The finite model theory toolbox of a database theoretician.

In Paredaens, J. and Su, J., editors, *Proceedings of the Twenty-Eighth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS 2009, June 19 - July 1, 2009, Providence, Rhode Island, USA*, pages 65–76. ACM.



Ullman, J. D. (1980).

Principles of Database Systems, 1st Edition.

Computer Science Press.



Vardi, M. Y. (1982).

The complexity of relational query languages (extended abstract).

In Lewis, H. R., Simons, B. B., Burkhard, W. A., and Landweber, L. H., editors, *Proceedings of the 14th Annual ACM Symposium on Theory of Computing, May 5-7, 1982, San Francisco, California, USA*, pages 137–146. ACM.