

FOUNDATIONS OF DATABASES AND QUERY LANGUAGES

Lecture 4: Complexity of FO Query Answering

Markus Krötzsch

TU Dresden, 4 May 2015

Overview

1. Introduction | Relational data model
2. First-order queries
3. Complexity of query answering
4. Complexity of FO query answering
5. Query optimization
6. Conjunctive queries
7. Limits of first-order query expressiveness
8. Introduction to Datalog
9. Implementation techniques for Datalog
10. Path queries
11. Constraints (1)
12. Constraints (2)
13. “Buffer time”
14. Outlook: database theory in practice

Markus Krötzsch, 4 May 2015

Foundations of Databases and Query Languages

slide 2 of 28

How to Measure Query Answering Complexity

Query answering as decision problem

↪ consider Boolean queries

Various notions of complexity:

- Combined complexity (complexity w.r.t. size of query and database instance)
- Data complexity (worst case complexity for any fixed query)
- Query complexity (worst case complexity for any fixed database instance)

Various common complexity classes:

$$L \subseteq NL \subseteq P \subseteq NP \subseteq PSPACE \subseteq EXPTIME$$

An Algorithm for Evaluating FO Queries

function Eval(φ, \mathcal{I})

```

01  switch ( $\varphi$ ) {
02      case  $p(c_1, \dots, c_n)$  : return  $\langle c_1, \dots, c_n \rangle \in p^{\mathcal{I}}$ 
03      case  $\neg\psi$  : return  $\neg$ Eval( $\psi, \mathcal{I}$ )
04      case  $\psi_1 \wedge \psi_2$  : return Eval( $\psi_1, \mathcal{I}$ )  $\wedge$  Eval( $\psi_2, \mathcal{I}$ )
05      case  $\exists x.\psi$  :
06          for  $c \in \Delta^{\mathcal{I}}$  {
07              if Eval( $\psi[x \mapsto c], \mathcal{I}$ ) then return true
08          }
09      return false
10  }
```

FO Algorithm Worst-Case Runtime

Let m be the size of φ , and let $n = |\mathcal{I}|$ (total table sizes)

- How many recursive calls of Eval are there?
 \leadsto one per subexpression: at most m
- Maximum depth of recursion?
 \leadsto bounded by total number of calls: at most m
- Maximum number of iterations of **for** loop?
 $\leadsto |\Delta^{\mathcal{I}}| \leq n$ per recursion level
 \leadsto at most n^m iterations
- Checking $\langle c_1, \dots, c_n \rangle \in p^{\mathcal{I}}$ can be done in linear time w.r.t. n

Runtime in $m \cdot n^m \cdot n = m \cdot n^{m+1}$

Time Complexity of FO Algorithm

Let m be the size of φ , and let $n = |\mathcal{I}|$ (total table sizes)

Runtime in $m \cdot n^{m+1}$

Time complexity of FO query evaluation

- Combined complexity: in EXPTIME
- Data complexity (m is constant): in P
- Query complexity (n is constant): in EXPTIME

FO Algorithm Worst-Case Memory Usage

We can get better complexity bounds by looking at memory

Let m be the size of φ , and let $n = |\mathcal{I}|$ (total table sizes)

- For each (recursive) call, store pointer to current subexpression of φ : $\log m$
- For each variable in φ (at most m), store current constant assignment (as a pointer): $m \cdot \log n$
- Checking $\langle c_1, \dots, c_n \rangle \in p^{\mathcal{I}}$ can be done in logarithmic space w.r.t. n

Memory in $m \log m + m \log n + \log n = m \log m + (m + 1) \log n$

Space Complexity of FO Algorithm

Let m be the size of φ , and let $n = |\mathcal{I}|$ (total table sizes)

Memory in $m \log m + (m + 1) \log n$

Space complexity of FO query evaluation

- Combined complexity: in PSPACE
- Data complexity (m is constant): in L
- Query complexity (n is constant): in PSPACE

FO Combined Complexity

The algorithm shows that FO query evaluation is in PSPACE.
Is this the best we can get?

Hardness proof: reduce a known PSPACE-hard problem to FO query evaluation

~> QBF satisfiability

Let $Q_1X_1.Q_2X_2.\dots.Q_nX_n.\varphi[X_1, \dots, X_n]$ be a QBF (with $Q_i \in \{\forall, \exists\}$)

- Database instance \mathcal{I} with $\Delta^{\mathcal{I}} = \{0, 1\}$
- One table with one row: true(1)
- Transform input QBF into Boolean FO query

$$Q_1x_1.Q_2x_2.\dots.Q_nx_n.\varphi[X_1 \mapsto \text{true}(x_1), \dots, X_n \mapsto \text{true}(x_n)]$$

Combined Complexity of FO Query Answering

Theorem

The evaluation of FO queries is PSPACE-complete with respect to combined complexity.

We have actually shown something stronger:

Theorem

The evaluation of FO queries is PSPACE-complete with respect to query complexity.

PSPACE-hardness for DI Queries

The previous reduction from QBF may lead to a query that is not domain independent

Example: QBF $\exists p. \neg p$ leads to FO query $\exists x. \neg \text{true}(x)$

Better approach:

- Consider QBF $Q_1X_1.Q_2X_2.\dots.Q_nX_n.\varphi[X_1, \dots, X_n]$ with φ in negation normal form: negations only occur directly before variables X_i (still PSPACE-complete: exercise)
- Database instance \mathcal{I} with $\Delta^{\mathcal{I}} = \{0, 1\}$
- Two tables with one row each: true(1) and false(0)
- Transform input QBF into Boolean FO query

$$Q_1x_1.Q_2x_2.\dots.Q_nx_n.\varphi'$$

where φ' is obtained by replacing each negated variable $\neg X_i$ with false(x_i) and each non-negated variable X_i with true(x_i).

Data Complexity of FO Query Answering

The algorithm showed that FO query evaluation is in L
~> can we do any better?

What could be better than L?

$$? \subseteq L \subseteq NL \subseteq P \subseteq \dots$$

~> we need to define circuit complexities first

Boolean Circuits

Definition

A **Boolean circuit** is a finite, directed, acyclic graph where

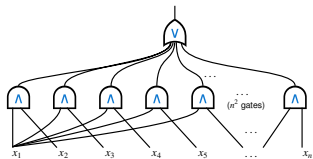
- each node that has no predecessors is an **input node**
- each node that is not an input node is one of the following types of **logical gate**: AND, OR, NOT
- one or more nodes are designated **output nodes**

→ we will only consider Boolean circuits with exactly one output

→ propositional logic formulae are Boolean circuits with one output and gates of fanout ≤ 1

Circuits as a Model for Parallel Computation

Previous example:



→ n^2 processors working in parallel

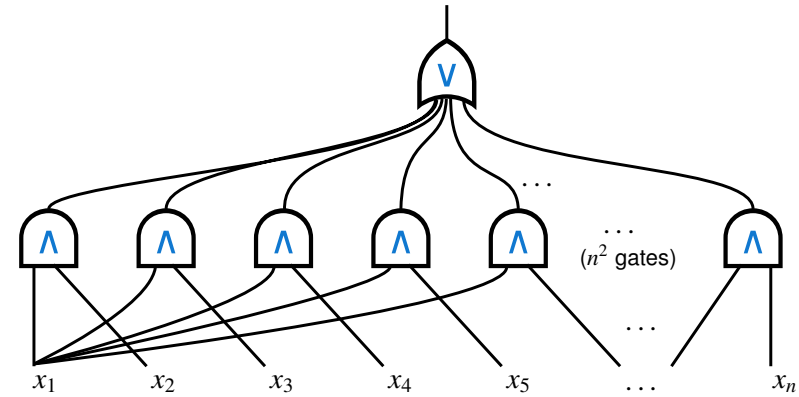
→ computation finishes in 2 steps

- **size**: number of gates = total number of computing steps
- **depth**: longest path of gates = time for parallel computation

→ refinement of polynomial time taking parallelizability into account

Example

A Boolean circuit over an input string $x_1x_2 \dots x_n$ of length n



Corresponds to formula $(x_1 \wedge x_2) \vee (x_1 \wedge x_3) \vee \dots \vee (x_{n-1} \wedge x_n)$

→ accepts all strings with at least two 1s

Solving Problems With Circuits

Observation: the input size is “hard-wired” in circuits

→ each circuit only has a finite number of different inputs

→ not a computationally interesting problem

How can we solve interesting problems with Boolean circuits?

Definition

A **uniform family** of Boolean circuits is a set of circuits C_n ($n \geq 0$) that can be computed from n (usually in logarithmic space or time; we don't discuss the details here).

A language $\mathcal{L} \subseteq \{0, 1\}^*$ is **decided by** a uniform family $(C_n)_{n \geq 0}$ of Boolean circuits if for each word w of length $|w|$:

$$w \in \mathcal{L} \quad \text{if and only if} \quad C_{|w|}(w) = 1$$

Measuring Complexity with Boolean Circuits

How to measure the computing power of Boolean circuits?

Relevant metrics:

- **size** of the circuit: overall number of gates (as function of input size)
- **depth** of the circuit: longest path of gates (as function of input size)
- **fan in**: two inputs per gate or any number of inputs per gate?

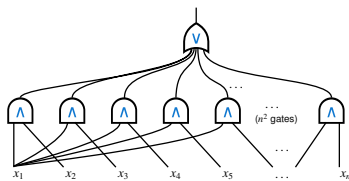
Important classes of circuits: small-depth circuits

Definition

$(C_n)_{n \geq 0}$ is a family of **small-depth circuits** if

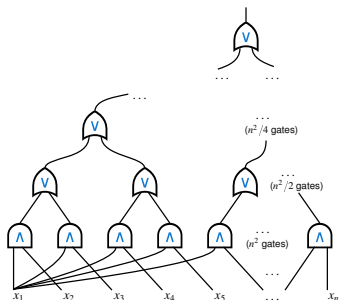
- the size of C_n is polynomial in n ,
- the depth of C_n is poly-logarithmic in n , that is, $O(\log^k n)$.

Example



family of polynomial size,
constant depth,
arbitrary fan-in circuits
 \leadsto in AC^0

We can eliminate arbitrary fan-ins by using more layers of gates:



family of polynomial size,
logarithmic depth,
bounded fan-in circuits
 \leadsto in NC^1

The Complexity Classes NC and AC

Two important types of small-depth circuits

Definition

NC^k is the class of problems that can be solved by uniform families of circuits $(C_n)_{n \geq 0}$ of fan-in ≤ 2 , size polynomial in n , and depth in $O(\log^k n)$.

The class NC is defined as $NC = \bigcup_{k \geq 0} NC^k$.

(“Nick’s Class” named after Nicholas Pippenger by Stephen Cook)

Definition

AC^k and AC are defined like NC^k and NC , respectively, but for circuits with arbitrary fan-in.

(A is for “Alternating”: AND-OR gates alternate in such circuits)

Relationships of Circuit Complexity Classes

The previous sketch can be generalised:

$$NC^0 \subseteq AC^0 \subseteq NC^1 \subseteq AC^1 \subseteq \dots \subseteq AC^k \subseteq NC^{k+1} \subseteq \dots$$

Only few inclusions are known to be proper: $NC^0 \subset AC^0 \subset NC^1$

Direct consequence of above hierarchy: $NC = AC$

Interesting relations to other classes:

$$NC^0 \subset AC^0 \subset NC^1 \subseteq L \subseteq NL \subseteq AC^1 \subseteq \dots \subseteq NC \subseteq P$$

Intuition:

- Problems in NC are parallelisable
- Problems in $P \setminus NC$ are inherently sequential

However: it is not known if $NC \neq P$

Theorem

The evaluation of FO queries is complete for (logtime uniform) AC^0 with respect to data complexity.

Proof:

- **Membership:** For a fixed Boolean FO query, provide a uniform construction for a small-depth circuit based on the size of a database
- **Hardness:** Show that circuits can be transformed into Boolean FO queries in logarithmic time (not on a standard TM ... not in this lecture)

Example

We consider the formula

$$\exists z. (\exists x. \exists y. R(x, y) \wedge S(y, z)) \wedge \neg R(a, z)$$

Over the database instance:

R:

a	a
a	b

S:

b	b
b	c

Active domain: {a, b, c}

From Query to Circuit

Assumption:

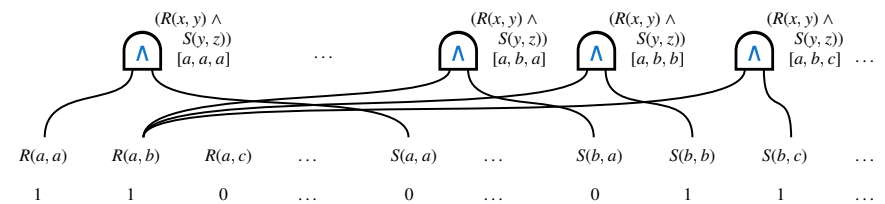
- query and database schema is fixed
- database instance (and thus active domain) are variable

Construct circuit uniformly based on size of active domain

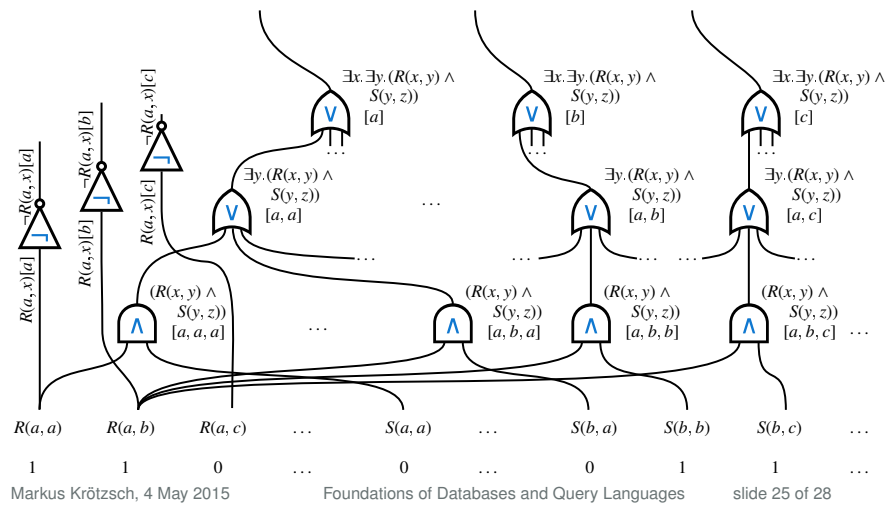
Sketch of construction:

- one input node for each possible database tuple (over given schema and active domain)
 - ~> true or false depending on whether tuple is present or not
- Recursively, for each subformula, introduce a gate for each possible tuple (instantiation) of this formula
 - ~> true or false depending on whether the subformula holds for this tuple or not
- Logical operators correspond to gate types: basic operators obvious, \forall as generalised conjunction, \exists as generalised disjunction
- subformula with n free variables ~> $|\mathbf{adom}|^n$ gates
 - ~> especially: $|\mathbf{adom}|^0 = 1$ output gate for Boolean query

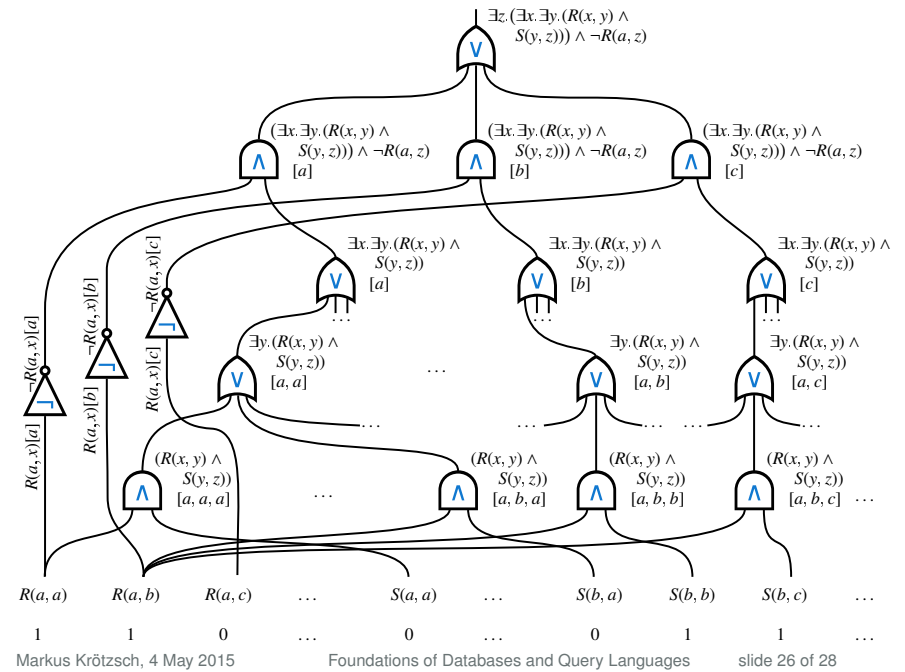
Example: $\exists z. (\exists x. \exists y. R(x, y) \wedge S(y, z)) \wedge \neg R(a, z)$



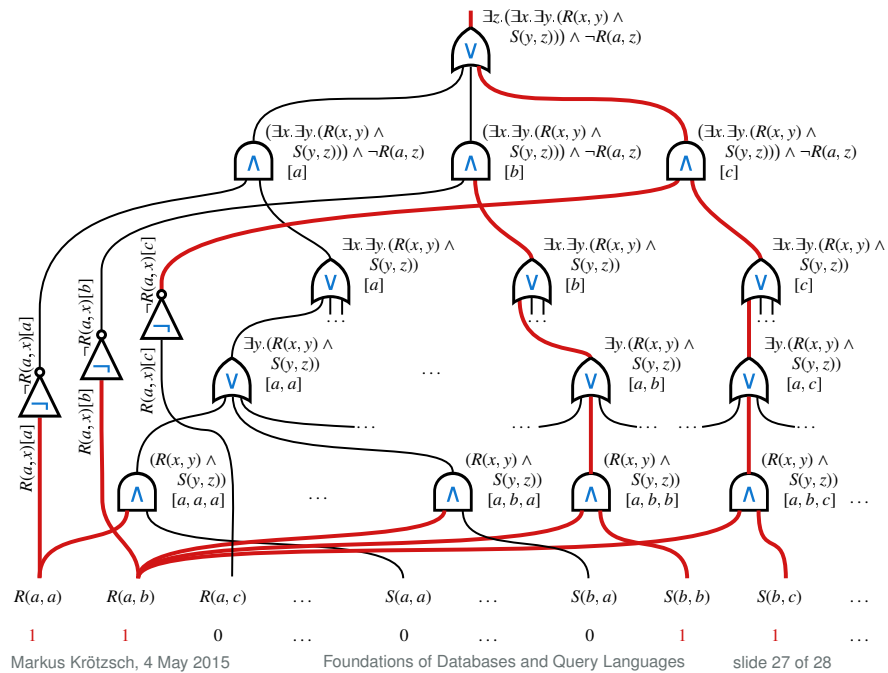
Example: $\exists z. (\exists x. \exists y. R(x, y) \wedge S(y, z)) \wedge \neg R(a, z)$



Example: $\exists z. (\exists x. \exists y. R(x, y) \wedge S(y, z)) \wedge \neg R(a, z)$



Example: $\exists z. (\exists x. \exists y. R(x, y) \wedge S(y, z)) \wedge \neg R(a, z)$



Summary and Outlook

The evaluation of FO queries is

- PSPACE-complete for combined complexity
- PSPACE-complete for query complexity
- AC⁰-complete for data complexity

Circuit complexities help to identify highly parallelisable problems in P

Open questions:

- Which other computing problems are interesting? (next lecture)
- Are there query languages with lower complexities?
- How can we study the expressiveness of query languages?