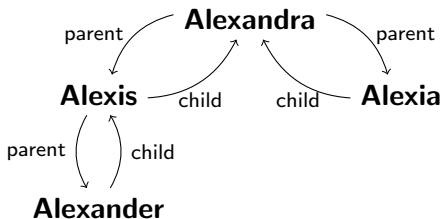


Conjunctive Context-Free Path Queries

Jelle Hellings

Hasselt University and transnational University of Limburg

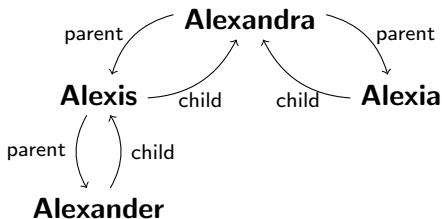
Motivation: graph data and path queries



Example (using CRPQ: conjunctive regular path queries)

- ▶ All pairs of relatives: $Q(n, m) \leftarrow (\text{parent} + \text{child})^*(n, m)$
- ▶ All pairs (ancestor, descendants): $Q(n, m) \leftarrow \text{parent}^+(n, m)$
- ▶ All pairs of n -th generation descendants of an ancestor?

Motivation: limitations of CRPQ

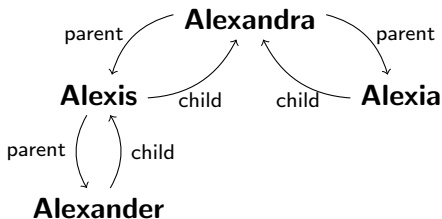


All pairs of n -th generation descendants of an ancestor?

Solution 1: path variables and regular relations [Barceló et al.]

$$Q(n, m) \leftarrow \exists \pi \exists \rho \exists z \ n \pi z \wedge m \rho z \wedge \text{child}^+(\pi) \wedge \text{child}^+(\rho) \wedge |\pi| = |\rho|$$

Motivation: limitations of CRPQ



All pairs of n -th generation descendants of an ancestor?

Solution 1: path variables and regular relations [Barceló et al.]

$$Q(n, m) \leftarrow \exists \pi \exists \rho \exists z \ n \pi z \wedge m \rho z \wedge \text{child}^+(\pi) \wedge \text{child}^+(\rho) \wedge |\pi| = |\rho|$$

Solution 2: context-free grammars

Production rules: $N \rightarrow \text{parent child}$, $N \rightarrow \text{parent } N \text{ child}$

$$Q(n, m) \leftarrow N(n, m)$$

Overview

1. Motivation
2. Introducing CCFPQ
3. Extending CCFPQ
 - 3.1 Adding path variables
 - 3.2 Grammars over tuples of paths
 - 3.3 Negation
4. Conclusions and future work

Context-free grammars

Definition (context-free grammar: $G = (N, A, P)$)

- ▶ N : a finite set of non-terminals
- ▶ A : a finite set of symbols (the alphabet)
- ▶ P : a finite set of production rules

Definition (language of context-free grammar G)

$$\mathcal{L}(G_N) = \{s_1 \dots s_j \text{ a finite string over } A \mid N \rightarrow_G s_1 \dots s_j\}$$

Definition (normal form)

All productions are of the form $A \rightarrow BC$, $A \rightarrow \sigma$, or $A \rightarrow \lambda$

CCFPQ - syntax

Query $Q(\vec{v}) \leftarrow \exists \vec{\mu} \bigwedge N_i(n_i, m_i)$

- ▶ \vec{v} : a tuple of node variables
- ▶ $\vec{\mu}$: a tuple of distinct node variables that do not occur in \vec{v}
- ▶ N_i : a non-terminal from a grammar
- ▶ n_i and m_i : are node variables taken from \vec{v} or $\vec{\mu}$

CCFPQ - semantics

Query $Q(\vec{\nu}) \leftarrow \exists \vec{\mu} \bigwedge N_i(n_i, m_i)$

- ▶ The *trace* of path π is denoted by $\mathcal{T}(\pi)$
- ▶ Define $\mathcal{R}_{N_i} \subseteq \mathcal{V} \times \mathcal{V}$ as:

$$\mathcal{R}_{N_i} = \{(n, m) \mid \exists n \pi m (\mathcal{T}(\pi) \in \mathcal{L}(G_{N_i}))\}$$

- ▶ Interpret atoms $N_i(n_i, m_i)$ as relational atoms $\mathcal{R}_{N_i}(n_i, m_i)$
- ▶ Threat result as normal conjunctive query

CCFPQ - query evaluation

Consider graph $\mathcal{D} = (\mathcal{V}, \mathcal{E})$ and query $Q(\vec{v}) \leftarrow \exists \vec{\mu} \bigwedge N_i(n_i, m_i)$

- ▶ \mathcal{R}_{N_i} is a binary relation
- ▶ Compute \mathcal{R}_{N_i} : reduces CCFPQ to conjunctive queries
- ▶ We have $\mathcal{R}_{N_i} \subseteq \mathcal{V} \times \mathcal{V}$

Context-free recognizer for graphs - algorithm

Input: Edge-labeled directed graph $\mathcal{D} = (\mathcal{V}, \mathcal{E})$

Context-free grammar $G = (N, A, P)$

Output: $\{(N, n, m) \mid (n, m) \in \mathcal{R}_N\}$

- 1: $r := \{(N, n, n) \mid (n \in \mathcal{V}) \wedge (N \rightarrow \lambda \in P)\} \cup$
 $\{(N, n, m) \mid ((n, l, m) \in \mathcal{E}) \wedge (N \rightarrow l \in P)\}$
- 2: $new := r$
- 3: **while** $new \neq \emptyset$ **do**
- 4: pick and remove a (N, n, m) from new
- 5: **for all** $(M, n', n) \in r$ **do**
- 6: **for all** $(N' \rightarrow MN \in P) \wedge ((N', n', m) \notin r)$ **do**
- 7: $new := new \cup \{(N', n', m)\}$
- 8: $r := r \cup \{(N', n', m)\}$
- 9: **for all** $(M, m, m') \in r$ **do**
- 10: ...
- 11: **return** r

Context-free recognizer for graphs - analysis

Input: Edge-labeled directed graph $\mathcal{D} = (\mathcal{V}, \mathcal{E})$

Context-free grammar $G = (N, A, P)$

Output: $\{(N, n, m) \mid (n, m) \in \mathcal{R}_N\}$

Proposition (correctness)

We have $(N, n, m) \in r$ if and only if $(n, m) \in \mathcal{R}_N$.

Theorem (worst-case complexity)

Running time is $\mathcal{O}(|N||\mathcal{E}| + (|N||\mathcal{V}|)^3)$.

Corollary

For CCFPQ queries, the query evaluation problem has

- 1. Polynomial time data complexity,*
- 2. NP-complete combined complexity.*

Extending CCFPQ: add path variables

Theorem

Consider $Q() \leftarrow \exists nm \exists n\pi m N(\pi) \wedge M(\pi)$, we have:

1. For a fixed graph \mathcal{D} , query evaluation is undecidable
2. For a fixed query Q , query evaluation is undecidable

Extending CCFPQ: CRPQ and ECRPQ

All pairs of n -th generation descendants of an ancestor?

Solution 1: path variables and regular relations [Barceló et al.]

$Q(n, m) \leftarrow \exists \pi \exists \rho \exists z \ n \pi z \wedge m \rho z \wedge \text{child}^+(\pi) \wedge \text{child}^+(\rho) \wedge |\pi| = |\rho|$

How to express $|\pi| = |\rho|$?

$|\pi| = |\rho|$ as a regular expression over *tuples of edge labels*:

$$\left[\sum_{\sigma_1, \sigma_2 \in A} \langle \sigma_1, \sigma_2 \rangle \right]^* (\pi, \rho)$$

j -ary context-free grammar

Definition

A j -ary context-free grammar over alphabet A is a context-free grammar of the form $G = (N, (A \cup \varsigma)^j, P)$.

Example (subsequences)

'assert', 'liver' *subsequences* of 'Hasselt University'

$$\begin{aligned} \text{Sub} &\rightarrow \lambda, & \text{Sub} &\rightarrow \langle \sigma, \sigma \rangle \text{Sub} \quad (\forall \sigma \in A), \\ & & \text{Sub} &\rightarrow \langle \sigma, \varsigma \rangle \text{Sub} \quad (\forall \sigma \in A). \end{aligned}$$

ECCFPQ - syntax

Query $Q(\vec{v}) \leftarrow \exists \vec{\mu} \bigwedge N_i(\vec{p}_i)$

- ▶ \vec{v} : a tuple of node variables
- ▶ $\vec{\mu}$: a tuple of distinct node variables that do not occur in \vec{v}
- ▶ N_i : a non-terminal from a j -ary grammar
- ▶ \vec{p}_i : a j -ary tuple of pairs of node variables taken from \vec{v} or $\vec{\mu}$

Example (subsequences)

$$\begin{aligned} \text{Sub} &\rightarrow \lambda, & \text{Sub} &\rightarrow \langle \sigma, \sigma \rangle \text{Sub} \ (\forall \sigma \in A), \\ & & \text{Sub} &\rightarrow \langle \sigma, \varsigma \rangle \text{Sub} \ (\forall \sigma \in A). \end{aligned}$$
$$Q_{\text{Sub}}(n_1, m_1, n_2, m_2) \leftarrow \text{Sub}(\langle n_1, m_1 \rangle, \langle n_2, m_2 \rangle)$$

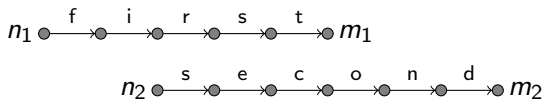
ECCFPQ - skip-traces

Definition

The *skip-traces* of path π , with $\mathcal{T}(\pi) = l_1 \dots l_j$, are defined as

$$\mathcal{T}_\zeta(\pi) = \{s_1 \uparrow l_1 \uparrow s_2 \uparrow \dots \uparrow s_i \uparrow l_i \uparrow s_{i+1} \mid s_1, \dots, s_{i+1} \text{ are finite sequences of } \zeta \text{ symbols}\}.$$

Example



$$\mathcal{T}_\zeta(n_1 \pi_1 m_1) = \{\text{first}, \zeta \text{first}, f \zeta \text{first}, \zeta f \zeta \text{first}, \dots\}$$

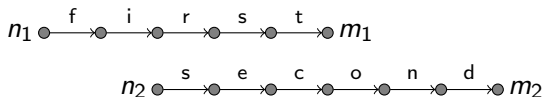
ECCFPQ - j -ary skip-traces

Definition

The *skip-traces* of a tuple of paths $\vec{\pi} = \langle \pi_1, \dots, \pi_j \rangle$ are defined as

$$\mathcal{T}_{\zeta j}(\vec{\pi}) = \{ \langle t_1^1, \dots, t_1^j \rangle \dots \langle t_l^1, \dots, t_l^j \rangle \mid (t^k = t_1^k \dots t_l^k) \wedge (t^k \in \mathcal{T}_{\zeta}(\pi_k)) \}.$$

Example



$$\mathcal{T}_{\zeta 2}(n_1 \pi_1 m_1, n_2 \pi_2 m_2) = \{ \langle \zeta \rangle \langle f \rangle \langle i \rangle \langle r \rangle \langle s \rangle \langle t \rangle, \langle \zeta \rangle \langle \zeta \rangle \langle f \rangle \langle i \rangle \langle r \rangle \langle s \rangle \langle t \rangle \langle \zeta \rangle, \dots \}$$

ECCFPQ - semantics

Query $Q(\vec{v}) \leftarrow \exists \vec{\mu} \bigwedge N_i(n_i, m_i)$

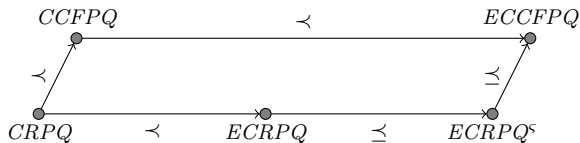
- ▶ Define $\mathcal{R}_{N_i} \subseteq (\mathcal{V} \times \mathcal{V})^j$ as:

$$\mathcal{R}_{N_i} = \{(\langle n_1, m_1 \rangle, \dots, \langle n_j, m_j \rangle) \mid \exists n_1 \pi_1 m_1 \dots \exists n_j \pi_j m_j \\ \exists t ((t \in \mathcal{T}_{s_j}(\langle \pi_1, \dots, \pi_j \rangle)) \wedge (t \in \mathcal{L}(G_{N_i})))\}$$

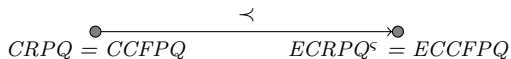
- ▶ Interpret atoms $N_i(\vec{p}_i)$ as relational atoms $\mathcal{R}_{N_i}(\vec{p}_i)$
- ▶ Threat result as normal conjunctive query

Expressive power of CCFPQ and variants

Theorem (on labeled graphs)



Theorem (on unlabeled graphs)



Extending CCFPQ: adding negation

Example (all pairs of *distinct* n -th generation descendants?)

Production rules: $N \rightarrow \text{parent child}$, $N \rightarrow \text{parent } N \text{ child}$

$\text{Equal} \rightarrow \lambda$

$$Q(n, m) \leftarrow N(n, m) \wedge \neg \text{Equal}(n, m)$$

Theorem (on unlabeled graphs)



Conclusions and Future Work

Conclusions

- ▶ Introduced context-free grammars for specifying paths
- ▶ Efficient query evaluation techniques

Future Work

- ▶ Paths as answers instead of nodes
- ▶ Determine all relations between expressive power of variants