

From relation algebra to semi-join algebra: an approach for graph query optimization

*Jelle Hellings*¹

Catherine L. Pilachowski² Dirk Van Gucht²

Marc Gyssens¹ Yuqing Wu³

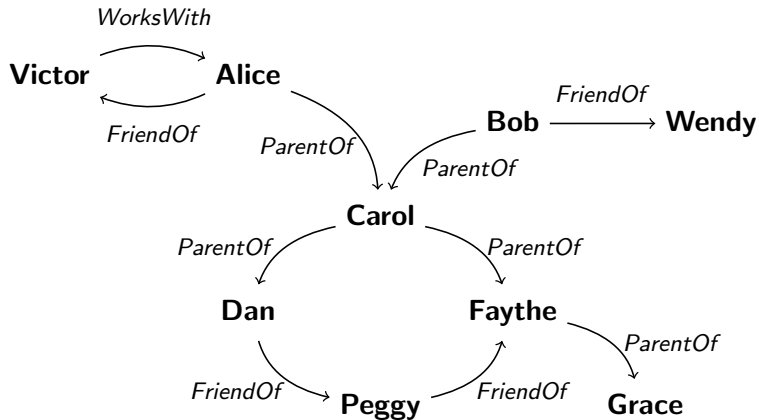
¹ Hasselt University

² Indiana University

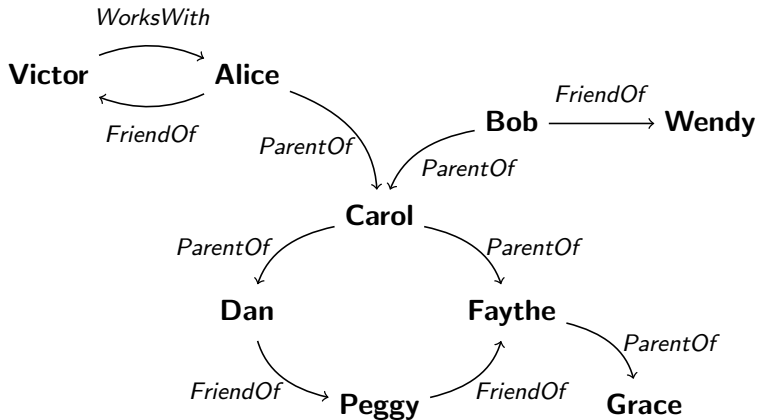
³ Pomona College



Graph queries: data model

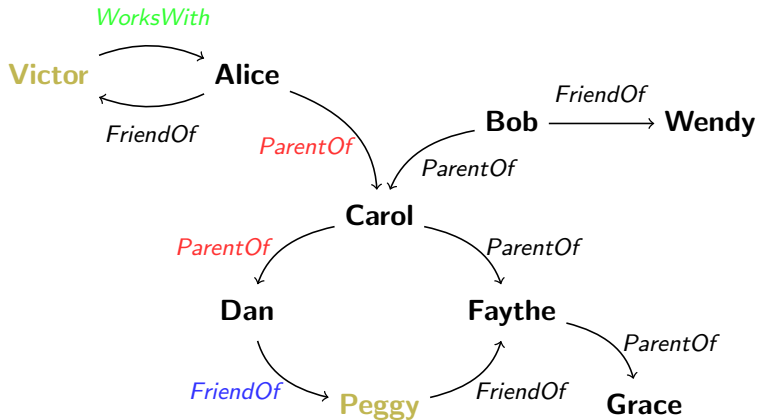


Graph queries: basic path queries



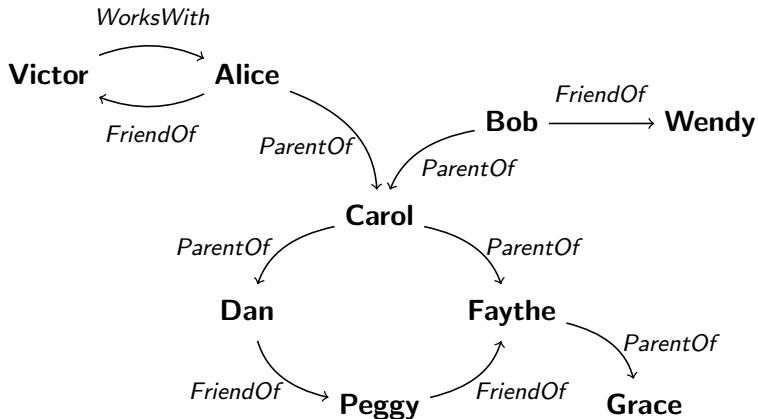
$(WorksWith \cup FriendOf) \circ [ParentOf]^+ \circ FriendOf$

Graph queries: basic path queries



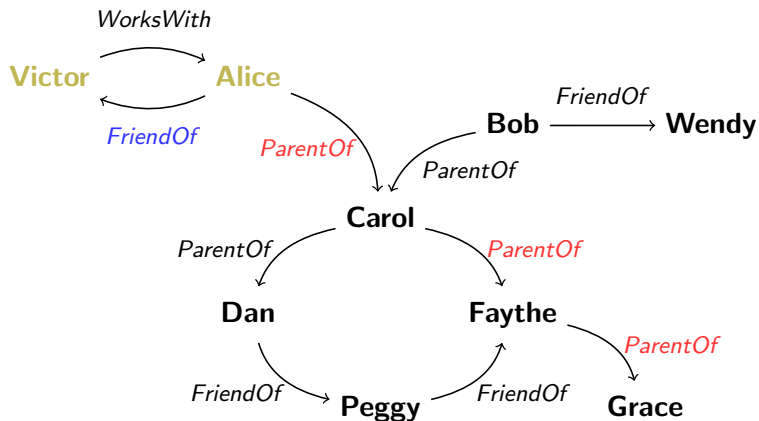
$$(\textit{WorksWith} \cup \textit{FriendOf}) \circ [\textit{ParentOf}]^+ \circ \textit{FriendOf}$$

Graph queries: node-tests and branching



$\pi_1[\text{ParentOf} \circ \text{ParentOf} \circ \text{ParentOf}] \circ \text{FriendOf}$

Graph queries: node-tests and branching



$\pi_1[\textit{ParentOf} \circ \textit{ParentOf} \circ \textit{ParentOf}] \circ \textit{FriendOf}$

Graph querying: relation algebra

id	\cup	\circ	$+$	\wedge	π	$\bar{\pi}$	\cap	$-$	di
RPQs									
2RPQs									
Nested RPQs									
Navigational XPath, Graph XPath									
FO[3] + transitive closure									

Relation algebra and query evaluation

id	\cup	\circ +	\wedge π	$\bar{\pi}$	\cap -	di
----	--------	-----------	----------------	-------------	----------	----

Cheap ($\cup, \wedge, \pi, \cap, -$).

Cost linearly upper bounded by operands

In between (id, $\bar{\pi}$).

Cost linearly upper bounded by #nodes

Expensive ($\circ, +, di$).

Worst-case quadratically lower bounded by #nodes

Naive query evaluation: an inefficient example

Return pairs of (great-grandparent, friend)

$$\pi_1[\text{ParentOf} \circ \text{ParentOf} \circ \text{ParentOf}] \circ \text{FriendOf}$$

1. Compute (grandparent, grandchild):

$$X = \text{ParentOf} \circ \text{ParentOf}$$

2. Compute (great-grandparent, great-grandchild):

$$Y = \text{ParentOf} \circ X$$

3. Throw away the great-grandchildren:

$$Z = \pi_1[Y]$$

4. Compute (great-grandparent, friend):

$$\text{Result} = Z \circ \text{FriendOf}$$

Optimize query evaluation: add specialized operators?

Return pairs of (great-grandparent, friend)

$$\pi_1[ParentOf \circ ParentOf \circ ParentOf] \circ FriendOf$$

1. Compute (grandparent, ???):

$$X = ParentOf \times ParentOf$$

2. Compute (great-grandparent, ???):

$$Y = ParentOf \times (X)$$

3. Throw away ???:

$$Z = \pi_1[Y]$$

4. Compute (great-grandparent, friend):

$$Result = Z \times FriendOf$$

$$\pi_1[ParentOf \times (ParentOf \times ParentOf)] \times FriendOf$$

Simple idea: automatic query rewriting

- ▶ Rewrite composition into semi-joins
- ▶ Rewrite transitive closure into fixpoints

In such a way that the rewritten query is equivalent

When are expressions equivalent?

Definition

Queries q_1 and q_2 are

path-equivalent if, for every graph \mathcal{G} , $\llbracket q_1 \rrbracket_{\mathcal{G}} = \llbracket q_2 \rrbracket_{\mathcal{G}}$
(denoted by $q_1 \equiv_{\text{path}} q_2$)

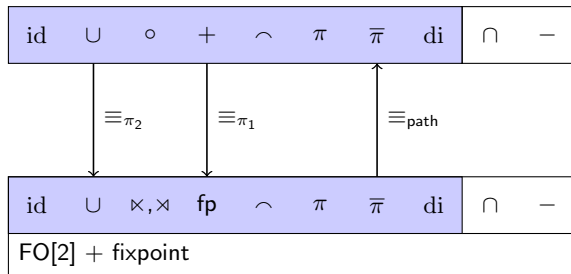
left-projection-equivalent if, for every graph \mathcal{G} , $\llbracket q_1 \rrbracket_{\mathcal{G}|_1} = \llbracket q_2 \rrbracket_{\mathcal{G}|_1}$
(denoted by $q_1 \equiv_{\pi_1} q_2$)

right-projection-equivalent if, for every graph \mathcal{G} , $\llbracket q_1 \rrbracket_{\mathcal{G}|_2} = \llbracket q_2 \rrbracket_{\mathcal{G}|_2}$
(denoted by $q_1 \equiv_{\pi_2} q_2$)

Example

- ▶ $R \cap S \equiv_{\text{path}} R - (R - S)$
- ▶ $R \circ S \equiv_{\pi_1} R \times S$
- ▶ $\pi_1[R \circ S] \equiv_{\text{path}} \pi_1[R \times S]$

The main result



- ▶ Collapse also holds for fragments (that include π)
- ▶ Example: Nested RPQs are projection-equivalent to expressions using only id , \cup , \bowtie , \bowtie , fp , \frown , and π

Intersection \cap and difference $-$

Issues when combining composition with \cap or $-$

$$(\text{FriendOf} \circ \text{FriendOf}) \cap \text{FriendOf}$$

- ▶ *Restricting*: use \cap and $-$ only on composition-free expressions
 - ▶ Exact syntactic fragment of $\text{FO}[3] + \text{TC}$ that is projection-equivalent to $\text{FO}[2] + \text{fixpoint}$.
- ▶ *Data models*: usage of \cap and $-$ is sometimes redundant
 - ▶ Sibling-ordered trees: $\text{FO}^{\text{tree}} \preceq_{\pi} \text{FO}[2] + \text{fixpoints}$.
 - ▶ Downward queries on trees [DBPL 2015]
 - ▶ ...
- ▶ *Partial rewriting*: keep compositions when necessary

The rewrite functions - partial rewriting

$$\begin{aligned}\tau(e) &\equiv_{\text{path}} e & \tau_{\pi_1}(e) &\equiv_{\pi_1} e & \tau_{\pi_2}(e) &\equiv_{\pi_2} e \\ \tau_{o_1}(e; \varepsilon) &\equiv_{\pi_1} e \times \varepsilon & \tau_{o_2}(e; \varepsilon) &\equiv_{\pi_2} \varepsilon \times e\end{aligned}$$

Example

$$\pi_1[(((WorksOn \circ WorksOn^{\wedge}) \cap FriendOf) \circ EditorOf) \circ StudentOf]$$

The rewrite functions - partial rewriting

$$\begin{aligned}\tau(e) &\equiv_{\text{path}} e & \tau_{\pi_1}(e) &\equiv_{\pi_1} e & \tau_{\pi_2}(e) &\equiv_{\pi_2} e \\ \tau_{o_1}(e; \varepsilon) &\equiv_{\pi_1} e \times \varepsilon & \tau_{o_2}(e; \varepsilon) &\equiv_{\pi_2} \varepsilon \times e\end{aligned}$$

Example

$\pi_1[(((\text{WorksOn} \circ \text{WorksOn}^\wedge) \cap \text{FriendOf}) \circ \text{EditorOf}) \circ \text{StudentOf}]$

$\tau(e)$

The rewrite functions - partial rewriting

$$\begin{aligned}\tau(e) &\equiv_{\text{path}} e & \tau_{\pi_1}(e) &\equiv_{\pi_1} e & \tau_{\pi_2}(e) &\equiv_{\pi_2} e \\ \tau_{o_1}(e; \varepsilon) &\equiv_{\pi_1} e \times \varepsilon & \tau_{o_2}(e; \varepsilon) &\equiv_{\pi_2} \varepsilon \times e\end{aligned}$$

Example

$\pi_1[(((WorksOn \circ WorksOn^\wedge) \cap FriendOf) \circ EditorOf) \circ StudentOf]$

$$\tau(e) = \tau_{\pi_2}(\pi_1[(((W \circ W^\wedge) \cap F) \circ E)]) \times \tau(S)$$

The rewrite functions - partial rewriting

$$\begin{aligned}\tau(e) &\equiv_{\text{path}} e & \tau_{\pi_1}(e) &\equiv_{\pi_1} e & \tau_{\pi_2}(e) &\equiv_{\pi_2} e \\ \tau_{o_1}(e; \varepsilon) &\equiv_{\pi_1} e \times \varepsilon & \tau_{o_2}(e; \varepsilon) &\equiv_{\pi_2} \varepsilon \times e\end{aligned}$$

Example

$\pi_1[(((WorksOn \circ WorksOn^\wedge) \cap FriendOf) \circ EditorOf) \circ StudentOf]$

$$\begin{aligned}\tau(e) &= \tau_{\pi_2}(\pi_1[(((W \circ W^\wedge) \cap F) \circ E)]) \times \tau(S) \\ &= \pi_1[\tau_{\pi_1}(((W \circ W^\wedge) \cap F) \circ E)] \times S\end{aligned}$$

The rewrite functions - partial rewriting

$$\begin{aligned}\tau(e) &\equiv_{\text{path}} e & \tau_{\pi_1}(e) &\equiv_{\pi_1} e & \tau_{\pi_2}(e) &\equiv_{\pi_2} e \\ \tau_{o_1}(e; \varepsilon) &\equiv_{\pi_1} e \times \varepsilon & \tau_{o_2}(e; \varepsilon) &\equiv_{\pi_2} \varepsilon \times e\end{aligned}$$

Example

$\pi_1[(((WorksOn \circ WorksOn^\wedge) \cap FriendOf) \circ EditorOf) \circ StudentOf]$

$$\begin{aligned}\tau(e) &= \tau_{\pi_2}(\pi_1[(((W \circ W^\wedge) \cap F) \circ E)]) \times \tau(S) \\ &= \pi_1[\tau_{\pi_1}(((W \circ W^\wedge) \cap F) \circ E)] \times S \\ &= \pi_1[\tau_{o_1}(((W \circ W^\wedge) \cap F; \tau_{\pi_1}(E)))] \times S\end{aligned}$$

The rewrite functions - partial rewriting

$$\begin{aligned}\tau(e) &\equiv_{\text{path}} e & \tau_{\pi_1}(e) &\equiv_{\pi_1} e & \tau_{\pi_2}(e) &\equiv_{\pi_2} e \\ \tau_{o_1}(e; \varepsilon) &\equiv_{\pi_1} e \times \varepsilon & \tau_{o_2}(e; \varepsilon) &\equiv_{\pi_2} \varepsilon \times e\end{aligned}$$

Example

$\pi_1[(((WorksOn \circ WorksOn^\wedge) \cap FriendOf) \circ EditorOf) \circ StudentOf]$

$$\begin{aligned}\tau(e) &= \tau_{\pi_2}(\pi_1[(((W \circ W^\wedge) \cap F) \circ E)]) \times \tau(S) \\ &= \pi_1[\tau_{\pi_1}(((W \circ W^\wedge) \cap F) \circ E)] \times S \\ &= \pi_1[\tau_{o_1}(((W \circ W^\wedge) \cap F; \tau_{\pi_1}(E)))] \times S \\ &= \pi_1[(\tau(W \circ W^\wedge) \cap \tau(F)) \times E] \times S\end{aligned}$$

The rewrite functions - partial rewriting

$$\begin{aligned}\tau(e) &\equiv_{\text{path}} e & \tau_{\pi_1}(e) &\equiv_{\pi_1} e & \tau_{\pi_2}(e) &\equiv_{\pi_2} e \\ \tau_{o_1}(e; \varepsilon) &\equiv_{\pi_1} e \times \varepsilon & \tau_{o_2}(e; \varepsilon) &\equiv_{\pi_2} \varepsilon \times e\end{aligned}$$

Example

$\pi_1[(((WorksOn \circ WorksOn^\wedge) \cap FriendOf) \circ EditorOf) \circ StudentOf]$

$$\begin{aligned}\tau(e) &= \tau_{\pi_2}(\pi_1[(((W \circ W^\wedge) \cap F) \circ E)]) \times \tau(S) \\ &= \pi_1[\tau_{\pi_1}(((W \circ W^\wedge) \cap F) \circ E)] \times S \\ &= \pi_1[\tau_{o_1}(((W \circ W^\wedge) \cap F; \tau_{\pi_1}(E)))] \times S \\ &= \pi_1[(\tau(W \circ W^\wedge) \cap \tau(F)) \times E] \times S \\ &= \pi_1[((\tau(W) \circ \tau(W^\wedge)) \cap F) \times E] \times S\end{aligned}$$

The rewrite functions - partial rewriting

$$\begin{aligned}\tau(e) &\equiv_{\text{path}} e & \tau_{\pi_1}(e) &\equiv_{\pi_1} e & \tau_{\pi_2}(e) &\equiv_{\pi_2} e \\ \tau_{o_1}(e; \varepsilon) &\equiv_{\pi_1} e \times \varepsilon & \tau_{o_2}(e; \varepsilon) &\equiv_{\pi_2} \varepsilon \times e\end{aligned}$$

Example

$\pi_1[(((WorksOn \circ WorksOn^\wedge) \cap FriendOf) \circ EditorOf) \circ StudentOf]$

$$\begin{aligned}\tau(e) &= \tau_{\pi_2}(\pi_1[(((W \circ W^\wedge) \cap F) \circ E)]) \times \tau(S) \\ &= \pi_1[\tau_{\pi_1}(((W \circ W^\wedge) \cap F) \circ E)] \times S \\ &= \pi_1[\tau_{o_1}(((W \circ W^\wedge) \cap F; \tau_{\pi_1}(E)))] \times S \\ &= \pi_1[(\tau(W \circ W^\wedge) \cap \tau(F)) \times E] \times S \\ &= \pi_1[((\tau(W) \circ \tau(W^\wedge)) \cap F) \times E] \times S \\ &= \pi_1[(((W \circ W^\wedge) \cap F) \times E)] \times S.\end{aligned}$$

Query optimization

- ▶ Cost of each operator ✓
- ▶ Input size of each operator

- ▶ Number of necessary evaluation steps

Query optimization

- ▶ Cost of each operator ✓
- ▶ Input size of each operator

Example

Let $R = \{(1, i) \mid 0 \leq i \leq m\}$. Consider

$$R \circ R^{\wedge} \equiv_{\pi_1} R \times R^{\wedge}.$$

- ▶ Number of necessary evaluation steps

Query optimization

- ▶ Cost of each operator ✓
- ▶ Input size of each operator ✓

Example

Let $R = \{(1, i) \mid 0 \leq i \leq m\}$. Consider

$$R \circ R^{\wedge} \equiv_{\pi_1} R \times R^{\wedge}.$$

Solution: use single-column evaluation algorithms

- ▶ Number of necessary evaluation steps

Query optimization

- ▶ Cost of each operator ✓
- ▶ Input size of each operator ✓

Example

Let $R = \{(1, i) \mid 0 \leq i \leq m\}$. Consider

$$R \circ R^{\wedge} \equiv_{\pi_1} R \times R^{\wedge}.$$

Solution: use single-column evaluation algorithms

- ▶ Number of necessary evaluation steps ✗

Expressions and evaluation steps

Expression size we denote the *expression size* of e by $\|e\|$.

Evaluation size we denote the *evaluation size* of e by $\text{eval-steps}(e)$.

Example

$$e_1 = ((R \circ R) \circ (R \circ R)) \circ ((R \circ R) \circ (R \circ R))$$

$$e_2 = R \times (R \times (R \times (R \times (R \times (R \times (R \times R))))))$$

- ▶ $e_1 \equiv_{\pi_1} e_2$
- ▶ We have $\|e_1\| = 7$ and $\text{eval-steps}(e_1) = 3$:
 1. $X = R \circ R$
 2. $Y = X \circ X$
 3. $\text{Result} = Y \circ Y$
- ▶ We have $\|e_2\| = 7$ and $\text{eval-steps}(e_2) = 7$.

Evaluation size and unions

Example

$$e_1 = (A \cup B) \circ (C \cup D) \circ (E \cup F)$$

$$e_2 = A \times (C \times E) \cup A \times (C \times F) \cup \dots$$

- ▶ $e_1 \equiv_{\pi_1} e_2$
- ▶ We have $\|e_1\| = \text{eval-steps}(e_1) = 5$.
- ▶ We have $\|e_2\| = \text{eval-steps}(e_2) = 23$.

Evaluation size and unions

Example

$$e_1 = (A \cup B) \circ (C \cup D) \circ (E \cup F)$$

$$e_2 = A \times (C \times E) \cup A \times (C \times F) \cup \dots$$

- ▶ $e_1 \equiv_{\pi_1} e_2$
- ▶ We have $\|e_1\| = \text{eval-steps}(e_1) = 5$.
- ▶ We have $\|e_2\| = \text{eval-steps}(e_2) = 23$.

$$e_3 = (A \times X) \cup (B \times X),$$

$$X = (C \times Y) \cup (D \times Y), Y = (E \cup F)$$

- ▶ $e_1 \equiv_{\pi_1} e_3$
- ▶ We have $\|e_2'\| = 13$ and $\text{eval-steps}(e_2') = 7$.

Evaluation size and unions

Example

$$e_1 = (A \cup B) \circ (C \cup D) \circ (E \cup F)$$

$$e_2 = A \times (C \times E) \cup A \times (C \times F) \cup \dots$$

- ▶ $e_1 \equiv_{\pi_1} e_2$
- ▶ We have $\|e_1\| = \text{eval-steps}(e_1) = 5$.
- ▶ We have $\|e_2\| = \text{eval-steps}(e_2) = 23$.

$$e_3 = (A \times X) \cup (B \times X),$$

$$X = (C \times Y) \cup (D \times Y), Y = (E \cup F)$$

- ▶ $e_1 \equiv_{\pi_1} e_3'$
- ▶ We have $\|e_2'\| = 13$ and $\text{eval-steps}(e_2') = 7$.
- ▶ $\tau_{o_i}(e; \varepsilon)$ does this! ✓

The main results (revised)

Theorem

Let e be an expression. We have $\tau(e) \equiv_{\text{path}} e$, $\tau_{\pi_i}(e) \equiv_{\pi_i} e$, and

1. $\text{eval-steps}(\tau(e)) \leq u + \|e\|$;
2. $\text{eval-steps}(\tau_{\pi_i}(e)) \leq u + \|e\|$;
3. $\|\tau(e)\| = \Theta(\|e\| \cdot 2^u)$ in the worst case;
4. $\|\tau_{\pi_i}(e)\| = \Theta(\|e\| \cdot 2^u)$ in the worst case,

with u the number of rewrite steps involving $\tau_{o_i}(e_1 \cup e_2; \varepsilon)$.

Conclusion and future work

1. Real-life systems
2. Relational databases
3. Intersection and difference elimination
4. Extending FO[3] (e.g. counting)

The FO[2] fixpoint we use

- ▶ Notation $\text{fp}_{i,\mathfrak{N}}$ [iterative case union base case]
- ▶ i specifies output column
- ▶ \mathfrak{N} is a variable representing the growing output (node-test)
- ▶ Subset of traditional inflationary fixpoints

Example

The query $\pi_1[[\text{ParentOf}]^+ \circ \text{OwnsPet}]$ returns ancestors of pet-owners. We rewrite this into

$$\pi_1[\text{fp}_{1,\mathfrak{N}}[\text{ParentOf} \times \mathfrak{N} \text{ union } \text{ParentOf} \times \text{OwnsPet}]]$$