

On Tarski's Relation Algebra

QUERYING TREES AND CHAINS

and

THE SEMI-JOIN ALGEBRA

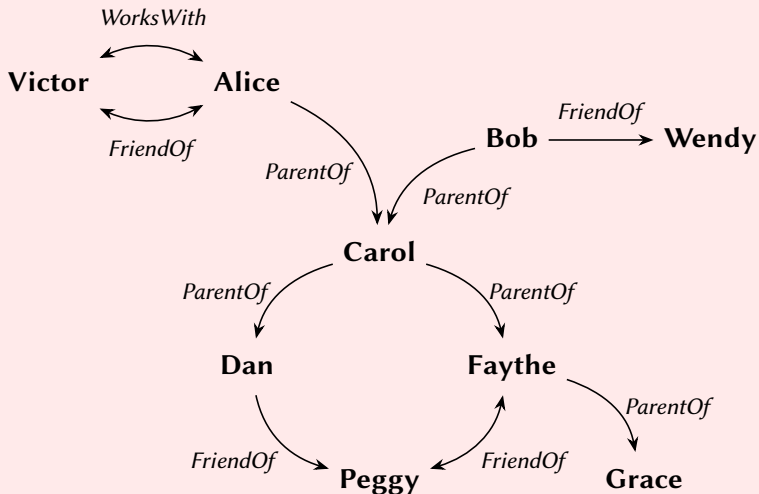
Jelle Hellings

Part I

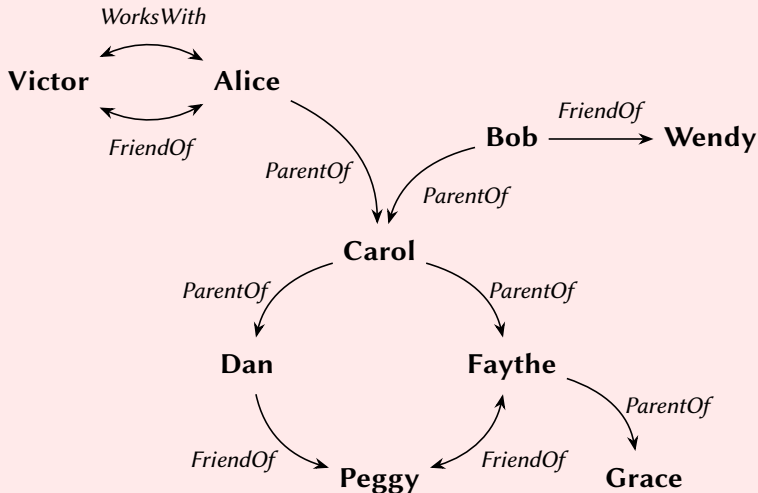
On Tarski's Relation Algebra

INTRODUCTION

Relation algebra: graphs and queries

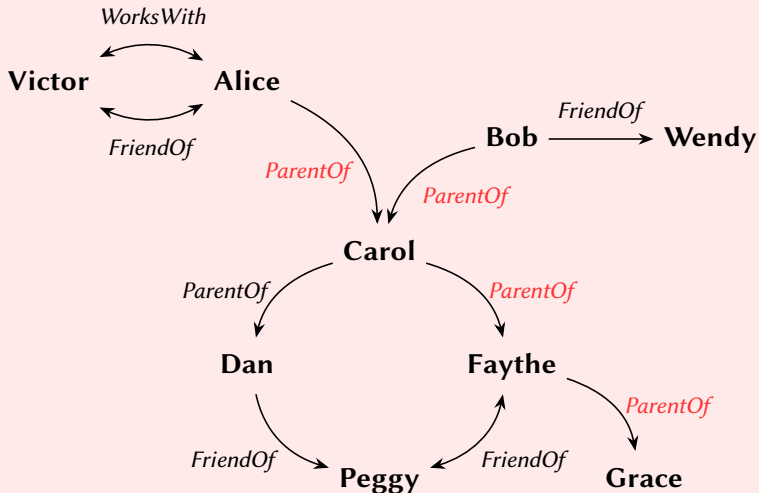


Relation algebra: graphs and queries



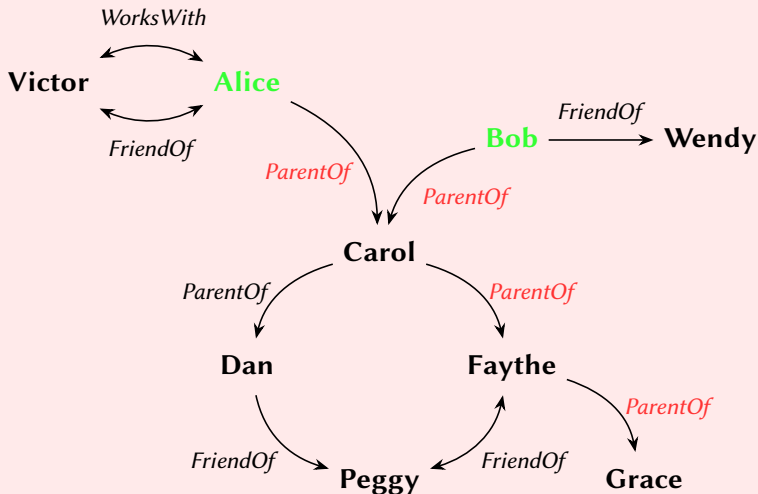
$$\text{FriendsOfGgp} = \pi_1[\text{ParentOf} \circ \text{ParentOf} \circ \text{ParentOf}] \circ \text{FriendOf}$$

Relation algebra: graphs and queries



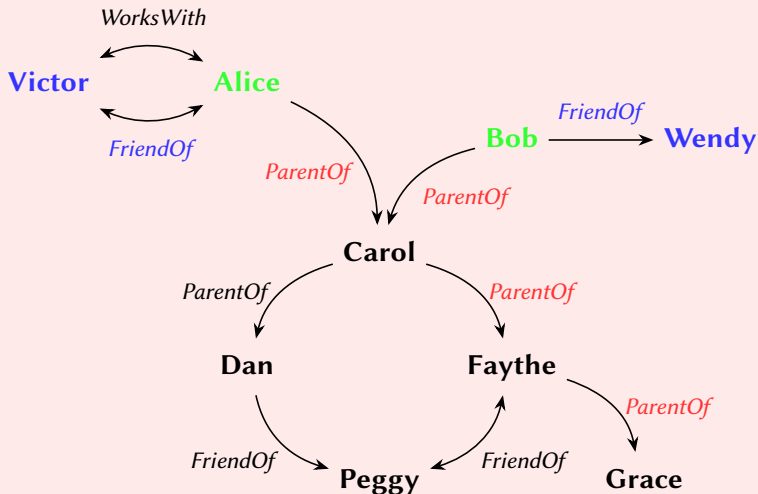
$$\text{FriendsOfGgp} = \pi_1[\text{ParentOf} \circ \text{ParentOf} \circ \text{ParentOf}] \circ \text{FriendOf}$$

Relation algebra: graphs and queries



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

Relation algebra: graphs and queries



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

Relation algebra: other operators

$$\text{ChildOf} = \text{ParentOf}^{\sim}$$

$$\text{AcquaintanceOf} = \text{FriendOf} \cup \text{WorksWith}$$

$$\text{WorkFriendOf} = \text{FriendOf} \cap \text{WorksWith}$$

$$\text{NonWorkFriend} = \text{FriendOf} - \text{WorksWith}$$

$$\text{Parent} = \pi_1[\text{ParentOf}]$$

$$\text{NonParent} = \overline{\pi}_1[\text{ParentOf}]$$

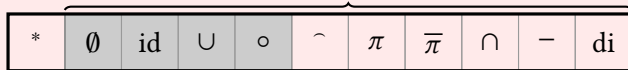
$$\text{GrandParentOf} = \text{ParentOf} \circ \text{ParentOf}$$

$$\text{AncestorOf} = [\text{ParentOf}]^+$$

Also: constants \emptyset , id, di.

Graph query languages

Tarski's Relation Algebra, FO[3]



RPQs

2RPQs

Nested RPQs

Navigational XPath, GXPath

Questions about the relation algebra

Questions

- ▶ Are all these operators necessary?
- ▶ What does each operator add?
- ▶ When can we replace complex operators by simpler ones?

Questions about the relation algebra

Questions

- ▶ Are all these operators necessary?
- ▶ What does each operator add?
- ▶ When can we replace complex operators by simpler ones?

What is the expressive power of fragments of the relation algebra?

Related work and our results

Previous work by Fletcher et al.

Relative expressive power of relation algebra fragments on graphs.

This work

Improve understanding of the relation algebra

- ▶ when querying trees or chains;
- ▶ in comparison with the semi-join algebra.

Background: when are queries equivalent?

Two types of queries and equivalences

Path-queries. The exact query result is important:

$$\text{FriendOf} \cap \text{WorksWith} \equiv_{\text{path}} \text{FriendOf} - (\text{FriendOf} - \text{WorksWith}).$$

Boolean-queries. The existence of a query result is important:

$$\text{FriendOf} \hat{\cap} \text{ParentOf} \equiv_{\text{bool}} \pi_1[\text{FriendOf}] \cap \pi_1[\text{ParentOf}].$$

Definition

Language \mathcal{L}_1 is **z-subsumed** by \mathcal{L}_2 if every query in \mathcal{L}_1 is z-equivalent to a query in \mathcal{L}_2 (denoted by $\mathcal{L}_1 \preceq_z \mathcal{L}_2$).

Background: query fragments

Let $\mathcal{F} \subseteq \{\text{di}, \wedge, \pi, \bar{\pi}, \cap, -, *\}$.

- ▶ We write $\mathcal{N}(\mathcal{F})$: only allows $\emptyset, \text{id}, \ell, \circ, \cup$, and all operators in \mathcal{F} .
- ▶ We write $\mathcal{N}(\underline{\mathcal{F}})$ to represent all operators expressible in $\mathcal{N}(\mathcal{F})$ using the following basic rewrite rules:

$$\pi_1[e] = \bar{\pi}_j[\bar{\pi}_1[e]] = (e \circ [e]^{-1}) \cap \text{id} = (e \circ (\text{id} \cup \text{di})) \cap \text{id};$$

$$\pi_2[e] = \bar{\pi}_j[\bar{\pi}_2[e]] = ([e]^{-1} \circ e) \cap \text{id} = ((\text{id} \cup \text{di}) \circ e) \cap \text{id};$$

$$\bar{\pi}_j[e] = \text{id} - \pi_j[e];$$

$$e_1 \cap e_2 = e_1 - (e_1 - e_2).$$

Examples

- ▶ $\mathcal{N}(\underline{\bar{\pi}}) = \mathcal{N}(\pi, \bar{\pi})$.
- ▶ $\mathcal{N}(\underline{\wedge}, -) = \mathcal{N}(\wedge, \pi, \bar{\pi}, \cap, -)$.

Part II

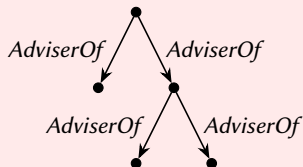
On Tarski's Relation Algebra

QUERYING TREES AND CHAINS

Why studying trees

Definition

A **tree** is an acyclic graph in which exactly one node, the *root*, has no incoming edges, and all other nodes have exactly one incoming edge.

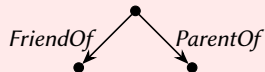
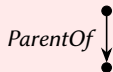


- ▶ Hierarchical relations: taxonomies, corporate structures.
- ▶ XML and JSON data.
- ▶ Nested relational data.

Initial classification

		downward		non-downward			
non- [*] -free		$\mathcal{N}(\pi, \bar{\pi}, \cap, -, *)$	$\mathcal{N}(\cap, \pi, \bar{\pi}, \cap, *)$	$\mathcal{N}(\text{di}, \cap, \pi, \bar{\pi}, *)$	\mathcal{N}	non-monotone	
		$\mathcal{N}(\pi, \cap, *)$	$\mathcal{N}(\cap, \pi, \cap, *)$	$\mathcal{N}(\text{di}, \cap, \pi, *)$	$\mathcal{N}(\text{di}, \cap, \pi, \cap, *)$	monotone	
non-local, [*] -free				$\mathcal{N}(\text{di}, \cap, \pi, \bar{\pi})$	$\mathcal{N}(\text{di}, \cap, \pi, \bar{\pi}, \cap, -)$	non-monotone	
				$\mathcal{N}(\text{di}, \cap, \pi)$	$\mathcal{N}(\text{di}, \cap, \pi, \cap)$	monotone	
local		$\mathcal{N}(\pi, \bar{\pi}, \cap, -)$	$\mathcal{N}(\cap, \pi, \bar{\pi}, \cap)$		$\mathcal{N}(\cap, \pi, \bar{\pi}, \cap, -)$	non-monotone	
		$\mathcal{N}(\cap, -)$ $\mathcal{N}(\pi, \cap)$	$\mathcal{N}(\cap, \pi, \cap)$			monotone	
		1-subtree reducible		2-subtree reducible		3-subtree reducible	

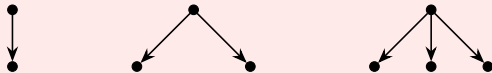
Recognizing branches: labeled trees



Only select the three on the right

- ▶ Not possible in \mathcal{N}^* .
- ▶ Using projection: $\pi_1[\textit{FriendOf}] \circ \pi_1[\textit{ParentOf}]$.
- ▶ Using converse: $\textit{FriendOf}^\smile \circ \textit{ParentOf}$.

Recognizing branches: unlabeled trees



Definition

A k -subtree-reduction step on tree \mathcal{T} consists of removing a subtree rooted at a node n with parent m , given that parent m has at least k other children isomorphic to n .

Theorem

1. $\mathcal{N}(\text{di}, \cap)$ and $\mathcal{N}(\frown, -)$ are closed under 3-subtree-reductions.
2. $\mathcal{N}(\text{di}, \frown, \pi, \bar{\pi}, *)$ is closed under 2-subtree-reductions.
3. $\mathcal{N}(\frown, \pi, \bar{\pi}, \cap, *)$ is closed under 1-subtree-reductions.

Monotonicity

- ▶ Without negation you cannot ‘forbid’ structures.
- ▶ $\bar{\pi}$ always provides negation: $NonParent = \bar{\pi}_1[ParentOf]$.
- ▶ Without negation, only a few Boolean queries are possible.

Theorem

1. *On unlabeled chains, we have $\mathcal{N}(di, \cap, \pi, \cap, *) \leq_{\text{bool}} \mathcal{N}()$.*
2. *On unlabeled trees, we have $\mathcal{N}(\cap, \pi, \cap, *) \leq_{\text{bool}} \mathcal{N}()$.*

Theorem

*Already on unlabeled chains, we have $\mathcal{N}(\bar{\pi}) \not\leq_{\text{bool}} \mathcal{N}(di, \cap, \pi, \cap, *)$.*

The Kleene-star

- ▶ Is not expressible in first-order logic.
- ▶ Path queries: always adds expressive power.
- ▶ Labeled structures: always adds expressive power.

Theorem

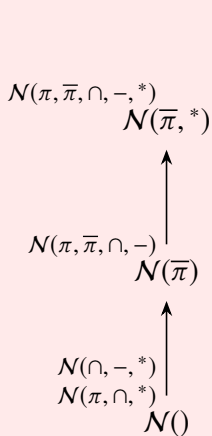
1. *On unlabeled chains, we have $\mathcal{N}(\text{di}, \frown, \pi, \cap, *) \leq_{\text{bool}} \mathcal{N}()$.*
2. *On unlabeled trees, we have $\mathcal{N}(\frown, \pi, \cap, *) \leq_{\text{bool}} \mathcal{N}()$.*

Theorem

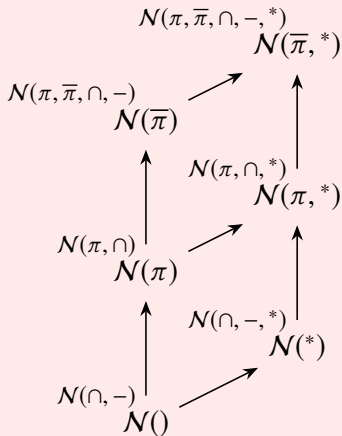
Let $\mathcal{F} \subseteq \{\text{di}, \frown, \pi\}$. On unlabeled trees, we have $\mathcal{N}(\mathcal{F} \cup \{\}) \leq_{\text{bool}} \mathcal{N}(\mathcal{F})$.*

Downward queries: Boolean semantics

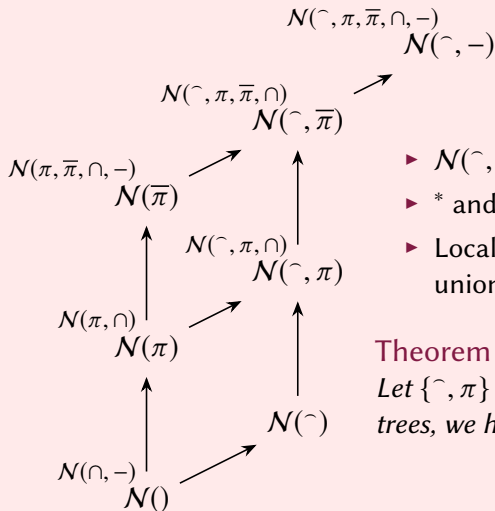
Unlabeled Trees



Labeled Trees



Local queries: path semantics



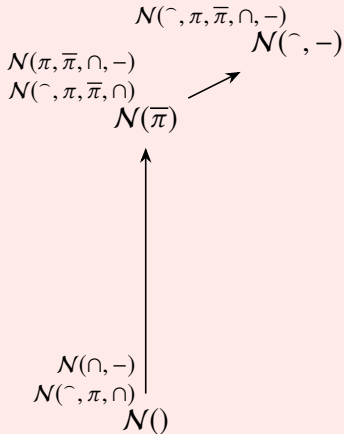
- ▶ $\mathcal{N}(\cap, \pi, \bar{\pi}, \cap, -)$ is local.
- ▶ * and di are not local.
- ▶ Local queries can be expressed by unions of **condition tree queries**.

Theorem

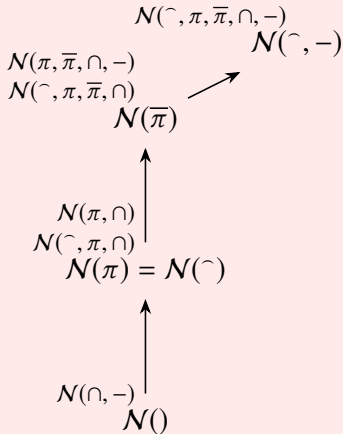
Let $\{\cap, \pi\} \subseteq \mathcal{F} \subseteq \{\cap, \pi, \bar{\pi}, \cap\}$. On labeled trees, we have $\mathcal{N}(\mathcal{F}) \leq_{\text{path}} \mathcal{N}(\mathcal{F} - \{\cap\})$.

Local queries: Boolean semantics

Unlabeled Trees



Labeled Trees



Conclusion

- ▶ We fully established relationships between:
 - ▶ downward fragments;
 - ▶ local fragments.
- ▶ On trees: solved most cases not involving $*$.
- ▶ On chains: solved most cases not involving di or $*$.

Main open problem

Let $\{di, \bar{\pi}, \cap\} \subseteq \mathcal{F} \subseteq \{di, \wedge, \pi, \bar{\pi}, \cap, *\}$ and let $z \in \{\text{bool}, \text{path}\}$. With respect to either labeled trees, unlabeled trees, labeled chains, or unlabeled chains, do we have a collapse $\mathcal{N}(\mathcal{F} \cup \{-\}) \leq_z \mathcal{N}(\mathcal{F})$ or not?

Part III

On Tarski's Relation Algebra

THE SEMI-JOIN ALGEBRA

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):

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

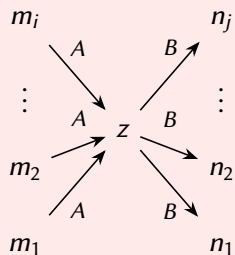
3. Throw away the great-grandchildren:

$$Z = \pi_1[Y]$$

4. Compute (great-grandparent, friend):

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

Introducing the semi-join algebra

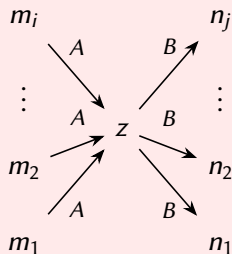


Composition: querying for paths

Consider $A \circ B$:

- ▶ yields (begin, end)-nodes connected by AB .
- ▶ yields $i \cdot j$ node pairs in total.

Introducing the semi-join algebra



Composition: querying for paths

Consider $A \circ B$:

- ▶ yields (begin, end)-nodes connected by AB .
- ▶ yields $i \cdot j$ node pairs in total.

Checking for existence of paths

Consider the semi-join $A \bowtie B$:

- ▶ yields the edges in A that connect to B .
- ▶ yields i node pairs in total.

Optimize query evaluation: using semi-joins

Return pairs of (great-grandparent, friend)

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

1. Compute (grandparent, ???):

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

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

$$Y = \text{ParentOf} \bowtie (X)$$

3. Throw away ???:

$$Z = \pi_1[Y]$$

4. Compute (great-grandparent, friend):

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

$$\pi_1[\text{ParentOf} \bowtie (\text{ParentOf} \bowtie \text{ParentOf})] \bowtie \text{FriendOf}.$$

Projection-equivalence

Requiring the rewrite to guarantee *path-equivalence* is too strong!

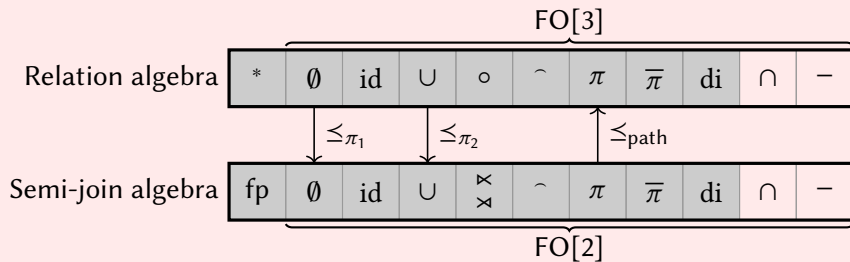
Left-projection-equivalence. We only need the first projection.

$$\begin{aligned} \text{ParentOf} \circ \text{ParentOf} &\equiv_{\pi_1} \text{ParentOf} \times \text{ParentOf} \\ \pi_1[\text{ParentOf} \circ \text{ParentOf}] &\equiv_{\text{path}} \pi_1[\text{ParentOf} \times \text{ParentOf}] \end{aligned}$$

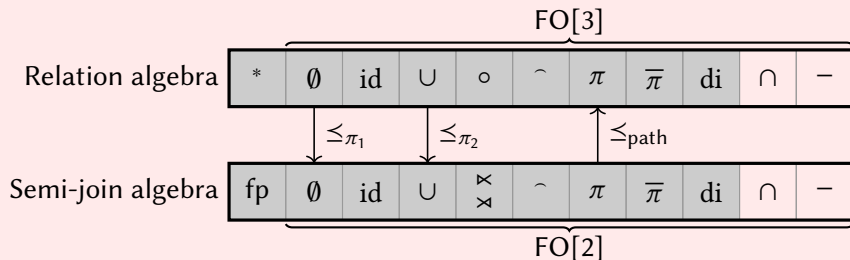
Right-projection-equivalence. We only need the second projection.

$$\begin{aligned} \text{ParentOf} \circ \text{ParentOf} &\equiv_{\pi_2} \text{ParentOf} \times \text{ParentOf} \\ \pi_2[\text{ParentOf} \circ \text{ParentOf}] &\equiv_{\text{path}} \pi_2[\text{ParentOf} \times \text{ParentOf}] \end{aligned}$$

The main result

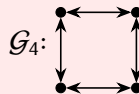
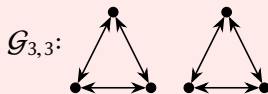


The main result



Intersection and difference

- ▶ Consider $(\mathcal{E} \circ \mathcal{E}) \cap \mathcal{E}$:



- ▶ We can allow \cap and $-$ on basic expressions (edges, \emptyset , id, di).

Partial rewriting to the semi-join algebra

Rewrite functions $\tau(e) \equiv_{\text{path}} e$, $\tau_{\pi_1}(e) \equiv_{\pi_1} e$ and $\tau_{\pi_2}(e) \equiv_{\pi_2} e$

Helper functions $\tau_{o_1}(e; \varepsilon) \equiv_{\pi_1} e \bowtie \varepsilon$ and $\tau_{o_2}(e; \varepsilon) \equiv_{\pi_2} \varepsilon \bowtie e$

Example

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

Rewriting e using $\tau(e)$:

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

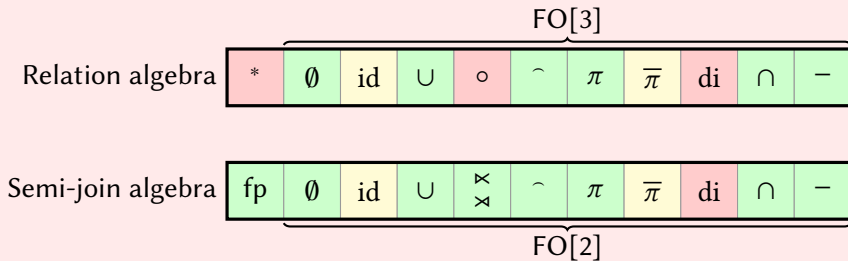
Query rewriting and optimization

- ▶ Cost of each operator χ .
- ▶ Input size of each operator χ .
- ▶ Number of evaluation steps χ .

Query rewriting and optimization

- ▶ Cost of each operator ✓.
- ▶ Input size of each operator ✗.
- ▶ Number of evaluation steps ✗.

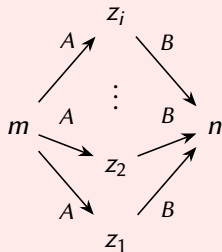
Cost of each operator



Query rewriting and optimization

- ▶ Cost of each operator ✓.
- ▶ Input size of each operator ✗.
- ▶ Number of evaluation steps ✗.

Input size of each operator



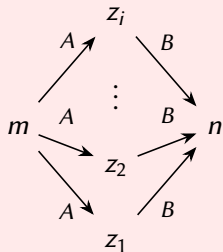
$$A \circ B \equiv_{\pi_1} A \times B$$

- ▶ $A \circ B$ yields $\{(m, n)\}$.
- ▶ $A \times B$ yields $\{(m, z_1), (m, z_2), \dots, (m, z_i)\}$.

Query rewriting and optimization

- ▶ Cost of each operator ✓.
- ▶ Input size of each operator ✓.
- ▶ Number of evaluation steps ✗.

Input size of each operator



$$A \circ B \equiv_{\pi_1} A \times B$$

- ▶ $A \circ B$ yields $\{(m, n)\}$.
- ▶ $A \times B$ yields $\{(m, z_1), (m, z_2), \dots, (m, z_i)\}$.
- ▶ Fix: add projection-steps into algorithms.

Query rewriting and optimization

- ▶ Cost of each operator ✓.
- ▶ Input size of each operator ✓.
- ▶ Number of evaluation steps \sim .

Number of evaluation steps

- ▶ Number of evaluation steps can increase.
- ▶ Complexity of each individual step can significantly decrease.

Example

Consider $e = \pi_1[(\ell \circ \ell) \circ (\ell \circ \ell)]$ and $\tau(e) = \pi_1[\ell \bowtie (\ell \bowtie (\ell \bowtie \ell))]$.

We can evaluate e in three steps:

1. compute $X = \ell \circ \ell$;
2. compute $Y = X \circ X$;
3. compute $\pi_1[Y]$.

Other expensive constructs: id, di, and $\overline{\pi}$

Idea: replace common use cases

- ▶ Replace id by equality-selections

$$(FriendOf \circ WorksWith) \cap id \equiv_{\text{path}} \sigma_{=}(FriendOf \circ WorksWith).$$

- ▶ Replace di by inequality-selections

$$(FriendOf \circ FriendOf) \cap di \equiv_{\text{path}} \sigma_{\neq}(FriendOf \circ FriendOf).$$

- ▶ Replace $\overline{\pi}$ by anti-semi-joins

$$FriendOf \circ \overline{\pi}_1[ParentOf] \equiv_{\text{path}} FriendOf \overline{\bowtie} ParentOf.$$

Further optimization: filter steps in practical queries

Find friend suggestions for *Alice*:

1. Compute all friends-of-friends (excluding friends and oneself):

$$\textit{FriendSuggestions} = (\textit{FriendOf} \circ \textit{FriendOf}) - (\textit{FriendOf} \cup \textit{id}).$$

2. Filter the first column on *Alice*.
3. Keep the second column.

Further optimization: filter steps in practical queries

Find friend suggestions for **Alice**:

1. Compute all friends-of-friends (excluding friends and oneself):

$$\text{FriendSuggestions} = (\text{FriendOf} \circ \text{FriendOf}) - (\text{FriendOf} \cup \text{id}).$$

2. Filter the first column on Alice.
3. Keep the second column.

Incorporate filter-step in the relation algebra

$$\pi_2[(\langle \text{Alice} \rangle \times \text{FriendOf}) \times (\text{FriendOf} \bar{\times} \pi_2[(\langle \text{Alice} \rangle \times \text{FriendOf}) \cup \langle \text{Alice} \rangle])]$$

Conclusion

Queries in the relation algebra can be partially rewritten into semi-join algebra queries that are easier to evaluate.

Future work

- ▶ Implementation in real systems.
- ▶ Interactions with other optimization techniques.
- ▶ Elimination of intersection and difference.

Part IV

On Tarski's Relation Algebra

CONCLUSION

Conclusion

Questions

- ▶ Are all these operators necessary?
- ▶ What does each operator add?
- ▶ When can we replace complex operators by simpler ones?

This work

Improve understanding of the relation algebra

- ▶ when querying trees or chains;
- ▶ in comparison with the semi-join algebra.