# Type Inference for Datalog with Complex Type Hierarchies

*Max Schäfer*    Oege de Moor

semmle **/**

# Example Program

```
reports_to(x,y) ← ∃ g.(in_group(x,g) ∧ manages(y,g))
                ∨ ∃ g,d.(manages(x,g) ∧ part_of(g,d)
                                     ∧ leads(y,d))
                ∨ senior_mgr(x) ∧ x≐y.
```

```
bonus(x,y)        ← ∃ s,f.(¬parttime(x) ∧ salary(x,s)
                          ∧ factor(x,f) ∧ y≐f*s)
                  ∨ parttime(x) ∧ y≐50.0.

factor(x,f)       ← senior_mgr(x) ∧ f≐0.5
                  ∨ ¬senior_mgr(x) ∧ f≐0.2.

query(x,y)        ← bonus(x,y) ∧ manager(x).
```

# Example Program

```
reports_to(x,y) ← ∃ g.(in_group(x,g) ∧ manages(y,g))
                ∨ ∃ g,d.(manages(x,g) ∧ part_of(g,d)
                                       ∧ leads(y,d))
                ∨ senior_mgr(x) ∧ x≐y.
```

```
bonus(x,y)      ← ∃ s,f.(¬parttime(x) ∧ salary(x,s)
                         ∧ factor(x,f) ∧ y≐f*s)
                ∨ parttime(x) ∧ y≐50.0.

factor(x,f)     ← senior_mgr(x) ∧ f≐0.5
                ∨ ¬senior_mgr(x) ∧ f≐0.2.

query(x,y)      ← bonus(x,y) ∧ manager(x).
```

# Example Program

```
reports_to(x,y) ← ∃ g.(in_group(x,g) ∧ manages(y,g))
                ∨ ∃ g,d.(manages(x,g) ∧ part_of(g,d)
                                     ∧ leads(y,d))
                ∨ senior_mgr(x) ∧ x≐y.
```

```
bonus(x,y)       ← ∃ s,f.(¬parttime(x) ∧ salary(x,s)
                           ∧ factor(x,f) ∧ float(y))
                 ∨ parttime(x) ∧ float(y).

factor(x,f)      ← senior_mgr(x) ∧ float(f)
                 ∨ ¬senior_mgr(x) ∧ float(f).

query(x,y)       ← bonus(x,y) ∧ manager(x).
```

- Schema $\mathscr{S}$ assigns entity types to columns of extensionals

$$(\forall x, g.\texttt{in\_group}(x, g) \rightarrow \texttt{developer}(x) \land \texttt{group}(g))$$

$$(\forall x, g.\texttt{manages}(x, g) \rightarrow \texttt{junior\_mgr}(x) \land \texttt{group}(g))$$

$$(\forall x, s.\texttt{salary}(x, s) \rightarrow \texttt{employee}(x) \land \texttt{float}(s))$$

$$\ldots$$

**semmle/**

- Type Hierarchy $\mathcal{H}$ relates entity types

$$
\begin{aligned}
\forall x. \quad & (\texttt{parttime}(x) \rightarrow \texttt{developer}(x)) \\
\wedge \quad & (\texttt{developer}(x) \rightarrow \texttt{employee}(x)) \\
\wedge \quad & (\texttt{manager}(x) \rightarrow \texttt{employee}(x)) \\
\wedge \quad & (\texttt{junior\_mgr}(x) \vee \texttt{senior\_mgr}(x) \leftrightarrow \texttt{manager}(x)) \\
\wedge \quad & \neg(\texttt{developer}(x) \wedge \texttt{manager}(x)) \\
\wedge \quad & \neg(\texttt{junior\_mgr}(x) \wedge \texttt{senior\_mgr}(x))
\end{aligned}
$$

## Type-based Optimisation

Example query:

```
bonus(x, y) ∧ manager(x)
```

Unfolding definitions:

```
(∃ s,f.(¬parttime(x) ∧ salary(x, s)
        ∧ factor(x,f) ∧ y≐f*s)
 ∨ parttime(x) ∧ y≐50.0)
∧ manager(x)
```

# Type-based Optimisation

semmle**/**

Example query:

```
bonus(x, y) ∧ manager(x)
```

Unfolding definitions:

```
(∃ s,f.(¬parttime(x) ∧ salary(x, s)
        ∧ factor(x,f) ∧ y≐f*s)
 ∨ parttime(x) ∧ y=50.0)
∧ manager(x)
```
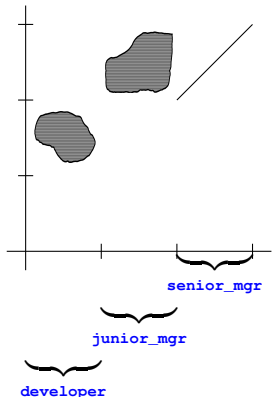
Type specialisation:

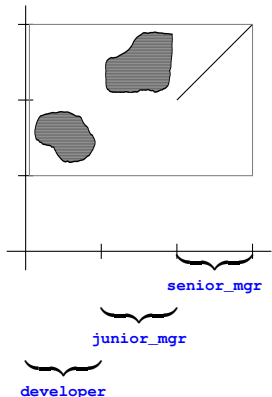$$\mathbf{parttime}(x) \wedge \mathbf{manager}(x) \models_{\mathscr{H}} \bot$$

- want to know whether

$$\mathscr{H}, \mathscr{S}, \varphi \models \bot$$

- undecidable for arbitrary programs
- classic approach: find upper envelope

```
reports_to(x,y) ←
   ∃ g.(in_group(x,g) ∧ manages(y,g))
∨ ∃ g,d.(manages(x,g) ∧ part_of(g,d)
                       ∧ leads(y,d))
∨ senior_mgr(x) ∧ x=y.
```

senior_mgr

junior_mgr

developer

# Cartesian Types

semmle**/**



- one type per variable

$$\lceil \texttt{reports\_to}(x, y) \rceil$$
$$= \textbf{employee}(x) \land \textbf{manager}(y)$$

- erroneous query:

**parttime**(y) $\land$ reports_to(x, y)

- disjunctive form

$$\lceil \textbf{reports\_to}(x, y) \rceil$$

$$= \quad \textbf{developer}(x) \wedge \textbf{junior\_mgr}(y)$$
$$\vee \quad \textbf{manager}(x) \wedge \textbf{senior\_mgr}(y)$$

- erroneous query:

```
parttime(x) ∧
reports_to(x, y) ∧ senior_mgr(y)
```

- disjunctive form

$$\lceil \mathtt{reports\_to}(x, y) \rceil$$

$=\quad \mathtt{developer}(x) \land \mathtt{junior\_mgr}(y)$
$\lor\quad \mathtt{junior\_mgr}(x) \land \mathtt{senior\_mgr}(y)$
$\lor\quad \mathtt{senior\_mgr}(x) \land \mathtt{senior\_mgr}(y)$
$$\land\; x = y$$

senior_mgr

junior_mgr

developer

```
reports_to(x,y)    ← ∃ g.(in_group(x,g) ∧ manages(y,g))
                   ∨ ∃ g,d.(manages(x,g) ∧ part_of(g,d)
                                          ∧ leads(y,d))
                   ∨ senior_mgr(x) ∧ x=̇y.
```
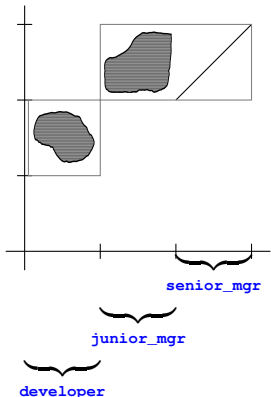
$$\begin{aligned}
\text{reports\_to(x,y)} \leftarrow\ & \exists\ g.(\textbf{in\_group}(x,g) \wedge \textbf{manages}(y,g)) \\
& \vee\ \exists\ g,d.(\textbf{manages}(x,g) \wedge \textbf{part\_of}(g,d) \\
& \qquad\qquad\qquad\qquad\quad \wedge \textbf{leads}(y,d)) \\
& \vee\ \textbf{senior\_mgr}(x) \wedge x\dot{=}y.
\end{aligned}$$

$$\Downarrow$$

$$\begin{aligned}
\lceil \text{reports\_to(x,y)} \rceil \leftarrow\ & \exists\ g.(\textbf{developer}(x) \wedge \textbf{group}(g) \\
& \qquad\qquad\qquad\qquad\quad \wedge \textbf{junior\_mgr}(y)) \\
& \vee\ \exists\ g,d.(\textbf{junior\_mgr}(x) \wedge \textbf{group}(g) \\
& \qquad\qquad \wedge \textbf{department}(d) \wedge \textbf{senior\_mgr}(y)) \\
& \vee\ \textbf{senior\_mgr}(x) \wedge x\dot{=}y.
\end{aligned}$$

```
reports_to(x,y)    ← ∃ g.(in_group(x,g) ∧ manages(y,g))
                   ∨ ∃ g,d.(manages(x,g) ∧ part_of(g,d)
                                         ∧ leads(y,d))
                   ∨ senior_mgr(x) ∧ x≐y.
```

$$\Downarrow$$

```
⌈reports_to(x,y)⌉ ← developer(x) ∧ junior_mgr(y)
                     ∧ ∃ g.(group(g))
                   ∨ junior_mgr(x) ∧ senior_mgr(y)
                     ∧ ∃ g.(group(g)) ∧ ∃ d.(department(d))
                   ∨ senior_mgr(x) ∧ x≐y.
```

# Types as Upper Envelopes

- type inference: $\lceil \cdot \rceil : P \to T$ where
  - $T \subset P$
  - containment, emptiness decidable on $T$
  - $\mathscr{H}, \mathscr{S}, \varphi \models \lceil \varphi \rceil$

# semmle /

- types are existential programs with monadic extensionals
  - polyadic intensionals
  - recursion, existentials, equality
  - no negation
- $\lceil e \rceil := \mathscr{S}(e)$ for extensionals
- $\lceil \neg \varphi \rceil := \top$

semmle **/**

### Soundness and Optimality

- for any $\varphi \in P$:
$$\mathscr{S}, \varphi \models \lceil \varphi \rceil$$

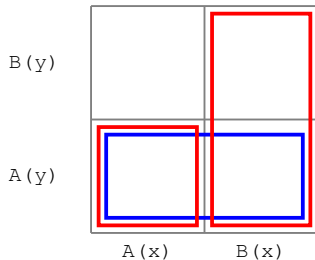- for negation-free $\varphi \in P$ and $\vartheta \in T$:

$$\text{if} \quad \mathscr{S}, \mathscr{H}, \varphi \models \vartheta$$
$$\text{then} \quad \mathscr{S}, \mathscr{H}, \lceil \varphi \rceil \models \vartheta$$

- intensional relations in $T$ can always be eliminated:
- can be written as disjunction of conjunctions
- every conjunct of the form $t(x)$, $x \doteq y$ or $\exists z.t(z)$ for propositional $t$
  $\Rightarrow$ compact representation using BDDs
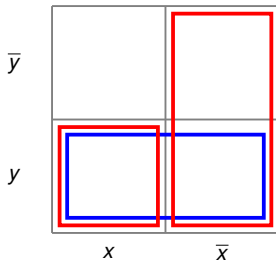- containment and emptiness decidable

# Conclusions

semmle /

- natural, expressive language of types for Datalog
- support for complex type hierarchies
- simple, sound, optimal type inference
- compact representation of types

- $S <: T$ if every disjunct of $S$ implies some disjunct of $T$
- obviously sound, $S <: T$ implies $S \models T$
- but incomplete!

**semmle/**



$$\underbrace{\overset{\sigma_1}{\overbrace{\text{C}(x) \wedge \text{A}(y)}}}_{S} \models \underbrace{\overset{\tau_1}{\overbrace{\text{A}(x) \wedge \text{A}(y)}} \vee \overset{\tau_2}{\overbrace{\text{B}(x) \wedge \text{C}(y)}}}_{T}$$

yet $\sigma_1 \not\models \tau_1$ and $\sigma_1 \not\models \tau_2$

# Prime Implicants

**semmle/**



$$y \models x \wedge y \vee \neg x$$

- prime implicants of $x \wedge y \vee \neg x$: $\{y, \neg x\}$
- obtained by consensus formation: $(x \wedge y) \oplus_x \neg x = y$

**semmle/**

- new interpretation of consensus:

$$(x \wedge y) \oplus_x (\neg x \wedge (y \vee \neg y)) = (x \vee \neg x) \wedge y = y$$

"disjunction on $x$, conjunction on $y$"
- can be generalised to types, need to do consensus on <u>sets</u> of variables
- $\mathtt{sat}(T)$: all prime implicants of $T$

> ### Complete Containment Check
> If $S \models_{\mathcal{H}} T$, then $S <: \mathtt{sat}(T)$.