# Specialized Static Analysis Framework: Type Inference

Woosuk Lee

CSE 6049 Program Analysis

Hanyang University, Korea

# Goal of This Lecture

- Learn practical alternatives to the aforementioned general, abstract interpretation framework

- For simple languages and properties, there are frameworks that are simple yet powerful enough

- But with several limitations

# Static Analysis by Proof Construction

- Static analysis = proof construction in a finite proof system

- Finite proof system = a finite set of inference rules for a predefined set of judgements

- The soundness corresponds to the soundness of the proof system

  - The input program is provable => the program satisfies the proven judgement.

# Example: Type Inference

- A simple ML-like language

$$
\begin{aligned}
E \quad \rightarrow \quad & n \\
\mid \quad & x \\
\mid \quad & E + E \\
\mid \quad & E - E \\
\mid \quad & \texttt{iszero } E \\
\mid \quad & \texttt{if } E \texttt{ then } E \texttt{ else } E \\
\mid \quad & \texttt{let } x = E \texttt{ in } E \\
\mid \quad & \texttt{proc } x\ E \\
\mid \quad & E\ E
\end{aligned}
$$

# Types

Types are defined inductively:

$$
\begin{aligned}
T \quad \rightarrow \quad & \text{int} \\
| \quad & \text{bool} \\
| \quad & T \rightarrow T
\end{aligned}
$$

Examples:

- int
- bool
- int $\rightarrow$ int
- bool $\rightarrow$ int
- int $\rightarrow$ (int $\rightarrow$ bool)
- (int $\rightarrow$ int) $\rightarrow$ (bool $\rightarrow$ bool)
- (int $\rightarrow$ int) $\rightarrow$ (bool $\rightarrow$ (bool $\rightarrow$ int))

# Types of Expressions

- Judgement that says expression E has type t is written as

$$\Gamma \vdash e : t$$

- $\Gamma$ is a set of type assumptions for the free variables in E (called *type environment*)

$$\Gamma : Var \rightarrow T$$

# Examples

- $[] \vdash \mathbf{3} : \mathsf{int}$
- $[x \mapsto \mathsf{int}] \vdash x : \mathsf{int}$
- $[] \vdash \mathbf{4} - \mathbf{3} :$
- $[x \mapsto \mathsf{int}] \vdash x - \mathbf{3} :$
- $[] \vdash \mathtt{iszero}\ \mathbf{11} :$
- $[] \vdash \mathtt{proc}\ (x)\ (x - \mathbf{11}) :$
- $[] \vdash \mathtt{proc}\ (x)\ (\mathtt{let}\ y = x - \mathbf{11}\ \mathtt{in}\ (x - y)) :$
- $[] \vdash \mathtt{proc}\ (x)\ (\mathtt{if}\ x\ \mathtt{then}\ \mathbf{11}\ \mathtt{else}\ \mathbf{22}) :$
- $[] \vdash \mathtt{proc}\ (x)\ (\mathtt{proc}\ (y)\ \mathtt{if}\ y\ \mathtt{then}\ x\ \mathtt{else}\ \mathbf{11}) :$
- $[] \vdash \mathtt{proc}\ (f)\ (\mathtt{if}\ (f\ \mathbf{3})\ \mathtt{then}\ \mathbf{11}\ \mathtt{else}\ \mathbf{22}) :$
- $[] \vdash (\mathtt{proc}\ (x)\ x)\ \mathbf{1} :$
- $[f \mapsto \mathsf{int} \rightarrow \mathsf{int}] \vdash (f\ (f\ \mathbf{1})) :$

# Type System

Inductive rules for assigning types to expressions:

$$\overline{\Gamma \vdash n : \mathsf{int}} \qquad \overline{\Gamma \vdash x : \Gamma(x)}$$

$$\frac{\Gamma \vdash E_1 : \mathsf{int} \qquad \Gamma \vdash E_2 : \mathsf{int}}{\Gamma \vdash E_1 + E_2 : \mathsf{int}} \qquad \frac{\Gamma \vdash E_1 : \mathsf{int} \qquad \Gamma \vdash E_2 : \mathsf{int}}{\Gamma \vdash E_1 - E_2 : \mathsf{int}}$$

$$\frac{\Gamma \vdash E : \mathsf{int}}{\Gamma \vdash \mathtt{iszero}\ E : \mathsf{bool}} \qquad \frac{\Gamma \vdash E_1 : \mathsf{bool} \qquad \Gamma \vdash E_2 : t \qquad \Gamma \vdash E_3 : t}{\Gamma \vdash \mathtt{if}\ E_1\ \mathtt{then}\ E_2\ \mathtt{else}\ E_3 : t}$$

$$\frac{\Gamma \vdash E_1 : t_1 \qquad [x \mapsto t_1]\Gamma \vdash E_2 : t_2}{\Gamma \vdash \mathtt{let}\ x = E_1\ \mathtt{in}\ E_2 : t_2} \qquad \frac{\Gamma \vdash E_1 : t_1 \to t_2 \qquad \Gamma \vdash E_2 : t_1}{\Gamma \vdash E_1\ E_2 : t_2}$$

$$\frac{[x \mapsto t_1]\Gamma \vdash E : t_2}{\Gamma \vdash \mathtt{proc}\ x\ E : t_1 \to t_2}$$

We say that a closed expression $E$ has type $t$ iff we can derive $[] \vdash E : t$.

# Example

Program $\texttt{proc} \ (x) \ (x - 11)$ is typed $\text{int} \to \text{int}$

because we can prove $[] \vdash \texttt{proc} \ (x) \ (x - 11) : \text{int} \to \text{int}$

as follows:

$$\frac{}{[] \vdash \texttt{proc} \ (x) \ (x - 11) : \text{int} \to \text{int}}$$

# Soundness of Type System

**Theorem (Soundness of the proof rules)**

*Let $E$ be a program, an expression without free variables. If $\emptyset \vdash E : \tau$, then the program runs without a type error and returns a value of type $\tau$ if it terminates.*

# Automatic Type Inference

- A static analysis algorithm that automatically figures out types of expressions by observing how they are used.
- The algorithm is *sound and complete* with respect to the type system design.
  - ▶ (Sound) If the analysis finds a type for an expression, the expression is well-typed with the type according to the type system.
  - ▶ (Complete) If an expression has a type according to the type system, the analysis is guaranteed to find the type.
- The algorithm consists of two steps:
  1. Generate type equations from the program text.
  2. Solve the equations.

# Generating Type Equations

For every subexpression and variable, introduce type variables and derive equations between the type variables.

# Type Equations

- Type equations are conjunctions of "type equalities": e.g.,

$$
\begin{aligned}
t_0 &= t_f \rightarrow t_1 \\
t_1 &= t_x \rightarrow t_4 \\
t_3 &= \text{int} \\
t_4 &= \text{int} \\
t_2 &= \text{int} \\
t_f &= \text{int} \rightarrow t_3 \\
t_f &= t_x \rightarrow t_4
\end{aligned}
$$

- Type equations ($\textit{TyEqn}$) are defined inductively:

$$
\begin{aligned}
\textit{TyEqn} \quad &\rightarrow \quad \emptyset \\
&\mid \quad T \doteq T \;\wedge\; \textit{TyEqn}
\end{aligned}
$$

# Example

$$\text{proc} \left( \underbrace{f}_{t_f} \right) \text{proc} \left( \underbrace{x}_{t_x} \right) (\underbrace{(\underbrace{(f \ 3)}_{t_3} - \underbrace{(f \ x)}_{t_4})}_{t_2})$$

$t_1$

$t_0$

$$
\begin{aligned}
t_0 &= t_f \rightarrow t_1 \\
t_1 &= t_x \rightarrow t_4 \\
t_3 &= \text{int} \\
t_4 &= \text{int} \\
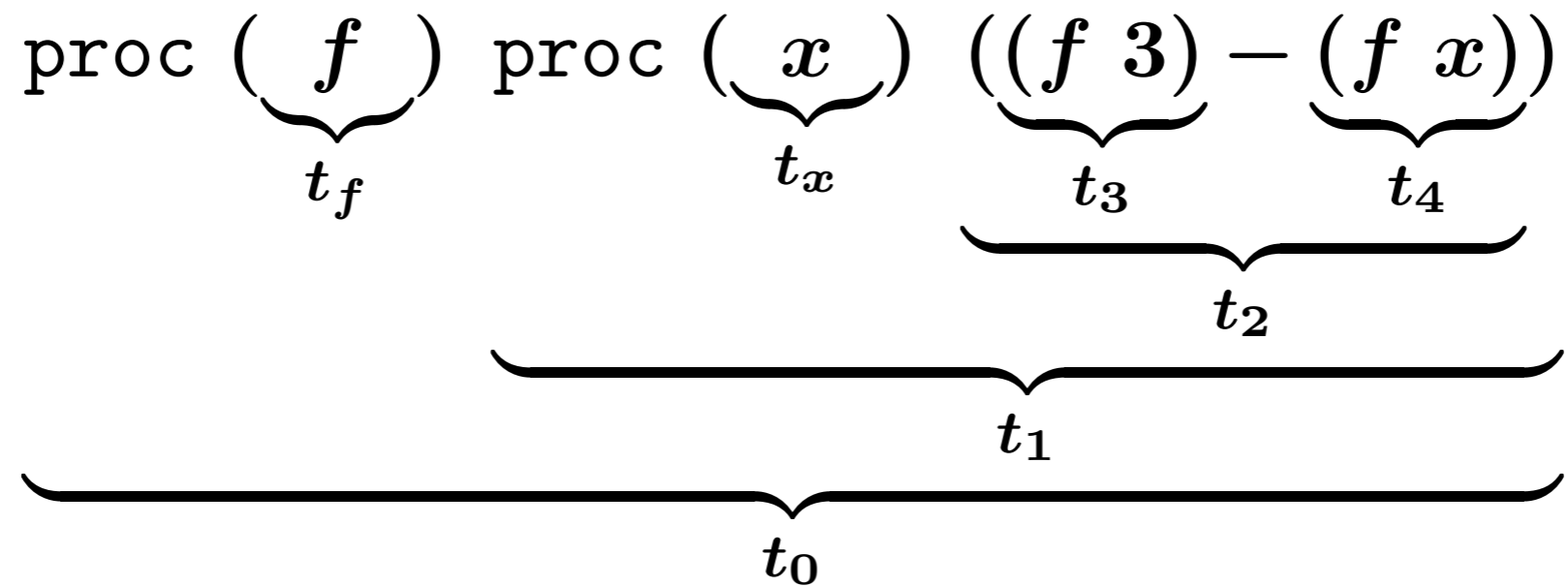t_2 &= \text{int} \\
t_f &= \text{int} \rightarrow t_3 \\
t_f &= t_x \rightarrow t_4
\end{aligned}
$$

# Example

$$\texttt{proc}\ (\ \underbrace{f}_{t_f}\ )\ (\underbrace{f\ 11}_{t_1})$$

$$\underbrace{\phantom{\texttt{proc}\ (\ f\ )\ (f\ 11)}}_{t_0}$$

$$t_0 = t_f \rightarrow t_1$$
$$t_f = \text{int} \rightarrow t_1$$

# Example

$$\underbrace{\text{if } \underbrace{x}_{t_x} \text{ then } \underbrace{(x - 1)}_{t_1} \text{ else } 0}_{t_0}$$

$$
\begin{aligned}
t_x &= \text{bool} \\
t_1 &= t_0 \\
\text{int} &= t_0 \\
t_x &= \text{int} \\
t_1 &= \text{int}
\end{aligned}
$$

# Generating Type Equations

- Algorithm for generating equations:

$$\mathcal{V} : (\textit{Var} \rightarrow T) \times E \times T \rightarrow \textit{TyEqn}$$

- $\mathcal{V}(\Gamma, e, t)$ generates the condition for $e$ to have type $t$ in $\Gamma$:

$$\Gamma \vdash e : t \text{ iff } \mathcal{V}(\Gamma, e, t) \text{ is satisfied.}$$

- Examples:
  - $\mathcal{V}([x \mapsto \text{int}], \texttt{x+1}, \alpha) = \alpha \doteq \text{int}$
  - $\mathcal{V}(\emptyset, \texttt{proc } (x) \ (\texttt{if } x \texttt{ then } 1 \texttt{ else } 2), \alpha \rightarrow \beta) =$
    $\alpha \doteq \text{bool} \wedge \beta \doteq \text{int}$

- To derive type equations for closed expression $E$, we call $\mathcal{V}(\emptyset, E, \alpha)$, where $\alpha$ is a fresh type variable.

# Generating Type Equations

$$\mathcal{V}(\Gamma, n, t) \;=\; t \doteq \mathsf{int}$$

$$\mathcal{V}(\Gamma, x, t) \;=\; t \doteq \Gamma(x)$$

$$\mathcal{V}(\Gamma, e_1 + e_2, t) \;=\; t \doteq \mathsf{int} \;\wedge\; \mathcal{V}(\Gamma, e_1, \mathsf{int}) \;\wedge\; \mathcal{V}(\Gamma, e_2, \mathsf{int})$$

$$\mathcal{V}(\Gamma, \mathtt{iszero}\; e, t) \;=\; t \doteq \mathsf{bool} \;\wedge\; \mathcal{V}(\Gamma, e, \mathsf{int})$$

$$\mathcal{V}(\Gamma, \mathtt{if}\; e_1\; e_2\; e_3, t) \;=\; \mathcal{V}(\Gamma, e_1, \mathsf{bool}) \;\wedge\; \mathcal{V}(\Gamma, e_2, t) \;\wedge\; \mathcal{V}(\Gamma, e_3, t)$$

$$\mathcal{V}(\Gamma, \mathtt{let}\; x = e_1\; \mathtt{in}\; e_2, t) \;=\; \mathcal{V}(\Gamma, e_1, \alpha) \;\wedge\; \mathcal{V}([x \mapsto \alpha]\Gamma, e_2, t) \;(\text{new } \alpha)$$

$$\mathcal{V}(\Gamma, \mathtt{proc}\; (x)\; e, t) \;=\; t \doteq \alpha_1 \to \alpha_2 \;\wedge\; \mathcal{V}([x \mapsto \alpha_1]\Gamma, e, \alpha_2)$$
$$(\text{new } \alpha_1, \alpha_2)$$

$$\mathcal{V}(\Gamma, e_1\; e_2, t) \;=\; \mathcal{V}(\Gamma, e_1, \alpha \to t) \;\wedge\; \mathcal{V}(\Gamma, e_2, \alpha) \;(\text{new } \alpha)$$
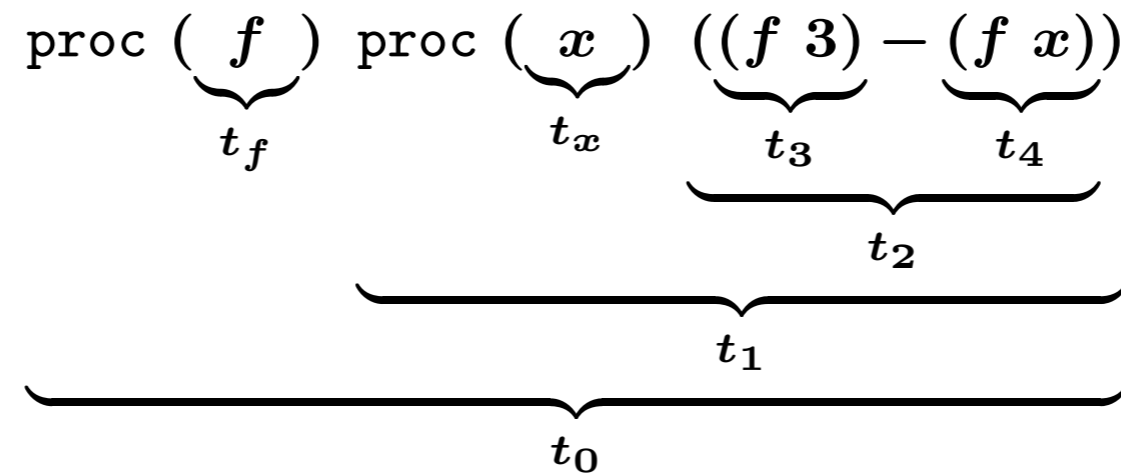
# Example

$$\mathcal{V}(\emptyset, (\texttt{proc } (x) \ (x)) \ 1, \alpha)$$
$$= \mathcal{V}(\emptyset, \texttt{proc } (x) \ (x), \alpha_1 \to \alpha) \wedge \mathcal{V}(\emptyset, 1, \alpha_1) \qquad \text{new } \alpha_1$$
$$= \alpha_1 \to \alpha \doteq \alpha_2 \to \alpha_3 \wedge \mathcal{V}([x \mapsto \alpha_2], x, \alpha_3) \wedge \alpha_1 \doteq \text{int} \quad \text{new } \alpha_2, \alpha_3$$
$$= \alpha_1 \to \alpha \doteq \alpha_2 \to \alpha_3 \wedge \alpha_2 \doteq \alpha_3 \wedge \alpha_1 \doteq \text{int}$$

# Finding a Solution of Type Equations

Find the values of type variables that make all the equations true.

$$\underbrace{\text{proc} \; (\; \underbrace{f}_{t_f}\;) \; \text{proc} \; (\; \underbrace{x}_{t_x}\;) \; \underbrace{(\underbrace{(f\;3)}_{t_3} - \underbrace{(f\;x)}_{t_4})}_{\substack{t_2 \\ t_1 \\ t_0}}}$$

|  | Equations |  |  |  | Solution |  |
|---|---|---|---|---|---|---|
| $t_0$ | $=$ | $t_f \to t_1$ | | $t_0$ | $=$ | $(\text{int} \to \text{int}) \to (\text{int} \to \text{int})$ |
| $t_1$ | $=$ | $t_x \to t_2$ | | $t_1$ | $=$ | $\text{int} \to \text{int}$ |
| $t_3$ | $=$ | $\text{int}$ | | $t_2$ | $=$ | $\text{int}$ |
| $t_4$ | $=$ | $\text{int}$ | | $t_3$ | $=$ | $\text{int}$ |
| $t_2$ | $=$ | $\text{int}$ | | $t_4$ | $=$ | $\text{int}$ |
| $t_f$ | $=$ | $\text{int} \to t_3$ | | $t_f$ | $=$ | $\text{int} \to \text{int}$ |
| $t_f$ | $=$ | $t_x \to t_4$ | | $t_x$ | $=$ | $\text{int}$ |

Static type systems find such a solution using *unification algorithm*.

# Finding a Solution of Type Equations

The calculation is split into equations to be solved and substitution found so far. Initially, the substitution is empty:

| | | Equations | Substitution |
|---|---|---|---|
| $t_0$ | $=$ | $t_f \rightarrow t_1$ | |
| $t_1$ | $=$ | $t_x \rightarrow t_2$ | |
| $t_3$ | $=$ | int | |
| $t_4$ | $=$ | int | |
| $t_2$ | $=$ | int | |
| $t_f$ | $=$ | int $\rightarrow t_3$ | |
| $t_f$ | $=$ | $t_x \rightarrow t_4$ | |

# Finding a Solution of Type Equations

Consider each equation in turn. If the equation's left-hand side is a variable, move it to the substitution:

| | Equations | | Substitution | |
|---|---|---|---|---|
| $t_1$ | $=$ | $t_x \to t_2$ | $t_0 \;=\; t_f \to t_1$ | |
| $t_3$ | $=$ | int | | |
| $t_4$ | $=$ | int | | |
| $t_2$ | $=$ | int | | |
| $t_f$ | $=$ | int $\to t_3$ | | |
| $t_f$ | $=$ | $t_x \to t_4$ | | |

# Finding a Solution of Type Equations

Move the next equation to the substitution and propagate the information to the existing substitution (i.e., substitute the right-hand side for each occurrence of $t_1$):

| Equations | | | Substitution | | |
|---|---|---|---|---|---|
| $t_3$ | $=$ | int | $t_0$ | $=$ | $t_f \to (t_x \to t_2)$ |
| $t_4$ | $=$ | int | $t_1$ | $=$ | $t_x \to t_2$ |
| $t_2$ | $=$ | int | | | |
| $t_f$ | $=$ | int $\to t_3$ | | | |
| $t_f$ | $=$ | $t_x \to t_4$ | | | |

# Finding a Solution of Type Equations

Same for the next three equations:

| Equations | | | Substitution | | |
|---|---|---|---|---|---|
| $t_4$ | $=$ | int | $t_0$ | $=$ | $t_f \to (t_x \to t_2)$ |
| $t_2$ | $=$ | int | $t_1$ | $=$ | $t_x \to t_2$ |
| $t_f$ | $=$ | int $\to t_3$ | $t_3$ | $=$ | int |
| $t_f$ | $=$ | $t_x \to t_4$ | | | |

| Equations | | | Substitution | | |
|---|---|---|---|---|---|
| $t_2$ | $=$ | int | $t_0$ | $=$ | $t_f \to (t_x \to t_2)$ |
| $t_f$ | $=$ | int $\to t_3$ | $t_1$ | $=$ | $t_x \to t_2$ |
| $t_f$ | $=$ | $t_x \to t_4$ | $t_3$ | $=$ | int |
| | | | $t_4$ | $=$ | int |

| Equations | | | Substitution | | |
|---|---|---|---|---|---|
| $t_f$ | $=$ | int $\to t_3$ | $t_0$ | $=$ | $t_f \to (t_x \to \text{int})$ |
| $t_f$ | $=$ | $t_x \to t_4$ | $t_1$ | $=$ | $t_x \to \text{int}$ |
| | | | $t_3$ | $=$ | int |
| | | | $t_4$ | $=$ | int |
| | | | $t_2$ | $=$ | int |

# Finding a Solution of Type Equations

Consider the next equation $t_f = \text{int} \to t_3$. The equation contains $t_3$, which is already bound to int in the substitution. Substitute int for $t_3$ in the equation. This is called *applying* the substitution to the equation.

| Equations | | | Substitution | | |
|---|---|---|---|---|---|
| $t_f$ | $=$ | $\text{int} \to \text{int}$ | $t_0$ | $=$ | $t_f \to (t_x \to \text{int})$ |
| $t_f$ | $=$ | $t_x \to t_4$ | $t_1$ | $=$ | $t_x \to \text{int}$ |
| | | | $t_3$ | $=$ | $\text{int}$ |
| | | | $t_4$ | $=$ | $\text{int}$ |
| | | | $t_2$ | $=$ | $\text{int}$ |

Move the resulting equation to the substitution and update it.

| Equations | | | Substitution | | |
|---|---|---|---|---|---|
| $t_f$ | $=$ | $t_x \to t_4$ | $t_0$ | $=$ | $(\text{int} \to \text{int}) \to (t_x \to \text{int})$ |
| | | | $t_1$ | $=$ | $t_x \to \text{int}$ |
| | | | $t_3$ | $=$ | $\text{int}$ |
| | | | $t_4$ | $=$ | $\text{int}$ |
| | | | $t_2$ | $=$ | $\text{int}$ |
| | | | $t_f$ | $=$ | $\text{int} \to \text{int}$ |

# Finding a Solution of Type Equations

Apply the substitution to the equation:

| Equations | Substitution |
|---|---|
| $\text{int} \rightarrow \text{int} \;=\; t_x \rightarrow \text{int}$ | $\begin{aligned} t_0 &= (\text{int} \rightarrow \text{int}) \rightarrow (t_x \rightarrow \text{int}) \\ t_1 &= t_x \rightarrow \text{int} \\ t_3 &= \text{int} \\ t_4 &= \text{int} \\ t_2 &= \text{int} \\ t_f &= \text{int} \rightarrow \text{int} \end{aligned}$ |

If neither side of the equation is a variable, simplify the equation by yielding two new equations:

| Equations | Substitution |
|---|---|
| $\begin{aligned} \text{int} &= t_x \\ \text{int} &= \text{int} \end{aligned}$ | $\begin{aligned} t_0 &= (\text{int} \rightarrow \text{int}) \rightarrow (t_x \rightarrow \text{int}) \\ t_1 &= t_x \rightarrow \text{int} \\ t_3 &= \text{int} \\ t_4 &= \text{int} \\ t_2 &= \text{int} \\ t_f &= \text{int} \rightarrow \text{int} \end{aligned}$ |

# Finding a Solution of Type Equations

Switch the sides of the first equation and move it to the substitution:

| Equations | | | Substitution | | |
|---|---|---|---|---|---|
| int | $=$ | int | $t_0$ | $=$ | $(\text{int} \rightarrow \text{int}) \rightarrow (\text{int} \rightarrow \text{int})$ |
| | | | $t_1$ | $=$ | $\text{int} \rightarrow \text{int}$ |
| | | | $t_3$ | $=$ | int |
| | | | $t_4$ | $=$ | int |
| | | | $t_2$ | $=$ | int |
| | | | $t_f$ | $=$ | $\text{int} \rightarrow \text{int}$ |
| | | | $t_x$ | $=$ | int |

The final substitution is the solution of the original equations.

# Finding a Solution of Type Equations

$$\texttt{proc (} \underbrace{f}_{t_f} \texttt{) (iszero } \underbrace{\overbrace{(f\ f)}^{t_2}}_{\underbrace{t_1}}\texttt{))}$$

$$\underbrace{\phantom{\texttt{proc (}f\texttt{) (iszero }(f\ f))}}_{t_0}$$

$$
\begin{aligned}
t_0 &= t_f \rightarrow t_1 \\
t_1 &= \texttt{bool} \\
t_2 &= \texttt{int} \\
t_f &= t_f \rightarrow t_2
\end{aligned}
$$

# Finding a Solution of Type Equations

Solving as usual, we encounter a problem:

| Equations | Substitution |
|-----------|--------------|
| $t_f = t_f \rightarrow \text{int}$ | $t_0 = t_f \rightarrow \text{bool}$ |
| | $t_1 = \text{bool}$ |
| | $t_2 = \text{int}$ |

- There is no type $t_f$ that satisfies the equation, because the right-hand side of the equation is always larger than the left.

- If we ever deduce an equation of the form $t = \ldots t \ldots$ where the type variable $t$ occurs in the right-hand side, we must conclude that there is no solution. This is called *occurrence check*.

# Unification Algorithm

For each equation in turn,

- Apply the current substitution to the equation.

- If the equation is always true (e.g. int $=$ int), discard it.

- If the left- and right-hand sides are contradictory (e.g. bool $=$ int), the algorithm fails.

- If neither side is a variable (e.g. int $\rightarrow t_1 = t_2 \rightarrow$ bool), simplify the equation, which eventually generates an equation whose left- or right-hand side is a variable.

- If the left-hand side is not a variable, switch the sides.

- If the left-hand side variable occurs in the right-hand side, the algorithm fails.

- Otherwise, move it to the substitution and substitute the right-hand side for each occurrence of the variable in the substitution.

# Substitutions

Solutions of type equations are represented by substitution:

$$S \in Subst = TyVar \rightarrow T$$

Applying a substitution to a type:

$$
\begin{aligned}
S(\text{int}) &= \text{int} \\
S(\text{bool}) &= \text{bool} \\
S(\alpha) &= \begin{cases} t & \text{if } \alpha \mapsto t \in S \\ \alpha & \text{otherwise} \end{cases} \\
S(T_1 \rightarrow T_2) &= S(T_1) \rightarrow S(T_2)
\end{aligned}
$$

# Example

Applying the substitution

$$S = \{t_1 \mapsto \mathsf{int}, t_2 \mapsto \mathsf{int} \to \mathsf{int}\}$$

to to the type $(t_1 \to t_2) \to (t_3 \to \mathsf{int})$:

$$S((t_1 \to t_2) \to (t_3 \to \mathsf{int}))$$
$$= S(t_1 \to t_2) \to S(t_3 \to \mathsf{int})$$
$$= (S(t_1) \to S(t_2)) \to (S(t_3) \to S(\mathsf{int}))$$
$$= (\mathsf{int} \to (\mathsf{int} \to \mathsf{int})) \to (t_3 \to \mathsf{int})$$

# Unification Algorithm

Update the current substitution with equality $t_1 \doteq t_2$.

$$\textbf{unify} : T \times T \times Subst \rightarrow Subst$$

$$
\begin{aligned}
\textbf{unify}(\mathsf{int}, \mathsf{int}, S) &= S \\
\textbf{unify}(\mathsf{bool}, \mathsf{bool}, S) &= S \\
\textbf{unify}(\alpha, \alpha, S) &= S \\
\textbf{unify}(\alpha, t, S) &= \begin{cases} \mathsf{fail} & \alpha \text{ occurs in } t \\ \text{extend } S \text{ with } \alpha \doteq t & \text{otherwise} \end{cases} \\
\textbf{unify}(t, \alpha, S) &= \textbf{unify}(\alpha, t, S) \\
\textbf{unify}(t_1 \rightarrow t_2, t_1' \rightarrow t_2', S) &= \text{let } S' = \textbf{unify}(t_1, t_1', S) \text{ in} \\
&\quad \text{let } S'' = \textbf{unify}(S'(t_2), S'(t_2'), S') \text{ in} \\
&\quad\quad S'' \\
\textbf{unify}(\_, \_, \_) &= \mathsf{fail}
\end{aligned}
$$

# Unification Algorithm

$$\mathbf{unifyall} : TyEqn \rightarrow Subst \rightarrow Subst$$

$$
\begin{aligned}
\mathbf{unifyall}(\emptyset, S) \;\; &= \;\; S \\
\mathbf{unifyall}((t_1 \doteq t_2) \;\wedge\; u, S) \;\; &= \;\; \text{let } S' = \mathbf{unify}(S(t_1), S(t_2), S) \\
&\quad\;\; \text{in } \mathbf{unifyall}(u, S')
\end{aligned}
$$

Let $\mathcal{U}$ be the final unification algorithm:

$$\mathcal{U}(u) = \mathbf{unifyall}(u, \emptyset)$$

# Automatic Type Inference

## typeof $: E \rightarrow T$

The final type inference algorithm that composes equation derivation ($\mathcal{V}$) and equation solving ($\mathcal{U}$):

$$
\begin{aligned}
&\textbf{typeof}(E) = \\
&\quad \textbf{let } S = \mathcal{U}(\mathcal{V}(\emptyset, E, \alpha)) \quad (\text{new } \alpha) \\
&\quad \textbf{in } S(\alpha)
\end{aligned}
$$

# Example

**typeof**$((\texttt{proc } (x)\ x)\ 1)$

# Correctness of Automatic Type Inference

**Theorem (Correctness of the algorithm)**

*Solving the equations $\equiv$ proving in the simple type system*

# Limitations

- For target languages that lack a sound static type system, we have to invent it.

  - Design a finite proof system

  - Prove the soundness of the proof system

  - Design its algorithm that automates proving

  - Prove the correctness of the algorithm

- What if the unification procedure is not enough?

  - For some properties, the algorithm can generate constraints that are unsolvable by the unification procedure

- For some conventional imperative language, sound and precise-enough static type systems are elusive.