

Semantic Analysis

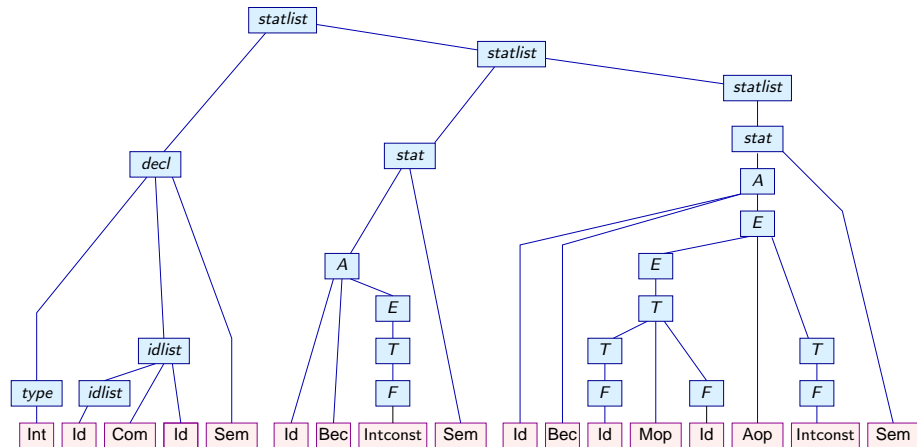
Sebastian Hack
Saarland University

W2017

Saarland University, Computer Science

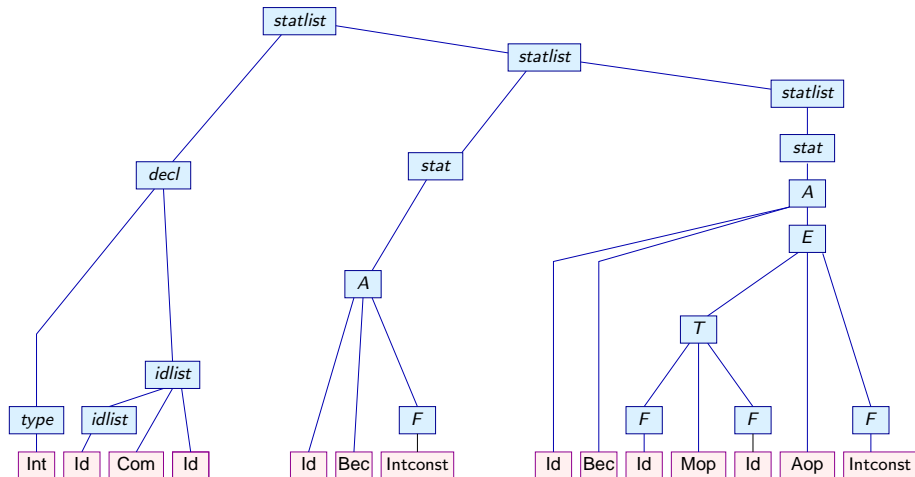
Abstract Syntax Trees (AST)

- Delete Terminals that were only there to make the grammar unambiguous
- Remove Nonterminals in chain productions (typical if precedence and associativity were expressed in grammar)



Abstract Syntax Trees (AST)

- Delete Terminals that were only there to make the grammar unambiguous
- Remove Nonterminals in chain productions (typical if precedence and associativity were expressed in grammar)



Semantic Analysis

Many programming languages have a **static** semantics. Interpretation of the syntax without knowing the inputs of the program.

Purpose:

- Refute (more) meaningless programs
- Infer information about program helpful for implementation

Program analyses with respect to the static semantics are often called **semantic analyses**. Commonly, we have two of them:

- Name Analysis
- Type Checking

Do we need semantic analysis?

Consider the following Python program:

```
def f(x):  
    if x = 0:  
        return x  
    else:  
        return y
```

What happens on $f(0)$ and $f(1)$?

Do we need semantic analysis?

Python does not have a static semantics, hence it has to “rule out” many invalid programs at runtime:

- When an identifier is declared, enter it into a definition table
- When it is used, look up, if there is a value for that identifier
- Incurs significant runtime overhead
- Allows more flexible (really really really?) code:
i.e. can add fields to objects at runtime

Similar things happen because of “duck typing”:

- Resolve operator overloading at runtime: $a + b$
Do we add two ints? Floats? Append to a string? Or is a an object of a class that redefines $+$?
- Similarly for methods: $a.b(c)$
Does the class of a really provide a method b or a superclass?

To avoid runtime overhead, we use **static semantics** to make these questions decidable **statically**.

Limits of static semantics

- C leaves semantics of many programs (under certain inputs) undefined. Result: runtime errors, e.g. division by zero, dereference of non-allocated memory
- Static semantics cannot capture them: in general undecidable
- Make language more restrictive to avoid some of them. E.g. Java:
 - Cannot take address of a field or array cell
 - Garbage collection policy avoids dangling pointers
- “Totalize” behavior, e.g.
 - Throw exception on `null` dereference
 - Array out-of-bounds check
 - Incurs runtime overhead!

Name Analysis

- Every entity (variable, function, data type, etc.) has a defining occurrence (definition) that gives it a name and describes its properties with respect to static semantics (e.g. type)
- Every defining occurrence has a scope in which it is **valid**
- Scope is the static pendant to **lifetime** of referred object (at runtime)
- An applied occurrence refers to a definition

Purpose of name analysis:

- Ensure that every applied occurrence refers to a unique defining occurrence
- Relate applied occurrence to defining occurrence in AST

Issues: Validity, Visibility, Qualification, Namespaces, Overloading

Validity

The **scope** (of validity) of a definition of an identifier x is the part of the program (nodes in the AST) in which an applied occurrence of x refers to it.

Inside a scope, each identifier can only be defined once.

Example

```
{  
    int y;  
    // ...  
    int x = 1;  
    // ...  
    y = x + 1;  
}
```

The definition of x is valid inside the entire block. But only visible after its definition.

Visibility

Inside its scope, a definition of an identifier x can be **invisible** (hidden, overwritten). Then, it is still forbidden to redefine x in the same scope (but in a nested one!) but not possible to refer to that definition:

```
for (int i = 0; i < n; i++)
    for (int i = 0; i < m; i++)
        A[i] = ...; // app occ of i refers
                    // to innermost definition
```

Possible in C. Java forbids overwriting definitions of local variables but allows

```
class X {
    int x;
    void foo(int x) {
        x = 1;
    }
}
```

Qualification

Some entities are composed of other entities (e.g. modules, classes, structs). Their definition opens a scope in which their constituents are defined:

```
struct vec3_t {  
    float x, y, z;  
};
```

The scope of `x`, `y`, and `z` is limited by `{ }`.

Use **qualification** to refer to constituents:

```
vec3_t v;  
// ...  
v.x = 1.0f;
```

In the static semantics, the `.` makes the identifier on the right-hand side refer to its definition inside the scope of the entity described by the left-hand side.

Multiple Name Spaces

Many programming languages allow multiple name spaces. More than one definition of the same identifier can be valid and visible at the same program point. Must be clear from the context which entity it is referred to.

Example (C)

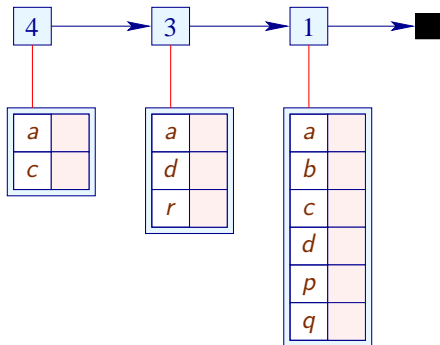
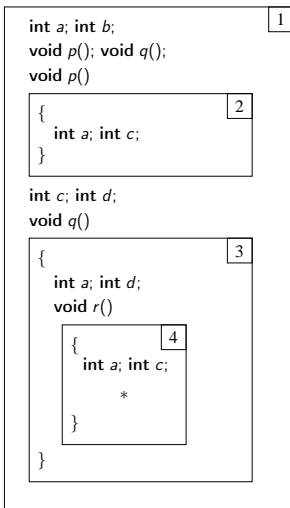
labels, structs/unions, and functions/variables are in disjoint name spaces.

```
struct x {
    int x;
};

int foo(int x) {
    struct x y; // x refers to struct
    y.x = x;    // refers to parameter
    goto x;    // refers to label
x:
    return y.x; // refers to parameter
}
```

Symbol (Definition) Tables

To detect the valid defining occurrence for an applied occurrence and to handle visibility one typically uses a **stack of environments** ($\text{Env} = \text{Id} \rightarrow \text{Type}$).



Pink boxes are pointers to the AST node of the identifier's declaration.

Symbol Table API

```
public class Scope {  
  
    private final Map<String, Declaration> decls;  
    private final Scope parent;  
  
    private Scope(Scope parent) {  
        this.parent = parent;  
    }  
  
    private Scope() {  
        this(null);  
    }  
  
    public boolean isRoot()    { return parent == null; }  
    public Scope enterBlock() { return new Scope(this); }  
    public Scope leaveBlock() { return parent; }  
}
```

Symbol Table API

```
public static class NotFound extends Exception { }
public static class AlreadyDefined extends Exception {

public Declaration lookup(String name) throws NotFound {
    if (decls.containsKey(name))
        return decls.get(name);
    else if (!isRoot())
        return parent.lookup(name);
    else
        throw new NotFound();
}

public void add(Declaration decl) throws AlreadyDefined {
    String name = decl.getToken().getText();
    if (decls.containsKey(name))
        throw new AlreadyDefined();
    decls.put(name, decl);
}
}
```

Type Checking C

Type checking C is done by a bottom-up traversal of the AST.

Resolve overloading based on types of operands and implicit type cast rules:

```
public Binary extends Expression {
    protected abstract Type applyOperatorTypeRules(Type left, Type right);

    public Type getType() {
        if (this.type == null) this.type = computeType();
        return this.type;
    }

    protected Type computeType() {
        Type lt = left.getType();
        Type rt = right.getType();
        if (lt.isError() || rt.isError())
            return Type.ERROR;
        return applyOperatorTypeRules(lt, rt);
    }
}
```


Type Inference

Often types of variables are uniquely identified by the way the variables are used. Use this information to avoid declarations and **infer** types automatically.

This is common in functional language and modern functional/imperative languages such as Scala, Rust, OCaml, ...

Consider following example:

```
let rec fac = fun → if x ≤ 0 then 1
                    else x * fac (x - 1)
```

Because one branch uses the int constant 1, the term $x * \text{fac } (x - 1)$ must be an int as well. Hence x and $\text{fac } (x - 1)$ must be ints as well. And finally fac is a function from $\text{int} \rightarrow \text{int}$.

A Toy Functional Language: Syntax

Syntax:

$$\begin{aligned} \text{Expr} \ni e & ::= c \mid x \mid \oplus \bar{e} \\ & \mid \text{if } e_0 \text{ then } e_1 \text{ else } e_2 \\ & \mid (e_1, \dots, e_k) \mid [] \mid e_1 :: e_2 \\ & \mid (e_1 \ e_2) \mid (\text{fun } x \rightarrow e) \\ & \mid \text{let } x_1 = e_1 \text{ in } e_0 \\ & \mid \text{let rec } x_1 = e_1 \text{ and } \dots \text{ and } x_n = e_n \text{ in } e_0 \end{aligned}$$

Types:

$$\text{Type} \ni t ::= \text{int} \mid \text{bool} \mid (t_1 * \dots * t_n) \mid t \ \text{list} \mid t_1 \rightarrow t_n$$

A Toy Functional Language: Static Semantics

Define a relation

$$\Gamma \vdash e : t \subseteq (Id \rightarrow Type) \times Expr \times Type$$

to indicate that under the environment Γ , expression e has type t .

Define typing rules inductively over the syntax.

Axioms:

$\Gamma \vdash c : t_c$	$(Const)$	t_c type of constant c
$\Gamma \vdash [] : t$ list	(Nil)	$\forall t$
$\Gamma \vdash x : (\Gamma x)$	(Var)	look up type of x in environment

A Toy Functional Language: Static Semantics

$$\text{Op: } \frac{\Gamma \vdash e_1 : \mathbf{int} \quad \Gamma \vdash e_2 : \mathbf{int}}{\Gamma \vdash e_1 + e_2 : \mathbf{int}}$$

$$\text{Comp: } \frac{\Gamma \vdash e_1 : t \quad \Gamma \vdash e_2 : t}{\Gamma \vdash e_1 = e_2 : \mathbf{bool}}$$

$$\text{If: } \frac{\Gamma \vdash e_0 : \mathbf{bool} \quad \Gamma \vdash e_1 : t \quad \Gamma \vdash e_2 : t}{\Gamma \vdash (\mathbf{if } e_0 \mathbf{ then } e_1 \mathbf{ else } e_2) : t}$$

$$\text{Tuple: } \frac{\Gamma \vdash e_1 : t_1 \quad \dots \quad \Gamma \vdash e_m : t_m}{\Gamma \vdash (e_1, \dots, e_m) : (t_1 * \dots * t_m)}$$

$$\text{Cons: } \frac{\Gamma \vdash e_1 : t \quad \Gamma \vdash e_2 : t \text{ list}}{\Gamma \vdash (e_1 :: e_2) : t \text{ list}}$$

$$\text{App: } \frac{\Gamma \vdash e_1 : t_1 \rightarrow t_2 \quad \Gamma \vdash e_2 : t_1}{\Gamma \vdash (e_1 e_2) : t_2}$$

$$\text{Fun: } \frac{\Gamma \oplus \{x \mapsto t_1\} \vdash e : t_2}{\Gamma \vdash \mathbf{fun } x \rightarrow e : t_1 \rightarrow t_2}$$

$$\text{Let: } \frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \oplus \{x_1 \mapsto t_1\} \vdash e_0 : t}{\Gamma \vdash (\mathbf{let } x_1 = e_1 \mathbf{ in } e_0) : t}$$

$$\text{Letrec: } \frac{\Gamma' \vdash e_1 : t_1 \quad \dots \quad \Gamma' \vdash e_m : t_m \quad \Gamma' \vdash e_0 : t}{\Gamma \vdash (\mathbf{let rec } x_1 = e_1 \mathbf{ and } \dots \mathbf{ and } x_m = e_m \mathbf{ in } e_0) : t}$$

with $\Gamma' = \Gamma \oplus \{x_1 \mapsto t_1, \dots, x_m \mapsto t_m\}$

An Example Derivation

$$\frac{\Gamma \vdash x : \mathbf{int} \quad \Gamma \vdash 0 : \mathbf{int}}{\Gamma \vdash x \leq 0 : \mathbf{bool}} \dots$$

$$\dots \frac{\Gamma \vdash 1 : \mathbf{int} \quad \frac{\Gamma \vdash x : \mathbf{int} \quad \frac{\Gamma \vdash \mathbf{fac} : \mathbf{int} \rightarrow \mathbf{int} \quad \frac{\Gamma \vdash x : \mathbf{int} \quad \Gamma \vdash 1 : \mathbf{int}}{\Gamma \vdash x - 1 : \mathbf{int}}}{\Gamma \vdash \mathbf{fac}(x - 1) : \mathbf{int}}}{\Gamma \vdash x \cdot \mathbf{fac}(x - 1) : \mathbf{int}}}{\Gamma \vdash \mathbf{if } x \leq 0 \mathbf{ then } 1 \mathbf{ else } x \cdot \mathbf{fac}(x - 1) : \mathbf{int}}$$

Remarks

- An important property of a type system is **preservation**:
The type of an expression does not change under evaluation.

If $\Gamma \vdash e_1 : \tau$ and $e_1 \Rightarrow^* e_2$ then $\Gamma \vdash e_2 : \tau$

Our type system exhibits preservation.

- preservation + progress = type safety
“Well-typed programs don’t go wrong”
- Given an environment, the rules can be used to **check** if the environment leads to a correct typing
- Note that an expression can have multiple types. For every Γ and every t , we can derive

$\Gamma \vdash \text{fun } x \rightarrow x : t \rightarrow t$

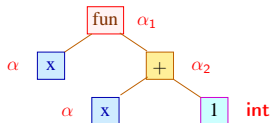
- How do we characterize all valid types for an expression?

An Equation System

- Build an equation system to describe all valid types for an expression
- Introduce **type variables** for the types of expressions with unknown type
- Technicality: To have one equation system for an expression, name variable apart to make their name unique
- Build equations based on typing rules

Consider `fun x → x + 1`. Based on the typing rules, we come up with the following equations:

$$\begin{array}{llll} \text{Fun} & : & \alpha_1 & = \alpha \rightarrow \alpha_2 \\ \text{Op} & : & \alpha_2 & = \mathbf{int} \\ & & \alpha & = \mathbf{int} \\ & & \mathbf{int} & = \mathbf{int} \end{array}$$



and hence

$$\alpha = \mathbf{int} \quad \alpha_1 = \mathbf{int} \rightarrow \mathbf{int} \quad \alpha_2 = \mathbf{int}$$

An Equation System

Let $\alpha[e]$ be the type variable for expression e .

Const:	$e \equiv b$	$\alpha[e] = t_b$
Nil:	$e \equiv []$	$\alpha[e] = \alpha \text{ list}$ (α fresh)
Op:	$e \equiv e_1 + e_2$	$\alpha[e] = \mathbf{int}$ $\alpha[e_1] = \mathbf{int}$ $\alpha[e_2] = \mathbf{int}$
Comp:	$e \equiv e_1 = e_2$	$\alpha[e_1] = \alpha[e_2]$ $\alpha[e] = \mathbf{bool}$
Tuple:	$e \equiv (e_1, \dots, e_m)$	$\alpha[e] = (\alpha[e_1] * \dots * \alpha[e_m])$
Cons:	$e \equiv e_1 :: e_2$	$\alpha[e_2] = \alpha[e_1] \text{ list}$ $\alpha[e] = \alpha[e_1] \text{ list}$
If:	$e \equiv \mathbf{if } e_0 \mathbf{ then } e_1 \mathbf{ else } e_2$	$\alpha[e_0] = \mathbf{bool}$ $\alpha[e] = \alpha[e_1]$ $\alpha[e] = \alpha[e_2]$
Fun:	$e \equiv \mathbf{fun } x \rightarrow e_1$	$\alpha[e] = \alpha[x] \rightarrow \alpha[e_1]$
App:	$e \equiv e_1 e_2$	$\alpha[e_1] = \alpha[e_2] \rightarrow \alpha[e]$
Letrec:	$e \equiv \mathbf{let rec } x_1 = e_1 \mathbf{ and } \dots \mathbf{ and } x_m = e_m \mathbf{ in } e_0$	$\alpha[x_1] = \alpha[e_1] \quad \dots$ $\alpha[x_m] = \alpha[e_m]$ $\alpha[e] = \alpha[e_0]$

A solution σ is a **unifier**, a **substitution** (= map from type variables to types) such that $\sigma s_i \equiv \sigma t_i$ for every equation $s_i = t_i$ in the equation system.

Correctness

To make use of the equation system, we have to show the following theorem:

Theorem

Let e be an expression, V be the set of variables in e and E be the equation system for e . Then,

1. If σ is a solution of E , then there is a derivation of $\Gamma \vdash e : t$ with

$$\Gamma = \{x \mapsto \sigma(\alpha[x]) \mid x \in V\} \quad \text{and} \quad t = \sigma(\alpha[e])$$

2. Consider a derivation of the judgement $\Gamma \vdash e : t$ that contains judgements $\Gamma \vdash e' : t_{e'}$ for all subterms e' of e . Then, the substitution

$$\sigma(\alpha[e']) = \begin{cases} t_{e'} & e' \text{ subterm of } e \\ \Gamma x & e' \equiv x \in V \end{cases}$$

is a solution of E .

Substitutions: Examples

1. Consider the equation

$$Y = X \rightarrow X$$

The set of solutions is given by the substitution

$$\{X \mapsto t, Y \mapsto (t \rightarrow t)\}$$

for every type t .

2. The equation

$$X \rightarrow \text{int} = \text{bool} \rightarrow Z$$

has exactly one solution:

$$\{X \mapsto \text{bool}, Z \mapsto \text{int}\}$$

3. The equation

$$\text{bool} = X \rightarrow Y$$

has no solution.

Substitutions: Definitions

- σ is **idempotent** if $\sigma \circ \sigma = \sigma$
Means that there is no variable Y with $\sigma(Y) = X$ with $\sigma(X) \neq X$.
- An idempotent σ is **most general** if for every other idempotent unifier τ , there is an appropriate substitution τ' with $\tau = \tau' \circ \sigma$. E.g.

$$\{Y \mapsto (X \rightarrow X), X \mapsto X\}$$

is the most general unifier of the first example on the last slide.

- A set of term equations $s_i = t_i$, $1 \leq i \leq m$ has either no solution or a most general idempotent unifier.

Computing the most general unifier

Given the equation system $s_i = t_i$, $1 \leq i \leq m$, start with
 $\text{unifyList} [(s_1, t_1), \dots, (s_m, t_m)] \emptyset$

```
let rec unify (s, t)  $\theta$  = if  $\theta s \equiv \theta t$  then  $\theta$ 
  else match ( $\theta s, \theta t$ )
    with (X, t)  $\rightarrow$  if occurs (X, t) then Fail
      else  $\{X \mapsto t\} \circ \theta$ 
    | (t, X)  $\rightarrow$  if occurs (X, t) then Fail
      else  $\{X \mapsto t\} \circ \theta$ 
    | (f(s1, ..., sk), f(t1, ..., tk))  $\rightarrow$ 
      unifyList [(s1, t1), ..., (sk, tk)]  $\theta$ 
    | (a, a)  $\rightarrow \theta$ 
    | _  $\rightarrow$  Fail
and unifyList list  $\theta$  = match list
  with []  $\rightarrow \theta$ 
    | ((s, t) :: rest)  $\rightarrow$  let  $\theta = \text{unify } (s, t) \theta$ 
      in if  $\theta = \text{Fail}$  then Fail
        else unifyList rest  $\theta$ 
```

Syntax-Directed Version (Algorithm W)

Solving the equation system as a whole does not give precise error messages (we don't know where an error was). Instead of collecting equations, solve equations in a syntax-directed way:

```
let rec  $\mathcal{W} e (\Gamma, \theta) = \text{match } e$ 
  with
  |  $c$             $\rightarrow (t_c, \theta)$ 
  |  $[]$           $\rightarrow \text{let } \alpha = \text{new}() \text{ in } (\alpha \text{ list}, \theta)$ 
  |  $x$            $\rightarrow (\Gamma(x), \theta)$ 
  |  $(e_1 :: e_2)$   $\rightarrow \text{let } (t_1, \theta) = \mathcal{W} e_1 (\Gamma, \theta) \text{ in}$ 
                        $\text{let } (t_2, \theta) = \mathcal{W} e_2 (\Gamma, \theta) \text{ in}$ 
                        $(t_2, \text{unify } (t_1 \text{ list}, t_2) \theta)$ 
  |  $(e_1' + e_2)$   $\rightarrow \text{let } (t_1, \theta) = \mathcal{W} e_1 (\Gamma, \theta) \text{ in}$ 
                        $\text{let } (t_2, \theta) = \mathcal{W} e_2 (\Gamma, \theta) \text{ in}$ 
                        $\text{let } \theta = \text{unify } (\text{int}, t_1) \theta \text{ in}$ 
                        $\text{let } \theta = \text{unify } (\text{int}, t_2) \theta \text{ in}$ 
                        $(\text{int}, \theta)$ 
  |  $(e_1 e_2)$     $\rightarrow \text{let } (t_1, \theta) = \mathcal{W} e_1 (\Gamma, \theta) \text{ in}$ 
                        $\text{let } (t_2, \theta) = \mathcal{W} e_2 (\Gamma, \theta) \text{ in}$ 
                        $\text{let } \alpha = \text{new}() \text{ in}$ 
                        $\text{let } \theta = \text{unify } (t_1, t_2 \rightarrow \alpha) \theta \text{ in}$ 
                        $(\alpha, \theta)$ 
```

...

Syntax-Directed Version (Algorithm W)

```
| (fun x → e)
  → let α = new() in
     let (t, θ) = W e (Γ ⊕ {x ↦ α}, θ) in
       (α → t, θ)

| (let x1 = e1 in e0)
  → let (t1, θ) = W e1 (Γ, θ) in
     in let Γ = Γ ⊕ {x1 ↦ t1} in
        in let (t0, θ) = W e0 (Γ, θ) in
           (t0, θ)

| (let rec x1 = e1 and ... and xm = em in e0)
  → let α1 = new()
     ...
     let αm = new() in
     let Γ = Γ ⊕ {x1 ↦ α1, ..., xm ↦ αm}} in
     let (t1, θ) = W e1 (Γ, θ) in
     let θ = unify (α1, t1) θ in
     ...
     let (tm, θ) = W em (Γ, θ) in
     let θ = unify (αm, tm) θ in
     let (t0, θ) = W e0 (Γ, θ) in
     (t0, θ)
```

(Let-) Polymorphism

Consider the expression

```
let single = fun y → [y]
in single (single 1)
```

We derive $\gamma \rightarrow \gamma \text{ list}$ as the type for `single`.

Because of `(single 1)` we unify γ with `int`.

The outer expression then wants to unify `int` with `int list` which fails.

Reason:

Type variables have to be instantiated **once** for the expression.

Solution:

Introduce **type schemas** and allow type variables to be instantiated **differently** in subexpressions.

$$\forall \alpha_1, \dots, \alpha_n. t$$

generalizes type variables $\alpha_1, \dots, \alpha_n$ in type expression t .

(Let-) Polymorphism

Introduce new typing rules for let and let rec that generalize type variables that do not occur in environment:

$$\text{Inst: } \frac{\Gamma(x) = \forall \alpha_1, \dots, \alpha_k. t}{\Gamma \vdash x : t[t_1/\alpha_1, \dots, t_k/\alpha_k]} \quad (\text{for all } t_1, \dots, t_k)$$
$$\text{Let: } \frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \oplus \{x \mapsto \text{close } t_1 \Gamma\} \vdash e_0 : t_0}{\Gamma \vdash (\text{let } x_1 = e_1 \text{ in } e_0) : t_0}$$

The function $\text{close } t \Gamma$ generalizes all type variables in t that do not occur in Γ . We modify Algorithm W accordingly:

$$\begin{array}{l} \text{fun inst } (\forall \alpha_1, \dots, \alpha_k. t) = \\ \quad \text{let } \beta_1 = \text{new}() \\ \quad \dots \\ \quad \text{in let } \beta_k = \text{new}() \\ \quad \text{in } t[\beta_1/\alpha_1, \dots, \beta_k/\alpha_k] \end{array} \quad \begin{array}{l} \text{let rec } \mathcal{W} e (\Gamma, \theta) = \dots \\ | \quad x \quad \rightarrow \quad (\text{inst } (\theta(\Gamma(x))), \theta) \\ | \quad (\text{let } x_1 = e_1 \text{ in } e_0) \\ \quad \quad \rightarrow \quad \text{let } (t_1, \theta) = \mathcal{W} e_1 (\Gamma, \theta) \\ \quad \quad \text{in let } s_1 = \text{close } (\theta t_1) (\theta \circ \Gamma) \\ \quad \quad \text{in let } \Gamma = \Gamma \oplus \{x_1 \mapsto s_1\} \\ \quad \quad \text{in let } (t_0, \theta) = \mathcal{W} e_0 (\Gamma, \theta) \\ \quad \quad \text{in } (t_0, \theta) \end{array}$$