

CS510
Semantics & Types
Winter'08 Lecture Notes
Type Inference

Type Inference

Type inference is the problem of producing a well-typed term from an **untyped term**. Having a type inference mechanism for a language can be very useful, since it saves the programmer from having to declare the types of all identifiers (though it does not **prevent** such declarations).

Type inference is **not** possible for every language and typing system, but there **is** a type inference algorithm relating the untyped and simply-typed λ -calculi.

Formally, we define a function **erase**, which takes a typed term, removes all the type annotations from the λ -bindings, and returns the resulting untyped term.

An untyped term M is **typable** if there exists a (simply-)typed term M' , a type τ and an environment E (with $FV(M) \in \text{Dom}(E)$) such that

$$E \vdash M' : \tau \text{ and } M = \text{erase}(M')$$

A (typable) untyped term may be the erasure of more than one typed term. Example: $\lambda x. x$ is the erasure of both $\lambda x : \text{int}. x$ and $\lambda x : \text{bool}. x$.

Type Schemes

To make the type inference algorithm precise, it is helpful to extend the language of type expressions to include **type variables**, denoted by $\{\alpha, \beta, \gamma, \dots\}$.

The resulting grammar for **type schemes** is:

$$\tau ::= b \mid \alpha \mid \tau_1 \rightarrow \tau_2 \quad (b \in B)$$

They are called schemes because each expression can be thought of as a schematic representation of a whole family of (ordinary) types, produced by **instantiating** the type variables with ordinary types.

Examples

α
 $\alpha \rightarrow \beta$
 $\alpha \rightarrow (\text{int} \rightarrow \alpha)$
 bool

An ordinary type is just a type scheme with no type variables. They are also called **ground types** or **monotypes**.

We modify our typing system so that typing rules, environments, λ -binding annotations, assertions, etc., all refer to type **schemes** where they used to refer to ground types. All this requires is a suitable change in the interpretation of the metavariables τ .

Type Inference

Our type inference algorithm will produce a type **scheme** for a given expression (in a given environment). Working with type schemes allows us to concentrate on the **structure** of the term being typed (and the corresponding typing derivation), which is the same regardless of the specific types in question.

Recall that the form of a type derivation for a given term is a tree isomorphic to the syntax tree for the term. To infer a type for the term, we start with a derivation tree of the correct form in which

- each λ -binding type scheme annotation is a fresh type variable
- the type scheme of each term is a fresh type variable.

Example $\lambda y. y \ 3$

$$\frac{\frac{\frac{}{\{y:\alpha_0\} \vdash y:\alpha_3} \text{VAR}}{\{y:\alpha_0\} \vdash (y \ 3):\alpha_2} \text{APP}}{\emptyset \vdash (\lambda y:\alpha_0. y \ 3):\alpha_1} \text{ABS}}{\frac{}{\{y:\alpha_0\} \vdash 3:\alpha_4} \text{CONST}}{\emptyset \vdash (\lambda y:\alpha_0. y \ 3):\alpha_1} \text{APP}}$$

To make this a **valid** derivation, certain relationships must hold between the type variables mentioned in adjacent rules. We can express these in the form of **equations between type schemes**. These equations represent **constraints** on the possible (valid) instantiations of the type variables.

Example continued

$$\frac{\frac{\frac{\frac{\emptyset \vdash y:\alpha_0}{\emptyset \vdash y:\alpha_3} VAR}{\emptyset \vdash (y\ 3):\alpha_2} APP}{\emptyset \vdash (\lambda y:\alpha_0. y\ 3):\alpha_1} ABS}{\frac{\frac{\emptyset \vdash y:\alpha_0}{\emptyset \vdash y:\alpha_3} VAR}{\emptyset \vdash 3:\alpha_4} CONST} APP$$

Each use of an axiom or rule produces one constraint equation. In this example, we get the following

- ABS rule: $\alpha_1 = \alpha_0 \rightarrow \alpha_2$
- APP rule: $\alpha_3 = \alpha_4 \rightarrow \alpha_2$
- VAR rule: $\alpha_3 = \alpha_0$
- CONST rule: $\alpha_4 = \text{int}$

To produce a valid type derivation, we must solve all the constraint equations **simultaneously**.

Here's one way to do this. We eliminate α_3 and α_4 using the (CONST) and (VAR) constraints, leaving

$$\begin{aligned} \alpha_1 &= \alpha_0 \rightarrow \alpha_2 \\ \alpha_0 &= \text{int} \rightarrow \alpha_2 \end{aligned}$$

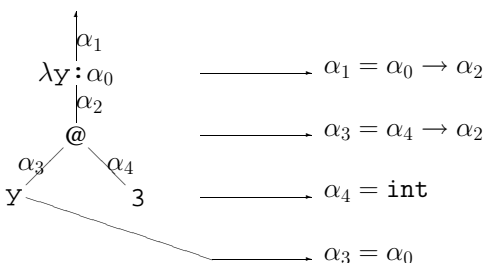
Substituting the second equation into the first, we get:

$$\alpha_1 = (\text{int} \rightarrow \alpha_2) \rightarrow \alpha_2$$

Constraints

$$\begin{aligned} \text{(CONST)} \quad E \vdash c:\tau \quad (\text{TypeOf}(c) = \tau) & \quad \begin{array}{c} \uparrow \\ \tau' \\ \hline c \end{array} \quad \tau = \tau' \\ \\ \text{(VAR)} \quad E \vdash x:\tau \quad (E(x) = \tau) & \quad \begin{array}{c} \uparrow \\ \tau' \\ \hline x \end{array} \quad \tau = \tau' \\ \\ \text{(ABS)} \quad \frac{E + \{x:\tau\} \vdash M:\tau'}{E \vdash (\lambda x:\tau. M):\tau''} & \quad \begin{array}{c} \uparrow \\ \tau'' \\ \hline \lambda x:\tau \\ \downarrow \\ \tau' \end{array} \quad \tau'' = \tau \rightarrow \tau' \\ \\ \text{(APP)} \quad \frac{E \vdash M:\tau \quad E \vdash N:\tau'}{E \vdash (MN):\tau''} & \quad \begin{array}{c} \uparrow \\ \tau'' \\ \hline @ \\ \swarrow \quad \searrow \\ \tau \quad \tau' \end{array} \quad \tau = \tau' \rightarrow \tau'' \end{aligned}$$

Given a complete expression, we patch the patterns together to get constraint equations, e.g.:



Example continued again

Instantiating the original typing tree with this solution, we have the following derivation:

$$\frac{\frac{\frac{\frac{\emptyset \vdash y:\text{int} \rightarrow \alpha_2}{\emptyset \vdash y:\text{int} \rightarrow \alpha_2} VAR}{\emptyset \vdash (y\ 3):\alpha_2} APP}{\emptyset \vdash (\lambda y:\text{int} \rightarrow \alpha_2. y\ 3):(\text{int} \rightarrow \alpha_2) \rightarrow \alpha_2} ABS}{\frac{\frac{\emptyset \vdash y:\text{int} \rightarrow \alpha_2}{\emptyset \vdash y:\text{int} \rightarrow \alpha_2} VAR}{\emptyset \vdash 3:\text{int}} CONST} APP$$

It is easy to see that every **instance** of

$$\emptyset \vdash (\lambda y:\text{int} \rightarrow \alpha_2. y\ 3):(\text{int} \rightarrow \alpha_2) \rightarrow \alpha_2$$

i.e., every assertion resulting from instantiation of the remaining type variable (α_2), will also be derivable. Moreover, since this scheme was constrained as minimally as possible, **every** derivable scheme must be an instance of it.

Such a type scheme is called **principal**. The approach outlined here always generates a principal type scheme, (if the given term is typable).

The inference algorithm consists of two steps:

- The term's structure is analyzed to derive a set of constraint equations. This can be described simply in terms of the syntax tree.
- The equations are solved simultaneously. This is an instance of the **unification** problem.

Constraint Equations

We now have a collection of equations between type schemes

$$\begin{aligned} \tau_{11} &= \tau_{12} \\ \tau_{21} &= \tau_{22} \\ \dots & \\ \tau_{n1} &= \tau_{n2} \end{aligned}$$

to be solved simultaneously.

This is a special case of the general problem called **unification**. A solution to a unification problem is a **substitution** from variables to terms that makes all the equations true. Such a substitution is called a **unifier**.

Example: the set of equations

$$\begin{aligned} \alpha_1 &= \alpha_0 \rightarrow \alpha_2 \\ \alpha_3 &= \alpha_4 \rightarrow \alpha_2 \\ \alpha_4 &= \text{int} \\ \alpha_3 &= \alpha_0 \end{aligned}$$

has the unifier

$$S = \{ \alpha_1 \mapsto (\text{int} \rightarrow \alpha_2) \rightarrow \alpha_2, \alpha_0 \mapsto \alpha_3 \mapsto \text{int} \rightarrow \alpha_2, \alpha_4 \mapsto \text{int} \}$$

Moreover, S is a **most general unifier**, meaning that every unifying substitution S' can be obtained from S by a further substitution.

Unification

In general, the unification problem is as follows:

Given a set of equations between terms built from variables and constructors of (fixed) arity ≥ 0 , either:

- compute a most general unifier that solves the equations; or
- report that the equations have no solution.

(It can be proved that if the set of equations have any unifier, they have a most general unifier.)

Unification is a widely applicable technique; e.g., it is the fundamental computational process in logic programming. There are several efficient algorithms known for computing unifiers.

For type inference, we specialize the problem as follows:

- the variables are type variables;
- there is one constructor (\rightarrow) of arity 2;
- each base type is treated as a constructor of arity 0.

If there is a most general unifier for the equations, it yields a **principal** schematic typing for the original λ -expression. If there is no solution, the original expression is not typable.

Example Again

Use the unification algorithm to solve the set of constraints generated from $\lambda x. x \ 3$. To make the role and arity of the constructors clearer, we rewrite the infix \rightarrow constructor as the prefix arity-2 constructor C_{\rightarrow} , and the base type `int` as the arity-0 constructor C_{int} .

$$\begin{aligned}\alpha_1 &= C_{\rightarrow}(\alpha_0, \alpha_2) \\ \alpha_3 &= C_{\rightarrow}(\alpha_4, \alpha_2) \\ \alpha_4 &= C_{\text{int}} \\ \alpha_3 &= \alpha_0\end{aligned}$$

We calculate the m.g.u. by calling \mathcal{U} on each equation in turn, passing the result substitution of each call as the initial substitution for the next:

$$\begin{aligned}S_1 &= \mathcal{U}(\emptyset, \alpha_1, C_{\rightarrow}(\alpha_0, \alpha_2)) \\ S_2 &= \mathcal{U}(S_1, \alpha_3, C_{\rightarrow}(\alpha_4, \alpha_2)) \\ S_3 &= \mathcal{U}(S_2, \alpha_4, C_{\text{int}}) \\ \text{m.g.u.} &= S_4 = \mathcal{U}(S_3, \alpha_3, \alpha_0)\end{aligned}$$

(Notice that the result substitution grows steadily; the algorithm never backtracks. An efficient implementation can therefore maintain a single substitution in a global variable rather than threading it through each call to \mathcal{U} .)

A Unification Algorithm

The essence of unification is the recursive traversal of two terms to be equated (“unified”).

We denote a substitution by S , and extend its domain from variables to arbitrary terms in the obvious way. Substitutions can be **composed**; we write $S_1 \circ S_2$ to mean the result of applying S_2 and then S_1 .

Function \mathcal{U} takes an existing substitution and two terms, and returns an extension of the original substitution that unifies the terms (if possible):

$$\begin{aligned}\mathcal{U}(S, t_1, t_2) &= \\ \text{case } (S(t_1), S(t_2)) \text{ of} & \\ (v_1, v_2) : (v_1 \mapsto v_2) \circ S & \\ (v_1, t'_2) : \text{if } v_1 \text{ occurs in } t'_2 \text{ then error “No Unifier”} & \\ \text{else } (v_1 \mapsto t'_2) \circ S & \\ (t'_1, v_2) : \text{if } v_2 \text{ occurs in } t'_1 \text{ then error “No Unifier”} & \\ \text{else } (v_2 \mapsto t'_1) \circ S & \\ (C_1(t_{11}, t_{12}, \dots, t_{1m}), C_2(t_{21}, t_{22}, \dots, t_{2n})) : & \\ \text{if } C_1 \neq C_2 & \\ \text{or } m \neq n \text{ then error “No Unifier”} & \\ \text{else } \mathcal{U}(\dots \mathcal{U}(\mathcal{U}(S, t_{11}, t_{21}), t_{12}, t_{22}) \dots, t_{1m}, t_{2m}) &\end{aligned}$$

To find the m.g.u. of a set of equations

$$t_{11} = t_{21}, t_{12} = t_{22}, \dots, t_{1n} = t_{2n}$$

we calculate $\mathcal{U}(\dots \mathcal{U}(\mathcal{U}(\emptyset, t_{11}, t_{21}), t_{12}, t_{22}) \dots, t_{1n}, t_{2n})$

Example continued

To pursue the calculation:

$$\begin{aligned}S_1 &= \{ \alpha_1 \mapsto C_{\rightarrow}(\alpha_0, \alpha_2) \} \\ S_2 &= \{ \alpha_3 \mapsto C_{\rightarrow}(\alpha_4, \alpha_2) \} \circ S_1 \\ &= \{ \alpha_1 \mapsto C_{\rightarrow}(\alpha_0, \alpha_2), \\ &\quad \alpha_3 \mapsto C_{\rightarrow}(\alpha_4, \alpha_2) \} \\ S_3 &= \{ \alpha_4 \mapsto C_{\text{int}} \} \circ S_2 \\ &= \{ \alpha_1 \mapsto C_{\rightarrow}(\alpha_0, \alpha_2), \\ &\quad \alpha_3 \mapsto C_{\rightarrow}(C_{\text{int}}, \alpha_2), \\ &\quad \alpha_4 \mapsto C_{\text{int}} \} \\ S_4 &= \{ \alpha_0 \mapsto C_{\rightarrow}(C_{\text{int}}, \alpha_2) \} \circ S_3 \\ &= \{ \alpha_1 \mapsto C_{\rightarrow}(C_{\rightarrow}(C_{\text{int}}, \alpha_2), \alpha_2), \\ &\quad \alpha_3 \mapsto C_{\rightarrow}(C_{\text{int}}, \alpha_2), \\ &\quad \alpha_4 \mapsto C_{\text{int}} \} \\ &\quad \alpha_0 \mapsto C_{\rightarrow}(C_{\text{int}}, \alpha_2)\end{aligned}$$

This is the same m.g.u. we derived before by ad-hoc substitution.

Note that the compositions performed in the calculation of S_1 and S_2 are trivial, whereas those in S_3 and S_4 are less so.

It is crucial to remember that the substitution argument to \mathcal{U} must be applied to both the term arguments **before** performing the case dispatch. This is illustrated in the calculation of the new substitution component of S_4 .

Polymorphism

Polymorphism means using the same code to operate on values of different types at different points in the same program. We are mainly interested in so-called **parametric polymorphism**, which refers to functions that are “oblivious” to the types of their arguments. An example is the function

$$\text{pair} \equiv \lambda x. (x, x)$$

which forms pairs of any type. Obviously, we need some conventions for uniform data representation in order to actually compile such functions.

We distinguish this from **ad-hoc polymorphism**, a.k.a. **overloading**, where the same source code function must have different runtime behavior depending on its argument’s types, e.g.,

$$\text{double} \equiv \lambda x. x + x$$

which has different behavior for $x:\text{real}$ and $x:\text{int}$.

We want a type system for polymorphism that is (efficiently) decidable, sound, and has an (efficient decidable) inference algorithm.

2nd-order Polymorphic λ -calculus

One way to express polymorphic programs of this kind is to use an **explicitly-typed** calculus which allows **abstraction** and **application** of type variables as well as ordinary variables.

Examples

$$(\lambda f:\forall\alpha.\alpha \rightarrow (\alpha \times \alpha). (f \text{ [int] } 3, f \text{ [bool] } \text{true}))$$

$$(\Lambda\alpha.\lambda x:\alpha. (x, x))$$

$$(\lambda a:\text{int}.\lambda b:\text{bool}.\lambda f:\forall\alpha.\alpha \rightarrow \alpha. (f \text{ [int] } a, f \text{ [bool] } b))$$

$$1 \text{ true } (\Lambda\alpha.\lambda x:\alpha.x)$$

In general, expressions are:

M	$:=$	c	Built-in constants
		v	Variable names
		$(M_1 M_2)$	Function applications
		$(\lambda v:\sigma.M)$	Function abstractions
		$(\Lambda\alpha.M)$	Type abstractions
		$(M[\sigma])$	Type applications

where polymorphic type schemes are:

$$\sigma := b \mid \alpha \mid \sigma_1 \rightarrow \sigma_2 \mid \forall\alpha.\sigma$$

Limitations of schemes

We already are part-way there through our use of type schemes.

When we derive a schematic type for `pair`, e.g.,

$$\emptyset \vdash \lambda x. (x, x) : \alpha \rightarrow (\alpha \times \alpha)$$

this implies that the application of `pair` to an expression of **any** type τ will type-check, with α being unified to τ .

Unfortunately, we have no way to express the idea that `pair` may be used with **different** types in different parts of the same expression, because this would require unifying α with two conflicting types.

For example, the following will not type-check:

$$(\lambda f. (f \ 3, f \ \text{true})) (\lambda x. (x, x))$$

since α would need to unify with both `int` and `bool`.

What we need is a way to **quantify** type variables over sub-expressions, so they can be **instantiated** to different types at different points in the expression.

ML Polymorphism

It is perfectly possible to give (sound) typing rules for the 2nd-order polymorphic calculus. We’ll investigate this in the next chapter. But meanwhile we want a type inference mechanism, so that the user doesn’t have to write down these complex types.

Unfortunately, this calculus has no principal types, and typability of untyped terms is actually undecidable.

So nearly all FL’s use a more restrictive style of polymorphism (called **ML-polymorphism**, **let-style polymorphism**, **Hindley-Milner polymorphism**, etc.) for which there is an inference algorithm.

The key idea is that we can do inference successfully if we have a visible **definition** for every polymorphic function.

Syntax

We use the following **untyped** syntax for expressions:

$M ::= c$	Built-in constants
v	Variable names
$(M_1 M_2)$	Function applications
$(\lambda v. M)$	Function abstractions
$(\text{let } v = M_1 \text{ in } M_2)$	let abstractions

Note that `let` is now a basic expression form in its own right. `let x = a in b` still has the same **operational** semantics as $(\lambda x. b) a$, but the two expressions will be **different** for typing purposes. Specifically, `let` is used to define expressions that are to be used **polymorphically**.

We would rewrite our first earlier example as:

```
let f = λx.(x,x) in (f 3, f true)
```

The use of `let` is essential; we **can't** write the body here as an ordinary λ -abstraction

```
λf.(f 3, f true)
```

as this will still not type check.

We cannot express the second example at all, however; this style of polymorphism is strictly less powerful.

ML Polymorphism Type System

We can give rules similar to the ones we introduced earlier for the simply-typed λ -calculus.

In these rules, we use ordinary type schemes just as before

$$\tau ::= b \mid \alpha \mid \tau_1 \rightarrow \tau_2 \mid \tau_1 \times \tau_2 \mid \dots \quad (b \in B)$$

and also **polymorphic type schemes** of the following form

$$\sigma ::= \tau \mid \forall \alpha. \sigma$$

Note that in these polymorphic schemes any quantifiers always come **first**. For example,

$$\forall \alpha. (\alpha \rightarrow (\alpha \times \alpha))$$

is a legitimate polymorphic scheme;

$$(\forall \alpha. \alpha) \rightarrow (\forall \alpha. \alpha)$$

is not.

Top-level Expressions

Note that expressions entered at the “top-level” in CAML are treated like `let`-bound expressions; that is, for typing purposes the sequence

```
let a = exp1;;
let b = exp2;;
...;;
```

is equivalent to

```
let a = exp1 in let b = exp2 in ...;;
```

In particular, library functions are treated as if they were `let`-bound in the scope of ordinary programs.

Instantiation

To get an (ordinary) type scheme from a polymorphic scheme $\sigma \equiv \forall \alpha_1 \forall \alpha_2 \dots \forall \alpha_n. \tau$ we instantiate each of the quantified type variables with a type. The resulting scheme will be of the form $\tau' \equiv S(\tau)$, where S is substitution with domain $\{\alpha_1, \dots, \alpha_n\}$. In this case we write

$$\sigma > \tau'$$

Examples

$$\begin{aligned} \forall \alpha. \forall \beta. \alpha \rightarrow \gamma \rightarrow \beta > \text{int} \rightarrow \gamma \rightarrow \text{bool} \\ \text{with } S = \{\alpha \mapsto \text{int}, \beta \mapsto \text{bool}\} \\ \forall \alpha. \forall \beta. \alpha \rightarrow \gamma \rightarrow \beta > \text{int} \rightarrow \gamma \rightarrow (\delta \rightarrow \delta) \\ \text{with } S = \{\alpha \mapsto \text{int}, \beta \mapsto (\delta \rightarrow \delta)\} \\ \forall \alpha. \alpha \rightarrow \beta \rightarrow \alpha \not> \text{int} \rightarrow \text{bool} \rightarrow \text{int} \end{aligned}$$

Efficiency of Algorithm W

Interestingly, this is an example of an algorithm that is efficient in practice but very inefficient (deterministic exponential time complete) in theory. For example, the following nasty example

```
let pair = λx.λy.λz.zxy
  in let x1 = λy.pair y y
    in let x2 = λy.x1(x1(y))
      in ...
        in let xn = λy.xn-1(xn-1(y))
          in xn(λz.z)
```

produces a principal type of length $2^{2^{cn}}$ (as a string). For $n = 5$ this produces 173 printed pages of output after several hours.

Side-effects

We have to be very careful in extending these typing rules to programs with side-effects. For example, we might naively add the following polymorphic type schemes for ML-style references to `TypeOf`:

$$\left\{ \begin{array}{l} \text{ref} : \forall \alpha. \alpha \rightarrow \text{aref}, \\ := : \forall \alpha. \alpha \text{ ref} \rightarrow \alpha \rightarrow \text{unit}, \\ ! : \forall \alpha. \alpha \text{ ref} \rightarrow \alpha \end{array} \right\}$$

where `ref` is a new type constructor.

Unfortunately, under these rules the following expression is typable, though it surely shouldn't be!

```
let r = ref(λx.x) in
  (r := λx.x+1; (!r) true)
```

Intuitively, the problem is that we infer the type $(\alpha \rightarrow \alpha)$ `ref` for `r`'s defining expression, and then **abstract** over α to form a polymorphic scheme for `r`. This abstraction is bogus, since α really refers to the particular type of the value currently held by the reference.

Similar problems arise with other imperative features, such as exceptions. Various solutions to these problems have been proposed and implemented. The cleanest is to put a restriction on the form of `let` definitions, namely the defining expression must be a λ -abstraction, a constant, or a simple variable. In all these cases, the **cbv evaluation** of the expression is essentially a no-op, which turns out to avoid the typing problems.

Recursion

As before, we can't define the Y combinator directly in an ML-polymorphism type system, so we need to add explicit support for recursion. Here's a possible definition for a LETREC rule:

$$\frac{E + \{f : \tau_2 \rightarrow \tau_1, x : \tau_2\} \vdash M : \tau_1 \quad E + \{f : \text{Clos}(\tau_2 \rightarrow \tau_1, E)\} \vdash N : \tau}{E \vdash \text{letrec } f = \lambda x. M \text{ in } N : \tau}$$

Example

```
letrec map = λf.λx.if null x
  then nil
  else cons(f (hd x),map f (tail x))
in map negate [1,2,3]; map not [true,false]
```

will typecheck correctly as expected, with `map` being given a polymorphic type scheme

$$\forall \alpha. \forall \beta. (\alpha \rightarrow \beta) \rightarrow \alpha \text{ list} \rightarrow \beta \text{ list}.$$

Note that this typing rule doesn't allow us to define a recursive function that invokes itself in a polymorphic way:

```
letrec f = λn.λx.if n = 0
  then x
  else f (n-1) (x,x)
```

Unfortunately, if we extend the typing system to permit expressions of this kind to be typable, type inference becomes undecidable!