

## Lecture 5: Type Inference

The main point of type annotations is to document your code and catch programming errors. However, it can be quite annoying to have to write down every single type in a program. If you've used Java Generics or something similar, you've probably encountered the annoyance of type annotations in real code. The type-checker you wrote earlier also has the same property.

The type inference problem is to take a program with no type annotations, such as this version of factorial:

```
fix f.λn.if n = 0 then 1 else n × f (n - 1)
```

and produce a program with type annotations:

```
fix f:ℕ→ℕ.λn:ℕ.if n = 0 then 1 else n × f (n - 1)
```

We know how to type-check the latter.

It is easy to code up with heuristics for inferring types. However, we're going to cover an approach known as the *Hindley-Milner algorithm*, which has two key properties:

- It is *sound*, so it always produces correct types. Therefore, you can trust the output of the algorithm and there is no need to type-check after inference (assuming you implement it correctly!)
- It is *complete*, so it always produces a type, if there is a type to be found. Therefore, if the algorithm fails to find a type, then the program is truly un-typable.

The algorithm has one other important property, which we'll look at later.

**An Informal Example** Suppose we were given the untyped factorial function and were asked to infer its type in our head. We'll run through a detailed analysis below. For brevity, we'll use the following notation:

- We write  $\llbracket e \rrbracket$  to mean “the type of the expression  $e$ ”.
- We use Greek letters to denote metavariables that stand for types that are unknown.

We might reason through the types as follows:

	Claim	Justification
1	$\llbracket \text{fix } f.\lambda n.\dots \rrbracket = \llbracket \lambda n.\dots \rrbracket$	<code>fix</code> evaluates to the body.
2	$\llbracket f \rrbracket = \llbracket \text{fix } f.\lambda n.\dots \rrbracket$	$f$ is bound to the whole expression during evaluation.
3	$\llbracket \lambda n.\dots \rrbracket = \alpha \rightarrow \beta$	The expression is a function.
4	$\llbracket n \rrbracket = \alpha$	$n$ is the argument of a function with type $\alpha \rightarrow \beta$ .
5	$\llbracket \text{if } n = 0 \text{ then } 1 \text{ else } n \times f(n - 1) \rrbracket = \beta$	The expression is the body of a function with type $\alpha \rightarrow \beta$ .
6	$\llbracket n = 0 \rrbracket = \text{Bool}$	The expression is used as the condition in an if-expression.
7	$\llbracket n \rrbracket = \text{Int}$	$n$ is used in a comparison. (We can only compare integers.)
8	$\llbracket 1 \rrbracket = \llbracket n \times f(n - 1) \rrbracket$	Both branches must have the same type.
9	$\llbracket 1 \rrbracket = \llbracket \text{if } n = 0 \text{ then } 1 \text{ else } n \times f(n - 1) \rrbracket$	Type of <code>if</code> is the type of either branch.
10	$\llbracket f \rrbracket = \gamma \rightarrow \delta$	$f$ appears in function position in an application.
11	$\llbracket \gamma \rrbracket = \llbracket (n - 1) \rrbracket$	The argument to $f$ , which has type $\gamma$ .
12	$\llbracket (n - 1) \rrbracket = \text{Int}$	Subtraction produces integers.
13	$\llbracket f(n - 1) \rrbracket = \delta$	The return type of $f$ is the type of the application.
14	$\llbracket \delta = \text{Int} \rrbracket$	The expression of type $\delta$ is used as an argument of $\times$ .

At this point, we can solve the constraints to figure out the types of the metavariables  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$ :

- $\alpha = \text{Int}$  by (4) and (7),
- $\beta = \text{Int}$  by (5), (9) (given that  $\llbracket 1 \rrbracket = \text{Int}$ ),
- $\gamma = \text{Int}$  by (11) and (12),
- $\delta = \text{Int}$  by (14).

Now, we have enough information to reconstruct the type annotations in the program.

## 1 Introduction

Type inference works with two syntaxes, shown in fig. 7.1. The *implicitly-typed syntax*, on the left-hand side, is the syntax in which the user writes the program. The *explicitly-typed syntax*, on the right-hand side, is the result of type-inference. The difference between these two syntaxes is that the explicitly-typed syntax has type annotations that can be used to type-check the program. In contrast, the implicitly-typed syntax has no type annotations. However, it is important to note that **the implicitly-typed language is typed**. The types aren't written down, but we could still write typing rules for this language.<sup>1</sup> The explicitly-typed syntax just makes types manifest, so that a simple type-checker can type-check the program.

Types in the explicitly-typed language include metavariables,  $\alpha$ ,  $\beta$ , etc. Metavariables are not types themselves, but are placeholders that stand for types. We need metavariables to describe type inference. However, the result of type inference will produce a program that has no metavariables.

Type inference has several steps:

<sup>1</sup>We would have to “guess” the type of binding identifiers when building typing derivations.

**Binary Operators** $op_2 ::= + \mid > \mid \dots$ **Constants**

$c ::= \text{true}$	True
$\text{false}$	False
$n$	Integers

**Expressions**

$e ::= c$	Constants
$x$	Identifiers
$op_2(e_1, e_2)$	Bin. Ops.
$\text{if } e_1 \text{ then } e_2 \text{ else } e_3$	Conditionals
$e_1 e_2$	Applications
$\lambda x. e$	Functions
$\text{fix } x. e$	Fixpoints

(a) Implicitly-Typed Syntax.

**Types**

$T ::= \text{Int}$	Integer type
$\text{Bool}$	Boolean type
$T_1 \rightarrow T_2$	Function types
$\alpha$	Type metavariables

**Binary Operators** $op_2 ::= + \mid > \mid \dots$ **Constants**

$c ::= \text{true}$	True
$\text{false}$	False
$n$	Integers

**Expressions**

$e ::= c$	Constants
$x$	Identifiers
$op_2(e_1, e_2)$	Bin. Ops.
$\text{if } e_1 \text{ then } e_2 \text{ else } e_3$	Conditionals
$e_1 e_2$	Applications
$\lambda x : T. e$	Functions
$\text{fix } x : T. e$	Fixpoints

(b) Explicitly-Typed Syntax.

Figure 7.1: Syntaxes for type inference.

1. We transform a program in the implicitly-typed syntax to a identical program in the explicitly-typed syntax, using unique metavariables for each type annotation. For example, the program  $\lambda x. \lambda y. x + y$  would be transformed to  $\lambda x : \alpha. \lambda y : \beta. x + y$ . These metavariables will allow us to refer to identifier's types without getting scope mixed up.
2. We generate a set of constraints by recursively processing the program. Each constraint equates two types to each other. To handle variables and scoping correctly, we use an environment that makes identifiers to metavariables. For example, the program  $\lambda x : \alpha. \lambda y : \beta. x + y$  would produce the constraints  $\{\alpha = \text{Int}, \beta = \text{Int}\}$ .
3. We solve the constraints, using a classic algorithm called *unification*. Constraint-solving produces a *substitution* from metavariables to types that have no metavariables within them. For example, solving the constraints  $\{\alpha = \text{Int}, \beta = \text{Int}\}$  would produce the substitution  $[\alpha \mapsto \text{Int}, \beta \mapsto \text{Int}]$ . Several things can go wrong during constraint-solving. For example, we may derive an unsatisfiable constraint, such as  $\text{Int} = \text{Bool}$ , which indicates that the program has a type error.
4. We annotate the program by applying the substitution to the type variables we generated in the first step. For example, the program  $\lambda x : \alpha. \lambda y : \beta. x + y$  and the substitution  $[\alpha \mapsto \text{Int}, \beta \mapsto \text{Int}]$  would produce the program  $\lambda x : \text{Int}. \lambda y : \text{Int}. x + y$ .
5. Finally, we can type-check the resulting program to verify the result of type-inference. However, if the steps above are perfectly correct, there is no need to verify the result.

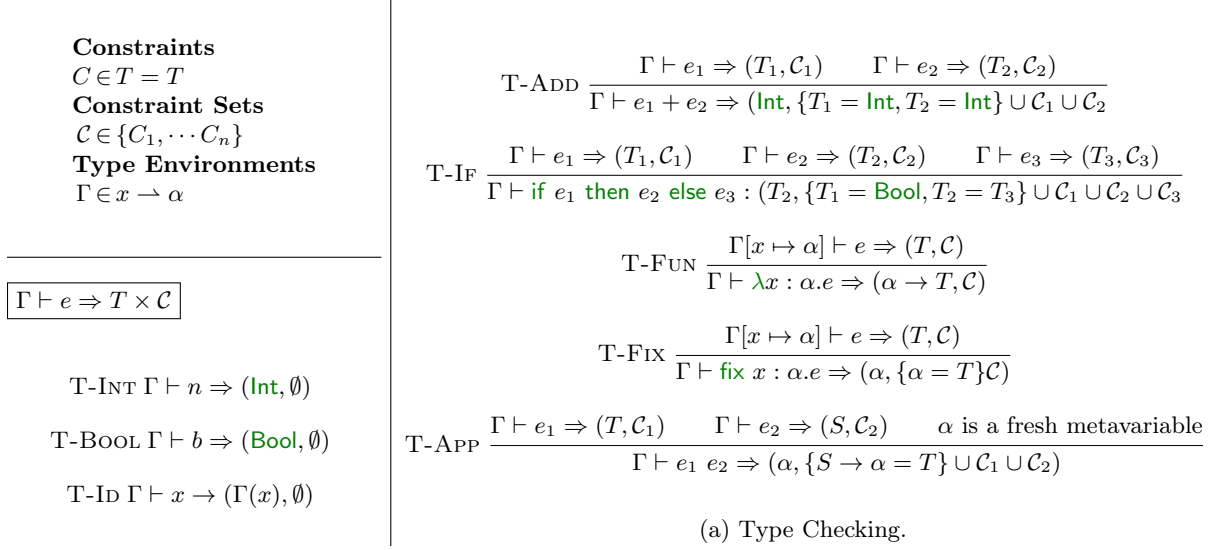


Figure 7.2: Constraint Generation

## 2 Constraint Generation

Figure 7.2 gives the constraint generation rules for our language. We can read the rule  $\Gamma \vdash e \Rightarrow T \times \mathcal{C}$  as *in the environment  $\Gamma$ ,  $e$  has type  $T$  and produces the set of constraints  $\mathcal{C}$* . There are two key differences between the type-inference rules and the type-checking rules:

- Unlike the type-checking rules, the type  $T$  may not be a complete type and may include metavariables ( $\alpha$ ,  $\beta$ , etc.).
- Constraint-generation does not catch type-errors. For example, the expression `true + 10` will fail to type-check. However, the constraint generation rules will produce the type `Int` and the set of constraints  $\{\mathbf{Bool} = \mathbf{Int}, \mathbf{Int} = \mathbf{Int}\}$ .

However, constraint-generation does catch unbound-identifier errors. If a program has a free variable, then it won't be bound to a metavariable in the environment.

Notice that the types produced by constraint-generation are used to further generate additional constraints. For example, in the expression  `$\lambda x : \alpha. x + 3$` , the bound identifier  $x$  produces the type  $\alpha$ , which is used to produce the constraint  $\alpha = \mathbf{Int}$  when generating constraints for  $x + 3$ .

## 3 Solving Type Constraints

After constraint-generation completes, we can solve constraints by *unification*. The UNIFY algorithm, shown in fig. 7.3, takes a single constraint and produces as output a *substitution*, which is a finite map from metavariables to types. A substitution can be applied to a type to replace metavariables with concrete types.

For example,  $\text{UNIFY}(\alpha \rightarrow \mathbf{Int} = \mathbf{Bool} \rightarrow \beta)$  produces the substitution  $[\alpha \mapsto \mathbf{Int}, \beta \mapsto \mathbf{Bool}]$ . If we apply this substitution to  $\alpha \rightarrow \mathbf{Int}$  or to  $\mathbf{Bool} \rightarrow \beta$ , we get the type  $\mathbf{Bool} \rightarrow \mathbf{Int}$ . Notice

### Substitutions

$$\begin{aligned} \pi &\in \alpha \rightarrow T \\ (\pi_1 \cdot \pi_2)(\alpha) &= \pi_2(\pi_1(\alpha)) \end{aligned}$$

### Unification

$$\begin{aligned} \text{UNIFY} &\in C \rightarrow \pi \\ \text{UNIFY}(T = T) &= \cdot \\ \text{UNIFY}(\alpha = T) &= [\alpha \mapsto T] \text{ if } \alpha \text{ does not occur in } T \\ \text{UNIFY}(T = \alpha) &= [\alpha \mapsto T] \text{ if } \alpha \text{ does not occur in } T \\ \text{UNIFY}(S_1 \rightarrow S_2 = T_1 \rightarrow T_2) &= \text{UNIFY}(\pi(S_2), \pi(T_2)) \cdot \pi \\ &\text{ where } \pi = \text{UNIFY}(S_1, T_1) \end{aligned}$$

Figure 7.3: Unification

that we get the same type on either side, which is a key property of the algorithm. However,  $\text{UNIFY}(\gamma \rightarrow \text{Int} = \text{Bool} \rightarrow \gamma)$  should fail, since there is no substitution that can be applied to either side to produce the same type.

UNIFY is a simple recursive algorithm with three cases:

- $\text{UNIFY}(\text{Int} = \text{Int})$  should produce the empty substitution, whereas  $\text{UNIFY}(\text{Int} = \text{Bool})$  should fail. In general, unifying two base types should succeed with the empty substitution if the types are the same or fail if the types are different.
- $\text{UNIFY}(\alpha = T)$  and  $\text{UNIFY}(T = \alpha)$  should produce the substitution  $[\alpha \mapsto T]$ . However, there is an important caveat discussed below, called the *occurs check*.
- $\text{UNIFY}(S_1 \rightarrow S_2 \mapsto T_1 \rightarrow T_2)$  needs to recursively unify the argument and the result types. However, the two substitutions have to be *composed* together.

**The Occurs Check** Consider the constraint  $\alpha = \text{Int} \rightarrow \alpha$ . If unification is done naively, we will produce the substitution  $[\alpha \mapsto \text{Int} \rightarrow \alpha]$ . If we apply this substitution to either side, we get the constraint  $\text{Int} \rightarrow \alpha = \text{Int} \rightarrow (\text{Int} \rightarrow \alpha)$ . We can unify this new constraint to get a larger type on either side. In fact, we can repeatedly apply unification to get larger and larger types. The real problem is that the original constraint is circular and contradictory. There is no substitution  $\pi$  such that  $\pi(\alpha) = \pi(\text{Int} \rightarrow \alpha)$ . To avoid this infinite regress, unification needs an *occurs check* to detect cyclic constraints. It is enough to run the occurs check when unifying metavariables with types, as shown in fig. 7.3.

**Unifying a set of constraints** Given UNIFY, it is straightforward to unify a set of constraints. After unifying each constraint in the set, we need to apply the intermediate substitution to the remaining constraints.

