Programming Languages Types

Dr Russ Ross

Utah Tech University—Department of Computing

Fall 2023

We will use the term type to refer to a static check, i.e., one that can be done purely with the program source. This means:

- types cannot refer to dynamic conditions
- they may suffer from either false-positive or false-negative errors

In return:

- they give us guarantees without ever having to run the program
- every part of the program can be checked, not just parts that have test coverage

Like objects, types are not well standardized. Many languages do not have them and those that do often disagree on their form. However, there are parts with widespread agreement.

Types are an abstraction of run-time values. We can have many distinct numbers, strings, etc., but when considering types we are concerned with the distinctions between categories and not the distinctions within them.

Let's start with a basic interpreter:

```
(define-type BinOp
 [plus])
(define-type Expr
 [binE (operator : BinOp)
               (left : Exp)
                    (right : Exp)]
 [numE (value : Number)])
```

```
(calc : (Exp -> Number))
(define (calc e)
  (type-case Exp e
    [(binE o l r)
    (type-case BinOp o
        [(plus) (+ (calc l) (calc r))])]
    [(numE v) v]))
```

(test (calc (binE (plus) (numE 5) (numE 6))) 11)

What needs to happen with a type-checker? The name "checker" suggests that the job of a type-checker is to *pass judgement* on a programs, i.e., to determine whether or not they are type-correct. Thus a natural type for the checker would be:

```
(tc : (Exp -> Boolean))
```

(In practice, of course, we would want more information in case the program is not type-correct, i.e., we'd like an error diagnostic. But we're ignoring human factors considerations here.) With this type, we can now rewrite the relevant parts of the interpreter above:

```
(define (tc e)
 (type-case Exp e
  [(binE o l r)
  (type-case BinOp o
      [(plus) (and (tc l) (tc r))])]
  [(numE v) #true]))
```

(test (tc (binE (plus) (numE 5) (numE 6))) #true)

Look at this for a moment. Given a number, it returns #true. In the recursive case, it checks the pieces and ands the result. There is no way to return #false, so every program is always type-correct.

The problem is that we only have one type (numbers) and one operation on numbers (plus) so there is nothing that *could* go wrong.

We need more types and operations:

(define-type BinOp
 [plus] [++])

```
(define-type Expr
[binE (operator : BinOp)
                          (left : Exp)
                               (right : Exp)]
[numE (value : Number)]
[strE (value : String)])
```

Various things break, and need to be fixed. How about this?

```
(define (tc e)
  (type-case Exp e
    [(binE o l r)
    (type-case BinOp o
        [(plus) (and (tc l) (tc r))]
        [(++) (and (tc l) (tc r))])]
    [(numE v) #true]
    [(strE v) #true]))
```

So this looks pretty good, right?

Here are two tests that demonstrate *desirable* behavior:

```
(test (tc (binE (++) (numE 5) (numE 6))) #false)
(test (tc (binE (plus) (strE "hello") (strE "world"))) #false)
```

The first string-concatenates two numbers, the second adds two strings. Therefore, both should be rejected by the type-checker. Yet both of them pass (i.e., the tests above fail).

What is the core problem here? Given an expression, we only know

- whether its sub-expressions typed correctly,
- but not *what* their types are.

This is insufficient to determine whether the current expression is type-correct. For instance:

- the ++ operator needs to check not only whether its two sub-expressions are well-typed, but also
- whether they produce strings;
- if they did not, then the concatenation is erroneous.

We need the type-checker to have a richer type: it must instead be

(tc : (Exp -> Type))

That is, the type "checker" must actually be a type *calculator*.

- it closely parallels the evaluator
- but over the universe of abstracted values (types)
- rather than concrete ones

Following convention, however, we will continue to call it a checker, because it *also* checks in the process of calculating types.

Type is a new (plait type) definition that records the possible types:

```
(define-type Type [numT] [strT])
```

With this, we can rewrite our type-"checker":

Now it passes some useful tests:

```
(test (tc (binE (plus) (numE 5) (numE 6))) (numT))
(test (tc (binE (++) (strE "hello") (strE "world"))) (strT))
```

```
(test/exn (tc (binE (++) (numE 5) (numE 6))) "strings")
(test/exn (tc (binE (plus) (strE "hello") (strE "world"))) "numbers")
```

There are three take-aways from this:

- 1. The type-checker follows the same implementation schema as the interpreter: an algebraic datatype to represent the AST, and structural recursion to process it. This is the schema we're calling SImPl.
- 2. A type-checker, unlike an interpreter, operates with "weak" values: note, for instance, how the numE case ignores the actual numeric values. Both the strengths and weaknesses of traditional type-checking arise from this ignorance.
- 3. In mathematical terms, the upgrade we performed in going from a type-checker to a type-calculator was a process of strengthening the inductive hypothesis: instead of returning only a Boolean, we had to return the actual type of each expression. This may not seem like a literal strengthening; but it is inasmuch as the former #true has been replaced by a Type and the #false by an error.

As we extend our type system, it is increasingly unwieldy to write everything out as code. Instead, we will adopt a notation commonly used in the world of types (though it can also be used for interpreters and other SImPI programs). We will write terms of the form

|- e : T

where the e are expressions, T are types, and : is pronounced as "has type": i.e., the notation above says "e has type T". For now we won't pronounce |- as anything at all; later, we will see that it should be read as "proves".

A Concise Notation

First, we can very concisely say that all numeric expressions have numeric type and all string expressions have string type:

|- n : Num |- s : Str

where n stands for all the syntactic terms with the syntax of numbers, and s likewise for strings.

(We can think of this as an infinite number of rules, one for each number and each string. We're in the realm of mathematics, so what's an infinite number of rules between friends?) The former is exactly equivalent to writing

[(numC n) (numT)]

but much more concisely.

A Concise Notation

When we get to Booleans, we have a choice: we can either write

|- b : Bool

where b stands for all the syntactic terms with the syntax of Booleans, or—because there are only two of them—just enumerate them explicitly:

```
|- true : Bool
|- false : Bool
```

That covers the base cases of the type-checker. These are called *axioms*. Next we will look at the *conditional* cases, which are called *(typing) rules*.

A Concise Notation

Remember our code for typing addition:

We can write it in this notation very concisely as follows:

|- e1 : Num |- e2 : Num |- (+ e1 e2) : Num

We read the line as "if (what's above) then (what's below)", and the space as "and".

So this says: "if e1 has type Num and e2 has type Num, then (+ e1 e2) has type Num". This is of course the exact same thing the code says, but with rather less noise.

Terminology: The part above is called the antecedent (that which goes before) and the part below is called the consequent (that which comes after). Don't call these the numerator and denominator!

Handling Division

Addition, multiplication, and subtraction are *total* functions over numbers: they consume two numbers and produce one. In contrast, division is a *partial* function: it isn't defined when the denominator is zero.

There are a few strategies for handling this:

- 1. Declare that division doesn't return a number but instead something else that captures its partiality.
 - Works, but every single use of division will need to check which kind of result it got
- 2. Declare that division only consumes non-zero numbers in its second argument.
 - A major change to the type system, because now numbers are not all a single type
 - Affects all callers of division, who must now prove that they are not dividing by zero. The type checker cannot prove this on its own in all cases (see Rice's Theorem)
- 3. Give it the same type as other binary numeric operations, but declare that division by zero will be handled by an exception or error
 - This puts the burden on the rest of the program, which must be away of this possibility and handle it

Most languages take the third option (exceptions/panics), but some are exploring the other two.

- They get around Rice's Theorem by proving non-zero-ness when they can, and putting the burden of proof on the programmer the rest of the time
- This creates more effort for the programmer, but increases the program's robustness

Another Perspective on Types

We have seen one way to think about types:

- An abstraction of values
- Type-checking as running a program over these abstract values

Another perspective is to think of types:

- As a static discipline
- A way of making judgements about programs

This is kind of like parsing:

- A parser statically (before the program runs)
- passes judgement (decides some programs are good and some bad)

Types can be viewed as an extension of this idea.

Aside: In computability theory terms, parsers are usually context-free, whereas types usually reflect context-sensitive constraints. Computability theory then helps us understand why we might separate these checks into two separate phases, and in particular why we might do one before the other. Essentially, the type-checker only needs to deal with programs that have already passed the parsing, i.e., context-free check, so it has much less complexity than if it had to do everything. We already saw this: our previous checker only consumed Exprs, which are produced by the parser.

From Axioms and Rules to Judgements

When we need to apply type rules to a program, we compose them recursively, just as the type-checker runs. Consider this program:

(+ 5 (+ 6 7))

To decide its type, we will use our current rules:

- It does not fit any axiom, because the program does not match the syntax of a single number or string.
- We have to use a conditional rule.

We have seen only one so far, and fortunately this term does match the consequent:

- it requires two terms, and we have two terms, so e1 is 5 and e2 is (+ 6 7)
- Therefore, applying this conditional rule, we get:

```
|- 5 : Num |- (+ 6 7) : Num
|- (+ 5 (+ 6 7)) : Num
```

From Axioms and Rules to Judgements

Let's look at the two terms in the antecedent. The first matches an axiom, so we are finished with it:

|-5:Num |-(+67):Num |-(+5(+67)):Num

For the other, we have to apply the same conditional rule again:

From Axioms and Rules to Judgements

These new terms also match the axiom for numbers:

|- 6 : Num |- 7 : Num ______ |- 5 : Num |- (+ 6 7) : Num ______ |- (+ 5 (+ 6 7)) : Num

Every part of the tree now terminates in an axiom, therefore:

- We consider this program to have successfully type-checked
- This tree is called a judegement, because it passes judgement on the initial term
- It is judged to have type-checked and to produce a value of type Num

This is the same pattern of execution we had with the type-checker program

- We were able to skip the tedious details of passing and returning things
- We just used pattern-matching
- This will save work going forward

Judgements and Errors

Let's see another example:

(+ 5 (+ 6 "hi"))

This starts out similar to the previous example. We attempt a judgement:

|- 6 : Num |- "hi" : Num ______ |- 5 : Num |- (+ 6 "hi") : Num ______ |- (+ 5 (+ 6 "hi")) : Num

But now we have a problem: we need to type-check

|- "hi" : Num

but we don't have a rule that matches. Therefore we cannot construct a successful tree.

Judgements and Errors

We cannot construct a successful tree:

|- 6 : Num |- "hi" : Num ___________ |- 5 : Num |- (+ 6 "hi") : ________ |- (+ 5 (+ 6 "hi")) :

Remember the "if ... and ... then" interpretation. Because we cannot satisfy all the antecedents, we cannot prove anything about the consequents, leaving the tree incomplete.

A type error is simply a failure to construct a judgment.

• It may not be the most satisfying user feedback, but our concern here is with a concise way of expressing ideas;

• going from this to an implementation is not too hard, and the user interface details can be added to the latter.

Judgements and Errors

This requires some clarification.

- We only call it a judgment if the tree is "checked off" completely:
 - every antecedent is generated using given rules, and
 - all the leaves are actual axioms.
- In this example, we are unable to check off the tree:
 - there is no available rule or axiom that lets us conclude that "hi" is a Num.
 - Therefore, we cannot "judge" the initial expression.

This is a technical meaning of the word "judgment", not to be confused with the broader use of this term. Imagine that we started with this program:

(+ 5 (- 6 7))

We would get this far:

|- 5 : Num |- (- 6 7) : Num |- (+ 5 (- 6 7)) : Num

Again we would fail, this time because we haven't provided a (conditional) rule for (- e1 e2). Obviously it's not difficult to define one; we just haven't done so yet, so our pattern-matcher would fail.

Now we're ready to add a rule for if. Languages have different rules for what can go in the conditional clause.

- The goal of a type-checker is to catch type errors, so it is common for languages with type-checkers to demand that the conditional be a Boolean (without a truthy/falsy set of Boolean values).
- Our goal here is not to make a value judgment but rather to illustrate how we would add a type rule for it.

By now, we can see that we will need a conditional rule (because we want to type-check more than just constants); following SImPI, and we will need the antecedent to say something about the sub-expressions. We need at least:

|- C : Bool ...

|- (if C T E) : ...

Okay, what now? What is the type of the entire conditional expression?

- Technically, it should be whatever type is returned by the branch that was executed.
- However, a type-checker can't know which branch will be executed; over time, both might.
- We have to somehow capture the uncertainty in this situation.

There are two common solutions:

- 1. Introduce a new kind of type that stands for "this type or that type" (a *union*). This is easy to introduce but creates a burden for every piece of code that will consume such a value.
- 2. Just rule that both branches should have the same type. The latter is a very elegant solution, because it eliminates the uncertainty entirely.

Okay, so we need to do the following things:

- Compute the type of T.
- Compute the type of E.
- Make sure T and E have the same type.
- Make this (same) type the result of the conditional.

That seems like a lot: how will we express all that? Very easily, actually:

|- C : Bool |- T : U |- E : U |- (if C T E) : U

Here, U is a placeholder:

- it isn't a concrete type but rather stands for whatever type might go in that place.
- The repeated use of U accomplishes all of our goals above.

```
|- C : Bool |- T : U |- E : U
|- (if C T E) : U
```

Read this as:

- if C has type Bool and
- T has type U and
- E has [the same] type U,
- then (if C T E) has [the same] type U

Let's see this in action on the following program:

(if true 1 2)

We get:

Either of the axioms for the other two antecedents tells us what U must be, which lets us fill in the result of U everywhere:

|- true : Bool |- 1 : Num |- 2 : Num |- (if true 1 2) : Num

Fortunately, the other two antecedents are also axioms:

|- true : Bool |- 1 : Num |- 2 : Num |- (if true 1 2) : Num

This lets us conclude that the overall term is well-typed, and that it has type Num.

Now let's look at:

(if 4 1 2)

Applying the conditional rule gives us:

However, we do not have any axiom or conditional rule that lets us conclude that 4 has type Bool (because, in fact, it does not). Therefore, we cannot complete the judgment and the program is (rightly) judged to have a type error.

One last example:

(if true 1 "hi")

Again, applying the conditional rule and checking off the first antecedent:

But now we have a problem. If we apply the axiom for numbers, we replace all instances of U with Num to get:

Maybe we just tried the wrong axiom? We do have one more option!

However, it ends up with the same net effect:

|- true : Bool |- 1 : Str |- "hi" : Str -------|- (if true 1 "hi") : Str

Because there is no way to construct a judgment for this program, it too has a type error.

Exercise: Let's now add functions. We need two new constructs: one to introduce them (lambda) and one to use them (function application). Write down judgments for each.

Where Types Diverge from Evaluation

Something very important, and subtle, happened above.

Compare the type rule for a conditional with the evaluation process. If the rule is too abstract, just look at the example judgments (or failed judgments) above.

- The evaluator evaluates only one branch out of T and E; indeed, that is the entire point of a conditional.
- The type-checker, in contrast, traverses both branches! In other words, it looks at code that might evaluate, not only code that absolutely does evaluate.

In other words, the idea that a type-checker is like an "evaluator that runs over simple values" is a convenient starting analogy, but it is in fact false.

- An evaluator and type-checker follow different traversal strategies.
- That is why a program like (if true 1 "hi") might run without any difficulty but is rejected by a type-checker.

Where Types Diverge from Evaluation

While this particular example may make the type-checker look overly pedantic, what if the same program were

```
(if (is-full-moon) 1 "hi")
```

Should the type-checker pass the program every month? Should it consider the moon's phase at the time of type-checking or at execution? Unfortunately, the type-checker doesn't know when the program will run; indeed, the program is type-checked once but may run an arbitrary number of times. Therefore, a type-checker must necessarily be *conservative*.

Where Types Diverge from Evaluation

This also lets us relate type-checking to testing.

- In software testing, making sure that all branches are visited is called *branch coverage*, and making sure all branches have coverage is both important and very difficult (because each branch may have additional branches which in turn may have even more branches which...).
- In contrast, a type-checker effortlessly covers both branches. The trade-off is that it does so only at the *type* level (and indeed, the abstraction of values to types is precisely what enables it to do this).

Thus, testing and type-checking are complementary.

- Type-checking provides code coverage at a lightweight level;
- testing typically provides only partial coverage but at the deep level of specific values.

In recent years, people have invented a notion of *concolic*—i.e., "concrete" + "symbolic"—testing to try to create the best of both worlds.

Growing Types: Typing Functions

Let us grow our language further to include functions. Concepts like functions tend to come in pairs:

- a way to introduce them, a lambda form in our case
- a way to use ("eliminate") them, namely function application

We use syntactic sugar over functions to obtain forms like let, so typing functions covers those cases as well.

Typing Function Applications

A function application has two parts: the function and the arguments. We will work with single-argument functions.

Because functions are first-class values, the function position is itself an expression. We have to check each sub-expression before we can type the whole expression. Therefore, function applications are conditional rules with two terms in the antecedent:

|- F : ??? |- A : ??? |- (F A) : ???
Typing Function Applications

First, let's notice that functions are different kinds of values than other values:

- a function is not itself a number, or string, or Boolean
- it may produce one of those, but it is not itself one of those (an important distinction).
- Therefore, we need a different type for functions, which reflects what functions consume and what they produce.

A natural idea is to assume functions have some "function" type, here called Fun:

What do we know about the argument expression (the actual parameter)?

- It had better match the type demanded by the formal parameter.
- But how do we check that here?

We've collapsed all functions in the world into a single type, Fun. That's far too coarse.

Typing Function Applications

Instead, following convention, we'll use the "arrow" syntax for functions:

```
|- F : (??? → ???) |- A : ???
|- (F A) : ???
```

(Technically, the arrow is a constructor of function types. It's a two-place constructor, for reasons we will see below.)

With this, we can now say that the function's formal parameter's type had better match up with the type of the actual argument.

Which type, exactly? Functions could consume numbers, strings, even other functions...all we know is that these should be consistent.

This is very similar to the consistency we expected of the branches of a conditional. We can again encode this by using the same placeholder in both places:

```
|- F : (T -> ???) |- A : T
|- (F A) : ???
```

Typing Function Applications

Now, what about what the function returns? Again, it could return values of any type. Whatever that type is, that is what the entire application produces. Again, we use a common placeholder to reflect this:

So here's how we read this:

- Type-check the F position. Make sure it's a function type (->). Assuming it is, call the formal parameter's type T and the return type U.
- Type-check the actual parameter (the argument). Make sure it has the same type as what the function is expecting in its formal parameter.
- If both of those hold, then the function's return type is the type of the entire application.

This list of steps is what a conventional type-checker would implement. Observe that again, a type error is the result of a failure to construct a judgment. If, for instance, the actual argument's type doesn't match that of the formal parameter, then the conditional rule above doesn't apply (it applies only when we can write the same type for the T placeholder), which is how we learn that the program has a type error.

Typing Function Definitions

Now we're ready to type lambda. Here, we have to be careful about how many sub-expressions there are.

- Given (lambda V B), it is tempting to think that there are two:
 - 1. V (the formal parameter) and
 - 2. B (the body).
- This is wrong! The formal parameter is a literal name, not an expression: we can't replace that name with some larger expression, which is what it would mean for it to be an expression.
- Furthermore, we can't evaluate it: it would (most likely) produce an unbound variable error, because its whole job is to bind that variable, so it can't assume it has already been bound.

Therefore, there is only one sub-expression, the body.

Typing Function Definitions

Therefore, we expect to end up with a conditional rule that looks like this:

If we think about this for a moment, we can see that there's going to be a problem.

- We just said that the lambda introduces a binding for the variable in the V position.
- This is precisely so that the body, B, can make use of that variable.

So let's imagine the simplest function:

But we don't have any typing rule that covers variables! Furthermore, we have no way of knowing what the type of any old variable will be. So we have a problem.

Typing Variables

Remember how we addressed this problem in our interpreter:

• we had an environment for recording the value bound to each variable.

We will use this same idea again:

• we'll have a type environment for recording the type of each variable.

That is, just as our interpreter had the type

```
(interp : (Exp Env -> Value))
```

our type-checker will have the type

```
(tc : (Exp TEnv -> Type))
```

Typing Variables

In our type-checker notation, we will use a slightly different way of writing it, which will finally make make |- stop being silent and take it proper pronounciation, "proves": all type rules will have the form

Г |- е : Т

where Γ , the capital Greek letter gamma, is conventionally used for the environment.

We read this as "the environment Γ proves that e has type T". So in fact there's been an environment hiding in all our judgments, but we didn't have to worry about it when we didn't have variables.

But now we do, so from now on we have to make it explicit. Fortunately, in most cases the environment is unchanged, and just passes recursively to the sub-terms, as you would expect from writing the interpreter.

With this, we can write a type for variables. What is the type of a variable? It's whatever the environment says it is! We'll treat the environment as a function, so we can just write the following axiom (where v stands for all the syntactically valid variable names):

 $\Gamma \mid -v : \Gamma(v)$

Now we're in a position to fill in the holes. When we check the body of the function, we should do it in an extended environment:

```
Γ[V <- ???] |- B : ???
-----
Γ |- (lambda V B) : ???
```

where $\Gamma[V \leq -]$ is how we write " Γ is extended with V bound to _": this is the same environment-extension function that we've written before, for type environments instead of value environments, but operationally the same.

Okay, but two questions: extend which environment, and extend it with what?

Which is easy: it's the environment of the function definition (static scope!). The repetition of Γ in both the consequent and antecedent accomplishes that.

In terms of what: We need to provide a type for the variable so that, when we try to look up its type, the environment can return something.

But we don't know what to extend it with! The type-checker needs the *programmer to tell it* what type the function is expecting. This is one of the reasons why programming languages expect annotations in function and method definitions.

(Another-equally good-reason is because it better documents the function for people who have to use it and maintain it.)

Therefore, we have to extend the syntax of functions to include a type annotation:

(lambda V : T B)

which says that V is expecting to be bound to a value of type T in body B.

Once we accept this modification, we can make progress on the conditional rule:

```
Γ[V <- T] |- B : ???
-----
Γ |- (lambda V : T B) : ???
```

What type are we expecting for the function definition? Clearly it must be a function type:

```
Γ[V <- T] |- B : ???
------
Γ |- (lambda V : T B) : (??? -> ???)
```

Furthermore, we know that the type expected by the function must be T:

```
Γ[V <- T] |- B : ???
------
Γ |- (lambda V : T B) : (T -> ???)
```

Given a value of type T, the function will return whatever the body produces:

And that gives us our final rule for function definitions.

More Divergence Between Types and Evaluation

It is interesting to contrast the above pair of typing rules with the corresponding evaluation rules.

- In the evaluator, we visit the body of the function on every *application*—which could be as many as an infinite number of times in a program.
- In contrast, we visit the body of the function on *definition*, which happens only once.
- Therefore, even if the program runs forever, the type-checker is guaranteed to terminate!

Why can we get away with this?

- The evaluator has to run the body with the *specific* value it was given.
- The type-checker, however, has abstracted the concrete values away.
- Therefore, it only needs to make one pass through the body with the "abstract value", the type.

Aside: Earlier, when we proposed the type Fun, we said that it collapsed all functions in the world into one type. This was too coarse, and we had to refine the type of a function. However, we are *still* collapsing an infinite number of functions into each of those function types—just as we collapse an infinite number of strings into Str, and so on. Both the strength and weakness of type-checking lies in this collapsing.

For the same reason, observe that a function application rule only cares about the *type* of the function, not *which* specific function is being applied. Therefore, any function that has that type can be used. For that same reason, the type-checker *cannot* traverse the function's body at application time—it doesn't even know which function might be used! All communication between the function body and application must happen entirely through the type boundary.

Assume-Guarantee Reasoning

There is a delicate dance going on between these typing rules for application and definition (now updated to have the environment):

```
Γ |- F : (T -> U) Γ |- A : T
------
Γ |- (F A) : U
```

The rule for lambda *assumes* the parameter will be given a value of type T; the application rule *guarantees* that that the actual parameter will indeed have the expected type. The application rule *assumes* that the function, if given a T, will produce a U (because the type is $(T \rightarrow U)$); the lambda rule *guarantees* that the function will indeed perform that way.

Aside: The notation $(T \rightarrow U)$ is not chosen at random. The \rightarrow may remind you of the notation for implication in mathematics. That's intentional. We can read the type as "giving the function a T implies that it will produce a U" (not giving it a T implies nothing about what it will do...). It is that *implication* that is assumed in the application rule, and that is guaranteed by the rule for lambda.

This assume-guarantee reasoning shows up in many places, so look out for this pattern in other places as well.

Recursion and Infinite Loops

We alluded, earlier, to how we can desugar more interesting features into functions and application. Let's take a look at a very specific feature: an infinite loop. Let's first confirm that we can write an infinite loop. Here's a program that does it:

fun f(): f()

f()

But this assumes we already have recursion. Can we write it without recursion? Actually we can! We'll use historical names (ω is the lower-case Greek omega):

(let ([ω (lambda (x) (x x))]) (ω ω))

Run this in Racket and confirm that it runs forever!

Recursion and Infinite Loops

Now let's see what happens when we try to type this. We have to provide a type annotation:

```
(let ([ω (lambda (x : ???) (x x))])
 (ω ω))
```

Historically, the overall term is called Ω (the capital Greek omega).

Okay, so what is the annotation? To determine a type for x, we have to see how it's used. It's used twice. One use is in a function application position, so we know that the type must be of the form $(T \rightarrow U)$; now we have to determine what T and U are. Let's focus on the parameter type, T. But what are we passing in? We're passing in x, whose type is $(T \rightarrow U)$. So we need a solution to the equation

 $T = (T \rightarrow U)$

with one coming from the application position and the other from the argument position. Of course, there is no finite type that can fit this equation! Therefore, it appears that this program cannot be typed!

Of course, this is not a proof. However, there is a formal property associated with this programming language, which is called the Simply Typed Lambda Calculus (STLC): the property is called *strong normalization*, and it means that *all programs in this language terminate*.

Aside: If you have heard about the Halting Problem, how does that square with what you just read?

Recursion and Infinite Loops

It may seem rather useless to have a language in which all programs terminate—you can't write an operating system, or Web server, or many other programs in such a language. However, that misses two things.

- 1. There are many cases where we want programs to always terminate.
 - You don't want a network packet filter or a device driver or a compiler or a type-checker or ... to run forever. Of course we also want them to run quickly, but it would be nice if we had a guarantee that no matter what we did, we cannot create an infinite loop.
 - The STLC is very useful in some of these settings.
 - Another example of a place where we want guaranteed termination is in program linking, and the module language of Standard ML is therefore built atop the STLC: it lets you even write higher-order programs, but the type language guarantees that all module compositions (linkages) will terminate.
- 2. Second, many long-running programs are actually a composition of an infinite loop and a short-running program.
 - Think about an operating system with device drivers, a Web server with a Web application, a GUI with callbacks, etc. In each case, there is a "spine" of an infinite loop that simply keeps the program reactive, and "ribs" of short computations that do a little specific work and terminate.
 - In fact, on the Web these programs must terminate quickly, otherwise the Web browser thinks the server has hung and offers to kill the window!
 - These kinds of reactive systems are therefore a composition of a very generic infinite loop calling out to specific programs for which a termination guarantee will often be very useful.

Finally, observe that we've learned something profound. Until now, we have probably thought of types as just a convenience or as a way of eliminating basic errors. However, we have just now seen that adding a type system can change the expressive power of a language. That is, these types are "semantic".

What went wrong above? The problem is that each application "uses up an arrow" in a function type; because a program text must be finite, it can contain at most a finite number of "arrows", so eventually the program must terminate. To get around this, we need a way to effectively have an "infinite quiver".

We typically do this by adding a recursive function construct to the language, and create a custom type for it. Let's start with a type rule for the analogous, but simpler, let:

Note that we're going to expect an annotation in let for the same reason we do for function definitions. So this says that we'll check that E actually does have the type promised in the declaration, T; when we extend the type environment with the V having type T, if the body B produces type U, then that's the type of the whole expression.

Aside: Notice that there's an assume-guarantee pair in the antecedent: the first term is guaranteeing the annotation, which the second term is assuming.

Aside: Technically, the type of E could be calculated. Therefore, the T annotation is not strictly necessary.

Observe that this is basically the type rule we would get from expanding the syntactic sugar for let. Therefore, this still doesn't let us write a recursive definition. We need something more. Let's introduce a new construct, rec, for recursive definitions. An example of a rec (in an untyped setting) might be

```
(rec ([inf-loop (lambda (n) (inf-loop n))])
  (inf-loop 0))
(rec ([fact (lambda (n) (if (zero? n) 1 (* n (fact (- n 1))))])
  (fact 10))
```

In the typed world, we'll want rec to have the form

(rec V : T E B)

so we'd instead have to write

```
(rec inf-loop : (Number -> Number)
      (lambda (n) (inf-loop n))
  (inf-loop 0))
```

```
(rec fact : (Number -> Number)
        (lambda (n) (if (zero? n) 1 (* n (fact (- n 1)))))
    (fact 10))
```

where V is fact, T is (Number -> Number), E is the big lambda term, and B is (fact 10).

So this introduces a recursive definition, and then uses it. How might we type this?

```
???
-----
Γ |- (rec V : T E B) : U
```

Well, clearly one part of it must be the same: we have to type the body in the extended environment, and the environment must be extended with the annotated type:

We also know that we need to confirm that the annotation is correct:

But clearly, something needs to be different, otherwise we've just reproduced let.

Look at the example use of rec: the E term also needs to have V bound in it! In other words, both E and B are typed in the same environment:

```
Γ[V <- T] |- E : T Γ[V <- T] |- B : U
------
Γ |- (rec V : T E B) : U
```

From the type, we can read off how the recursion happens: the extended environment for B initiates the recursion, while that for E sustains it. Essentially, the environment of E enables arbitrary recursive depth.

In short, to obtain arbitrary recursion—and hence infinite loops—we have to add a special construct to the language and its type-checker; we cannot obtain it just through desugaring. Once we add rec to the STLC, however, we obtain a conventional programming language again.

Safety and Soundness

A critical component of SMoL is the concept of safety:

- some operations are partial over the set of all values
- a SMoL language enforces this by reporting violations

Typical examples of partiality may include + applying only to certain types of values. However, I intentionally write "operations" rather than, say, "functions", because these could be primitive operations like application (expecting the first position to be a function or method) as well. In fact, in some languages like JavaScript, there are very few violations, as the Wat talk shows.

How must these be enforced? It can be either statically or dynamically. In Python and JavaScript, for instance, all safety violations are reported dynamically. In Java or OCaml, most of them are reported statically. Either way, safety means that data have integrity: there is some notion of "what they are", and that identity is respected by operations. Put differently, data are not misinterpreted.

These are all very abstract statements, which we will soon concretize.

Revisiting the Basic Calculator

We will start with a very basic calculator that has two types, numbers and strings, and an operation (addition and concatenation, respectively) on them. Note that it helps to have more than one type if we want to talk about safety. We will skip most of the boilerplate code and focus on the core of the calculator:

```
(define-type Exp
                                                      (define (num+ lv rv)
  [num (n : Number)]
                                                        (type-case Value lv
  [str (s : String)]
                                                          ((numV ln)
  [plus (l : Exp) (r : Exp)]
                                                           (type-case Value rv
  [cat (1 : Exp) (r : Exp)])
                                                             ((numV rn) (numV (+ ln rn)))
                                                             (else (error '+ "right not a number"))))
                                                          (else (error '+ "left not a number"))))
(define-type Value
  (numV (n : Number))
                                                      (define (str++ lv rv)
  (strV (s : String)))
                                                        (type-case Value lv
(calc : (Exp -> Value))
                                                          ((strV ls)
                                                           (type-case Value rv
(define (calc e)
                                                             ((strV rs) (strV (string-append ls rs)))
  (type-case Exp e
                                                             (else (error '++ "right not a string"))))
    [(num n) (numV n)]
                                                          (else (error '++ "left not a string"))))
    [(str s) (strV s)]
    [(plus l r) (num+ (calc l) (calc r))]
    [(cat l r) (str++ (calc l) (calc r))]))
```

Revisiting the Basic Calculator

With a suitable parser, we can run tests such as the following:

(test (calc (plus (num 1) (num 2))) (numV 3)) (test (calc (plus (num 1) (plus (num 2) (num 3)))) (numV 6)) (test (calc (cat (str "hel") (str "lo"))) (strV "hello")) (test (calc (cat (cat (str "hel") (str "l")) (str "o"))) (strV "hello")) (test/exn (calc (cat (num 1) (str "hello"))) "left") (test/exn (calc (plus (num 1) (str "hello"))) "right")

The last two, in particular, show that the language is safe. The checks inside the primitives—in num+, for instance—are called *safety checks.

Now we're going to do something fun: we're going to make the memory allocation of values explicit. As we go through this, remember what we've said before: a value in SMoL is just a memory address.

Let's do this in stages. First, we'll use a vector to represent memory:

```
(define MEMORY (make-vector 100 -1))
```

The value -1 is useful for identifying parts of memory that have not yet been touched (assuming, of course, we don't write a program that produces -1—which we can avoid doing easily enough in this illustration).

It will be useful to have a helper to use the next available bit of memory:

```
(define next-addr 0)
(define (write-and-bump v)
 (let ([n next-addr])
    (begin
        (vector-set! MEMORY n v)
        (set! next-addr (add1 next-addr))
        n)))
```

Now let's say we want to store a number in memory. We put it in the next available memory place, and return the *address* of the place where the number was stored. Be careful here: the number we return is a memory address (which, here, is represented as an array index), which is not at all necessarily the same as the *numeric value* being stored.

```
(define (store-num n)
 (write-and-bump n))
```

Correspondingly, when we want to read a number, we simply return what is at the address corresponding to the number.

```
(define (read-num a)
  (vector-ref MEMORY a))
```

We want the property that when we read-num from the address where we store-num a number, we get back that same number: for all N,

(read-num (store-num N)) is N

Now let's look at strings. We are going to convert the string into a sequence of character codes, and store those codes explicitly:

In particular, the value stored at the address representing the string is the *length* of the string, followed by the individual characters. (Endless blood has been spent over whether strings should store their lengths at the front, or whether they should only be delimited by a special value, or both. The question is uninteresting here.)

Thus, suppose with a fresh memory we run

(store-str "hello")

this would return the address 0. The resulting value of MEMORY would be

'#(5

-1

That is, at address 0 we have the length of the string, followed by five character codes; these six memory entries together constitute the five-character string "hello". The rest of the memory remains untouched.

To read a string we have to reassemble it:

Once again, we want the result of reading a written string to give us the same string.

Now let's update the calculator. First, we're in for a surprise: we no longer need (*or want*) a fancy Racket datatype to track values, because values are just addresses (i.e., array indices)! So:

(define-type-alias Value Number)

The type of the calculator doesn't change; it still produces values. It's just that the representation of values has changed...dramatically. (Recall, again, that these Numbers are addresses, not numeric values *in* the interpreted language.)

The calculator remains the same. What has changed is in the helper functions. In the primitive value cases, we have to explicitly allocate them—which is what we were doing when we called the previous definitions of numV and strV (which store data on the heap), except it may not have been so evident. We will make it explicit as follows:

(define numV store-num)
(define strV store-str)

Okay, now to update the helper functions. Let's focus on num+. The core logic is currently

```
[(numV rn) (numV (+ ln rn))]
```

Observe that now we're calling it on the result of calling calc, i.e., on Values. That means num+ is going to get two addresses as arguments, and it needs to look up the corresponding numbers in memory, and then produce the resulting number:

```
(define (num+ la ra)
 (numV (+ (read-num la) (read-num ra))))
```

Analogously, we can define concatenation as well:

```
(define (str++ la ra)
  (strV (string-append (read-str la) (read-str ra))))
```

Finally, we have to update our tests as well. Because calc now returns *addresses*, all our answers appear to be incorrect. Instead, we have to obtain the corresponding numbers or strings at those addresses. Once we do so, calc passes the tests:

```
(test (read-num (calc (plus (num 1) (num 2)))) 3)
(test (read-num (calc (plus (num 1) (plus (num 2) (num 3))))) 6)
(test (read-str (calc (cat (str "hel") (str "lo")))) "hello")
(test (read-str (calc (cat (cat (str "hel") (str "l")) (str "o")))) "hello")
```

Except...does it? These two tests do not pass:

```
(test/exn (calc (cat (num 1) (str "hello"))) "left")
(test/exn (calc (plus (num 1) (str "hello"))) "right")
```

In fact, how can they? In all the above code, there are no errors left!

Rather, when we run

```
(calc (cat (num 1) (str "hello")))
```

we get an address back (maybe 69; it depends on what you ran earlier and hence what is in MEMORY). In fact, we can decide how we want to treat this: as a number?

```
> (read-num 69)
- Number
6
```

```
> (read-str 69)
```

```
- String
```

```
"\u0005hello"
```

How can something be both a number and a string? Well, actually, the situation is a bit more confusing than that: 69 above is just an address in memory from which we can read off whatever we want *however we want it* (i.e., the content of that address is *interpreted* by the function that reads from it), which can result in garbage.

It can get even worse:

```
> (read-num (calc (plus (num 1) (str "hello"))))
- Number
6
> (read-str (calc (plus (num 1) (str "hello"))))
- String
. . integer->char: contract violation
    expected: valid-unicode-scalar-value?
    given: -1
```

That is, we've tried to read "off the end of memory". It was dumb luck that we had a -1 as the initial value; the -1 triggered an error when we tried to convert it to a character *because Racket's primitives are safe*, which halted the program. If integer->char did not have a safety check, we would have gotten some garbled string instead.

In short, what we have created is an *unsafe* language. Data have no integrity. Any value can be treated as any kind of datum. This, in short, is the memory model of C, and it's largely proven to be a disaster for modern programming, which is why SMoL languages evolved.

Fortunately, it does not take too much work to make the language safe again. What we've just written holds the key:

- every value needs to record what kind of value it is.
- And any use of that value needs to check that it's the right kind of value.

This information is called a *tag*; it takes a fixed amount of space, and represents *metadata* about the subsequent datum. All subsequent values are interpreted in accordance with the tag. We need two tags for the two kinds of values. Let's use

(define NUMBER-TAG 1337) (define STRING-TAG 5712)

It's important that the two tags be different, so they are unambiguous. However, we don't need to worry about the tags themselves being confused with other data (e.g., numbers), because the tags will never be processed directly as program data (unless, of course, there is a bug in our implementation that accidentally does so...which is why language implementations need to be tested extensively).
Now, when we allocate a number, we write its tag into the first address, followed by the actual numeric value:

```
(define (store-num n)
 (let ([a0 (write-and-bump NUMBER-TAG)])
  (begin
    (write-and-bump n)
    a0)))
```

And when we try to read a number, we first check that it really is a number, and only then obtain the actual numeric value:

```
(define (safe-read-num a)
 (if (= (vector-ref MEMORY a) NUMBER-TAG)
        (vector-ref MEMORY (add1 a))
        (error 'number (number->string a))))
```

Strings are analogous:

```
(define (store-str s)
  (let ([a0 (write-and-bump STRING-TAG)])
    (begin
    (write-and-bump (string-length s))
    (map write-and-bump
         (map char->integer (string->list s)))
   a0)))
(define (safe-read-str a)
  (if (= (vector-ref MEMORY a) STRING-TAG)
      (letrec ([loop
                 (lambda (count a)
                   (if (zero? count)
                       empty
                       (cons (vector-ref MEMORY a)
                             (loop (sub1 count) (add1 a)))))])
        (list->string
         (map integer->char
              (loop (vector-ref MEMORY (add1 a)) (+ a 2)))))
      (error 'string (number->string a))))
```

So now, starting from a fresh memory, running

(store-str "hello")

still produces 0, but the content of MEMORY looks a bit different:

'#(5712

111

-1 -1

- 1

-1

...)

That is, at address 0 we first encounter the tag for strings. Only then do we get the string's length, followed by its contents.

Observe that now, storing the length up front makes even more sense:

- the first two locations contain the tag and the length,
- both of which are metadata that help us interpret what comes later,
- with the second (the length) refining the first (the tag).

With this change, the interpreter stays unchanged, and effectively so do the helpers, other than using the new names we've chosen:

```
(define (num+ la ra)
 (store-num (+ (safe-read-num la) (safe-read-num ra))))
```

```
(define (str++ la ra)
  (store-str (string-append (safe-read-str la) (safe-read-str ra))))
```

All our "good" tests still pass, but interestingly, our "bad" tests now fail:

```
(test/exn (calc (cat (num 1) (str "hello"))) "string")
(test/exn (calc (plus (num 1) (str "hello"))) "number")
```

Our safe evaluator has, however, come at a price relative to the unsafe evaluator:

- In terms of running time, we are now clearly paying for the overhead of safety checks.
- In terms of **space**, we are paying for the tags.
- Thus, we have had to get worse space and time.

Nevertheless, the price of unsafe languages is so high—e.g., in the form of security problems—and the cost of safety is often so low, that programmers gladly pay this price (or do so without even particularly noticing it).

Still, it would be nice if we didn't have to pay the price at all. And there is a way to accomplish that: types.

Look at our "bad" programs:

- These are programs that can *statically* be rejected by a type-checker.
- If we could reject all such programs, then—since no "bad" programs would be left—we can then run the program on the unsafe evaluator without worrying about negative consequences.
- This, in effect, is what most typed languages, like Java and OCaml, do.

Thus we find another use for types: to improve program performance. But this requires care.

Soundness

Running on an unsafe evaluator is, as the name suggests, dangerous. Therefore, we should only do it if we can be sure that nothing can go wrong. That means that our type system needs to come with a *guarantee*.

The way this guarantee is usually formulated is as follows. Suppose we have

e : t

and suppose we evaluate it and find that

e -> v

The latter—its value—is the ground truth. The type checker's job is to make sure it matches what the evaluator produces. That is, we would ideally like that

```
e : t if and only if e -> v and v : t
```

This says that the type checker's job is to perfectly mirror the evaluator: whatever type the program's result value has is the same type the type-checker says it has.

Soundness

Unfortunately, for a Turing-complete language, this full guarantee is impossible to obtain, because of Rice's Theorem. Instead, we have to compromise and see if we can get at least one of the two directions. When we think about it, we realize that, in a typed language, we're only really interested in programs that pass the type-checker (i.e., have a type). Therefore, we expect that

```
If e : t then
if e \rightarrow v, then v : t
```

This says that whatever type the type-checker predicted is exactly the type that the program has. That means we can rely on the type-checker's prediction. Which in turn means that we can be sure there are no type violations. Which tells us we can safely run the program atop an unsafe evaluator! This property is called *type soundness*.

Note that soundness is not a given: it's a property that must be formally, mathematically *proven* of a given type-checker and evaluator. The proof can be quite complex. This is because the "shape" of program evaluation and that of type-checking can be very different, as we have seen before for conditionals and functions.

And failure to prove it correctly—i.e., claiming it holds when in fact it doesn't—means we've allowed a vulnerability to slip through. This can manifest as uncaught exceptions, crashes, segmentation faults, etc. In addition, a clever attacker can construct a program that exploits the vulnerability, and our system can be subjected to a security or other attack. Thus, any soundness violations are emergencies and result in panic.

One of the consequences of our tagged representation is that when extracting a value from memory, we don't have to know whether to use safe-read-num or safe-read-str; the tag at the address can tell us which to use. That is, we can define

```
(define (generic-read a)
 (let ([tag (vector-ref MEMORY a)])
    (cond
     [(= tag NUMBER-TAG) (safe-read-num a)]
     [(= tag STRING-TAG) (safe-read-str a)]
      [else (error 'generic-print "invalid tag")])))
```

Unfortunately, this code can't be typed by plait because the two branches return different types.

Generic Printing

We can solve this in two ways:

- 1. We can use a hack: use #lang plait #:untyped, which provides the same syntactic language, features, and run-time behavior, but turns off the type-checker. (Curiously, we were using the type-checker to keep us disciplined: so that the only values we could store in MEMORY would be numbers! Therefore, it's good to not use the untyped version often.)
- 2. Notice that in the end, what printers do is essentially print a string.

Therefore, we just need to return a string in all cases:

```
(define (generic-read a)
 (let ([tag (vector-ref MEMORY a)])
  (cond
    [(= tag NUMBER-TAG) (number->string (safe-read-num a))]
    [(= tag STRING-TAG) (safe-read-str a)]
    [else (error 'generic-print "invalid tag")])))
```

A consequence of having this function is that we can rewrite our tests to be more proper:

```
(test (generic-read (calc (plus (num 1) (num 2)))) "3")
(test (generic-read (calc (plus (num 1) (plus (num 2) (num 3))))) "6")
(test (generic-read (calc (cat (str "hel") (str "lo")))) "hello")
(test (generic-read (calc (cat (cat (str "hel") (str "l")) (str "o"))))
"hello")
```

This is much closer to how we would write the test in the original interpreter; the only difference here is that the evaluator produces an address as the value, but we would like to inspect the value in a human-readable and -writable form, so we use generic-read.

Unannotated Programs and Types

Consider the following plait program:

(lambda (x y) (if x (+ y 1) (+ y 2)))

If we enter this program into plait, e.g., as follows, something remarkable happens:

```
> (lambda (x y)
    (if x
        (+ y 1)
        (+ y 2)))
- (Boolean Number -> Number)
#>
```

In response, plait *figures out* the type of this function *without* our having to provide any annotations. This is in contrast to the type-checker we just wrote, which required us to extend the syntax just to provide (required) type annotations. That tells us that something different—and more—must be happening under plait.

Unannotated Programs and Types

In contrast, consider another example:

(lambda (x) (if x (+ x 1) (+ x 2)))

This produces an error, observing that we are using x both in a position that requires it to be a Boolean (in if) and a number (in the two additions). Again, plait has figured this out without our having to write any annotations at all!

The algorithm that sits underneath plait is essentially the same algorithm under OCaml, Haskell, and several other programming languages. These languages provide *type inference*: figuring out (inferring) types automatically from the program source. Now we're going to see how this works.

Unannotated Programs and Types

The key idea is to break this seemingly very complex problem into two rather simple parts.

- In the first, we recursively visit each sub-expression of the program (following SImPI) and generate a set of *constraints* that formally do what we've been doing informally above.
- The second phase *solves* this set of constraints, using a process that is a generalization of the process you used for solving "systems of simultaneous equations" in school.

The solution is a type for each variable. That lets us fill in the annotations that the programmer left blank.

The process of generation will also have applied the type constraints, so there will be no further need to type-check the program;

- But we can use the annotations, for instance, in an IDE for tool-tips, in a compiler for optimization, etc.
- That is, with inference, we can program as if we're in a "scripty" language without annotations, yet achieve most of the benefits of types.
- (I say "most" because one of the benefits is documentation; leaving off all annotations makes programs harder to read and understand. For that reason, inference should be used sparingly.)

Imagining a Solution

Until now, our type checker has required us to annotate the parameter of every function. But let's imagine someone handed us a piece of code without annotations; can we figure out the type anyway? For instance, consider:

(+ 1 2)

We clearly know the type of this; even our type-checker can calculate it for us without any annotations. But of course that's not surprising: there are no variables to annotate. So now consider this expression:

```
(lambda (x : ___) (+ x 1))
```

With a moment's inspection, we can tell that the function has type ($Num \rightarrow Num$). But our type-checker couldn't have calculated that, because it would have tripped on the empty annotation. So how can *we* figure it out?

Imagining a Solution

Well, let's see. First we have to figure out the type of x. To determine its type, we should look for uses of x. There is only one, and it's used in an addition. But the rule for addition

Γ |- e1 : Num Γ |- e2 : Num -----Γ |- (+ e1 e2) : Num

tells us that the term in that position must have type Num.

- There is no additional information we have about x (this remark will become clearer in a moment).
- Therefore, we can determine that its type must be Num.
- Furthermore, we know that the result of an addition is also a Num.

From that, we can conclude that the function has type (Num \rightarrow Num).

Unique Variable Names

In what follows, we will assume that all variable names in the program are unique. That is, a given variable name is bound in at most one place in a program. This greatly simplifies the presentation below, because we can speak of the type of a variable and know *which* variable it refers to, instead of having to constantly qualify which variable of that name we mean.

This restriction does not actually preclude any programs in a language with static scope. Consider this program, which produces 7:

(let ([x 3]) (+ (let ([x 4]) x) x))

We can just as well *consistently rename* one of the xs to something else (heck, we can even use the DrRacket interface to have Racket do the renaming for us), and leave the program meaning exactly the same:

```
(let ([x 3])
(+ (let ([y 4])
y)
x))
```

This renaming process is called alpha conversion or alpha renaming.

With that important detail out of the way, let's return to our process of *inferring* or *reconstructing* the types of variables from the way they're used in a program. Here's another example with a two-parameter function:

Once again, we can't just calculate its type with our type-checker; instead, we must reconstruct the type from the function body. Let's do that. What can we tell?

Let's again refer to the conditional rule:

Γ |- C : Bool Γ |- T : U Γ |- E : U -----Γ |- (if C T E) : U

- This tells us that what's in the C position—here, x—must be a Bool.
- Furthermore, both branches (+ y 1) and (+ y 2) must have the same type.

That's all we can learn from the rule for if!

But now we can (and must) recur into the sub-expressions.

- Each one is an addition, and the addition rule tells us that both arguments must be Nums.
- Both of these indicate that the type of y must be Num.

Furthermore, both indicate that the overall addition returns a Num. From that we can tell that the entire expression must have the type

(Bool Num -> Num)

By this process,

- we can figure out what types to put in the missing annotations.
- More subtly, notice that by running through this process, we have effectively applied all the typing rules;
- therefore, if we have successfully reconstructed the type annotations, we need not bother type-checking the program with those annotations: it will have to type-check.

Now let's consider a slight variation on the above program:

```
(lambda (x : ___)
(if x
(+ x 1)
(+ x 2)))
```

Now let's figure out everything we can learn about x from the function's body:

- x is used in the conditional position of an if. Therefore, it must have type Bool.
- x is used as a parameter to +. Therefore, it must have type Num.
- x is again used as a parameter to +. Therefore, it must have type Num.

Notice that each of these conclusions is perfectly fine on its own. However, when we *put them together* (which is what we meant by "additional information" above), there's a problem:

- x cannot be both of those.
- That is, we are unable to find a single type for x.
- This inability to find a type for x means that the program has a *type error*.

And indeed, there is no type we could have given that would have enabled this program to execute safely.

Observe something subtle.

- While we can report that the program clearly has a type error, our error message must necessarily be much more ambiguous.
- Previously, when we had a type annotation on x, we could pinpoint where the error occurred.
- Now, all we can say is that the program is not type-*consistent*, but cannot blame one spot or the other without potentially misleading the programmer.

Instead, we must report all these locations and let the programmer decide where the error is based on their *unstated intent* (in the form of a type annotation).

The details of this algorithm—called *Hindley-Milner* inference—are fascinating, and worked out in detail in both the first and second editions of the book, PLAI (Chapter 30 in the first edition and Chapter 15.3.2 in the second edition).

For several worked examples of both constraint generation and constraint solving, refer to the first edition.

The first edition has a more algorithmic presentation, while the second provides code (it may be useful to compare the two). The prose in the second is different from that in the first, so different readers may prefer one over the other.