

Type Systems and Type Inference

Programming involves a wide range of computational constructs, such as data structures, functions, objects, communication channels, and threads of control. Because programming languages are designed to help programmers organize computational constructs and use them correctly, many programming languages organize data and computations into collections called types. In this chapter, we look at the reasons for using types in programming languages, methods for type checking, and some typing issues such as polymorphism, overloading, and type equality. A large section of this chapter is devoted to type inference, the process of determining the types of expressions based on the known types of some symbols that appear in them. Type inference is a generalization of type checking, with many characteristics in common, and a representative example of the kind of algorithms that are used in compilers and programming environments to determine properties of programs. Type inference also provides an introduction to polymorphism, which allows a single expression to have many types.

6.1 TYPES IN PROGRAMMING

In general, a *type* is a collection of computational entities that share some common property. Some examples of types are the type `int` of integers, the type `int→int` of functions from integers to integers, and the Pascal subrange type `[1 .. 100]` of integers between 1 and 100. In concurrent ML there is the type `int channel` of communication channels carrying integer values and, in Java, a hierarchy of types of exceptions.

There are three main uses of types in programming languages:

- naming and organizing concepts,
- making sure that bit sequences in computer memory are interpreted consistently,
- providing information to the compiler about data manipulated by the program.

These ideas are elaborated in the following subsections.

Although some programming language descriptions will say things like, “Lisp is an untyped language,” there is really no such thing as an untyped programming language. In Lisp, for example, lists and atoms are two different types: list operations

can be applied to lists but not to atoms. Programming languages do vary a great deal, however, in the ways that types are used in the syntax and semantics (implementation) of the language.

6.1.1 Program Organization and Documentation

A well-designed program uses concepts related to the problem being solved. For example, a banking program will be organized around concepts common to banks, such as accounts, customers, deposits, withdrawals, and transfers. In modern programming languages, customers and accounts, for example, can be represented as separate types. Type checking can then check to make sure that accounts and customers are treated separately, with account operations applied to accounts but not used to manipulate customers. Using types to organize a program makes it easier for someone to read, understand, and maintain the program. Types therefore serve an important purpose in documenting the design and intent of the program.

An important advantage of type information, in comparison with comments written by a programmer, is that types may be checked by the programming language compiler. Type checking guarantees that the types written into a program are correct. In contrast, many programs contain incorrect comments, either because the person writing the explanation was careless or because the program was later changed but the comments were not.

6.1.2 Type Errors

A *type error* occurs when a computational entity, such as a function or a data value, is used in a manner that is inconsistent with the concept it represents. For example, if an integer value is used as a function, this is a type error. A common type error is to apply an operation to an operand of the wrong type. For example, it is a type error to use integer addition to add a string to an integer. Although most programmers have a general understanding of type errors, there are some subtleties that are worth exploring.

Hardware Errors. The simplest kind of type error to understand is a machine instruction that results in a hardware error. For example, executing a “function call”

```
x()
```

is a type error if *x* is not a function. If *x* is an integer variable with value 256, for example, then executing *x()* will cause the machine to jump to location 256 and begin executing the instructions stored at that place in memory. If location 256 contains data that do not represent a valid machine instruction, this will cause a hardware interrupt. Another example of a hardware type error occurs in executing an operation

```
float_add(3, 4.5)
```

where the hardware floating-point unit is invoked on an integer argument 3. Because the bit pattern used to represent 3 does not represent a floating-point number in the form expected by the floating-point hardware, this instruction will cause a hardware interrupt.

Unintended Semantics. Some type errors do not cause a hardware fault or interrupt because compiled code does not contain the same information as the program source code does. For example, an operation

```
int_add(3, 4.5)
```

is a type error, as `int_add` is an integer operation and is applied here to a floating-point number. Most hardware would perform this operation. Because the bits used to represent the floating-point number 4.5 represent an integer that is not mathematically related to 4.5, the operation it is not meaningful. More specifically, `int_add` is intended to perform addition, but the result of `int_add(3, 4.5)` is not the arithmetic sum of the two operands.

The error associated with `int_add(3, 4.5)` may become clearer if we think about how a program might apply integer addition to a floating-point argument. To be concrete, suppose a program defines a function `f` that adds three to its argument,

```
fun f(x) = 3+x;
```

and someplace within the scope of this definition we also declare a floating-point value `z`:

```
float z = 4.5;
```

If the programming language compiler or interpreter allows the call `f(z)` and the language does not automatically convert floating-point numbers to integers in this situation, then the function call `f(z)` will cause a run-time type error because `int_add(3, 4.5)` will be executed. This is a type error because integer addition is applied to a noninteger argument.

The reason why many people find the concept of type error confusing is that type errors generally depend on the concepts defined in a program or programming language, not the way that programs are executed on the underlying hardware. To be specific, it is just as much of a type error to apply an integer operation to a floating-point argument as it is to apply a floating-point operation to an integer argument. It does not matter which causes a hardware interrupt on any particular computer.

Inside a computer, all values are stored as sequences of bytes of bits. Because integers and floating-point numbers are stored as four bytes on many machines, some integers and floating-point numbers overlap; a single bit pattern may represent an integer when it is used one way and a floating-point number when it is used another. Nonetheless, a type error occurs when a pattern that is stored in the computer for the

purpose of representing one type of value is used as the representation of another type of value.

6.1.3 Types and Optimization

Type information in programs can be used for many kinds of optimizations. One example is finding components of records (as they are called in Pascal and ML) or structs (as they are called in C). The component-finding problem also arises in object-oriented languages. A record consists of a set of entries of different types. For example, a student record may contain a student name of type string and a student number of type integer. In a program that also has records for undergraduate students, these might be represented as related type that also contains a field for the year in school of the student. Both types are written here as ML-style type expressions:

```
Student = {name : string, number : int}
Undergrad = {name : string, number : int, year : int}
```

In a program that manipulates records, there might be an expression of the form `r.name`, meaning the name field of the record `r`. A compiler must generate machine code that, given the location of record `r` in memory at run time, finds the location of the field name of this record at run time. If the compiler can compute the type of the record at compile time, then this type information can be used to generate efficient code. More specifically, the type of `r` makes it is possible to compute the location of `r.name` relative to the location `r`, at compile time. For example, if the type of `r` is `Student`, then the compiler can build a little table storing the information that name occurs before number in each `Student` record. Using this table, the compiler can determine that name is in the first location allocated to the record `r`. In this case, the expression `r.name` is compiled to code that reads the value stored in location `r+1` (if location `r` is used for something else besides the first field). However, for records of a different type, the name field might appear second or third. Therefore, if the type of `r` is not known at compile time, the compiler must generate code to compute the location of name from the location of `r` at run time. This will make the program run more slowly. To summarize: Some operations can be computed more efficiently if the type of the operand is known at compile time.

In some object-oriented programming languages, the type of an object may be used to find the relative location of parts of the object. In other languages, however, the type system does not give this kind of information and run-time search must be used.

6.2 TYPE SAFETY AND TYPE CHECKING

6.2.1 Type Safety

A programming language is *type safe* if no program is allowed to violate its type distinctions. Sometimes it is not completely clear what the type distinctions are in a specific programming language. However, there are some type distinctions that are meaningful and important in all languages. For example, a function has a different

type from an integer. Therefore, any language that allows integers to be used as functions is not type safe. Another action that we always consider a type error is to access memory that is not allocated to the program.

The following table characterizes the type safety of some common programming languages. We will discuss each form of type error listed in the table in turn.

Safety	Example languages	Explanation
Not safe	C and C++	Type casts, pointer arithmetic
Almost safe	Pascal	Explicit deallocation; dangling pointers
Safe	Lisp, ML, Smalltalk, Java	Complete type checking

Type Casts. Type casts allow a value of one type to be used as another type. In C in particular, an integer can be cast to a function, allowing a jump to a location that does not contain the correct form of instructions to be a C function.

Pointer Arithmetic. C pointer arithmetic is not type safe. The expression $*(p+i)$ has type A if p is defined to have type A*. Because the value stored in location $p+i$ might have any type, an assignment like $x = *(p+i)$ may store a value of one type into a variable of another type and therefore may cause a type error.

Explicit Deallocation and Dangling Pointers. In Pascal, C, and some other languages, the location reached through a pointer may be deallocated (freed) by the programmer. This creates a *dangling pointer*, a pointer that points to a location that is not allocated to the program. If p is a pointer to an integer, for example, then after we deallocate the memory referenced by p, the program can allocate new memory to store another type of value. This new memory may be reachable through the old pointer p, as the storage allocation algorithm may reuse space that has been freed. The old pointer p allows us to treat the new memory as an integer value, as p still has type int. This violates type safety. Pascal is considered “mostly safe” because this is the only violation of type safety (after the variant record and other original type problems are repaired).

6.2.2 Compile-Time and Run-Time Checking

In many languages, type checking is used to prevent some or all type errors. Some languages use type constraints in the definition of legal program. Implementations of these languages check types at compile time, before a program is started. In these languages, a program that violates a type constraint is not compiled and cannot be run. In other languages, checks for type errors are made while the program is running.

Run-Time Checking. In programming languages with run-time type checking, the compiler generates code so that, when an operation is performed, the code checks to make sure that the operands have the correct type. For example, the Lisp language operation `car` returns the first element of a cons cell. Because it is a type error to apply `car` to something that is not a cons cell, Lisp programs are implemented so that, before `(car x)` is evaluated, a check is made to make sure that x is a cons cell. An advantage of run-time type checking is that it catches type errors. A disadvantage is the run-time cost associated with making these checks.

Compile-Time Checking. Many modern programming languages are designed so that it is possible to check expressions for potential type errors. In these

languages, it is common to reject programs that do not pass the compile-time type checks. An advantage of compile-time type checking is that it catches errors earlier than run-time checking does: A program developer is warned about the error before the program is given to other users or shipped as a product. Because compile-time checks may eliminate the need to check for certain errors at run time, compile-time checking can make it possible to produce more efficient code. For a specific example, compiled ML code is two to four times faster than Lisp code. The primary reason for this speed increase is that static type checking of ML programs greatly reduces the need for run-time tests.

Conservativity of Compile-Time Checking. A property of compile-time type checking is that the compiler must be conservative. This means that compile-time type checking will find all statements and expressions that produce run-time type errors, but also may flag statements or expressions as errors even if they do not produce run-time errors. To be more specific about it, most checkers are both sound and conservative. A type checker is sound if no programs with errors are considered correct. A type checker is conservative if some programs without errors are still considered to have errors.

There is a reason why most type checkers are conservative: For any Turing-complete programming language, the set of programs that may produce a run-time type error is undecidable. This follows from the undecidability of the halting problem. To see why, consider the following form of program expression:

```
if (complicated-expression-that-could-run-forever)
  then (expression-with-type-error)
  else (expression-with-type-error)
```

It is undecidable whether this expression causes a run-time type error, as the only way for expression-with-type-error to be evaluated is for complicated-expression-that-could-run-forever to halt. Therefore, deciding whether this expression causes a run-time type error involves deciding whether complicated-expression-that-could-run-forever halts.

Because the set of programs that have run-time type errors is undecidable, no compile-time type checker can find type errors exactly. Because the purpose of type checking is to prevent errors, type checkers for type-safe languages are conservative. It is useful that type checkers find type errors, and a consequence of the undecidability of the halting problem is that some programs that could execute without run-time error will fail the compile-time type-checking tests.

The main trade-offs between compile-time and run-time checking are summarized in the following table.

Form of Type Checking	Advantages	Disadvantages
Run-time	Prevents type errors	Slows program execution
Compile-time	Prevents type errors Eliminates run-time tests Finds type errors <i>before</i> execution and run-time tests	May restrict programming because tests are <i>conservative</i> .

Combining Compile-Time and Run-Time Checking. Most programming languages actually use some combination of compile-time and run-time type checking. In Java, for example, static type checking is used to distinguish arrays from integers, but array bounds errors (which are a form of type error) are checked at run time.

6.3 TYPE INFERENCE

Type inference is the process of determining the types of expressions based on the known types of some symbols that appear in them. The difference between type inference and compile-time type checking is really a matter of degree. A *type-checking* algorithm goes through the program to check that the types declared by the programmer agree with the language requirements. In *type inference*, the idea is that some information is not specified, and some form of logical inference is required for determining the types of identifiers from the way they are used. For example, identifiers in ML are not usually declared to have a specific type. The type system *infers* the types of ML identifiers and expressions that contain them from the operations that are used. Type inference was invented by Robin Milner (see the biographical sketch) for the ML programming language. Similar ideas were developed independently by Curry and Hindley in connection with the study of lambda calculus.

Although practical type inference was developed for ML, type inference can be applied to a variety of programming languages. For example, type inference could, in principle, be applied to C or other programming languages. We study type inference in some detail because it illustrates the central issues in type checking and because type inference illustrates some of the central issues in algorithms that find any kind of program errors.

In addition to providing a flexible form of compile-time type checking, ML type inference supports polymorphism. As we will see when we subsequently look at the type-inference algorithm, the type-inference algorithm uses *type variables* as placeholders for types that are not known. In some cases, the type-inference algorithm resolves all type variables and determines that they must be equal to specific types such as `int`, `bool`, or `string`. In other cases, the type of a function may contain type variables that are not constrained by the way the function is defined. In these cases, the function may be applied to any arguments whose types match the form given by a type expression containing type variables.

Although type inference and polymorphism are independent concepts, we discuss polymorphism in the context of type inference because polymorphism arises naturally from the way type variables are used in type inference.

6.3.1 First Examples of Type Inference

Here are two ML type-inference examples to give you some feel for how ML type inference works. The behavior of the type-inference algorithm is explained only superficially in these examples, just to give some of the main ideas. We will go through the type inference process in detail in Subsection 6.3.2.

Example 6.1

```

- fun f1(x) = x + 2;
val f1 = fn : int → int

```

The function `f1` adds 2 to its argument. In ML, 2 is an integer constant; the real number 2 would be written as 2.0. The operator `+` is overloaded; it can be either integer addition or real addition. In this function, however, it must be integer addition because 2 is an integer. Therefore, the function argument `x` must be an integer. Putting these observations together, we can see that `f1` must have type `int → int`.

Example 6.2

```

- fun f2(g,h) = g(h(0));
val f2 = fn : ('a → 'b) * (int → 'a) → 'b

```

The type `('a → 'b) * (int → 'a) → 'b` inferred by the compiler is parsed as `((('a → 'b) * (int → 'a)) → 'b)`.

The type-inference algorithm figures out that, because `h` is applied to an integer argument, `h` must be a function from `int` to something. The algorithm represents “something” by introducing a type variable, which is written as `'a`. (This is unrelated to Lisp `'a`, which would be syntax for a Lisp atom, not a variable.) The type-inference algorithm then deduces that `g` must be a function that takes whatever `h` returns (something of type `'a`) and then returns something else. Because `g` is not constrained to return the same type of value as `h`, the algorithm represents this second something by a new type variable, `'b`. Putting the types of `h` and `g` together, we can see that the first argument to `f2` has type `('a → 'b)` and the second has type `(int → 'a)`. Function `f2` takes the pair of these two functions as an argument and returns the same type of value as `g` returns. Therefore, the type of `f2` is `((('a → 'b) * (int → 'a)) → 'b)`, as shown in the preceding compiler output.

6.3.2 Type-Inference Algorithm

The ML type-inference algorithm uses the following three steps:

1. Assign a type to the expression and each subexpression. For any compound expression or variable, use a type variable. For known operations or constants, such as `+` or `3`, use the type that is known for this symbol.
2. Generate a set of constraints on types, using the parse tree of the expression. These constraints reflect the fact that if a function is applied to an argument, for example, then the type of the argument must equal the type of the domain of the function.
3. Solve these constraints by means of unification, which is a substitution-based algorithm for solving systems of equations. (More information on unification appears in the chapter on logic programming.)

The type-inference algorithm is explained by a series of examples. These examples present the following issues:

- explanation of the algorithm
- a polymorphic function definition
- application of a polymorphic function
- a recursive function
- a function with multiple clauses
- type inference indicates a program error

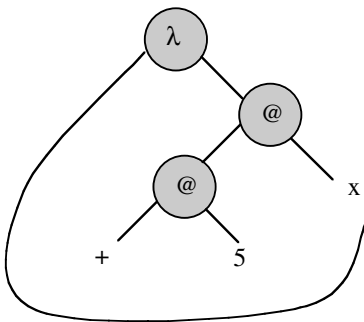
Altogether, these six examples should give you a good understanding of the type-inference algorithm, except for the interaction between type inference and overloading. The interaction between overloading and type inference is not covered in this book.

Example 6.3 Explanation of the Algorithm

We can see how the type-inference algorithm works by considering this example function:

```
– fun g(x) = 5 + x;
  val g = fn : int → int
```

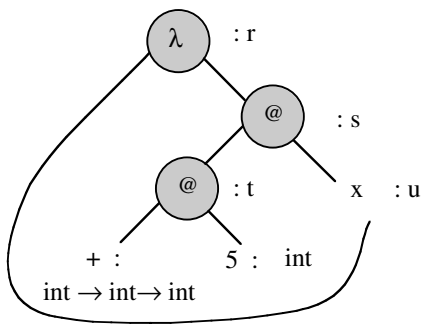
The easiest way to see how the algorithm works is by drawing the parse tree of the expression. We use an abbreviated form of parse tree that lists only the symbols that occur in the expression, together with the symbol @ for an application of a function to an argument. For the preceding expression we use the following graph. This is a form of parse tree, together with a special edge indicating the binding lambda for each bound variable. Here, the link from x to λ indicates that x is lambda bound at the beginning of the expression:



We use this graph to follow our type-inference steps:

1. *We assign a type to the expression and each subexpression:*

We illustrate this step by redrawing the graph, writing a type next to each node. To simplify notation, we use single letters r, s, t, u, v, \dots , for type variables instead of ML syntax 'a', 'b', and so on:



Recall that each node in a parse tree represents a subexpression. Repeating the information in the picture, the following table lists the subexpressions and their types:

Subexpression	Type
$\lambda x. ((+ 5) x)$	r
$((+ 5) x)$	s
$(+ 5)$	t
$+$	$\text{int} \rightarrow (\text{int} \rightarrow \text{int})$
5	int
x	u

Here we have written addition (+) as a Curried function and have chosen type $\text{int} \rightarrow \text{int} \rightarrow \text{int}$ for this operation. The prefix notation for addition is not ML syntax, of course, but it is used here to make the pictures simpler. As we saw earlier, + can either be integer addition or real-number addition. Here we can see from context that integer addition is needed. The actual ML type-inference algorithm will require a few steps to figure this out, but we are not concerned with the mechanics of overloading resolution here.

2. We generate a set of constraints on types, using the parse tree of the expression.

Constraints are equations between type expressions that must be solved. The way we generate them depends on the form of each subexpression. For each function application, constraints are generated according to the following rule.

Function Application: If the type of f is a , the type of e is b , and the type of fe is c , then we must have $a = b \rightarrow c$.

This typing rule for function application can be used twice in our expression.

- Subexpression $(+5)$, Constraint $\text{int} \rightarrow (\text{int} \rightarrow \text{int}) = \text{int} \rightarrow t$,
- Subexpression $(+5)x$, Constraint $t = u \rightarrow s$.

In the subexpression $(+ 5)$, the type of the function $+$ is $\text{int} \rightarrow (\text{int} \rightarrow \text{int})$, the type of the argument 5 is int and the type of the application is t . Therefore, we must have $\text{int} \rightarrow (\text{int} \rightarrow \text{int}), = \text{int} \rightarrow t$. The reasoning for subexpression $(+ 5)x$ is similar: In the subexpression $(+ 5) x$, the type of the function $(+ 5)$ is t , the type of the argument x is u , and the type of the application is s . Therefore, we must have $t = u \rightarrow s$.

Lambda Abstraction (Function Expression): If the type of x is a and the type of e is b , then the type of $\lambda x.e$ must equal $a \rightarrow b$.

For our example expression, we have one lambda abstraction. This gives us the following constraint:

Subexpression $\lambda x.((+5)x)$, Constraint $r = u \rightarrow s$.

In words, the type of the whole expression is r , the type of x is u , and the type of the subexpression $((+5)x)$ is s . This gives us the equation $r = u \rightarrow s$.

3. *We solve these constraints by means of unification.*

Unification is a standard algorithm for solving systems of equations by substitution. The general properties of this algorithm are not discussed here. Instead, the process is shown by example in enough detail that you should be able to figure out the types of simple expressions on your own.

For our example expression, we have the following constraints.

$$\begin{aligned} \text{int} \rightarrow (\text{int} \rightarrow \text{int}) &= \text{int} \rightarrow t, \\ t &= u \rightarrow s, \\ r &= u \rightarrow s. \end{aligned}$$

If there is a way of associating type expression to type variables that makes all of these equations true, then the expression is well typed. In this case, the type of the expression will be the type expression equal to the type variable r . If there is no way of associating type expression to type variables that makes all of these equations true, then there is no type for this expression. In this case, the type-inference algorithm will fail, resulting in an error message that says the expression is not well typed.

We can process these equations one at a time. The order is not very important, although it is convenient to put the equation involving the type of the entire expression last, as this is the output of the type-inference algorithm.

The first equation is true if $t = \text{int} \rightarrow \text{int}$. Because we need $t = \text{int} \rightarrow \text{int}$, we substitute $\text{int} \rightarrow \text{int}$ for t in the remaining equations. This gives us two equations to solve:

$$\begin{aligned} \text{int} \rightarrow \text{int} &= u \rightarrow s, \\ r &= u \rightarrow s. \end{aligned}$$

The only way to have $\text{int} \rightarrow \text{int} = u \rightarrow s$ is if $u = s = \text{int}$. Proceeding as before, we substitute int for both u and s in the remaining equation. This gives us

$$r = \text{int} \rightarrow \text{int},$$

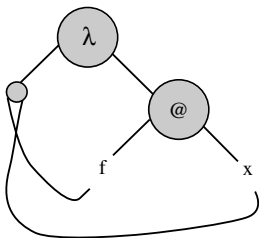
which tells us that the type r of the whole expression is $\text{int} \rightarrow \text{int}$. Because every constraint is solved, we know that the expression is typeable and we have computed a type for the expression.

Example 6.4 A Polymorphic Function Definition

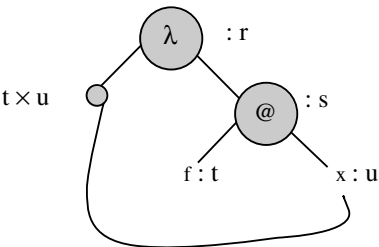
```
- fun apply(f,x) = f(x);
val apply = fn : ('a -> 'b) * 'a -> 'b
```

This is an example of a function whose type involves type variables. The type-inference algorithm begins by assigning a type to each subexpression. Because this makes it easiest to understand the algorithm, we write the function as a lambda expression with a pair $\langle f, x \rangle$ instead of a variable as a formal parameter: apply is defined by the lambda expression $\lambda \langle f, x \rangle. f\ x$ that maps a pair $\langle f, x \rangle$ to the result $f\ x$ of applying f to x .

Here is the parse graph of $\lambda \langle f, x \rangle. f\ x$, in which a pairing node is used on the left to indicate that the argument $\langle f, x \rangle$ of the function is a pair, with links to f and x .



The first step of the algorithm is to assign types to each node in the graph, as shown here:



The same information is repeated in the following table, showing the subexpressions represented by each node and their types.

Subexpression	Type
$\lambda \langle f, x \rangle. f\ x$	r
$\langle f, x \rangle.$	$t \times u$
$f\ x$	s
f	t
x	u

The second step of the algorithm is to generate a set of constraints. For this example, there is one constraint for the application and one for the lambda abstraction. The application gives us

$$t = u \rightarrow s.$$

In words, the type of the application $f\ x$ has type s , provided that the type of the function f is equal to $\langle \text{type of argument} \rangle \rightarrow s$. Because the type of the argument is u , this gives us the constraint $t = u \rightarrow s$.

The second constraint, from the lambda abstraction, is

$$r = t * u \rightarrow s.$$

In words, the type of $\lambda\langle f, x \rangle. fx$ is r , where r must be equal to $\langle \text{type of argument} \rangle \rightarrow s$, as s is the type of the subtree representing the function result. Because the argument is the pair $\langle f, x \rangle$, the type of the argument is $t * u$.

The constraints can be solved in order. The first requires $t = u \rightarrow s$, which we solve by substituting $u \rightarrow s$ for t in the remaining constraint. This gives us

$$r = (u \rightarrow s) * u \rightarrow s.$$

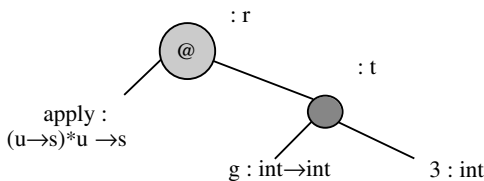
This tells us the type of the function. If we rewrite $(u \rightarrow s) * u \rightarrow s$ with 'a' and 'b' in place of u and s , then we get the compiler output previously shown. Because there are type variables in the type of the expression, the function may be used for many types of arguments. This is illustrated in the following example, which uses the type we have just computed for `apply`.

Example 6.5 Application of a Polymorphic Function

In the last example, we calculated the type of `apply`. The type of `apply` is $(a \rightarrow b) * a \rightarrow b$, which contains type variables. The type variables in this type mean that `apply` is a *polymorphic* function, a function that may be applied to different types of arguments. In the case of `apply`, the type $(a \rightarrow b) * a \rightarrow b$ means that `apply` may be applied to a pair of arguments of type $(a \rightarrow b) * a$ for any types 'a' and 'b'. In particular, recall that function `fun g(x) = 5 + x` from Example 6.3 has type $\text{int} \rightarrow \text{int}$. Therefore, the pair $(g, 3)$ has type $\text{int} \rightarrow \text{int} * \text{int}$, which matches the form $(a \rightarrow b) * a$ for function `apply`. In this example, we calculate the type of the application

```
apply(g,3);
```

Following the three steps of the type inference algorithm, we begin by assigning types to the nodes in the parse tree of this expression:



In this illustration, the smaller unlabeled circle is a pairing node. This node and the two below it represent the pair $(g, 3)$. In the previous example, we associated a product type with the pairing node. Here, to show that it is equivalent to use a type variable and constraint, we associate a type variable t with the pairing node and generate the constraint $t = (\text{int} \rightarrow \text{int}) * \text{int}$.

There are two constraints, one for the pairing node and one for the application node. For pairing, we have

$$t = (\text{int} \rightarrow \text{int}) * \text{int}.$$

For the application, we have

$$(u \rightarrow s)^* u \rightarrow s = t \rightarrow r.$$

In words, the type $(u \rightarrow s)^* u \rightarrow s$ of the function must equal $\langle \text{type of argument} \rangle \rightarrow r$, where r is the type of the application.

Now we must solve the constraints. The first constraint gives a type expression for t , which we can substitute for t in the second constraint. This gives us

$$(u \rightarrow s)^* u \rightarrow s = (\text{int} \rightarrow \text{int}) * \text{int} \rightarrow r.$$

This constraint has an expression on each side of the equal sign. To solve this constraint, corresponding parts of each expression must be equal. In other words, we can solve this constraint precisely by solving the following four constraints:

$$u = \text{int}, \quad s = \text{int}, \quad u = \text{int}, \quad s = r.$$

Because these require $u = s = \text{int}$, we have $r = \text{int}$. Because all of the constraints are solved, the expression $\text{apply}(g,3)$ is typeable in the ML type system. The type of $\text{apply}(g,3)$ is the solution for type variable r , namely int .

We can also apply apply to other types of arguments. If $\text{not} : \text{bool} \rightarrow \text{bool}$, then

```
apply(not, false)
```

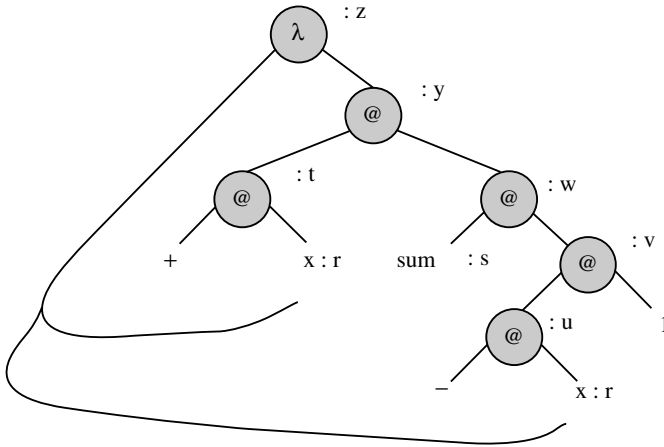
is a well-typed expression with type bool by exactly the same type-inference processes as those for $\text{apply}(g,3)$. This illustrates the polymorphism of apply : Because the type $(\text{'a} \rightarrow \text{'b}) * \text{'a} \rightarrow \text{'b}$ of apply contains type variables, the function may be applied to any type of arguments that can be obtained if the type variables in $(\text{'a} \rightarrow \text{'b}) * \text{'a} \rightarrow \text{'b}$ are replaced with type names or type expressions.

Example 6.6 A Recursive Function

When a function is defined recursively, we must determine the type of the function body without knowing the type of recursive function calls. To see how this works, consider this simple recursive function that sums the integers up to a given integer. This function does not terminate, but it does type check:

```
- fun sum(x) = x+sum(x-1);
val sum = fn : int -> int
```

Here is a parse graph of the function, with type variables associated with each of the nodes except for $+$, $-$, and 1 , as we ignore overloading and treat these as integer operations and integer constant. Because we are trying to determine the type of sum , we associate a type variable with sum and proceed:



Starting with the applications of $+$ and $-$ and proceeding from the lower right up, the constraints associated with the function applications and lambda abstraction in this expression are

$$\begin{aligned} \text{int} \rightarrow (\text{int} \rightarrow \text{int}) &= r \rightarrow t, \\ \text{int} \rightarrow (\text{int} \rightarrow \text{int}) &= r \rightarrow u, \end{aligned}$$

$$\begin{aligned} u &= \text{int} \rightarrow v, \\ s &= v \rightarrow w, \\ t &= w \rightarrow y, \\ z &= r \rightarrow y. \end{aligned}$$

To this list we add one more because the type of `sum` must also be the type of the entire expression:

$$s = z.$$

The constraint $s = z$ is the one additional constraint associated with the fact that this is a recursive declaration of a function. Solving these in order, we have

$$\begin{aligned} r &= \text{int}, \quad t = \text{int} \rightarrow \text{int}, \\ u &= \text{int} \rightarrow \text{int}, \\ v &= \text{int}, \\ s &= \text{int} \rightarrow w, \\ t &= w \rightarrow y, \\ z &= r \rightarrow y, \\ w &= \text{int}, \quad y = \text{int}, \\ z &= \text{int} \rightarrow \text{int}, \\ s &= \text{int} \rightarrow \text{int}. \end{aligned}$$

Because the constraints can be solved, the function is typeable. In the process of solving the constraints, we have calculated that the type of `sum` is $\text{int} \rightarrow \text{int}$.

Example 6.7 A Function with Multiple Clauses

Type inference for functions with several clauses may be done by a type check of each clause separately. Then, because all clauses define the same function, we impose the constraint that the types of all clauses must be equal. For example, consider the `append` function on lists, defined as follows:

```

- fun append(nil, l) = l
  | append(x::xs, l) = x :: append(xs, l);
val append = fn : 'a list * 'a list → 'a list.

```

As the type $: 'a \text{ list} * 'a \text{ list} \rightarrow 'a \text{ list}$ indicates, `append` can be applied to any pair of lists, as long as both lists contain the same type of list elements. Thus, `append` is a polymorphic function on lists.

We begin type inference for `append` by following the three-step algorithm for the first clause of the definition, then repeating the steps for the second clause. This gives us two types:

```

append : 'a list * 'b → 'b
append : 'a list * 'b → 'a list

```

Intuitively, the first clause has type $'a \text{ list} * 'b \rightarrow 'b$ because the first argument must match `nil`, but the second argument may be anything. The second clause has type $'a \text{ list} * 'b \rightarrow 'a \text{ list}$ because the return result is a list containing one element from the list passed as the first argument.

If we impose the constraint

```

'a list * 'b → 'b = 'a list * 'b → 'a list

```

then we must have $'b = a \text{ list}$. This gives us the final type for `append`:

```

append : 'a list * 'a list → 'a list

```

Example 6.8 Type Inference Indicates a Program Error

Here is an example that shows how type inference may produce output that indicates a programming error, even though the program may type correctly. Here is a sample (incorrect) declaration of a `reverse` function on lists:

```

- fun reverse (nil) = nil
  | reverse (x::lst) = reverse (lst);
val reverse = fn : 'a list → 'b list

```

As the compiler output shows, this function is typeable; there is no type error in this declaration. However, look carefully at the type of `reverse`. The type `'a list → 'b list` means that we can apply `reverse` to any type of list and obtain any type of list as a result. However, the type of the “reversed” list is not the same as the type of the list we started with!

Because it does not make sense for `reverse` to return a list that is a different type from its argument, there must be something wrong with this code. The problem is that, in the second clause, the first element `x` of the input list is not used as part of the output. Therefore, `reverse` always returns the empty list.

As this example illustrates, the type-inference algorithm may sometimes return a type that is more general than the one we expect. This does not indicate a type error. In this example, the faulty `reverse` can be used anywhere that a correct `reverse` function could be used. However, the type of `reverse` is useful because it tells the programmer that there is an error in the program.

6.4 POLYMORPHISM AND OVERLOADING

Polymorphism, which literally means “having multiple forms,” refers to constructs that can take on different types as needed. For example, a function that can compute the length of any type of list is polymorphic because it has type `'a list → int` for every type `'a`.

There are three forms of polymorphism in contemporary programming languages:

- *parametric polymorphism*, in which a function may be applied to any arguments whose types match a type expression involving type variables;
- *ad hoc polymorphism*, another term for overloading, in which two or more implementations with different types are referred to by the same name;
- *subtype polymorphism*, in which the subtype relation between types allows an expression to have many possible types.

Parametric and ad hoc polymorphism (overloading) are discussed in this section. Subtype polymorphism is considered in later chapters in connection with object-oriented programming.

6.4.1 Parametric Polymorphism

The main characteristic of parametric polymorphism is that the set of types associated with a function or other value is given by a type expression that contains type variables. For example, an ML function that sorts lists might have the ML type

```
sort : ('a * 'a → bool) * 'a list → 'a list
```

In words, `sort` can be applied to any pair consisting of a function and a list, as long as the function has a type of the form `'a * 'a → bool`, in which the type `'a` must also be the type of the elements of the list. The function argument is a less-than operation used to determine the order of elements in the sorted list.

In parametric polymorphism, a function may have infinitely many types, as there are infinitely many ways of replacing type variables with actual types. The sort

function, for example, may be used to sort lists of integers, lists of lists of integers, lists of lists of lists of integers, and so on.

Parametric polymorphism may be implicit or explicit. In *explicit parametric polymorphism*, the program text contains type variables that determine the way that a function or other value may be treated polymorphically. In addition, explicit polymorphism often involves explicit instantiation or type application to indicate how type variables are replaced with specific types in the use of a polymorphic value. C++ templates are a well-known example of explicit polymorphism. ML polymorphism is called *implicit parametric polymorphism* because programs that declare and use polymorphic functions do not need to contain types – the type-inference algorithm computes when a function is polymorphic and computes the instantiation of type variables as needed.

C++ Function Templates

For many readers, the most familiar type parameterization mechanism is the C++ template mechanism. Although some C++ programmers associate templates with classes and object-oriented programming, function templates are also useful for programs that do not declare any classes.

As an illustrative example, suppose you write a simple function to swap the values of two integer variables:

```
void swap(int& x, int& y){
    int tmp = x; x = y; y = tmp;
}
```

Although this code is useful for exchanging values of integer variables, the sequence of instructions also works for other types of variables. If we wish to swap values of variables of other types, then we can define a function template that uses a type variable *T* in place of the type name *int*:

```
template <typename T>
void swap(T& x, T& y){
    T tmp = x; x = y; y = tmp;
}
```

For those who are not familiar with templates, the main idea is to think of the type name *T* as a parameter to a function from types to functions. When applied to, or *instantiated* to, a specific type, the result is a version of *swap* that has *int* replaced with another type. In other words, *swap* is a general function that would work perfectly well for many types of arguments. Templates allow us to treat *swap* as a function with a type argument.

In C++, function templates are instantiated automatically as needed, with the types of the function arguments used to determine which instantiation is needed. This is illustrated in the following example lines of code.

```

int i,j; ... swap(i,j); // Use swap with T replaced with int
float a,b; ... swap(a,b); // Use swap with T replaced with float
String s,t; ... swap(s,t); // Use swap with T replaced with String

```

Comparison with ML Polymorphism

In ML polymorphism, the type-inference algorithm infers the type of a function and the type of a function application (as explained in Section 6.3). When a function is polymorphic, the actions of the type-inference algorithm can be understood as automatically inserting “template declarations” and “template instantiation” into the program. We can see how this works by considering an ML sorting function that is analogous to the C++ sort function previously declared:

```

fun insert(less, x, nil) = [x]
| insert(less, x, y::ys) = if less(x,y) then x::y::ys
                           else y::insert(less,x,ys);
fun sort(less, nil) = nil
| sort(less, x::xs) = insert(less, x, sort(less,xs));

```

For sort to be polymorphic, a less-than operation must be passed as a function argument to sort.

The types of insert and sort, as inferred by the type-inference algorithm, are

```

val insert = fn : ('a * 'a -> bool) * 'a * 'a list -> 'a list
val sort = fn : ('a * 'a -> bool) * 'a list -> 'a list

```

In these types, the type variable 'a can be instantiated to any type, as needed. In effect, the functions are treated as if they were “templates.” By use of a combination of C++ template, ML function, and type syntax, the functions previously defined could also be written as

```

template <type 'a>
fun insert(less : 'a * 'a -> bool, x : 'a, nil : 'a list) = [x]
| insert(less, x, y::ys) = if less(x,y) then x::y::ys
                           else y::insert(less,x,ys);
template <type 'a>
fun sort(less : 'a * 'a -> bool, nil : 'a list) = nil
| sort(less, x::xs) = insert(less, x, sort(less,xs));

```

These declarations are the explicitly typed versions of the implicitly polymorphic ML functions. In other words, the ML type-inference algorithm may be understood as a program preprocessor that converts ML expressions without type information

into expressions in some explicitly typed intermediate language with templates. From this point of view, the difference between explicit and implicit polymorphism is that a programming language processor (such as the ML compiler) takes the simpler implicit syntax and automatically inserts explicit type information, converting from implicit to explicit form, before programs are compiled and executed.

Finishing this example, suppose we declare a less-than function on integers:

```
- fun less(x,y) = x < y;
val less = fn : int * int -> bool
```

In the following application of the polymorphic sort function, the sort template is automatically instantiated to type `int`, so `sort` can be used to sort an integer list:

```
- sort (less, [1,4,5,3,2]);
val it = [1,2,3,4,5] : int list
```

6.4.2 Implementation of Parametric Polymorphism

C++ templates and ML polymorphic functions are implemented differently. The reason for the difference is not related to the difference between explicitly polymorphic syntax and implicitly polymorphic syntax. The need for different implementation techniques arises from the difference between data representation in C and data representation in ML.

C++ Implementation

C++ templates are instantiated at program link time. More specifically, suppose that the swap function template is stored in one file and compiled and a program calling swap is stored in another file and compiled separately. The so-called relocatable object files produced by compilation of the calling program will include information indicating that the compiled code calls a function swap of a certain type. The program linker is designed to combine the two program parts by linking the calls to swap in the calling program to the definition of swap in a separate compilation unit. It does so by instantiating the compiled code for swap in a form that produces code appropriate for the calls to swap.

If a program calls swap with several different types, then several different instantiated copies of swap will be produced. One reason that a different copy is needed for each type of call is that function swap declares a local variable `tmp` of type `T`. Space for `tmp` must be allocated in the activation record for swap. Therefore the compiled code for swap must be modified according to the size of a variable of type `T`. If `T` is a structure or object, for example, then the size might be fairly large. On the other hand, if `T` is `int`, the size will be small. In either case, the compiled code for swap must “know” the size of the datum so that addressing into the activation record can be done properly.

The linking process for C++ is relatively complex. We will not study it in detail.

However, it is worth noting that if `<` is an overloaded operator, then the correct

version of `<` must be identified when the compiled code for `sort` is linked with a calling program. For example, consider the following generic `sort` function:

```
template <typename T>
void sort( int count, T * A[count] ) {
    for (int i=0; i <count-1; i++)
        for (int j=i+1; j<count-1; j++)
            if (A[j] < A[i]) swap(A[i],A[j]);
}
```

If `A` is an array of type `T`, then `sort(n, A)` will work only if operator `<` is defined on type `T`. This requirement of `sort` is not declared anywhere in the C++ code. However, when the function template is instantiated, the actual type `T` must have an operator `<` defined or a link-time error will be reported and no executable object code will be produced.

ML Implementation

In ML, there is one sequence of compiled instructions for each polymorphic function. There is no need to produce different copies of the code for different types of arguments because related types of data are represented in similar ways. More specifically, pointers are used in parameter passing and in the representation of data structures such as lists so that when a function is polymorphic, it can access all necessary data in the same way, regardless of its type. This property of ML is called *uniform data representation*.

The simplest example of uniform data representation is the ML version of the polymorphic `swap` function:

```
- fun swap(x,y) = let val tmp = x in x := !y; y := !tmp end;
val swap = fn : 'a ref * 'a ref -> unit
```

As the type indicates, this `swap` function can be applied to any two references of the same type. Although ML references generally work like assignable variables in other languages, the value of a reference is a pointer to a cell that contains a value. Therefore, when a pair of references is passed to the `swap` function, the `swap` function receives a pair of pointers. The two pointers are the same size (typically 32 bits), regardless of what type of value is contained in the locations to which they point. In fact, we can complete the entire computation by doing only pointer assignment. As a result, none of the compiled code for `swap` depends on the size of the data referred to by arguments `x` and `y`.

Uniform data representation has its advantages and disadvantages. Because there is no need to duplicate code for different argument types, uniform data representation leads to smaller code size and avoids complications associated with C++-style linking. On the other hand, the resulting code can be less efficient, as uniform data representation often involves using pointers to data instead of storing data directly in structures.

For polymorphic list functions to work properly, all lists must be represented in exactly the same way. Because of this uniformity requirement, small values that would fit directly into the car part of a list cons cell cannot be placed there because large values do not fit. Hence we must store pointers to small values in lists, just as we store pointers to large values. ML programmers and compiler writers call the process of making all data look the same by means of pointers *boxing*.

Comparison

Two important points of comparison are efficiency and reporting of error messages. As far as efficiency, the C++ implementation requires more effort at link time and produces a larger code, as instantiating a template several times will result in several copies of the code. The ML implementation will run more slowly unless special optimizations are applied; uniform data representation involves more extensive use of pointers and these pointers must be stored and followed.

As a general programming principle, it is more convenient to have program errors reported at compile time than at link time. One reason is that separate program modules are compiled independently, but are linked together only when the entire system is assembled. Therefore, compilation is a “local” process that can be carried out by the designer or implementer of a single component. In contrast, link-time errors represent global system properties that are not known until the entire system is assembled. For this reason, C++ link-time errors associated with operations in templates can be irritating and a source of frustration.

Somewhat better error reporting for C++ templates could be achieved if the template syntax included a description of the operations needed on type parameters. However, this is relatively complicated in C++, because of overloading and other properties of the language. In contrast, the ML has simpler overloading and includes more information in parameterized constructs, allowing all type errors to be reported as a program unit is compiled.

6.4.3 Overloading

Parametric polymorphism can be contrasted with overloading. A symbol is *overloaded* if it has two (or more) meanings, distinguished by type, and resolved at compile time. In an influential historical paper, Christopher Strachey referred to ML-style polymorphism as *parametric polymorphism* (although ML had not been invented yet) and overloading as *ad hoc polymorphism*.

Example. In standard ML, as in many other languages, the operator $+$ has two distinct implementations associated with it, one of type $\text{int} * \text{int} \rightarrow \text{int}$, the other of type $\text{real} * \text{real} \rightarrow \text{real}$. The reason that both of these operations are given the name $+$ is that both compute numeric addition. However, at the implementation level, the two operations are really very different. Because integers are represented in one way (as binary numbers) and real numbers in another (as exponent and mantissa, following scientific notation), the way that integer addition combines the bits of its arguments to produce the bits of its result is very different from the way this is done in real addition.

An important difference between parametric polymorphism and overloading is that parameter polymorphic functions use one algorithm to operate on arguments

of many different types, whereas overloaded functions may use a different algorithm for each type of argument.

A characteristic of overloading is that overloading is *resolved* at compile time. If a function is overloaded, then the compiler must choose between the possible algorithms at compile time. Choosing one algorithm from among the possible algorithms associated with an overloaded function is called resolving the overloading. In many languages, if a function is overloaded, then only the function arguments are used to resolve overloading. For example, consider the following two expressions:

```
3 + 2;      /* add two integers */
3.0 + 2.0; /* add two real (floating point) numbers */
```

Here is how the compiler will produce code for evaluating each expression:

- **3 + 2:** The parsing phase of the compiler will build the parse tree of this expression, and the type-checking phase will compute a type for each symbol. Because the type-checking phase will determine that $+$ must have type $\text{int} * \text{int} \rightarrow \text{int}$, the code-generation phase of the compiler will produce machine instructions that perform integer addition.
- **3.0 + 2.0:** The parsing phase of the compiler will build the parse tree of this expression, and the type-checking phase will compute a type for each symbol. Because the type-checking phase will determine that $+$ must have type $\text{real} * \text{real} \rightarrow \text{real}$, the code-generation phase of the compiler will produce machine instructions that perform integer addition.

Automatic conversion is a separate mechanism that may be combined with overloading. However, it is possible for a language to have overloading and not to have automatic conversion. ML, for example, does not do automatic conversion.

6.5 TYPE DECLARATIONS AND TYPE EQUALITY

Many kinds of type declarations and many kinds of type equality have appeared in programming languages of the past. The reason why type declarations and type equality are related is that, when a type name is declared, it is important to decide whether this is a “new” type that is different from all other types or a new name whose meaning is equal to some other type that may be used elsewhere in the program.

Some programming languages have used fairly complicated forms of type equality, leading to confusing forms of type declarations. Instead of discussing many of the historical forms, we simply look at a few rational possibilities.

6.5.1 Transparent Type Declarations

There are two basic forms of type declaration:

- *transparent*, meaning an alternative name is given to a type that can also be expressed without this name,
- *opaque*, meaning a new type is introduced into the program that is not equal to any other type.

In the ML form of transparent type declaration,

```
type <identifier> = <type_expression>
```

the identifier becomes a synonym for the type expression. For example,

```
type Celsius = real;
type Fahrenheit = real;
```

declare two type names, Celsius and Fahrenheit, whose meaning is the type real, just the way that the two value declarations

```
val x = 3;
val y = 3;
```

declare two identifiers whose value is 3. (Remember that ML identifiers are not assignable variables; the identifier `x` will have value 3 wherever it is used.) If we declare an ML function to convert Fahrenheit to Celsius, this function will have real \rightarrow real:

```
- fun toCelsius(x) = ((x-32.0)* 0.555556);
val toCelsius = fn : real  $\rightarrow$  real
```

This should not be surprising because there is no indication that the function argument or return value has any type other than real. If we want to specify the types of the argument and result, we can add them to the function declaration. This produces a function of type Fahrenheit \rightarrow Celsius:

```
- fun toCelsius(x : Fahrenheit) = ((x-32.0)* 0.555556) : Celsius;
val toCelsius = fn : Fahrenheit  $\rightarrow$  Celsius
```

This version of the toCelsius function is more informative to read, as the types indicate the intended purpose of the function. However, because Fahrenheit and Celsius are synonyms for real, this function can be applied to a real argument:

```
- toCelsius(74.5);
val it = 23.61113 : Celsius
```

The ML type checker gives the result type `Celsius`, but because `Celsius=real`, the result can be used in real expressions.

ML abstract types are an example of an opaque type declaration. These are discussed in Chapter 9. ML data type is another form of opaque type declaration; the opacity of ML data type declarations is discussed in Subsection 6.5.3.

Two historical names for two forms of type equality are *name type equality* and *structural type equality*. Intuitively, name equality means that two type names are considered equal in type checking only if they are the same name. Structural equality means that two type names are the same if the types they name are the same (i.e., have the same structure). Although these terms may seem simple and innocuous, there are lots of confusing phenomena associated with the use of name and structural type equivalence in programming languages.

6.5.2 C Declarations and Structs

The basic type declaration construct in C is `typedef`. Here are some simple examples:

```
typedef char byte;
typedef byte ten_bytes[10];
```

the first declaring a type `byte` that is equal to `char` and the second an array type `ten_bytes` that is equal to arrays of 10 bytes. Generally speaking, C `typedef` works similarly to the transparent ML type declaration discussed in the preceding subsection. However, when structs are involved, the C type checker considers separately declared type names to be unequal, even if they are declared to name the same struct type. Here is a short program example showing how this works:

```
typedef struct {int m;} A;
typedef struct {int m;} B;
A x;
B y;
x=y; /* incompatible types in assignment */
```

Here, although the two struct types used in the two declarations are the same, the C type checker does not treat `A` and `B` as equal types. However, if we replace the two declarations with `typedef int A; typedef int B;`, using `int` in place of structs, then the assignment is considered type correct.

6.5.3 ML Data-Type Declaration

The ML data-type declaration, discussed in Subsection 5.4.4, is a form of type declaration that simultaneously defines a new type name and operations for building and making use of elements of the type. Because all of the examples in Subsection 5.4.4 were monomorphic, we take a quick look at a polymorphic declaration before discussing type equality.

Here is an example of a polymorphic data type of trees. You may wish to compare this with the monomorphic (nonpolymorphic) example in Subsection 5.4.4:

```
datatype 'a tree = LEAF of 'a | NODE of ('a tree * 'a tree);
```

This declaration defines a polymorphic type `'a tree`, with instances `int tree`, `string tree`, and so on, together with polymorphic constructors `LEAF` and `NODE`:

```
- LEAF;
val it = fn : 'a -> 'a tree
- NODE;
val it = fn : 'a tree * 'a tree -> 'a tree
```

The following function checks to see if an element appears in a tree. The function uses an exception, discussed in Section 8.2, when the element cannot be found. ML requires an exception to be declared before it is used:

```
- exception NotFound;
exception NotFound

- fun inTree(x, EMPTY) = raise NotFound
  |   inTree(x, LEAF(y)) = x = y
  |   inTree(x, NODE(y,z)) = inTree(x, y) orelse inTree(x, z);

val inTree = fn : "a * "a tree -> bool
```

Each ML data-type declaration is considered to define a new type different from all other types. Even if two data types have the same structure, they are not considered equivalent. The design of ML makes it hard to declare similar data types, as each constructor has only one type. For example, the two declarations

```
datatype A = C of int;
datatype B = C of int;
```

declare distinct types `A` and `B`. Because the second declaration follows the first and ML considers each declaration to start a new scope, the constructor `C` has type `int → B` after both declarations have been processed. However, we can see that `A` and `B` are considered different by writing a function that attempts to treat a value of one type as the other,

```
fun f(x:A) = x : B;
```

which leads to the message `Error: expression doesn't match constraint [tycon mismatch]`

6.6 CHAPTER SUMMARY

In this chapter, we studied reasons for using types in programming languages, methods for type checking, and some typing issues such as polymorphism, overloading, and type equality.

Reasons for Using Types

There are three main uses of types in programming languages:

- *Naming and organizing concepts:* Functions and data structures can be given types that reflect the way these computational constructs are used in a program. This helps the programmers and anyone else reading a program figure out how the program works and why it is written a certain way.
- *Making sure that bit sequences in computer memory are interpreted consistently:* Type checking keeps operations from being applied to operands in incorrect ways. This prevents a floating-point operation from being applied to a sequence of bits that represents a string, for example.
- *Providing information to the compiler about data manipulated by the program:* In languages in which the compiler can determine the type of a data structure, for example, the type information can be used to determine the relative location of a part of this structure. This compile-time type information can be used to generate efficient code for indexing into the data structure at run time.

Type Inference

Type inference is the process of determining the types of expressions based on the known types of some of the symbols that appear in them. For example, we saw how to infer that the function `g` declared by

```
fun g(x) = 5+x;
```

has type `int → int`. The difference between type inference and compile-time type checking is a matter of degree. A type-checking algorithm goes through the program to check that the types declared by the programmer agree with the language requirements. In type inference, the idea is that some information is not specified and some form of logical inference is required for determining the types of identifiers from the way they are used.

The following steps are used for type inference:

1. Assign a type to the expression and each subexpression by using the known type of a symbol of a type variable.
2. Generate a set of constraints on types by using the parse tree of the expression.
3. Solve these constraints by using unification, which is a substitution-based algorithm for solving systems of equations.

In a series of examples, we saw how to apply this algorithm to a variety of expressions. Type inference has many characteristics in common with the kind of algorithms that are used in compilers and programming environments to determine properties of programs. For example, some useful alias analysis algorithms that try to determine

whether two pointers might point to the same location have the same general outline as that of type inference.

Polymorphism and Overloading

There are three forms of polymorphism: parametric polymorphism, ad hoc polymorphism (another term for overloading), and subtype polymorphism. The first two were examined in this chapter, with subtype polymorphism left for later chapters on object-oriented languages. Parametric polymorphism can be either implicit, as in ML, or explicit, as with C++ templates. There are also two ways of implementing parametric polymorphism, one in which the same data representation is used for many types of data and the other in which explicit instantiation of parametric code is used to match each different data representation.

The difference between parametric polymorphism and overloading is that parametric polymorphism allows one algorithm to be given many types, whereas overloading involves different algorithms. For example, the function `+` is overloaded in many languages. In an expression adding two integers, the integer addition algorithm is used. In adding two floating-point numbers, a completely different algorithm is used for computing the sum.

Type Declarations and Type Equality

We discussed opaque and transparent type declarations. In opaque type declarations, the type name stands for a distinct type different from all other types. In transparent type declarations, the declared name is a synonym for another type. Both forms are used in many programming languages.

EXERCISES

6.1 ML Types

Explain the ML type for each of the following declarations:

- (a) `fun a(x,y) = x+2*y;`
- (b) `fun b(x,y) = x+y/2.0;`
- (c) `fun c(f) = fn y => f(y);`
- (d) `fun d(f,x) = f(f(x));`
- (e) `fun e(x,y,b) = if b(y) then x else y;`

Because you can simply type these expressions into an ML interpreter to determine the type, be sure to write a short explanation to show that you understand why the function has the type you give.

6.2 Polymorphic Sorting

This function performing insertion sort on a list takes as arguments a comparison function `less` and a list `l` of elements to be sorted. The code compiles and runs correctly:

```
fun sort(less, nil) = nil |
  sort(less, a : : l) =
  let
```

```

fun insert(a, nil) = a :: nil |
    insert(a, b :: l) = if less(a,b) then a :: (b :: l)
                        else b :: insert(a, l)
in
    insert(a, sort(less, l))
end;

```

What is the type of this sort function? Explain briefly, including the type of the subsidiary function *insert*. You do not have to run the ML algorithm on this code; just explain why an ordinary ML programmer would expect the code to have this type.

6.3 Types and Garbage Collection

Language *D* allows a form of “cast” in which an expression of one type can be treated as an expression of any other. For example, if *x* is a variable of type integer, then *(string)x* is an expression of type string. No conversion is done. Explain how this might affect garbage collection for language *D*.

For simplicity, assume that *D* is a conventional imperative language with integers, reals (floating-point numbers), pairs, and pointers. You do not need to consider other language features.

6.4 Polymorphic Fixed Point

A *fixed point* of a function *f* is some value *x* such that $x = f(x)$. There is a connection between recursion and fixed points that is illustrated by this ML definition of the factorial function $\text{factorial} : \text{int} \rightarrow \text{int}$:

```

fun Y f x = f (Y f) x;
fun F f x = if x=0 then 1 else x*f(x-1);
val factorial = Y F;

```

The first function, *Y*, is a fixed-point operator. The second function, *F*, is a function on functions whose fixed point is factorial. Both of these are curried functions; using the ML syntax $\text{fn } x \Rightarrow \dots \text{ for } \lambda x \dots$, we could also write the function *F* as

```

fun F(f) = fn x =>
    if x=0 then 1 else x*f(x-1)

```

This *F* is a function that, when applied to argument *f*, returns a function that, when applied to argument *x*, has the value given by the expression $\text{if } x=0 \text{ then } 1 \text{ else } x*f(x-1)$.

- What type will the ML compiler deduce for *F*?
- What type will the ML compiler deduce for *Y*?

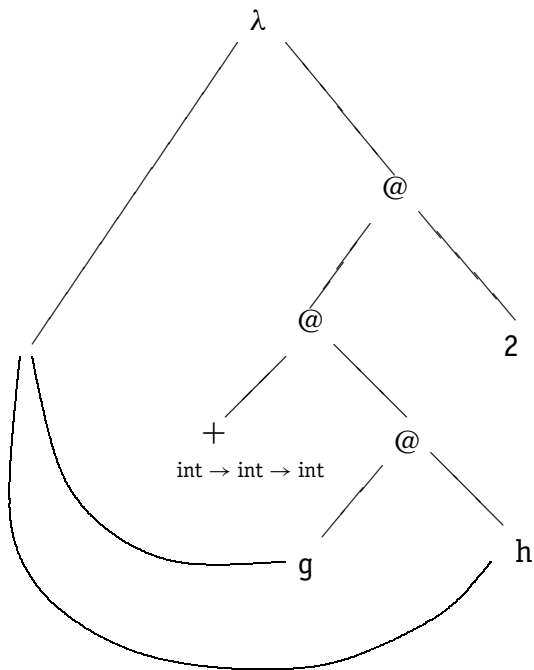
6.5 Parse Graph

Use the following parse graph to calculate the ML type for the function

```

fun f(g,h) = g(h) + 2;

```

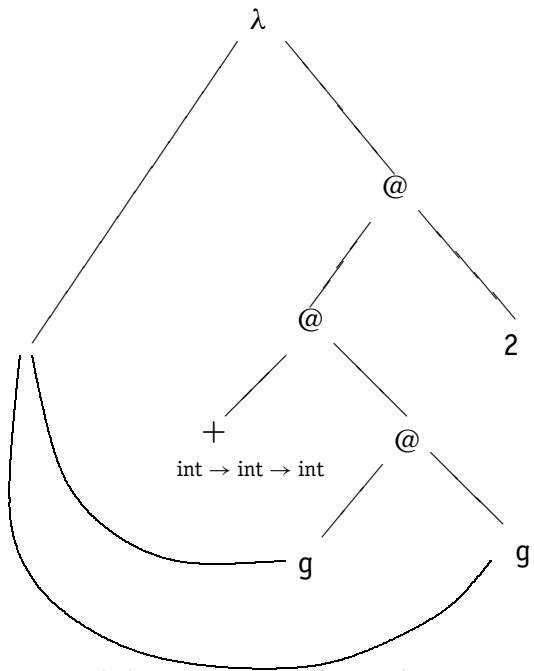


6.6 Parse Graph

Use the following parse graph to follow the steps of the ML type-inference algorithm on the function declaration

```
fun f(g) = g(g) + 2;
```

What is the output of the type checker?



6.7 Type Inference and Bugs

What is the type of the following ML function?

```
fun append(nil, l) = l
|   append(x : l, m) = append(l, m);
```

Write one or two sentences to explain succinctly and informally why `append` has the type you give. This function is intended to append one list onto another. However, it has a bug. How might knowing the type of this function help the programmer to find the bug?

6.8 Type Inference and Debugging

The `reduce` function takes a binary operation, in the form of a function `f`, and a list, and produces the result of combining all elements in the list by using the binary operation. For example;

```
reduce plus [1,2,3] = 1 + 2 + 3 = 6
```

if `plus` is defined by

```
fun plus (x, y : int) = x+y
```

A friend of yours is trying to learn ML and tries to write a `reduce` function. Here is his incorrect definition:

```
fun reduce(f, x) = x
|   reduce(f, (x : t)) = f(x, reduce(f, y));
```

He tells you that he does not know what to return for an empty list, but this should work for a nonempty list: If the list has one element, then the first clause returns it. If the list has more than one element, then the second clause of the definition uses the function `f`. This sounds like a reasonable explanation, but the type checker gives you the following output:

```
val reduce = fn : (('a * 'a list) -> 'a list) * 'a list -> 'a list
```

How can you use this type to explain to your friend that his code is wrong?

6.9 Polymorphism in C

In the following C `min` function, the type `void` is used in the types of two arguments. However, the function makes sense and can be applied to a list of arguments in which `void` has been replaced with another type. In other words, although the C type of this function is not polymorphic, the function could be given a polymorphic type if C had a polymorphic type system. Using ML notation for types, write a type for this `min` function that captures the way that `min` could be meaningfully applied to arguments of various types. Explain why you believe the function has the type you have written.

```
int min (
    void *a[],          /* a is an array of pointers to data of unknown type */
    int n,              /* n is the length of the array */
    int (*less)(void*, void*) /* parameter less is a pointer to function */
)                      /* that is used to compare array elements */
{
    int i;
    int m;
    m=0;
```

```

    for (i=1; i < n; i++)
        if (less(a[i], a[m])) m=i;
    return(m);
}

```

6.10 Typing and Run-Time Behavior

The following ML functions have essentially identical computational behavior,

```

fun f(x) = not f(x);
fun g(y) = g(y) * 2;

```

because except for typing differences, we could replace one function with the other in any program without changing the observable behavior of the program. In more detail, suppose we turn off the ML type checker and compile a program of the form $\mathcal{P}[\text{fun } f(x) = \text{not } f(x)]$. Whatever this program does, the program $\mathcal{P}[\text{fun } f(y) = f(y) * 2]$ we obtain by replacing one function definition with the other will do exactly the same thing. In particular, if the first does not lead to a run-time type error such as adding an integer to a string, neither will the second.

- (a) What is the ML type for f ?
- (b) What is the ML type for g ?
- (c) Give an informal explanation of why these two functions have the same run-time behavior.
- (d) Because the two functions are equivalent, it might be better to give them the same type. Why do you think the designers of the ML typing algorithm did not work harder to make it do this? Do you think they made a mistake?

6.11 Dynamic Typing in ML

Many programmers believe that a run-time typed programming language like Lisp or Scheme is more expressive than a compile-time typed language like ML, as there is no type system to “get in your way.” Although there are some situations in which the flexibility of Lisp or Scheme is a tremendous advantage, we can also make the opposite argument. Specifically, ML is more expressive than Lisp or Scheme because we can define an ML data type for Lisp or Scheme expressions.

Here is a type declaration for pure historical Lisp:

```

datatype LISP = Nil
    | Symbol of string
    | Number of int
    | Cons of LISP * LISP
    | Function of (LISP -> LISP)

```

Although we could have used (Symbol “nil”) instead of a primitive Nil, it seems convenient to treat nil separately.

- (a) Write an ML declaration for the Lisp function *atom* that tests whether its argument is an atom. (Everything except a cons cell is an atom – The word *atom* comes from the Greek word *atomos*, meaning indivisible. In Lisp, symbols, numbers, nil, and functions cannot be divided into smaller pieces, so they are considered to be atoms.) Your function should have type $\text{LISP} \rightarrow \text{LISP}$, returning atoms Symbol(“T”) or Nil.

- (b) Write an ML declaration for the Lisp function `islist` that tests whether its argument is a *proper* list. A proper list is either `nil` or a cons cell whose `cdr` is a proper list. Note that not all listlike structures built from cons cells are proper lists. For instance, `(Cons (Symbol("A"), Symbol("B")))` is not a proper list (it is instead what is known as a dotted list), and so `(islist (Cons (Symbol("A"), Symbol("B"))))` should evaluate to `Nil`. On the other hand, `(Cons (Symbol("A"), (Cons (Symbol("B"), Nil))))` is a proper list, and so your function should evaluate to `Symbol("T")`. Your function should have type $LISP \rightarrow LISP$, as before.
- (c) Write an ML declaration for Lisp `car` function and explain briefly. The function should have type $LISP \rightarrow LISP$.
- (d) Write Lisp expression `(lambda (x) (cons x 'A))` as an ML expression of type $LISP \rightarrow LISP$. Note that 'A means something completely different in Lisp and ML. The 'A here is part of a Lisp expression, not an ML expression. Explain briefly.