

Principles of Programming Languages

Achim Blumensath

Spring Semester 2017

Contents

1	Expressions and Functions	1
1.1	Arithmetic expressions	1
1.2	Local definitions	1
1.3	Functions	2
1.4	Static and dynamic scoping	3
1.5	Higher-order and first-class functions	5
1.6	Function parameters	6
1.7	Conditionals	7
1.8	Constructors and pattern matching	8
1.9	Recursion	10
1.10	Lazy evaluation	14
1.11	Programming examples	15
2	Types	19
2.1	Static and dynamic typing	19
2.2	Type annotations	20
2.3	Common Types	21
2.4	Type checking	24
2.5	Polymorphism	26
2.6	Type inference	27
3	State and Side-Effects	31
3.1	Assignments	31
3.2	Ramifications	32
3.3	Parameter passing	35
3.4	Memory management	39
3.5	Loops	41
3.6	Programming Examples	43
4	Modules	45
4.1	Simple modules	45
4.2	Encapsulation	45
4.3	Abstract Data Types	46
4.4	Module expressions	48
5	Control-Flow	51
5.1	Continuation passing style	51
5.2	Continuations	54

Contents

5.3	Generators	54
5.4	Exceptions	56
6	Constraints	59
6.1	Single-assignment variables	59
6.2	Unification	60
6.3	Backtracking	61
6.4	Programming examples	63
7	Objects	67
7.1	Dynamic dispatch	67
7.2	Subtyping	73
7.3	Encapsulated state	75
7.4	Inheritance	77
7.5	Discussion	83
8	Concurrency	85
8.1	Fibres	85
8.2	Ramifications	89
8.3	Message passing	89
8.4	Shared-memory	94

1 Expressions and Functions

1.1 Arithmetic expressions

In one form or other *expressions* are present in nearly every programming language. In the abstract, an expression can be defined as a program construct that computes a value. A prototypical example are arithmetical expressions in mathematics. In so-called *functional* languages, expressions form the central construct around which the whole language is built. In this chapter we will introduce a functional kernel language, starting with arithmetical expressions.

$$\langle expr \rangle ::= \langle num \rangle \mid (\langle expr \rangle) \mid \langle expr \rangle + \langle expr \rangle \mid \langle expr \rangle * \langle expr \rangle$$

The evaluation strategy for such expressions is obvious: we recursively evaluate all subexpressions and then combine their results using the operation at the current position. For instance,

```
1  1+2*3
2  => 7
```

In this case we can add the new operation as *syntactic sugar*, i.e., we express it in terms of the operations the core language already provides.

$$expr_1 - expr_2 \implies expr_1 + (-1) * expr_2 .$$

Hence, after parsing, but before any further analysis, we replace every occurrence of subtraction by its definition. Of course, this only works if we can express the new feature in terms of the old ones.

1.2 Local definitions

One central mechanism to improve the readability and maintainability of code is the ability to *name* a given program construct like an expression, a type, etc. and to refer to that construct using its new name. This process is called *abstraction*. It is vital in breaking a program into smaller, easier to understand parts. We can use it to hide unimportant details and thereby decrease complexity of our code. In addition abstraction also facilitates code reuse.

At the moment our kernel language has only one construct: expressions. To name an expression, we introduce *local definitions*.

$$\langle expr \rangle ::= \dots \mid \langle id \rangle \mid \mathbf{let} \langle id \rangle = \langle expr \rangle ; \langle expr \rangle$$

A remark on terminology: in our current setting without side-effects, we refer to these names as *identifiers*, instead of using the more common term ‘variables’. We reserve the latter for *mutable identifiers*, which we will introduce in Chapter 3.

1 Expressions and Functions

Examples:

```
1  let x = 1;                let pi = 3; // the integer version ;-)  
2  let y = 2;                2*pi*5  
3  x + 2*y                  => 30  
4  => 5  
5  
6  let x = (let y = 2; 2*y);  (let x = 2; x * x) - (let x = 1; x+4)  
7  x + 3                    => -1  
8  => 7
```

Apart from making code more readable and easier to write, local definitions can also be used to improve performance. If a complicated expression is used in several places, we can use a let-binding to evaluate the expression only once and then refer to its value via the corresponding identifier. For instance, if we want to rotate a vector, we only need to compute the sine and cosine once.

```
1  let s = ... compute the sine ...  
2  let c = ... compute the cosine ...  
3  u = c * x - s * y;  
4  v = s * x + c * y;
```

When introducing let-bindings a new phenomenon arises called *scoping*. The problem is, when we try to evaluate an expression and come upon an identifier x , which of the possibly several definitions for x do we use? The part of the code where a particular definition of x is in effect is called the *scope* of the definition. In our case, the scope of a definition $\mathbf{let\ } x = e; e'$ is the expression e' . That is, every occurrence of x inside e' refers to the value e . Other occurrences of x (for instance, those in e or in other parts of the program) refer to other definitions. We also say that this definition of x is *local* to e' and that the variable x is *bound* (to the value e) by this definition. In general the association of *names* in a program with the *entities* they refer to is called *binding*. The characteristic property of a local variable is that it can be *renamed* without changing the meaning of the program. (The technical term for such a renaming is α -conversion.)

$\mathbf{let\ } x = 2; x*x \iff \mathbf{let\ } y = 2; y*y$

In most languages scopes can be nested, but they cannot partially overlap. Therefore, they are usually implemented using a stack.

```
1  let x = 2;  
2    let y = x-1;  
3    x+y          } scope of y    } scope of x
```

1.3 Functions

Next we add function definitions to our language. Function definitions are one of the main mechanisms for *control abstraction* in programming languages. They facilitate *code reuse* and they can increase the *readability* of code by hiding low-level details and thereby revealing the logical

structures of the code. But note that overuse of this feature can degrade readability again, if functionality is split over too many places of the code. (This is a common problem with inheritance in object-oriented programming.)

For efficiency reasons, many languages (like C++ and Java) only support *non-nested functions*. In this case, a program is of the form

```

1  let f1(x) { expr };
2  ...
3  let fn(x) { expr };
4  expr

```

As we have seen when implementing non-nested functions, we can evaluate the body of a globally defined function in the empty environment. When allowing nested functions, we have to use the environment of the function definition instead. This complicates the implementation since we have to store this environment somewhere. We extend our language as follows.

$$\langle expr \rangle ::= \dots \mid \langle id \rangle (\langle expr \rangle) \mid \mathbf{let} \langle id \rangle (\langle id \rangle) \{ \langle expr \rangle \}; \langle expr \rangle$$

1.4 Static and dynamic scoping

When invoking a function, *static scoping* evaluates the function body in the scope of the function's *definition*, while *dynamic scoping* uses the scope of the function's *caller*.

Examples:

1	let x = 1;	let x = 1;	let x = 1;
2	let f(y) { x };	let g(y) { x };	let g(y) { x };
3	let x = 2;	let f(y) { g(y) };	let f(x) { g(0) };
4	f(3)	let x = 2;	let x = 2;
5		f(3)	f(3)

Dynamic scoping Examples of languages using dynamic scope are: the original Lisp, Emacs Lisp, TeX, Perl, and many scripting languages including early versions of Python and JavaScript.

Today, dynamic scoping is generally considered to be a mistake. The main problem is that dynamic scoping is not robust: changing local variables in some part of the program can have drastic influences on other parts. Hence, with dynamic scoping bound variables are not local in the sense defined above since renaming them *can* change the meaning of the program. This means that understanding code with dynamic scoping requires *global reasoning* about the program, which violates one of the fundamental principles of code readability.

For instance, consider a GUI library that provides an event-loop where the program can install call-backs to react to user input. If the event-loop and the user code happen to share a variable, the call-back will get the event-loops variable instead of its own. Problems of dynamic scoping include:

1 Expressions and Functions

- As seen in the above example, with dynamic scoping, every 3rd party library needs to document the names of all local variables it uses. This makes library maintenance more difficult, as new versions cannot introduce new local variables.
- Dynamic scoping also presents a security risk as it enables other parts of the code to access and modify sensitive information stored in local variables.

Let us conclude with an example of a programming idiom that is enabled by dynamic scope: one can use it to simulated *default parameters* for functions. If a certain function parameter has nearly always the same value, one can use a variable instead. For example, if we write a function converting numbers to strings we might want to support other bases than decimal. The code could look like this.

```
1  let base = 10;
2  let num_to_string(n) {
3    ... convert n into base base ...
4  };
5
6  let f(x) {
7    ...
8    let base = 16;
9    let str = num_to_string(137);
10   ...
11  };
```

Of course, with languages supporting default parameters one could simply write:

```
1  let num_to_string(n, base=10) {
2    ... convert n into base base ...
3  };
4
5  let f(x) {
6    ...
7    let str = num_to_string(137,16);
8    ...
9  };
```

Static scoping While static scoping is clearly superior to dynamic scoping, it is not without its problems. The way it is usually implemented, the scoping structure of a program is determined by its syntactic structure. This is a very simple way to specify scoping rules, which is not always adequate. Sometimes one would like to have more fine-grained control over scoping, say, by specifying which parts of the program are allowed to see a given identifier. Some languages have therefore tried to untie scoping from the syntactic structure by making it explicit. One example is the concept of *namespaces* in C++, which allows complete control over scoping. Slightly less general are *modules* or *packages* which are supported by most modern languages.

1.5 Higher-order and first-class functions

In many languages, functions are not values. You cannot assign them to variables or pass them as arguments to functions. Some languages, like C and C++, allow passing functions as arguments, but not returning them as results. This can be used for example to implement call-backs in GUI frameworks. Such languages support what is called *higher-order functions*.

```

1  let f(x) { x+1 };
2  let g(s) { s(1) };
3  g(f)
4  => 2

```

In some languages, like Lisp, ML, or JavaScript, functions are values like any others. In this case we speak of *first-class function*. In such languages, we need an operation to create new functions. This is called a *lambda abstraction*.

$$\langle expr \rangle ::= \dots \mid \mathbf{fun} \ (\langle id \rangle) \ \{ \langle expr \rangle \}$$

Example

```

1  let adder(n) { fun(x) { x + n } };
2  let add3 = adder(3);
3  add3(4)
4  => 7

```

First-class functions are frequently used, for instance, in GUI frameworks where they are called *call-backs*.

```

1  let mouse_button(button, x, y) {
2    ...
3    react to a mouse button being pressed
4    ...
5  };
6  register_call_back(MouseDown, mouse_button);

```

In functional programming, first-class functions are one of the main concepts used for abstraction. They allow the separation of the *action* to be performed on some data structure from the *traversal* of said data structure. For instance,

```

1  map(update, lst)    applies update to every element of lst
2  fold(sum, 0, lst)  adds all elements of lst

```

When using dynamic scoping first-class functions cause additional problems. Traditionally there are two possible ways to implement dynamic scoping for such functions: *shallow binding* and *deep binding*. The question is, which environment is used when calling a function value. With deep binding it is the environment where the function was declared, i.e., the same as when using static scoping. With shallow binding it is the environment of the function call instead, which is more in the spirit of dynamic binding.

1 Expressions and Functions

```
1  let f(x) { fun(y) { x } };
2  let x = 3;
3  f(1)(2)
```

Finally, note that, once we have first-class functions, we can simplify our kernel language by removing the let-construct and implementing it as syntactic sugar instead.

```
1  let x =  $expr_1$ ;  $expr_2$    $\implies$   (fun (x) {  $expr_2$  })( $expr_1$ )
```

1.6 Function parameters

Multiple parameters As function application is one of the most frequently used mechanism in programming, many languages provide features making it more convenient. The first such feature we consider are functions with multiple arguments. There are two ways to add such functions to our language. The first one implements them as syntactic sugar in terms of first-class functions. This is called *currying*. It is present in many functional languages like OCaml or Haskell. The idea is simple. We view a function with two parameters as a function that take the first argument and returns a function taking the second argument and returning the result.

```
fun (x,y) { x*x + y*y }
```

```
 $\implies$  fun (x) { fun(y) { x*x + y*y } }
```

In our kernel language, we can implement currying as syntactic sugar. We translate

```
let f(x,y,...,z) { body };  $expr$ 
```

```
 $\implies$  let f(x) { fun (y) { ... fun (z) { body } ... } };
 $expr$ 
```

and

```
f(a,b,...,c)  $\implies$  f(a)(b)...(c)
```

Note that this syntactic sugar allows us to use partially applied functions, that is, expressions like the following one.

```
1  let f(x,y) { ... };
2  f(1)
```

If we do not want to allow this, we have to add a pass for arity checks before doing the desugaring.

The other way of implementing functions with multiple parameters does not require first-class functions, but uses a tuple datatype instead. Instead of passing several arguments to a function, we pass a single tuple containing them. This is done for example in Standard ML.

```
fun (x,y) { x*x + y*y }
```

```
 $\implies$  fun (p) { p.x * p.x + p.y * p.y }
```

Keyword parameters One such feature are *named parameters* or *keyword parameters*. Ordinarily, arguments are passed to a function by *position*, that is, the i -th argument will be bound to the i -th formal parameter of the function. If a function takes many arguments, it becomes hard to remember the correct order of the parameters. In many languages it is therefore possible to assign names to the parameters and use these names when invoking a function. In this case, the arguments can be given in any order.

```

1  let f(serial_number, price, weight) { ... };
2
3  f(serial_number = 83927, weight = 60, price = 120);

```

To avoid ambiguities, if both positional and keyword parameters are used in the same function call, one usually requires all positional parameters to be listed first.

Default arguments Another feature are *default arguments*. When a function is frequently called with a fixed value for some argument, one can specify this value as the default and allow the programmer to omit the argument from a function call.

```

1  let int_to_string(num, base = 10) { ... };
2
3  int_to_string(17)

```

To avoid ambiguities, if positional parameters are used in a function call where some default arguments is omitted, one usually requires all arguments after the omitted one to be keyword parameters.

Variable number of arguments Some languages like C allow the definition of functions where the number of arguments is not fixed. There is a minimal number of arguments, but every function invocation can use more if needed.

```

1  let printf(format, ...) { ... };
2
3  printf("f(%d) = %d", x, f(x));

```

Conceptually what happens is that the first arguments are passed to the function as usual and the remaining arguments are passed in an array which the function body can inspect.

1.7 Conditionals

In preparation for adding recursion, let us implement conditionals first. These are needed to add a termination condition to a recursive function call. For simplicity, we only support equality predicates.

$$\langle expr \rangle ::= \dots \mid \text{if } \langle expr \rangle == \langle expr \rangle \text{ then } \langle expr \rangle \text{ else } \langle expr \rangle$$

Example

```

1  let f(n) {
2    if n == 0 then
3      0
4    else
5      n-1
6  };

```

There are two approaches to boolean values in programming languages. Languages with a strict type discipline define a type for boolean values and demand that the condition in an if-statement is of that type. Languages with a looser type discipline allow the condition to be of a different type and automatically coerce it to a boolean values. For such languages one uses the terminology of *truthy* values (those that are treated as *true*) and *falsy* values (those that are treated as *false*).

Automatic coercions are more convenient, but also more error-prone and make the type system much more complicated. (In general, coercions also make the code harder to understand, but for booleans in conditionals that is not the issue.) While for languages with a simple type system like, say, C, the rules for boolean conversions are easily understood and remembered. But for languages with a richer type system like JavaScript, Python, or Ruby, these rules become very complicated. (Is the empty array *false*? What about the empty string, or the string "0"? Are "00" and "0.0" treated the same as "0"?) What makes matters worse is that none of these languages agree on the precise rules.

1.8 Constructors and pattern matching

So far, our kernel language does not support any composite data structures. We only have numbers and functions. To add composite types like records and arrays, we need operations that create new data objects. In imperative languages like Java there are usually two different kind of such operations. There is an operation like `new` that allocates a piece of memory that then has to be initialised by the programmer. Furthermore, some of the types allow the programmer to write down values of the type directly. These constructs are called *literals* or (*data*) *constructors*. In a language without side-effects, we cannot initialise a data structure after it has already been allocated. We have to do both in one step. Therefore such languages usually only have constructors.

For our kernel language we will provide several built-in constructors and also allow user-defined ones. Each number literal is treated as a constructor. Furthermore, we have constructors for records, two constructors `True` and `False` for the boolean values, constructors `()` `Pair(y, y)` for the empty tuple and pairs, and two constructors `Cons(x, y)` and `Nil` to build lists. Using these last two constructors, we can represent a list like `[1, 2, 3]` in the form

$$\text{Cons}(1, \text{Cons}(2, \text{Cons}(3, \text{Nil}))) .$$

The more convenient notation `[1, 2, 3]` will be provided as syntactic sugar.

We add three new constructs to the language. We can define new constructors, we can call a constructor to create a data structure, and we can match a given data structure with a template to

extract its fields.

$$\begin{aligned}
 \langle expr \rangle ::= & \dots \mid \mathbf{type} \langle id \rangle = \mid \langle variant \rangle \dots \mid \langle variant \rangle ; \langle expr \rangle \\
 & \mid \mathbf{type} \langle id \rangle = [\langle id \rangle = \langle id \rangle , \dots , \langle id \rangle = \langle id \rangle] ; \langle expr \rangle \\
 & \mid \langle ctor \rangle (\langle expr \rangle , \dots , \langle expr \rangle) \\
 & \mid [\langle id \rangle = \langle expr \rangle , \dots , \langle id \rangle = \langle expr \rangle] \\
 & \mid \langle expr \rangle . \langle id \rangle \\
 & \mid \mathbf{case} \langle expr \rangle \mid \langle pattern \rangle \Rightarrow \langle expr \rangle \mid \dots \mid \langle pattern \rangle \Rightarrow \langle expr \rangle \\
 \langle pattern \rangle ::= & \langle id \rangle \mid \langle num \rangle \mid \langle ctor \rangle (\langle id \rangle , \dots , \langle id \rangle) \mid \mathbf{else} \\
 \langle variant \rangle ::= & \langle id \rangle \mid \langle id \rangle (\langle id \rangle , \dots , \langle id \rangle)
 \end{aligned}$$

For instance, we can create a pair and extract its two components again using the following definitions. (The arguments a and b in the definition of the constructor P are only used to specify the arity. Later on when we add a type system, these parameters will specify the types of the constructor arguments.)

```

1  type int_pair = | P(int, int);           type int_pair = [ x : int, y : int ];
2
3  let make_pair(x,y) { P(x,y) };          let make_pair(x,y) { [ x = x, y = y ] };
4
5  let fst(p) {                               let fst(p) { p.x };
6    case p
7    | P(x,y) => x
8  };
9
10 let snd(p) {                               let snd(p) { p.y };
11   case p
12   | P(x,y) => y
13 };

```

Similarly, we can define the following functions to create and traverse lists.

```

1  let empty_list      = Nil;
2  let add_to_list(x, lst) = Cons(x, lst);
3
4  let is_nil(lst) { case lst | Nil      => True | else => False };
5  let is_cons(lst) { case lst | Cons(x, xs) => True | else => False };
6  let head(lst) { case lst | Cons(x, xs) => x };
7  let tail(lst) { case lst | Cons(x, xs) => xs };

```

With the case-construct we can now implement if- and (non-recursive) let-statements as syntactic sugar.

$$\begin{aligned}
 \mathbf{if} \ c_0 == c_1 \ \mathbf{then} \ t \ \mathbf{else} \ e & \implies \mathbf{case} \ c_0 - c_1 \mid \emptyset \Rightarrow t \mid \mathbf{else} \Rightarrow e \\
 \mathbf{let} \ x = e; \ e' & \implies \mathbf{case} \ e \mid x \Rightarrow e'
 \end{aligned}$$

1 Expressions and Functions

In fact, we can now define the equality predicate explicitly and introduce a version of the if-statement that uses an arbitrary predicate.

$$\begin{aligned} e == e' & \implies \text{case } e - e' \mid 0 \Rightarrow \text{True} \mid \text{else} \Rightarrow \text{False} \\ \text{if } c \text{ then } t \text{ else } e & \implies \text{case } c \mid \text{True} \Rightarrow t \mid \text{False} \Rightarrow e \end{aligned}$$

Exercise Use case statements to define syntactic sugar for and and or operations that evaluate their arguments only as needed (short-circuit evaluation).

```
1  e1 and e2  -> case e1 | True => e2   | False => False
2  e1 or e2   -> case e1 | True => True  | False => e2
```

Exercise Introduce syntactic sugar for lists.

```
1  [e1, ..., en]  -> Cons(e1, Cons(..., Cons(en, Nil)...))
2  [e1, ..., en | e] -> Cons(e1, Cons(..., Cons(en, e)...))
```

Introduction and elimination forms Let us conclude this section with a remark an *introduction* and *elimination* constructs. In many aspects of a programming language, we can observe a duality between constructs introducing a certain object and ones eliminating it again. For instance, with data types we have (I) constructors to assemble a structure and (E) the case-statement to disassemble it into its components again. Similarly, for functions we have (I) lambda abstractions which create new functions and (E) function applications which turn functions into their return value.

1.9 Recursion

Every serious programming language needs a mechanism for *unbounded recursion* or iteration. For instance, we would like to define recursive functions like the following one.

```
let fac(n) { if n == 0 then 1 else n * fac(n-1) };
```

(Which is, in fact, a *bounded* (by n) recursion.) Implementing recursion is rather straightforward: in let-bindings `let x = e; e'` we just have to extend the scope of x to include e .

Note that, while straightforward, the addition of recursion *does* change the language considerably. In particular, it is now very easy to write *non-terminating programs*. (Strictly speaking, this is also possible in our old language with non-recursive let-bindings, but it requires some tricks and a lot of effort to do so, see below.)

From a theoretical perspective, this addition is much more involved, and books on programming language theory usually devote quite some space to the topic. The problem is how to implement recursion without using recursion. (We cheated in our implementation by using the built-in recursion of Haskell.) There are two ways to get around this problem.

The first one requires mutable state. When defining a recursive function f , we first allocate a variable for it (initialised with some dummy value), then we write the actual function into the variable using an assignment.

```

1  let f = fun (x) { x };          // dummy value
2  let f' = fun (x) { ... body using f ... }
3  f := f'

```

This is what most real language implementations do.

The second solution is much cleaner from a theoretical point of view. We add a *recursion operator* (also called a *fixed-point operator*) to the language which is defined by

$$\mathbf{rec}(f) = f(\mathbf{rec}(f))$$

(Of course, this is a recursive definition itself.) Then we can write

```

1  let fac_body(f) {
2    fun (n) { if n == 0 then 1 else n * f(n-1) }
3  };
4  let fac = rec(fac_body);

```

`fac_body(f)` is the body of the factorial function where we have replaced the recursive call by a call to the function `f`. Then we tie the knot by defining

$$\mathbf{fac} = \mathbf{rec}(\mathbf{fac_body}) = \mathbf{fac_body}(\mathbf{rec}(\mathbf{fac_body})) = \mathbf{fac_body}(\mathbf{fac})$$

Intuitively, the `rec` operator provides a marker indicating that 'at this position there is a recursive call'. Whenever such a marker is evaluated, we insert the body of the corresponding function (where all recursive calls are marked by `rec` again).

If our language is untyped (or if the type system supports recursive types), we can actually define `rec` as syntactic sugar.

$$\mathbf{let\ rec}(f) = (\mathbf{fun\ (x)\ \{ f(x(x)) \}})(\mathbf{fun\ (x)\ \{ f(x(x)) \}});$$

Then `rec(f)` evaluates to `f(rec(f))` (try it).

Simultaneous recursion Our `let`-construct only allows the definition of a single recursive function. Sometimes one would like to define several mutually recursive functions like

```

1  let f(x) { if x = 0 then 1 else g(x-1) };
2  let g(x) { if x = 0 then 1 else 1+f(x-1) };

```

There are three ways to implement such definitions. The first one is to extend the syntax of `let`-bindings to allow for the simultaneous definition of several identifiers. This is the most practical solution and implemented in all serious programming languages. In our kernel language we will not take this approach (just to make our life easier, at the cost of making the programmers' life harder). We can do so because simultaneous recursion can be implemented using single recursion. Suppose we have a definition like

```

1  let x = f(x,y) and y = g(x,y);
2  h(x,y);

```

We can either transform it into

```

1  let x = f(x, (let y = g(x,y); y));
2  let y = g(x,y);
3  h(x,y)

```

1 Expressions and Functions

or, if the language supports tuple or records (see the next section), we can use them to write

```
1  let (x,y) = (f(x,y), g(x,y));
2  h(x,y);
```

The first solution duplicates some code (which is then executed twice), the second one has to allocate memory for the tuple. In most cases this overhead is negligible.

Recursive data structures Since we are already talking about data structures, let us also mention the related problem of creating mutually recursive data structures. The most practical solution is again to use mutable data structures. Then we can (i) first allocate all the memory and then (ii) initialise it. For instance, to create two pairs

```
let p = (1, q) and q = (2, p);
```

we can write

```
1  let p = (1, 0);
2  let q = (2, 0);
3  p.2 := q;
4  q.2 := p;
```

Tail calls Finally, let us mention an important implementation detail. In a programming language where the only mechanism for unbounded iteration is recursion, it is essential that this feature is usable even if the number of iterations is large. For every recursive call, we have to allocate memory to store parameters and local variables. In a naive implementation this memory will only be freed once all recursive calls have returned. This leads to a memory consumption that is linear in the number of recursive calls, which is problematic if this number is large. There is an important situation where we can free this memory earlier: if the recursive call is the last expression of our function, i.e., the return value of the function is the value returned by the recursive call.

```
1  let find_next_prime(n) {
2    if n is prime then
3      n
4    else
5      find_next_prime(n+1)
6  };
1  let fac(n) {
2    if n == 0 then
3      1
4    else
5      n * fac(n-1)
6  };
```

In the situation on the left, we will not need the parameters and local variables after the recursive call has returned. Hence, we can free the memory containing them before the call to `find_next_prime` instead of after it. This is called a *tail-call optimisation*. It amounts to replacing the recursive call by a jump to the beginning of the function.

```
1  let find_next_prime(n) {
2    label start;
3    if n is prime then
4      n
```



```

5   else
6     (n := n+1; goto start)
7   };

```

After this transformation the function will use a constant amount of memory and is as efficient as an imperative solution using a while-loop.

Frequently, it is possible to transform a recursive definition that uses non-tail calls into a tail-call one by using an *accumulator*. For instance, we can define the factorial function as

```

1  let fac(n) {
2    let multiply(n, acc) {
3      if n == 0 then
4        acc
5      else
6        multiply(n-1, n*acc)
7    };
8    multiply(n,1);
9  };

```

After tail-call optimisation, this looks like

```

1  let fac(n) {
2    let multiply(n, acc) {
3      label start;
4      if n == 0 then
5        acc
6      else (
7        new_n   := n-1;
8        new_acc := n*acc;
9        n       := new_n;
10       acc     := new_acc;
11       goto start;
12     );
13     acc := 1;
14     goto start;
15   };

```

which (after some trivial optimisations) is equivalent to the imperative code

```

1  let fac(n) {
2    let acc = 1;
3    while n > 0 {
4      acc := n * acc;
5      n  := n - 1;
6    };
7    return acc;
8  };

```

1.10 Lazy evaluation

Since our language does not support side-effects, the order in which we evaluate expressions does not matter. Any order we choose produces the same result.

1	fun (x) {1+x*x} (1+1)	fun (x) {1+x*x} (1+1)
2	=> fun (x) {1+x*x} 2	=> 1+(1+1)*(1+1)
3	=> 1+2*2	=> 1+2*(1+1)
4	=> 1+4	=> 1+2*2
5	=> 5	=> 1+4
6		=> 5

There are two canonical orders in which we can evaluate expressions:

- *eager evaluation* evaluates an expression starting with the left-most, inner-most operations, while
- *lazy evaluation* starts with the left-most, outer-most operation.

Advantages of lazy evaluation It can be shown that lazy evaluation is more powerful than eager evaluation in the following sense: every computation that terminates using an *arbitrary* evaluation order also terminates with lazy evaluation and produces the same result. On the other hand, there are expressions that terminate with a result using lazy evaluation but not with eager evaluation.

It has turned out that there are two main areas where this property of lazy evaluation makes it superior to eager evaluation:

- (1) processing of infinite data structures and
- (2) evaluations of (mutually) recursive definitions.

(1) Using data constructors with lazy evaluation, it is very simple to define and process infinite data structures like infinite lists.

```

1  let ones = [1 | ones];
2  ones
3
4  let numbers i = [i | number(i+1)];
5  numbers
6
7  let add(x,y) { x+y };
8  let fib = [0, 1 | map2(add, fib, (tail(fib)))];
9  fib

```

But note that this means that there are no inductive lazy datatypes. For example

```
type nat = Zero | Succ(nat)
```

does not define the natural numbers since

```
let omega = Succ(omega);
```

defines the infinite number $\text{Succ}(\text{Succ}(\text{Succ}(\dots)))$. Hence, in order to be able to define inductive datatypes like natural numbers, finite lists, or finite trees, a lazy language must have support for eagerly evaluated data constructors.

(2) The definition of the Fibonacci sequence above is also an example of a recursive definition, that is very easy to write down using lazy evaluation, but much more involved when using eager evaluation. (say, we want to compute the list of the first n Fibonacci numbers.

```

1 // lazy                                // eager
2 let fib_list(n) = take(n, fib);        let fib_list(n) {
3                                         let iter(i,a,b) {
4                                             if i == n then
5                                                 []
6                                             else
7                                                 [a+b | iter (i+1,b,a+b)]
8                                         };
9                                         [0, 1 |iter(2,0,1)]
10                                        };

```

Disadvantages of lazy evaluation On the flip side, lazy evaluation has also severe disadvantages. The most prominent one is that it cannot be combined with side-effects as it obscures the order in which expressions are evaluated, which is of paramount importance in computations with side-effects.

Furthermore, it turned out that it is very hard to predict the memory consumption of programs using lazy data structures since one is never quite sure when a structure will be constructed and when the program is done processing it, so the garbage collector can free it again.

1.11 Programming examples

Let us conclude this chapter with several examples of programs in a functional style. We concentrate on functions for list processing.

```

1 type list =
2   | Nil
3   | Cons(a, b);
4
5 let nth(lst,i) {
6   if i == 0 then
7     head(lst)
8   else
9     nth(tail(lst), i-1)
10 };
11
12 let length(lst) {
13   case lst

```

1 Expressions and Functions

```
14   | Nil          => 0
15   | Cons(x,xs) => 1 + length(xs)
16 };
17
18 let sum(lst) {
19   case lst
20   | Nil          => 0
21   | Cons(x,xs) => x + sum(xs)
22 };
23
24 let map(f, lst) {
25   case lst
26   | Nil          => Nil
27   | Cons(x,xs) => Cons(f(x), map(f, xs))
28 };
29
30 let fold(f, acc, lst) {
31   case lst
32   | Nil          => acc
33   | Cons(x,xs) => fold(f, f(acc, x), xs)
34 };
35
36 let foldr(f, acc, lst) {
37   case lst
38   | Nil          => acc
39   | Cons(x,xs) => f(x, foldr(f, acc, xs))
40 };
41
42 let reverse(lst) {
43   let iter(lst, result) {
44     case lst
45     | Nil          => result
46     | Cons(x,xs) => iter(xs, Cons(x,result))
47   };
48   iter(lst, Nil)
49 };
50
51 // tail recursive version of foldr
52
53 let foldr(f, acc, lst) {
54   let g(x,y) { f(y,x) };
55   fold(g, acc, reverse(lst))
56 };
```

Exercise Write an implementation of balanced binary search trees in the kernel language as it is defined so far.

2 Types

2.1 Static and dynamic typing

In most languages there are operations that cannot be performed on every kind of input. For instance, division might be defined for numbers, but not for strings. For this reason one distinguishes several *types* of data. In some languages such as Haskell, Scala, or Rust, the type system is extremely sophisticated and subject to active research, other languages make do with rather impoverished type systems. For instance, the original Fortran had only two types: integers and floating point numbers.

Traditionally, there are two radically different ways of implementing types: *static typing* and *dynamic typing*. In static typing, every identifier of the program is associated with some type and the compiler ensures that the value of the identifier will always be of that type. In dynamic typing on the other hand, the types are not associated with the identifiers but with the values themselves. That means that every value in memory is *tagged* with its type and these tags are checked by all operations performed on the value. Each choice has its advantages and disadvantages.

Dynamic typing

- is slow: every operation performs runtime checks of the types,
- catches only type errors in those parts of the program that are executed,
- is more permissive and more convenient: no type annotations or other kinds of red tape.

For these reasons, dynamic typing is mainly useful in scripting languages, but not for writing non-trivial programs.

Static typing

- is stricter and catches more errors,
- the compiler can *prove* that the program is free of type errors,
- there is no runtime overhead,
- it can sometimes be inconvenient: the programmer has to write additional code in order to make correct code actually compile,
- not all properties can be checked statically (e.g., array bounds),
- with sophisticated type systems, the error messages from the type checker can be hard to understand,
- type annotations help document code,

- static type information can provide implicit context that changes the behaviour of a piece of code (e.g., with overloading).

Good static type systems try not to get in the way of the programmer. For instance, ML-like languages provide static type checking without requiring any kind of type annotations. Unfortunately, other languages are much less successful in this respect, think for example of template code in C++.

For serious software development, static type checking has turned out to be indispensable. First of all, we can use it as a means for the compiler to automatically *prove* that the program does not contain certain kinds of errors. The more expressive the type system is, the more kinds of errors we can catch.

Secondly, types also help with program design. When tasked with writing a certain submodule of a program, many programmers first design the types and data structures of the data involved. Then they use these types as a guide to write the actual code.

Thirdly, experience has shown that a good type system helps with refactoring large programs: after a change in one place of the program, the type checker can tell you all the other places you have to change as well.

Finally, let us note that the advantages of a type system apply much more to symbolic computations, than to numeric code (e.g., it doesn't catch sign errors).

2.2 Type annotations

To implement static typing in our kernel language, we add type annotations to every declaration. (This is not strictly necessary as there exist algorithms to *automatically infer* the types from a program without annotations. We will discuss such algorithms below.)

$$\begin{aligned} \langle expr \rangle ::= & \dots \mid \mathbf{let} \langle id \rangle : \langle type \rangle = \langle expr \rangle ; \langle expr \rangle \\ & \mid \mathbf{let} \langle id \rangle (\langle id \rangle : \langle type \rangle) : \langle type \rangle \{ \langle expr \rangle \}; \langle expr \rangle \\ & \mid \mathbf{fun} (\langle id \rangle : \langle type \rangle) : \langle type \rangle \{ \langle expr \rangle \} \end{aligned}$$

We also need to define which types the language provides. In our case we have the base type `int` for integers, function types $a \rightarrow b$, and one type `foo` for every type declaration

```
type foo = | A(a,b,...) | ... | Z(c,d,...);
```

in the program. So far, the parameters `a, b, c, d, ...` in constructor declarations served only to denote the arity of the constructor. Now we require them to be type expressions specifying the type of the constructors arguments. For instance, if we want `A` to take an integer and a boolean, we write `A(int, bool)`.

Examples

```
1  let fac(n: int) : int {
2    if n == 0 then 1 else n * fac(n-1)
3  };
```



```

4
5  let compose(f: int -> int, g: int -> int): int -> int {
6    fun (x: int) { f(g(x)) }
7  };
8  let twice(f : int -> int): int -> int {
9    compose(f, f)
10 };
11 let apply : (int -> int) -> int -> int =
12   fun (f: int -> int) {
13     fun (x: int) {
14       f(x)
15     }
16   };
17
18 type int_list = IntNil | IntCons(int, int_list);
19
20 let sum(l : int_list) : int { ... };

```

Exercise What could be the type of the following function?

```
let f = fun (x) { x(x) };
```

(Note that $f(f)$ evaluates to $f(f)$.)

2.3 Common Types

Let us give a short overview of types that are commonly found in real programming languages. Generally, we distinguish between *basic* and *composite* types. The former are atomic and built into the language, while that latter are composed out of one or several simpler types. Hence, basic types are types such as

- integers, signed and unsigned, of various precisions, including arbitrary precision integers,
- floating point numbers of various precisions, decimal numbers (0000.00), and arbitrary precision rational numbers,
- integer ranges (1..100),
- enumerations (enum colours { Red, Green, Blue, Yellow }),
- booleans,
- characters,
- strings,
- the empty type and the unit type,

while the composite types include

2 Types

- arrays,
- pointers and references,
- functions and procedures,
- records and tuples,
- unions and variants,
- lists and maps or dictionaries.

Arrays Arrays are homogeneous (all elements have the same type) collections of values. Some languages come with a very elaborate support for arrays. The language FORTRAN shines in this area, as it was specifically designed for numeric computations where arrays play an important rôle. In particular, FORTRAN supports higher-dimensional arrays and efficient array *slices*, which are (not necessarily contiguous) subsets of an array. For instance, one can define a slice consisting of the first 16 elements of every other row of an array. The important aspect of a slice is that it does not make a copy of the array, but only provides a new way to index the elements of the old array.

From the point of view of the type system, a point of concern is the fact that it is not possible to statically check that all array accesses are within bounds. This would be a very desirable thing to have, as array overflows are a very common source of bugs and security problems. Therefore, modern languages usually add dynamic bounds checks to each array access. (Usually these can be turned off selectively at places where efficiency matters.)

Product types Products are similar to arrays, but they are inhomogeneous (the elements can have different types) and their size is fixed. Set-theoretically they correspond to a cartesian product. Commonly the components of a product are labelled and can be referred to by their name. In this case such types are usually called *records* or *structures*.

```
1  type triple = int * int * int;  
2  type vector = [ x : float, y : float, z : float ];
```

Languages supporting product types come in two flavours depending on how the components of a product are accessed. If the language has *first-order fielding* the component is fixed at compile time, while a language with *first-class fielding* allows the runtime computation of the component. For instance, if we write $r.x$ to access the field named x of a record r , we know the field at compile time. But if we can write $r.(e)$ with an arbitrary expression e , it is only known at compile time which field we are accessing. Clearly, first-class fielding is much more expressive than first-order fielding, but it is unfortunately not possible to combine it with static type checking. For example, in

```
1  type foo = [ x : int, y : bool ];  
2  
3  r.(if i = 0 then x else y)
```

we cannot say, whether the expression evaluates to an integer or a boolean. For these reasons, first-class fielding is usually found only in dynamically typed scripting languages. One example, where it is rather useful is in writing serialisation and deserialisation code.

Sum types Sum types (also called *tagged unions*) are dual to products. Instead of storing several values at the same time, a sum type contains only a single value whose type may be one of a given list of types. Set-theoretically they correspond to a disjoint union.

```

1  type int_list = | Nil | Cons(int, int_list);
2  type expr = | Num(int) | Plus(expr, expr) | Times(expr, expr);
3  type nat = | Zero | Suc(nat)

```

In languages with sum types, one usually combines them with products, i.e., one allows the user to specify a type as a sum of products as in the example above. In this case one speaks of *variant types* or *algebraic types*.

Variant types are frequent in ML-like languages, but not well-supported by C-based or Pascal-based ones. C++ has *enums* which can be seen as sum types where the constructors have no arguments. It also has *untagged unions*, which can simulate sum types by adding the tag manually. Pascal supports a case-statement inside records which serves the same purpose as a sum type.

Note that sum types add a dynamic component to a type system. For instance, if we have a value of type

```
type either = | Left(int) | Right(bool);
```

it is unknown at compile time whether it is an integer or a boolean. Hence, we have to tag the value with its variant (Left or Right). Note that this is the same thing we do in set theory, where a disjoint union is usually defined as

$$A + B := \{0\} \times A \cup \{1\} \times B.$$

Here the first component (0 and 1) serves as a tag distinguishing the elements of A and B .

Unit and void type One has to distinguish between a *unit* type which has exactly one value, usually the empty tuple,

```
type unit = | Nothing;
```

and the *void* type which has no values at all.

```
type void = ;
```

If we want to treat procedures as functions with a special return value, this value must be of a unit type, since a function must return a value but we do not care which one it is. A function whose return type is void cannot return at all as it would have to produce a value of void type to do so.

Recursive types Most programming languages have at least some support for recursively defined types such as

```
type expr = | Num(int) | Plus(expr, expr) | Times(expr, expr);
```

Note that a value of the form, say, $\text{Plus}(e_1, e_2)$ is not stored in memory by having a memory segment consisting of a tag and two copies of the value e_1 and e_2 (which can be arbitrarily large). Instead, the memory segment contains the tag and *pointers* to the two argument values. In many languages, one is only allowed to define recursive types if the recursion is via such pointers. Some

2 Types

languages have full support for recursive types by allowing arbitrary recursive definitions. Unfolding such definitions produce a possibly infinite type expression. For instance,

```
type t = t -> t
```

is the type of all functions from t to t . It unfolds to

```
type t = (... -> ...) -> (... -> ... )
```

This is the type of the self-application function.

```
let f(x : t): t = { x(x) };
```

This means that with full support for recursive types, we can type the recursion operator as

```
1 type b = b -> a;  
2 let rec(f : a -> a) : a =  
3   (fun (x : b) : a { f(x(x)) })  
4   (fun (x : b) : a { f(x(x)) });
```

2.4 Type checking

Type equivalence Before being able to type check, we have to decide when we allow an argument of a given type a to be passed to a function expecting arguments of some, possibly different, type b . Clearly, this should be the case if the two types are *equivalent*. But what does being equivalent mean? There are basically two choices.

- With *name equivalence* two types are considered to be the same if they have the same name. Examples of languages using name equivalence are Pascal, C and their descendants.
- With *structural equivalence* two types are considered to be the same if they have the same structure, even if their names might be different. Languages in the ML-family typically use this kind of equivalence.

Example In C, which uses name equivalence for structures, all of the following types are considered to be distinct since they have different names.

```
1 type vector = [ x : int, y : int ];  
2 type pair   = [ x : int, y : int ];  
3 type pair2  = [ y : int, x : int ];
```

In ML the corresponding definitions would all define the same type, so we could pass a `pair2` to a function expecting a `pair`.

Example Suppose we want to use types to distinguish between measurements in metric units and in imperial units. How to do so depends on which kind of equivalence the type system uses.

```

1 // name equivalence                // structural equivalence
2 type metric = float;              type metric = M(float);
3 type imperial = float;            type imperial = I(float);
4
5 let f(x: metric): metric {        let f(x: metric) : metric {
6   ... 2*x ...                      ... 2 * unpack(x) ...
7 };                                  };
8 let x : imperial = 10;             let x : imperial = I(10);
9 f(x) // error                       f(x) // error

```

Type conversions There are cases where we can allow passing arguments to a function even if the types are not equivalent. For instance, this is the case when we can *convert* the argument to the expected type. For example, in C one can pass an integer to a function expecting a floating point argument and it will automatically converted into a floating point number. When such conversions are done automatically by the compiler, we speak of *type coercions*. Some languages like C, Perl, or JavaScript are very liberal with regard to type coercions, while other languages, like ML and Haskell do not allow coercions at all. Except for scripting languages, modern programming languages usually try to reduce the amount of coercions.

On the one hand, coercions are convenient since they make the code shorter and cleaner. On the other hand, they make the code harder to understand (implicit behaviour) and can hide type errors. This is the usual trade-off between an implicit effect and an explicit one. In moderation they can make the life of the programmer easier, but when overdone they easily create a mess.

Some languages with a permissive type system also allow *type casts*. A type cast is a command telling the compiler to regard a value as having a user-specified type instead of its real one.

There are several kinds of conversion between types (either in a coercion or a type cast). If every value of the first type has the same memory representation as the corresponding value of the second type, we can just change the type and there is no run-time overhead. If the memory representations differ (e.g., if we convert an integer to a floating-point number), we have to insert code that does the conversion. Some languages also support *non-converting type casts*. Such casts never change the memory representation, even if this does not make sense semantically. This feature makes the type system *unsound*, but it can be useful for system programming. For instance, in C one can cast from any pointer-type to any other in this way.

An additional complication arises if not every value of one type can be converted to the other type. In such cases one has to add a runtime check ensuring that the conversion is possible. For example, in object-oriented languages one sometimes what to cast an object of a superclass to one of its subclasses. In this case a runtime check is needed to make sure that the object is in fact of the required class.

Type checking After these preliminary remarks, let us finally turn to type checking itself. For the simple type system we have chosen for our kernel language, which is basically what the older mainstream languages like C and Pascal did provide, this is straightforward.

```

1 let fac(n : int) : int {
2   if n == 0 then 1 else n * fac(n-1)

```

```
3   };
```

2.5 Polymorphism

In the typing examples above, we have seen that, when adding type annotations to a program, we sometimes have to make arbitrary choices since some functions could be used with different types. For instance, the `identify` function

```
fun (x) { x }
```

can be given the type `int -> int`, or `bool -> bool`, or `(int -> bool) -> (int -> bool)`, and so on. It would be desirable, to use the same function definition for all suitable types instead of requiring a separate definition (with the same code) for each of them. This phenomenon is called *polymorphism*. Most modern programming languages support it in one form or other. One can broadly distinguish three different forms of polymorphism.

- (i) *ad-hoc polymorphism*, also called *overloading*,
- (ii) *parametric polymorphism* as can be found in ML-like languages, and
- (iii) *subtyping polymorphism* as is present in object-oriented languages.

Ad-hoc polymorphism In ad-hoc polymorphism the programmer can define several versions of a function. Which of them will be selected when the function is called will depend on the types of the supplied arguments. A typical example are arithmetic operations which in many languages are defined both for integers and floating point numbers.

```
1  + : int -> int -> int
2  + : float -> float -> float
3  + : string -> string -> string
```

Ad-hoc polymorphism is the most flexible form of polymorphism since it allows the programmer complete freedom. The disadvantage is that one has to write several different versions of each function which can become quite a chore. Furthermore, if ad-hoc polymorphism is used extensively the program can become hard to understand as it will be difficult to figure out which version of the function will be called at each call site.

Parametric polymorphism In parametric polymorphism we allow type expressions to contain *type variables*. For instance, we can specify the type of the `map` function as

```
map : (a -> b) -> list(a) -> list(b)
```

with two variables `a` and `b`. This is a simple and quite clean extension of the type system with few drawbacks. But it is less flexible than ad-hoc polymorphism. Most of the functional programming languages have adopted this version of polymorphism.

Subtyping polymorphism This kind of polymorphism is based on the *subtyping relation*. We say that a type a is a *subtype* of b if every value of type a can be used where a value of type b was expected. This is a situation that commonly arises in object-oriented languages where objects of a subclass automatically also belong to their superclasses. We will discuss subtyping polymorphism in more detail in Chapter 7. As far as the expressive power is concerned, there are things that subtyping can express, which parametric polymorphism cannot, and vice versa. Both approaches have their merits, but they have a very different *feel* to them. While parametric polymorphism is conceptually quite simple, subtyping makes a type system very complex.

Type inspection Some languages provide mechanisms to inspect the type of a polymorphic value at runtime. In this way a polymorphic function can behave differently depending on which type is supplied. A prominent example is serialisation, where an arbitrary value is converted to a string.

```

1  let serialise(value) {
2    case type_of(value)
3    | int    => int_to_string(value)
4    | bool  => bool_to_string(value)
5    | string => sanitise_string(value)
6    | cons  => "cons(" ++ serialise(fst(value)) ++ ","
7              ++ serialise(snd(value)) ++ ")"
8    | ...
9  };

```

Type inspection is a way to add the power of ad-hoc polymorphism back to a system based on parametric or subtyping polymorphism. It makes the type system much more powerful, but also less uniform and more complex.

Data polymorphism So far, we have looked at polymorphic functions. Data structures can also be polymorphic. For instance, the general list type and the types of the two constructors are

```

1  type list(a) = | Nil | Cons(a, list(a));
2  Nil  : list(a)
3  Cons : a -> list(a) -> list(a)

```

2.6 Type inference

Writing explicit type annotations at every declaration can become quite tedious, in particular, if we use a sophisticated type system where the type expressions are quite large (see, for instance, template code in C++). Many modern languages therefore implement a form of *type inference* where the types of expressions are automatically derived from the code without the help of type annotations. The amount to which this is possible strongly depends on the type system. In ML-like systems, type inference is possible without restrictions. In more complicated type systems, we need some type annotations but can infer others. The original type inference algorithm for

2 Types

ML was developed by Damas, Hindley, and Milner. Therefore, one frequently speaks of Hindley-Milner type inference. Given an expression the algorithm looks at every subexpression and extracts a list of equations between the types involved and solves them.

Examples

```
1  let twice(x) { 2 * x };
```

```
1  let fac(n) {  
2    if n == 0 then  
3      1  
4    else  
5      n * fac(n-1)  
6  };
```

```
1  let add(lst) {  
2    case lst  
3    | Nil      => 0  
4    | Cons(a,as) => a + add(as)  
5  };
```

```
1  let f(x) { x };
```

```
1  let map(f, lst) {  
2    case lst  
3    | Nil      => Nil  
4    | Cons(x,xs) => Cons(f(x), map(f, xs))  
5  };
```

Unification The process of solving a single type equation $s = t$ is called *unification*. Conceptually, the algorithm is very simple. If s or t is a type variable, we can set its value to be the other term. Otherwise, we check that the outermost operator of both type expressions is the same and recursively unify the arguments.

$$\begin{aligned}x = t & \rightsquigarrow x := t \\t = x & \rightsquigarrow x := t \\s \rightarrow s' = t \rightarrow t' & \rightsquigarrow s = t \wedge s' = t' \\c(s_1, \dots, s_n) = c(t_1, \dots, t_n) & \rightsquigarrow s_1 = t_1 \wedge \dots \wedge s_n = t_n \\s = t & \rightsquigarrow \text{failure}\end{aligned}$$

Type inference has its advantages and disadvantages. On the one hand, it is very convenient, relieving the programmer of the burden of having to annotate every declaration with a type. Furthermore, it will find the *most general type* for an expression and automatically introduce polymorphism. On the other hand, having explicit type annotations serves as a kind of documentation and improves the readability of the code (explicit vs. implicit information). Furthermore,

error messages from a type checker with type inference are usually more complicated and harder to read. One reason for this is that the equation-based approach of type inference obscures the *location* of the type error. The algorithm only determines that some equations are inconsistent, but it cannot deduce *which* of them is the cause.

3 State and Side-Effects

3.1 Assignments

In its current state our kernel language is purely functional, that is, running a program amounts to evaluating a mathematical expression that produce some value and nothing more. In this chapter we will add *side effects* to the picture. With side effects expressions do not only produce a value, but they can also modify the state of the world in certain ways, say, changing the contents of memory cells, drawing on the screen, or reading keystrokes from the keyboard. These are all essential features no serious general purpose programming language can do without. Even so-called purely functional languages must therefore support side effects, but they do it in a way which is separated from the rest of the program. For instance, a Haskell program consists of two phases. The first phase is pure and does not allow side effects. It computes a list of commands that *do* have side effects. This list is then executed in the second phase.

We start by extending our kernel language with two commands providing different kinds of side effects: an assignment statement to alter the contents of a memory cell and a print statement to produce screen output.

$$\langle expr \rangle ::= \dots \mid \mathbf{skip} \mid \mathbf{print} \langle msg \rangle \langle expr \rangle \mid \langle expr \rangle ; \langle expr \rangle \\ \mid \langle id \rangle := \langle expr \rangle$$

```
1  let x = 1;
2  print "x has value: " x;
3  x := 2;
4  print "now x has value: " x;
```

With these new statements we cannot regard an expression e anymore as a mathematical function $env \rightarrow val$ that, given values for the free identifiers in e , produces a value. Instead, we also have to specify its effect on the state of the world. That is, an expression now determines a function $env \times state \rightarrow val \times state$. In our case the state must contain the memory contents and also the produced output.

Note that, with assignments identifiers no longer represent constant values but *variables* instead. A variable in this context is an identifier associated with a *location* in memory which contains the value stored in the variable. This means that the notion of an environment is changed from a function mapping names to values to one mapping names to memory locations.

We have seen that in a language with assignments we must distinguish between expressions denoting values and those denoting memory locations. Only the latter can appear on the left-hand-side of an assignment, while the right-hand-side can contain both kinds of expressions. One frequently uses the terminology of *l-values* and *r-values* for locations and values, respectively. Here, the l and the r specify on which side of an assignment the expression can appear.

3 State and Side-Effects

In our kernel language, the only l-values are variables and expressions for structure access `r.m`. In real imperative languages like C several other kinds of expression can be l-values, for instance, expressions for array indexing `a[i]`.

3.2 Ramifications

The support for side effects has a drastic influence on all aspects of a programming language. Let us mention a few aspects.

Evaluation order First of all, with side effects the *order of evaluation* becomes important. Until now we could not have cared less about in which order subexpressions were evaluated (if we ignore termination issues for the moment), but with assignments and IO the order matters. For instance,

```
1  let x = 0;  
2  let y = (x := 1; 3) + (x := 2; 4);  
3  x + y
```

returns either 8 or 9 depending on which term in the definition of `y` is evaluated first. This means that with side effects we have to define an evaluation order, preferably one that can easily be read off from the syntactic structure of the code. This rules out lazy evaluation, where it is very hard for the programmer to determine in which order expressions are evaluated.

Uninitialised data structures A second new effect is that assignments allow for uninitialised or partially initialised data structures. Such things are not possible in a purely functional language since there is no way to retrospectively initialise objects. Partial initialisation is very convenient when creating mutually recursive data structures. We can first allocate the memory for all the structures and then fill in the pointers between them. Of course, having uninitialised data structures is also a source of bugs when such a structure escapes into a part of the program that expects its inputs to be fully initialised. (This is a common problem when writing constructors in, say, C++ as constructor code is executed while the objects in question is not yet fully initialised. So all methods called inside a constructor have to work with a partially initialised object.)

Aliasing A third effect is that with assignments we have to distinguish two notions of equality. To objects can have *the same value* or *the same memory location*. We can tell these two apart by changing the value of one object. If the value of the second object also changes, they share the same memory location, otherwise their locations are distinct. Having the same memory location accessible through several variables or expressions is called *aliasing*. For instance, consider the following code.

```
1  let x = 1;  
2  let y = x;  
3  x := 2;  
4  y
```

Depending on the semantics of our programming language y will or will not alias x and the code will return either 1 or 2.

When working with mutable data structures, aliasing has to be strictly controlled. If a piece of code wants to modify a data structure and it does not know whether there is aliasing involved, it has to make a copy of the structure before modification. In big programs written by a large team of programmers, it is not always clear at which places aliasing can occur. Therefore, one commonly makes copies of data ‘just to be sure’. This leads to a lot of unnecessary copying. For instance, in one version of the Chrome web browser, profiling revealed that every single keystroke in the URL field the browser resulted in several thousand memory allocations. This copying does not only slow down the program it also wastes space. When working with immutable data structures, one can have them share those parts they agree on. For instance, we can have two lists share a common tail or two search trees share common subtrees. This is a common situations for data structures in functional languages.

In light of the aliasing effect, a language designer has to decide what to do if a data structure gets assigned to a variable. The most efficient solution is to just let the variable point to the same object without making a copy. As we have discussed, this creates aliasing. If one wants to avoid aliasing, one has to make a copy of the data structure and, recursively, of all data structures reachable from the given one via pointers. This approach is called *deep copying*. It is quite slow and memory inefficient. There is also a compromise where only the first structure is copied, but not the pointed to structures. This approach, called *shallow copying*, is clearly inferior to the other two: it is less efficient than the first one, does not avoid aliasing, and it is also more complicated for the programmer. We will discuss these different strategies more below in the section on parameter passing.

Cleanup code Finally, since our code can now affect the state of the system, it needs to clean up when it is done by freeing the allocated resources like, say, memory, file handles, or windows. This means that we have to make sure that every code path leaving this part of the program calls the cleanup code. In practice, this can be a lot of work and rather a nuisance. It is also quite error prone as it is easy to forget to free one or two of the resources. Not that, in addition to direct returns we also have to check indirect ones like exceptions.

```

1  ...
2  let a = allocate_a();
3  if error then
4    return
5  ...
6  let b = allocate_b();
7  if error then {
8    free(a);
9    return
10 }
11 ...
12 let c = allocate_c();
13 if error then {
```

3 State and Side-Effects

```
14   free(b);
15   free(a);
16   return
17 }
18 ...
```

Many languages have added special constructs to help with cleanup. For instance, in Java a block can be annotated with a finally-statement which contains code that is always executed when control leaves the block.

```
1  let a = nil;
2  let b = nil;
3  let c = nil;
4  {
5    ...
6    a := allocate_a();
7    if error then return;
8    ...
9    b := allocate_b();
10   if error then return;
11   ...
12   c := allocate_c();
13   if error then return;
14   ...
15 }
16 finally {
17   if c then free(c);
18   if b then free(b);
19   if a then free(a);
20 }
```

A similar idea is to have a defer-statement which specifies commands to be executed when leaving the current block.

```
1  let a = nil;
2  let b = nil;
3  let c = nil;
4  {
5    ...
6    a := allocate_a();
7    if error then return;
8    defer free(a);
9    ...
10   b := allocate_b();
11   if error then return;
12   defer free(b);
```

```

13   ...
14   c := allocate_c();
15   if error then return;
16   defer free(c);
17   ...
18   }

```

Discussion Side effects drastically increase the power of a language. There are algorithmic problems that have very simple solutions using side effects, but where the corresponding side-effect free solutions are extremely cumbersome or inefficient. Furthermore, every serious language must support some form of IO, which is not possible without side effects.

On the flip side, side effects make the code much more complicated to reason about. They add implicit interactions between different parts of a program, for instance, via mutable global variables. This reduces encapsulation, makes the program harder to understand (non-local reasoning), and the coding more error prone.

So side effects are necessary but dangerous. Therefore it is desirable for a language to have some sort of separation between pure and impure code. This separation was already present in Algol which distinguishes between expressions and commands. A modern example is Haskell, which is particularly strict in this regard. Other languages are much more relaxed. For instance in ML or C++, one can declare variables to be constant (the default in ML) or mutable (the default in C++). This can be used to limiting side effects. So far, none of the solutions are really satisfactory. Either the separation is too lenient to offer real protection against side effects in places where they are not needed; or it is too strict making it very cumbersome. For instance, if during development one discovers that some part of a Haskell program would profit from a use of side effects, it is frequently necessary to rewrite large (and mostly unrelated) parts of the program to make the type system happy.

3.3 Parameter passing

Having introduced assignments and mutable state, we have to decide how it interacts with parameter passing. When we change a variable inside a function, does this effect become visible on the outside?

```

1  let f(x) { x := 1; };
2  let y = 0;
3  f(y);
4  y

```

Some languages allow the programmer to annotate function definitions with the desired behaviour for the parameters. For instance, Ada distinguishes between *in-mode*, *out-mode*, and *in/out-mode* parameters. In-mode parameters allow the passage of value from the call site to the function, out-mode parameters allow the passage in the opposite direction, and in/out-mode parameters can be used for both. Most other languages only provide in-mode parameters. Let us take a closer look at the various parameter passing mechanisms.

3 State and Side-Effects

Call-by-value is the standard mechanism for in-mode parameters. When calling a function, the argument values are passed as copies to the function body. Modifications of the copies do not affect the originals. This is a very safe method that avoids any confusion caused by unexpected modifications. The disadvantage is that it can be inefficient if large objects are passed in this way.

Call-by-result is the analog of call-by-value for out-mode parameters. No value is passed during the function call. Instead, when the function returns, the current contents of the variable corresponding to the parameter is copied back to the argument, which *must be an l-value*. Call-by-result has the same advantages and disadvantages as call-by-value. There are two additional problems that need to be addressed.

(i) What happens if the same variable is passed to two different out-mode parameters?

```
1  f(in x, out y, out z) {
2    y := x+1;
3    z := x+2;
4  };
5  let u = 0;
6  f(u,u,u);
```

(ii) At what time is the l-value passed as argument evaluated?

```
1  f(in x, out y, out z) {
2    y := x+1;
3    z := x+2;
4  };
5  let i = 0;
6  f(i, array[i], i);
```

Call-by-value-result/call-by-copy/call-by-copy-result combines call-by-value and call-by-result for in/out-mode parameters. The argument value is copied to the parameter when the function is called and copied back, when it returns.

```
1  let u = 1;
2  let f(x) {
3    print "u is " u;
4    x := 2;
5    print "u is now " u;
6  };
7  f(u);
8  print "u is now " u;
```

Call-by-reference is a more efficient version of call-by-value-result. Instead of copying the value back-and-forth, its address is passed to the function. Every Modification inside the function directly affects on the original l-value. This is very efficient, but can create aliasing problems.


```

1  let f(x) { x := 2 };
2  let u = 1;
3  f(u);
4  u

1  f(x, y) { x := 1; y := 2; }
2  g(x, y) { y := 2; x := 1; }
3  let u = 0;
4  f(u, u); print "after f:" u;
5  g(u, u); print "after g:" u;

```

Call-by-name is a radically different calling convention invented in Algol. Here the expression given as argument is substituted for the formal parameter in the function body using a *capture-avoiding substitution*, i.e., all local variables in the function will be renamed to avoid name clashes. In an implementation this amounts to passing the argument expression as a *thunk* (a suspended computation). This calling convention is the basis for implementing lazy evaluation. For code without side-effects, we have seen that call-by-name is nearly indistinguishable from call-by-value (except for issues of termination). In combination with side-effects, call-by-name is radically different from call-by-value.

```

1  let i = 0;
2  let array = [1,2];
3  let p(x) {
4    i := x;
5    x := 0;
6  };
7  p(array[i]);
8  print i array[0] array[1];

```

A famous example of call-by-name is what is called *Jensen's device*. The function

```

1  let sum(k, l, u, expr) {
2    let s = 0;
3    for k := l .. u {
4      s := s + expr;
5    };
6    s;
7  };

```

computes $\sum_{k=l}^u expr$ where the expression can be passed as an argument.

- `sum(i, 0, 99, array[i])` sums the first 100 entries of an array.
- `sum(i, 1, 100, i*i)` sums the first 100 square numbers.
- `sum(i, 0, 3, sum(j, 0, 3, m[i,j]))` sums the entries of a 4×4 matrix.

3 State and Side-Effects

Call-by-need is an optimised version of call-by-name useful for in-mode parameters. It is the standard calling convention used in lazy functional languages like Haskell. Here, after the first evaluation of a passed argument expression, the result is stored, so subsequent uses of the parameter do not need to evaluate the expression again. Of course, this only works if the argument expression has no side effects.

Call-by-macro-expansion is also similar to call-by-name but uses textual substitutions instead of capture-avoiding ones. Hence, the function works like a macro. This calling convention has its uses in a few limited cases, but it is clearly unsuited as the main calling convention of a language. Besides it being hard to implement efficiently (in particular, if recursion is involved), it also introduced non-local effects via unintended variable capturing. In particular, renaming local variables can change the behaviour of a program.

More examples

```
1  let f(x,y) { x := 2; y := 3; x };
2  let u = 1;
3  let v = 1;
4  let w = f(u,v);
5  print "u is now: " u;
6  print "v is now: " v;
7  print "w is now: " w;
8  let w = f(u,u);
9  print "u is now: " u;
10 print "v is now: " v;
11 print "w is now: " w;
12
13 let swap(x,y) { let tmp := x; x := y; y := tmp };
14 let a = 1;
15 let b = 2;
16 swap(a,b);
17 print "a is now: " a;
18 print "b is now: " b;
```

Discussion The consensus today is that one does want to have call-by-value in languages with side effects and call-by-need in languages without. The reason for call-by-value is to avoid *aliasing*, in particular *variable aliasing* where writing to one variable can change the contents of another one. There is also *data structure aliasing* where part of a data structure is accessible from different variables. If one wants to avoid this as well, we have to copy entire data structure when passing them to a function. This process is called *deep copying* as it involves following all pointers in the structure and recursively copying the memory pointed to. Since deep copying is very inefficient, it is implemented by only a few languages. Some languages prove a compromise where only the structure directly pointed to is copied, but no recursive copying occurs. This is called *shallow*

copying. Shallow copying has fallen out of favour, as it does not really solve the problem of aliasing and it is still no as efficient as the simple form of call-by-value. Therefore, most languages today use call-by-value where non-scalars, i.e., composite data structures, are passed by pointer. Some allow the simulation of call-by-reference by using explicit *reference* or *pointer types*. There is one exception: in a logical language call-by-reference is more natural, as it better fits the semantics expected by the user.

```
head([X|XS], X).
```

```
p(L) :- head(L, X), q(X).
```

In such languages, the problems of call-by-reference is reduced considerably as variables usually only support single-assignments (see Chapter 6), not multiple ones.

3.4 Memory management

When adding assignments we introduced the notion of a store. Our naive implementation added values to the store but never removed them again. In a real implementation this is of course unacceptable. Programs would run out of memory. So every real programming language must have some form of memory management that frees unused values in memory. There are three forms of memory management.

- In *manual memory management* the programmer is responsible for (nearly all) allocations and deallocations of memory blocks. (The exception is memory for local variables, which is usually managed automatically on a stack.)
- In *automatic memory management* the runtime system of the language performs allocations and deallocations automatically.
- In *type based memory management* the type system tells the compiler at which places it has to allocate and deallocate memory.

There are two types of problems memory management has to address.

- *Dangling pointers*, that is, pointers to already deallocated memory block. These can lead to program crashes and other undefined behaviour.
- *Unreachable objects*, that is, objects that are still allocated, but no longer reachable via pointers. These waste memory but are otherwise harmless.

Manual memory management For manual memory management, the language provides two operations to the programmer: one to allocated a certain amount of memory and one to deallocated it again. It is the responsibility of the programmer to make sure that memory that is not needed anymore is actually freed. Of course, this is quite tedious and error prone. It is easy to either forget to free memory, or to free it too soon. Both kinds of errors are hard to debug as the place where the error is made is usually not the place where its effects manifest.

Most implementations of manual memory management use a list (or several) of free memory blocks. If a certain amount of memory is to be allocated, this list is traversed until a block of

3 State and Side-Effects

suitable size is found. If later on the memory is freed again, it is simply added to the list. In actual implementations the picture is a bit more complicated as several techniques are added to increase efficiency. In particular, one should note that, in this scheme, allocating and freeing memory are both operations which take a non-negligible amount of time.

Automatic memory management With automatic memory management the programmer is relieved of the burden of managing memory herself. There are two main approaches. The first one is called *reference counting*. Here, every memory object has a counter which stores the number of pointer to this object. If, at some point, this counter reaches zero, the object is automatically deleted. The other approach is based on *garbage collection*. Here, objects are not freed right away. Instead, the program continues to allocate memory until the remaining amount of free memory drops below a certain threshold. Then the memory manager determines which part of the allocated memory is actually in use and frees the rest.

Reference counting is easy to implement, but very slow and it cannot deal with cyclic data structures. Garbage collection on the other hand, is very hard to implement well. But it has the advantage that allocations are very fast (usually just a pointer increment and a compare) and deallocations do not take any time at all. Of course there is also the collection phase, which can take quite some time. How much depends on the kind of collection being performed. We distinguish the following cases.

- During collection the whole program is stopped. This is the easiest to implement, but it causes latency problems.
- A collection is split into several pieces, which are interleaved with the program execution. This somewhat reduces the latency problem.
- The garbage collector and the main program run in parallel. This is very hard to implement well, but it completely eliminates the latency problem at the cost of further increasing the garbage collection overhead.

Type based memory management is a novel approach to memory management and still experimental. The only mainstream language implementing it is Rust. Here, one uses the type system to encode information about the lifetime of objects. Objects are deallocated when the type system says that they are dead. For instance, if an object is locally defined in some scope and no references to the object are passed out of the scope, we know that we can safely delete the object when the scope terminates. This approach tries to retain the safety guarantees of automatic memory management while avoiding its overhead. It remains to be seen how practical it will turn out to be.

Discussion Automatic memory management has clear advantages over manual management. It guarantees the absence of certain kinds of memory errors which historically have been the cause of many program crashes and security breaches. It also makes the code shorter and cleaner as the programmer does not need to write cleanup code. Finally, there are scenarios where it is even faster than manual memory management.

On the other hand, it also has several disadvantages. First of all, it is quite complex and hard to implement. In particular, if it is to be parallelised or if one wants to address real-time requirements. Furthermore, many of the faster garbage collection algorithms waste quite a bit of memory (frequently only half of the real memory is usable). And finally, even with all optimisations, there is still a considerable overhead associated with garbage collection. This makes it unusable for certain applications with strict performance requirements like, say, computer games.

3.5 Loops

The imperative analogue of recursion is a loop. We distinguish two kinds of loops: *bounded* and *unbounded* ones. A loop is bounded if the number of iterations is known beforehand. So for-loops are bounded and while-loops unbounded.

$$\langle expr \rangle ::= \dots \mid \mathbf{while} \langle expr \rangle \{ \langle expr \rangle \} \mid \mathbf{for} \langle id \rangle = \langle expr \rangle \dots \langle expr \rangle \{ \langle expr \rangle \}$$

There is a more fundamental primitive that can be used to implement loops: the goto-statement. A goto is an unconditional jump that transfers the program execution to the specified location.

$$\langle expr \rangle ::= \dots \mid \mathbf{label} \langle id \rangle \mid \mathbf{goto} \langle id \rangle$$

Using gotos we can replace a while-loop

```
while cond { expr }
```

by the following code:

```
1  label start;
2  if cond then (
3    expr;
4    goto start
5  )
6  else
7    skip
```

Similarly, a for-loop

```
for i = first to last { expr }
```

can be translated to

```
1  let i = first;
2  let l = last;
3  label start;
4  if i == l then
5    skip
6  else (
7    expr;
```

3 State and Side-Effects

```
8     i := i + 1;
9     goto start
10  )
```

Although `goto` is more expressive than `for`- and `while`-loops, it has the disadvantage that it can easily lead to unreadable code jumping willy-nilly from one location to another. The nesting imposed by loops prevents this kind of spaghetti code. There are several guidelines for the clean use of `goto`-statements. The simplest one is to only allow forward jumps in the code, but no backward ones. It can be shown that, if the language supports `while`-loops and `if`-statements, we can eliminate every `goto` by restructuring the code. For these reasons many modern programming languages have no `goto`-statements.

There are situations where the lack of a `goto`-statement leads to rather cumbersome code. The most common one is when one wants to jump out of the middle of a loop. Here a solution using an `if`-statement is rather inelegant, in particular if several such jumps are needed.

```
1  let terminate = False;
2
3  while ... and not(terminate) {
4    ...
5    if ... then
6      terminate := True
7    else {
8      ... rest of the loop ...
9    }
10 }
```

As this situation arises quite frequently, most languages provide specialised statements for them. A **break** statement terminates the innermost loop, a **continue** statement skips the rest of the loop's body and directly continues with the next iteration, and a **return** statement terminates the current function and returns to the caller.

$$\langle expr \rangle ::= \dots \mid \mathbf{break} \mid \mathbf{continue} \mid \mathbf{return} \langle expr \rangle$$

In some languages, it is also possible to jump out of nested loops by adding a label to the `break`- or `continue`-statement.

```
1  for i = 0 to 10 {
2    for k = 0 to 10 {
3      ...
4      continue i;
5      ...
6    }
7  }
```

3.6 Programming Examples

We have argued above that the use of side-effects can be problematic as it can make a program much harder to understand. On the other hand, judicious use of side-effects can also greatly simplify a program. Let us give some examples.

Recursive data structures As already explained in the section on recursion, we can use side-effects to create truly recursive data structures: first, we allocate all the memory needed for the various parts of the structure and before initialising it and creating all references between them.

Optimisation In certain cases using mutable variables makes an implementation more efficient. If we update some value and do not need the old value anymore, we can store the new value at the same memory location instead of allocating new memory. A typical example are accumulator variables used in loops. For instance, the list functions of Section 1.11 can be written using mutable variables in the following way.

```

1  let length(lst) {
2    let len = 0;
3    while is_cons(lst) {
4      len := len+1;
5      lst := snd(lst);
6    };
7    len
8  };
9
10 let sum(lst) {
11   let s = 0;
12   while is_cons(lst) {
13     s := s+1;
14     lst := snd(lst);
15   };
16   s
17 };
18
19 let map(f, lst) {
20   while is_cons(lst) {
21     fst lst := f(fst(lst));
22     lst := snd(lst);
23   }
24 };
25
26 let fold(f, acc, lst) {
27   while is_cons(lst) {
28     acc := f(acc, fst(lst));

```

3 State and Side-Effects

```
29     lst := snd(lst);
30   };
31   acc
32   };
```

Another common example are mutable data structures such as hash tables, search trees, etc. When programming in a functional style we have to create a new copy of the data structure whenever we update it. (Frequently, we do not need to copy the *whole* structure though, since we can share those parts that do not need to be modified with the old copy.) If we allow mutation, we can change the structure in place, which is usually more efficient. Of course, if we do so and we still need the old version of the structure, we have to manually make a copy first (which is less efficient as the functional implementation since in this case we usually cannot use sharing of parts of the structure).

Communication We can use mutable data structures to communicate between parts of the code. For example, if we want to implement a random number generator, we have to pass its state from one invocation to the next. In a functional implementation, the generator takes the form

```
random : state -> (state, int)
```

Hence, we have to pass the current state of the generator to every place where we want call this function and we have to pass the new state back to the next invocation. This is very tedious and decreases the readability of the code quite a bit. In an implementation with side-effect, we can store the current state in a mutable variable.

```
1  let state = ... some initial value ...;
2
3  let random(): int {
4    state := (1103515245 * state + 12345) mod 2147483647;
5    state
6  };
```

The problem with this use of side-effects is that it can make the program much harder to understand. Instead of explicitly passing values between the program parts in question, we do so implicitly by storing them in some shared variable. Hence, the programmer cannot understand one part of the program without the other, which violates the principle of local reasoning.

4 Modules

4.1 Simple modules

As programs grow larger it is necessary to divide them into manageable units commonly called *modules*, *packages*, or *program units*. A module is a part of a program with a well-defined interface that lists all the identifiers and their types defined within.

$$\begin{aligned} \langle expr \rangle ::= & \dots \mid \mathbf{module} \langle id \rangle \{ \langle declarations \rangle \} \mid \mathbf{module} \langle id \rangle = \langle module\text{-}expr \rangle \\ & \mid \langle module\text{-}expr \rangle . \langle id \rangle \mid \mathbf{import} \langle module\text{-}expr \rangle \\ \langle module\text{-}expr \rangle ::= & \langle id \rangle \mid \langle module\text{-}expr \rangle . \langle id \rangle \end{aligned}$$

Every module creates its own namespace. To access its elements, other parts of the program must prefix the identifier with the module name. Alternatively, one can use an **import**-statement to include the namespace of the module in the current one.

```
1  module Stack {
2    type stack(a) = list(a);
3
4    let empty = [];
5
6    let top(s)    { head(s) };
7    let pop(s)   { tail(s) };
8    let push(s, x) { [x|s] };
9  };
10
11 ...                               import Stack;
12 let s = Stack.empty;              let s = empty;
13 ...                               ...
14 Stack.push(s, 13);                push(s, 13);
15 ...                               ...
```

4.2 Encapsulation

The module mechanism addresses two ergonomic issues. Firstly, they help us manage namespaces and avoid name clashes between identifiers. Note that this could also be solved by adopting a strict coding style where, for instance, all identifier names in a given program unit start with a prefix indicating that unit. But this manual solution is cumbersome for the programmer (for instance, the convenience of import statements is lost) and not enforced by the compiler.

Secondly, they help to decompose the program into smaller, easier to understand parts (which themselves might be divided further into submodules). To understand such a program we only need to understand each component separately and then look at the way they interact. The second part is the easier the more limited the interaction between modules is. This is where declarative programming styles shine. If the modules are written declaratively, they can only interact via their specified inputs and outputs. We do not have to take further interactions into account, say, via mutable global state as is the case when using side effects.

This second use of modules is an example of a mechanism called *encapsulation*. Generally, encapsulation is the process of separating part of the program from the rest and allowing access only via a specified *interface*. This has several advantages.

First of all, as we have already explained above, it makes the program easier to understand since it reduces the amount of code a programmer must be aware of when looking at some part of the program. In particular, users of a module only need to know its interface, not the actual implementation. This is called *information hiding* and is the main way encapsulation contributes to program readability.

Secondly, it can be used to guarantee the *integrity* of data maintained by a module, since only code within the module is allowed to directly access the inner representation the data. In this way a module can enforce certain invariants a data structure must satisfy. (For instance, the requirement on red and black nodes in a red-black tree.)

Finally, it helps with program maintainability as one is free to change the inner representation of a modules data without affecting the rest of the program.

The simple module system defined in the previous section, supports the separating part of encapsulation, but not the interface part. For full encapsulation, we need to add a mechanism to restrict the access to the names defined in a module. There are basically two ways to do so. One is to allow definitions to be declared as either *public* or *private*. Only the public definitions are accessible from the outside. This is the method chosen by C++ and Java for class definitions, for example, and it is also what we will implement in our kernel language.

$$\langle decl \rangle ::= \dots \mid \mathbf{public} \mid \mathbf{private}$$

An alternative method, used in ML and also in C header files, for example, is to provide every module with a separate interface specification containing declarations of all identifiers visible to the outside. It requires more typing from the programmer, but it spacially separates the interface from the implementation. This makes it easier to read and also allows some more advanced mechanisms for module handling which we shall introduce below.

4.3 Abstract Data Types

An *abstract data type* is what we get when we apply the concept of encapsulation to the implementation of a data type. More concretely, an abstract data type is a data structure (usually defined inside a module) where the *representation* of the data, i.e., the *concrete implementation*, is hidden from the rest of the program (*information hiding*). The only access is via the operations defined in its *interface* (*encapsulation*). For instance, note that in most languages, built in types can be considered abstract, although this is a somewhat degenerate case.

Example Let us take a look at an abstract data type for stacks. We start with a functional version. The interface is

```

1  module Stack {
2    type stack(a);
3    let empty : stack(a)
4    let push  : stack(a) * a -> stack(a);
5    let top   : stack(a) -> a;
6    let pop   : stack(a) -> stack(a);
7  };

```

and the implementation is

```

1  module Stack {
2    type stack(a) = list(a);
3    let empty : stack(a) = nil;
4    let push(st : stack(a), x : a) : stack(a) {
5      pair(x, st)
6    };
7    let top(st : stack(a)) : a {
8      case st | pair(x,xs) => x
9    };
10   let pop(st : stack(a)) : stack(a) {
11     case st | nil => nil | pair(x,xs) => xs
12   };
13 };

```

The next version is imperative. The interface is

```

1  module Stack {
2    type stack(a);
3    let create : unit -> stack(a);
4    let empty  : stack(a) -> bool;
5    let push   : stack(a) * a -> unit;
6    let top    : stack(a) -> a;
7    let pop    : stack(a) -> unit;
8  };

```

and the implementation is

```

1  module Stack {
2    let create() : stack(a) {
3      [ elements = nil ]
4    };
5    let empty(st : stack(a)) : bool {
6      is_nil(st.elements)
7    };
8    let push(st : stack(a), x : a) : unit {

```

4 Modules

```
9     st.elements := [x|st.elements]
10 };
11 let top(st : stack(a)) : a {
12     head(st.elements)
13 };
14 let pop(st : stack(a)) : unit{
15     st.elements := tail(st.elements)
16 };
17 };
```

4.4 Module expressions

Most programming languages only offer a simple module system like the one presented above. A notable exception is ML where one can parametrise modules by other modules.

$$\begin{aligned} \langle expr \rangle ::= & \dots \mid \mathbf{module} \langle id \rangle (\langle id \rangle , \dots , \langle id \rangle) \{ \dots \} \\ & \mid \mathbf{module} \langle id \rangle = \langle module\text{-}expr \rangle \\ & \mid \langle module\text{-}expr \rangle . \langle id \rangle \mid \mathbf{import} \langle module\text{-}expr \rangle \\ \langle module\text{-}expr \rangle ::= & \langle id \rangle \mid \langle module\text{-}expr \rangle . \langle id \rangle \\ & \mid \langle module\text{-}expr \rangle (\langle module\text{-}expr \rangle , \dots , \langle module\text{-}expr \rangle) \end{aligned}$$

For instance, one way to define a map data type parametrised by the key type is as follows.

```
1  interface KEY {
2    type t;
3    type ord = | LT | EQ | GT;
4    let compare : t * t -> ord;
5  };
6  module Map(Key : KEY) {
7    type map(a) =
8    | Leaf
9    | Node(Key.t, a, map(a), map(a));
10
11    let empty : map(a) = Leaf;
12
13    let add(m : map(a), k : Key.t, v : a) : map(a) {
14      case m
15      | Leaf => Node(k, v, Leaf, Leaf)
16      | Node(k2, v2, l, r) => case compare(k, k2)
17          | LT => Node(k2, v2, add(l, k, v), r)
18          | EQ => Node(k2, v, l, r)
19          | GT => Node(k2, v2, l, add(r, k, v))
20    };
```

```

21     ...
22 };

```

First-class modules We can make the module system even more expressive by supporting *first-class modules*, i.e., adding the ability to pass modules around just like other values.

```

1  let add_two(M, x, y) {
2      let m = M.make();
3      M.add(m, x);
4      M.add(m, y);
5  };
6

```

```

1  let add_two(M, x, y) {
2      import M;
3      let m = make();
4      add(m, x);
5      add(m, y);
6  };

```

One can implement first-class modules by treating every module as a record containing the values of all identifiers defined within. Of course this means that referencing an element of a module now requires a memory lookup and cannot be done statically anymore (in general, the lookup can of course be optimised away in certain cases).

5 Control-Flow

5.1 Continuation passing style

In order to introduce the notion of a continuation let us take a look at the following example. Suppose we have a program that prompts the user for two inputs and then computes some result.

```
1  let f () {
2    let u = input("first: ");
3    let v = input("second: ");
4    process(u,v)
5  };
```

When we want to adapt this program to use a web-interface we face a problem with the way web-servers operate. Web-servers have the ability to call external programs to generate web-pages. But these programs are immediately terminated by the server after a web-page is produced. In our example we need three pages, two containing forms for the user to fill in the values of u and v , and one to display the computed result. As the program is terminated after each page, we have to figure out some way to pass the program state from one invocation to the next. Of course, web-sites with internal state are quite common, so web-servers do provide several mechanisms for doing so (cookies, hidden form fields, URL query string,...). What remains for us to do is to figure out, which information precisely to pass along. We need (i) a data structure storing at what place in the program we are and (ii) a way to use this data to resume the program at that point.

To resume the computation of the program from an arbitrary point we need to know

- where in the program we are, i.e., what the last evaluated expression was,
- what the result of this expression was, and
- what the values of the local variables were.

We can store this information as a function that, given the result of the last expression, continues the program from this point. Such a function is called a *continuation*. For instance, in the above example the continuation after having read the first input is

```
1  fun (u) {
2    let v = input("second: ");
3    process(u,v)
4  };
```

The continuation after the second input is

```
1  fun (v) {
2    process(u,v)
3  };
```

5 Control-Flow

In order to prepare our program for usage with a web-server, it is useful to translate it into a form that makes these continuations explicit. This form is called *continuation passing style*, CPS for short. In this form, every function takes an additional argument *k* that takes the continuation to be called when the function wants to return. Our example now looks as follows.

```
1  let f (k) {
2    input("first: ",
3      fun (u) {
4        input("second: ",
5          fun (v) {
6            process(u,v,k);
7          })
8      })
9  };
```

As a second example, let us take a look at the factorial function.

```
let fac(n) { if n == 0 then 1 else n * fac(n-1) };
```

We present two versions using continuation passing style. The first one is rather relaxed in the sense that we do not convert primitive operations.

```
1  let fac_cps(n,k) {
2    if n == 0 then
3      k(1)
4    else
5      fac_cps(n-1, fun (x) { k(n*x) })
6  };
```

If we also use continuation passing style for primitive operations like equality, subtraction, and multiplication, the code looks as follows.

```
1  let fac_cps(n,k) {
2    equal(n,0,
3      fun (c) {
4        if c then
5          k(1)
6        else
7          minus(n,1,
8            fun (a) {
9              times(n,x,
10             fun (b) { fac_cps(a, fun (x) { k(b) }) })
11            })
12      })
13  };
```

As we see in CPS a function never really returns, instead it calls its continuation. We can see CPS as a programming style where instead of using a call stack and we manually handle return

addresses by storing them in function closures, i.e., on the heap. This is of course a bit less efficient, since we removed the optimisation of using a stack, but as we will see below it offers more flexibility and allows for certain programming constructs not possible (or at least much harder to implement) with a stack discipline.

Let us conclude this section with a more involved example: a parsing function for regular expressions.

```

1  type regex =
2    | Char(char)
3    | Plus(regex, regex)
4    | Concat(regex, regex)
5    | Star(regex)
6
7  let parse_cps(str  : list(char),
8             regex : regex,
9             succ  : list(char) -> bool,
10            fail   : unit -> bool) : bool {
11    case regex
12    | Char(c) => if head(str) == c then
13                succ(tail(str))
14                else
15                  fail()
16    | Plus(r,s) =>
17      parse(str, r, succ,
18            fun () { parse(str, s, succ, fail) })
19    | Concat(r,s) =>
20      parse(str, r,
21            fun (str) { parse(str, s, succ, fail) },
22            fail)
23    | Star(r) =>
24      parse(str, r,
25            fun (str) { parse(str, Star(r), succ,
26                          fun () { succ(str) }) },
27            fun () { succ(str) })
28  };
29  let parse(str, regex) {
30    parse_cps(str, regex,
31              fun (s) { s == "" },
32              fun () { False })
33
34  };

```

5.2 Continuations

The problem with continuation passing style is that, in order to use it at one place in the program, we have to convert all functions used in that part to CPS. This is rather inconvenient and makes modifications of a program unnecessarily complicated. To avoid this overhead we can introduce a new construct into our language that makes the continuation at the current position available to the programmer.

$$\langle expr \rangle ::= \dots \mid \mathbf{letcc} \langle id \rangle \Rightarrow \langle expr \rangle$$

The statement

$$\mathbf{letcc} \ k \Rightarrow \ expr$$

evaluates the given expression after binding the current continuation to the identifier k . So when calling $k(a)$, the program behaves as if $expr$ returned the value a .

Examples

```

1  letcc k => 1
2  letcc k => k(1)
3  letcc k => k(1+2)
4  1 + letcc k => k(1)
5  1 + letcc k => k(1+1)
6  letcc k => (3 + k(1))
7  1 + letcc k => (3 + k(1))

```

There are two ways we can use the continuation k . We can call it within the expression $expr$, or we can store it somewhere and call it after the evaluation of the **letcc** statement is already finished. In the first case it acts like a return statement or an exception: we abort the evaluation of the $expr$ prematurely and return the specified result. In the second case, we perform some kind of backtracking: we restart the computation following the **letcc** statement with an alternative value for $expr$. We will see several examples where this can be used to good effect.

A word of warning A continuation can be seen as a function analogue of a goto statement. The only difference is that with continuations we can only jump to places we have already been, while a goto also allows jumps to unvisited program locations. As with gotos this flexibility can be misused. Many languages therefore try to replace arbitrary continuations with restricted versions like exception mechanisms (see below).

5.3 Generators

As a first application of continuations, let us implement what is sometimes called a generator. For-loops in our language are rather restricted. We can only iterate over the numbers between two given bounds. Many language designers thought it to be rather restrictive and tried to improve

for-loops to an imperative analogue of a fold function. For instance, there are languages where for-loops can also iterate over the elements of container types like arrays and lists. Instead of having built in support for a handful of such types, recent languages just like Python found a way to allow the user to define her own iterators for for-loops. Such a definition is called a *generator*. It is a function that produces, one after the other, all the values the loop should iterate over. The question is how we can pass these values to the loop construct. Using the return value of the function is cumbersome since we can return only one value at a time. So these languages introduced a new language construct **yield** that stops the evaluation of the current function, returns a value to the caller, and allows the caller to later resume the function at this position. Generally, such functions that can be stopped at an arbitrary position and resumed later on are called *coroutines*. For instance,

```

1  let gen() {
2    let n = 0;
3    while True {
4      yield n;
5      n := n+1;
6    }
7  };

```

generates the sequence 0, 1, 2, 3, 4, Similarly,

```

1  let gen(lst) {
2    while is_cons(lst) {
3      yield head(lst);
4      lst := tail(lst);
5    }
6  };

```

generates the sequence of all elements in the given list.

Looking at the definition of **yield**, we see that it looks a lot like a continuation, and we can in fact use continuations for the implementation.

```

1  let gen() {
2    let n = 0;
3    while True {
4      letcc k => {
5        gen := k;
6        return n;
7      };
8      n := n+1;
9    };
10 };
11 let gen(lst) {
12   while is_cons(lst) {
13     letcc k => {
14       gen := fun (x) { k() };

```

```

15     return head(lst);
16   };
17   lst := tail(lst);
18 }
19 };

```

5.4 Exceptions

Continuations can also be used to good effect for error handling. The problem with error handling is that, when an error occurs, we need to abort the current computation and go to an outer context where we can sensibly react to the failure. If we are using side effects, we also have to do the required clean up work required by the aborted computation. In the traditional way of error handling, error conditions are communicated via the return value of functions. This has the disadvantage that we have to surround every function call by an if-statement to test for the occurrence of an error, which is quite cumbersome, error prone (easy to forget), and clutters the code. Therefore programming languages have introduced a mechanism making the error checking implicit. This *exception* mechanism works similarly to the break-statement of imperative languages. But instead of jumping out of loops, i.e., out of a nested static scope, it allows the program to jump out of nested function calls.

$$\langle expr \rangle ::= \dots \mid \mathbf{try} \langle expr \rangle \mathbf{catch} \langle var \rangle \Rightarrow \langle expr \rangle$$

$$\mid \mathbf{throw} \langle expr \rangle$$

```

1  catch 2                                catch (2 + throw 4)
2  | x => x + 1                            | x => x + 1
3  => 2                                    => 5

```

```

1  type error = | EmptyList;
2
3  let head(lst) {
4    case lst
5    | []      => throw EmptyList
6    | [x|xs] => x
7  };
8
9  try head([])
10 catch x => 0

1  type error = | NotFound;
2
3  let lookup(lst : list([ key : a, val : b ]), k : a) : b {
4    case lst
5    | []      => throw NotFound

```

```

6   | [x|xs] => if x.key == k then
7           x.val
8           else
9           lookup(xs, k)
10  };

```

Exceptions can be implemented using continuations. Every function gets as an additional argument a continuation to call when raising an exception. A **catch** statement uses **letcc** to create such a continuation.

```

1  try e catch x => handler    ==>    letcc k => e(fun (x) { k(handler) })
2  throw e k                  ==>    k(e)

```

Exercise Implement exceptions using this translation.

Exceptions are not without problems. Although they can be considered as a generalised break-statement, the destination of an exception is determined dynamically by the sequence of function calls and not statically by the syntactic structure of the program. This makes reasoning about exceptions non-local. In particular when programming with side effects, it is important to know which function calls can throw exceptions, since we might need to perform some cleanup tasks if an exception occurs. Some languages therefore require programmers to annotate functions with a list of all exceptions they can throw. In practice, this has not turned out to be very successful, as many programmers consider this tedious and simply specify that there are no restrictions on the exceptions a function can throw.

Exceptions are uncritical for purely functional code, but problematic for code with side effects.

6 Constraints

6.1 Single-assignment variables

To support logic programming, we have to extend our language with two new constructs. The first is the concept of a *single-assignment variable*. Such variables may start uninitialised, but once a value is assigned it cannot be changed anymore. The only change in syntax to the previous chapter is that we allow to omit the value from let-bindings to introduced uninitialised variables.

$$\langle expr \rangle ::= \dots \mid \mathbf{let} \langle id \rangle ; \langle expr \rangle$$

But the semantics change. Assigning a value to a variable is only allowed if the variable is either uninitialised, or it already has the same value we are assigning.

```
1  let x;  
2  let y;  
3  x := 1;  
4  x := 1; // ok  
5  x := 2; // error  
6  y := x+1;
```

Furthermore, parameter-passing is now by reference.

```
1  let add(x,y,z) {  
2    z := x+y;  
3  };  
4  let u;  
5  add(1,2,u);  
  
1  let reverse(lst, ret) {  
2    case lst  
3    | [] => ret := []  
4    | [x|xs] => {  
5      let rev;  
6      ret := [x|rev];  
7      reverse(xs, rev);  
8    }  
9  };
```

Note that we made the function reverse tail recursive by putting the assignment to ret before the recursive call.

We can also use uninitialised variables in data structures. For instance, to implement lists that allow adding elements at the end.

6 Constraints

```
1  let make() {
2    let empty;
3    Pair(empty, empty)
4  };
5  let add_first(lst, x, ret) {
6    case lst
7    | Pair(first, last) =>
8      ret := Pair([x|first], last)
9  };
10 let add_tail(lst, x, ret) {
11   case lst
12   | Pair(first, last) =>
13     ( let empty;
14       last := [x|empty];
15       ret := Pair(first, empty) )
16   };
```

We can also use single-assignment variables to create cyclic data structures.

```
x := [1,2,3|x]
```

6.2 Unification

An assignment statement $x := e$ is asymmetric as we can only use l-values on the left-hand side, while arbitrary r-values are allowed on the right-hand side. When using single-assignment variables, we can define a symmetric version of the assignment statement which is called *unification*.

$$\langle expr \rangle ::= \dots \mid \langle expr \rangle ::= \langle expr \rangle$$

When unifying two values u and v , we try to assign values to all undefined variables in u and v in such a way that the resulting values become equal. Hence, a unification $u ::= v$ can be seen as solving the equation $u = v$ by substituting values for the variables.

```
1  1 ::= x           x := 1
2  x ::= y           identifies x and y
3  [x,2] ::= [1,y]  x := 1 and y := 2
```

Implementation We can solve equations of the form $u ::= v$ recursively as follows.

- If u is an uninitialised variable, we set it to v .
- If v is an uninitialised variable, we set it to u .
- If $u = m$ and $v = n$ are both numbers, we check that $m = n$. If this is not the case, unification fails.
- If $u = c(s_0, \dots, s_{m-1})$ and $v = d(t_0, \dots, t_{n-1})$ are both constructors, we check that $c = d$, $m = n$, and $s_i ::= t_i$, for all i .

- If $u = [l_0 = s_0, \dots, l_{m-1} = s_{m-1}]$ and $v = [k_0 = t_0, \dots, k_{n-1} = t_{n-1}]$ are both records, we check that there is some bijection $\varphi : m \rightarrow n$ such that $l_i = k_{\varphi(i)}$ and $s_i := t_{\varphi(i)}$, for all i .
- In all other cases, unification fails.

(Note in particular that we cannot unify function values.) There are a few things one has to keep in mind when implementing this procedure.

(1) We have to distinguish two kinds of uninitialised values. If we have just introduced an uninitialised variable x , we know nothing at all about its value. After a unification with another uninitialised variable y , we still do not know the value of x , but we already know that it is equal to that of y . So we need to distinguish between values for completely undefined variables and values for variables that are equal to other variables.

(2) A naïve recursive implementation of unification can go into an infinite loop if we unify cyclic data structures. For instance, the last unification in

```

1  let x; let y;
2  x := [1, x];
3  y := [1, 1, y];
4  x := y;

```

might not terminate. To prevent this, we need to remember during unification which equations we have already checked. If we try to check an equation $u := v$ for the second time, we do not need to recursively call the unification procedure, we can simply skip it and assume that it was successful.

6.3 Backtracking

What do we do if we use single-assignment variables and discover that we have assigned them the wrong value? *Backtracking* is a mechanism for reverting such choices by reverting the whole program to an earlier state. To implement it we add a nondeterministic choice operator to our language.

$$\langle expr \rangle ::= \dots \mid \mathbf{choose} \mid \langle expr \rangle \dots \mid \langle expr \rangle \mid \mathbf{fail}$$

Abstractly, a choice operator selects one of the given expressions that does not cause a failure and executes it. The actual implementation of course does not know which of the expressions will fail. What the operator therefore does is to create a kind of checkpoint and then executes the first expression. If, later on, a failure occurs, the program state saved at the last checkpoint is restored and the next alternative is tried instead. If all alternatives fail, the checkpoint is deleted again and the choose-statement itself fails. This means that only this last alternative executed (the one that succeeded) will have an effect on the program, those that have failed will not. It is if they never were executed. Not that the failure does not need to occur inside the expressions themselves, it may happen later on in the program.

```

1  choose
2  | { x := 1; y := 1; }
3  | { x := 2; y := 2; }

```

6 Constraints

tries first to set two variables to 1. If one of them already has a different value, the corresponding assignment fails and we try to set the variables to 2. If this fails as well, the whole choose-statement fails. In this case, none of the variables is modified. Note that a transaction can only undo memory changes, not other kinds of side-effects. So

```
1  choose
2  | { print "start..." 1; fail; print "stop..." 1; }
3  | print "start..." 2
```

will print out "start... 1start... 2" even though the first expression is aborted.

How do we implement the choose construct? We use two primitive operations `checkpoint k` and `rewind`. The first one takes a continuation as argument and creates a checkpoint storing the current machine state. If later on a `rewind` command is executed, we

- fetch the continuation associated with the last checkpoint,
- restore the machine state to its previous state (which deletes the last checkpoint),
- and call the fetched continuation.

Using these two operations we can translate a choose statement as follows.

```
1  choose |  $e_1$             $\implies$    $e_1$ 
2  choose |  $e_1$  |  $e_2$  ... |  $e_n$   $\implies$  letcc k =>
3                                     checkpoint
4                                     fun () { k(choose |  $e_2$  ... |  $e_n$ ) };
5                                      $e_1$ 
6  fail                                $\implies$   rewind
```

Of course, now we have to figure out how to implement `checkpoint` and `rewind`. Saving the whole machine state is very inefficient. What we will do instead is to record all memory changes and undo them when we `rewind`. As we only use single-assignment variables the only changes we need to undo are assignments of values to uninitialised variables. This can be done by simply marking those variables as uninitialised again. Hence, what we need to do is to store a list of all variables whose value has changed since the last checkpoint. Then `rewind` can traverse the list and undo those changes again.

This means our implementation looks as follows. We maintain a stack of checkpoints. Each entry of the stack contains a stored continuation and the list of variables whose value has changed. A `checkpoint k` command simply pushes a new entry on the stack consisting of `k` and the empty list. Each variable assignment now has to add the variable in question to the list in the top stack entry. Finally, a `rewind` command, retrieves the continuation from the top stack entry, walks the list of variables to mark them as uninitialised again, and then calls the retrieved continuation.

With single-assignment variables and backtracking, we can translate most Prolog programs (which do not use advanced features) into our kernel language. For instance,

```
1  edge(a,b).
2  edge(b,c).
3  trans(X,Y) :- edge(X,Y).
4  trans(X,Y) :- edge(X,Z), trans(Z,Y).
```

turns into

```

1  let edge(x,y) {
2    choose
3    | { x := a; y := b; }
4    | { x := b; y := c; }
5  }
6  let trans(x,y) {
7    choose
8    | edge(x,y)
9    | { let z; edge(x,z); trans(z,y); }
10 }

```

6.4 Programming examples

Let us write our standard list processing examples using single-assignment variables.

```

1  let nth(lst,n,z) {
2    choose
3    | { let t;
4        n := 0;
5        lst := [z|t] }
6    | { let h; let t;
7        lst := [h|t];
8        nth(t,n-1,z) }
9  };
10 let length(lst, n) {
11   choose
12   | { lst := []; n := 0 }
13   | { let h; let t; let m;
14       lst := [h|t];
15       length(t, m);
16       n := m+1 }
17 };
18 let sum(lst, n) {
19   choose
20   | { lst := []; n := 0 }
21   | { let h; let t; let m;
22       lst := [h|t];
23       sum(t,m);
24       n := m+h; }
25 };
26 let map(f, lst, z) {
27   choose

```

6 Constraints

```
28   | { lst := []; z := [] }
29   | { let h; let t; let y;
30       lst := [h|t];
31       z := [f(h)|y];
32       map(f, t, y); }
33 };
34 let fold(f, acc, lst, z) {
35   choose
36   | { lst := []; z := acc }
37   | { let h; let t;
38       lst := [h|t];
39       fold(f, f(acc, h), t, z) }
40 };
41 let foldr(f, acc, lst, z) {
42   choose
43   | { lst := []; z := acc }
44   | { let h; let t; let y;
45       lst := [h|t];
46       foldr(f, acc, t, y);
47       z := f(h, y) }
48 };
49 let append(x,y,z) {
50   choose
51   | { x := []; y := z }
52   | { let h; let t; let r;
53       x := [h|t];
54       z := [h|r];
55       append(t,y,r); }
56 };
57 let reverse(lst, z) {
58   let iter(lst, y, z) {
59     choose
60     | { lst := []; z := y }
61     | { let h; let t;
62         lst := [h|t]; iter(t, [h|y], z) }
63   };
64   iter(lst, [], z)
65 };
```

If we use a more Prolog-like syntax, the code becomes extremely clean.

```
1  nth([x|xs], 0, x).
2  nth([x|xs], i, y) :- nth(xs, i-1, y).
3
4  length([], 0).
```

```

5  length([x|xs], n) :- length(xs, n-1).
6
7  sum([], 0).
8  sum([x|xs], n) :- sum(xs, n-s).
9
10 map(f, [], []).
11 map(f, [x|xs], [f(x)|ys]) :- map(f, xs, ys);
12
13 fold(f, acc, [], acc).
14 fold(f, acc, [x|xs], z) :- foldr(f, f(acc, x), xs, z).
15
16 foldr(f, acc, [], acc).
17 foldr(f, acc, [x|xs], f(acc, z)) :- foldr(f, acc, xs, z).
18
19 append([], y, y).
20 append([x|xs], y, [x|z]) :- append(xs, y, z).
21
22 reverse(lst, rev) :- reverse_helper(lst, [], rev).
23
24 reverse_helper([], z, z).
25 reverse_helper([x|xs], y, z) :- reverse(xs, [x|ys], z).

```

We can use lists terminated by an unbound variable to efficiently implement queues.

```

1  type queue = | Queue(int,list(a),list(a));
2
3  let empty () {
4    let t;
5    Queue(0, t, t);
6  };
7  let is_empty(queue) {
8    case queue
9    | Queue(n,q,t) => n == 0
10 };
11 let insert(queue,x) {
12   case queue
13   | Queue(n,q,t) => { let s; t := [x|s]; Queue(n+1,q,s) }
14 };
15 let first(queue) {
16   case check(queue)
17   | Queue(n,q,t) => head(q)
18 };
19 let remove(queue) {
20   case check(queue)
21   | Queue(n,q,t) => Queue(n-1, tail(q), t)

```

6 Constraints

22 };

7 Objects

Object-oriented programming has created quite a hype after it became mainstream with the release of the first C++ compilers. It was seen as a panacea for all kinds of program design issues, mainly because of its clear advantages over the other mainstream languages of the time, most notably C. Fortunately, this hype is slowly fading over the last years, so a rational discussion of object-oriented programming is now possible.

Unfortunately, there is no standard definition of object-orientation as everybody uses his or her own version. The initial idea was to make the global state of a program more manageable by breaking it into smaller parts called *objects*. Each of these objects has its own local state which is not accessible to the outside. To communicate objects can pass *messages* between them. Thus, as a slogan we could say that object-orientation combines *encapsulated state* plus *message passing*.

At least that was the initial idea. Over time the meaning has changed slightly. Nowadays when introducing object-oriented programming one usually mentions as a key idea the concept of combining data and functions operating on it into a single data structure. According to this newer definition, an object is simply a record containing both functions and non-function values. In addition one usually considers a certain set of additional languages features (such as inheritance) to be an essential part of the definition. Which exactly depends on the person.

Still, the original definition is quite useful as it tells us *how to use* object-orientation, whereas the newer one simply tells us *how it is implemented*. In the following sections we will present several language features that can be used to implement object-oriented programming, or to make it more useful. In particular, we will consider

- dynamic dispatch,
- subtyping,
- encapsulated state,
- inheritance.

7.1 Dynamic dispatch

We will implement the features of object-oriented programming step by step starting with dynamic dispatch. If we want to send a certain message to an object, we do not know statically which function to call. Therefore, we have to store the function with the object and look it up at runtime. An easy way to do so is to represent the object as a record of functions, one for each message. For instance, a list object supporting the messages

```
1 next  : unit -> object
2 get   : unit -> int
3 iter  : (int -> unit) -> unit
```

7 Objects

```
4 length : unit -> int
```

would be represented by a record of type

```
1 [ next   : unit -> object,
2   get    : unit -> int,
3   iter   : (int -> unit) -> unit,
4   length : unit -> int ];
```

To send a message, we just call the corresponding function.

```
1 let new_empty() {
2   let n =
3     [ next   = fun () { n },
4     get     = fun () { error },
5     iter    = fun (f) { () },
6     length = fun() { 0 } ];
7   n
8 };
9 let new_node(val, next) {
10  [ next   = fun () { next },
11  get     = fun () { val },
12  iter    = fun (f) { f(val); next.iter(f) },
13  length = fun () { 1 + next.length() } ]
14 };
15
16 let n1 = new_empty();
17 let n2 = new_node(1,n1);
18 let n3 = new_node(2,n2);
19 n3.iter(fun (x) { print "value is " x });
```

This direct encoding via records quickly becomes cumbersome, but the right kind of syntactic sugar makes it usable.

```
1 object {  $m_1 : t_1 ; \dots ; m_k : t_k ;$  }
2  $\implies$  [  $m_1 : t_1 , \dots , m_k : t_k$  ]
3
4 object {
5    $m_1 ( a_1 ) \{ b_1 \} ;$ 
6   ...
7    $m_k ( a_k ) \{ b_k \} ;$ 
8 }
9  $\implies$ 
10 [  $m_1 = \mathbf{fun} ( a_1 ) \{ b_1 \},$ 
11   ...
12    $m_k = \mathbf{fun} ( a_k ) \{ b_k \}$  ]
```

With this notation we can write the above code as


```

1  type olist =
2  object {
3      next    : unit -> olist;
4      get     : unit -> int;
5      iter    : (int -> unit) -> unit;
6      length  : unit -> int
7  };
8
9  let new_empty() {
10     let n =
11         object {
12             next() { n };
13             get()  { error };
14             iter(f) { () };
15             length() { 0 };
16         };
17     n
18 };
19 let new_node(val, next) {
20     object {
21         next() { next };
22         get()  { val };
23         iter(f) { f(val); next.iter(f) };
24         length() { 1 + next.length() };
25     }
26 };
27
28 let n1 = new_empty();
29 let n2 = new_node(1,n1);
30 let n3 = new_node(2,n2);
31 n3.iter(fun (x) { print "value is " x });

```

Note that this approach of representing objects as record is based on structural type equivalence. Two objects (like empty and node above) have the same type if they support the same set of methods. Most object-oriented languages use name equivalence instead and would consider empty and node to have different types. In languages such as C++ there is nothing corresponding to object types like the type `olist` in the above example. But many of the modern object-oriented languages that are based on name equivalence have added such types as an additional concept. For instance, in Java they are called *interfaces*.

Example Our running example in this chapter will be a class hierarchy for geometric shapes. This is still simple enough to (mostly) fit on a single page but shares many properties with the more complicated hierarchies one finds in real-world programs, like class hierarchies for graphical user interfaces.

7 Objects

```
1  type shape = object {
2    draw      : unit -> unit;
3    move      : int -> int -> shape;
4    dimensions : unit -> [ min_x : int, min_y : int, max_x : int, max_y : int ];
5  };
6
7  let new_point(x : int, y : int) : shape =
8    object {
9      draw()      { draw_point(x,y)          };
10     move(dx, dy) { new_point(x + dx, y + dy) };
11     dimensions() { [ min_x = x, min_y = y, max_x = x, max_y = y ] };
12   };
13
14  let new_circle(x : int, y : int, r : int) : shape =
15    object {
16     draw()      { draw_circle(x,y,r)        };
17     move(dx, dy) { new_circle(x + dx, y + dy, r) };
18     dimensions() { [ min_x = x - r, min_y = y - r,
19                    max_x = x + r, max_y = y + r ] };
20   };
21
22  let new_rectangle(x : int, y : int, w : int, h : int) : shape =
23    object {
24     draw()      { draw_rectangle(x,y,w,h)    };
25     move(dx, dy) { new_rectangle(x + dx, y + dy, w, h) };
26     dimensions() { [ min_x = x, min_y = y,
27                    max_x = x + w, max_y = y + h ] };
28   };
29
30  let new_group(shapes : list(shape)) {
31    object {
32     draw()      { iter(fun (s) { s.draw() }, shapes) };
33     move(dx, dy) { new_group(map( fun (s) { s.move(dx, dy) }, shapes)) };
34     dimensions() {
35       case shapes
36       | []      => [ x = 0, y = 0, w = 0, h = 0 ]
37       | [s:ss] => let d = s.dimensions();
38                  fold(fun(d,s) {
39                    let d2 = s.dimensions();
40                    [ min_x = min(d.min_x, d2.min_x),
41                      min_y = min(d.min_y, d2.min_y),
42                      max_x = max(d.max_x, d2.max_x),
43                      max_y = max(d.max_y, d2.max_y) ]
44                  },
```

```

45         s.dimensions(),
46         ss)
47     };
48 };
49 };

```

Multi-methods One problem with dynamic dispatch as defined above is that it is asymmetric with respect to its arguments. The object we dispatch on is treated differently than the other arguments. Some languages have introduced the possibility to dispatch on the types of all arguments. This is called *multi-methods*.

```

1  let intersect(x : circle, y : circle) : shape { ... }
2  let intersect(x : circle, y : rectangle) : shape { ... }
3  let intersect(x : rectangle, y : circle) : shape { ... }
4  let intersect(x : rectangle, y : rectangle) : shape { ... }

```

The problem with multi-methods is that, as the number of classes grows, defining functions for all combinations quickly becomes unmanageable. While there are languages that support multi-methods, the approach has never really become popular.

Type classes An alternative approach to dynamic dispatch is provided by Haskell's *type classes*. A type class consists of a collection of function types associated with one or more parameter types. For each choice of parameter types, we can define an *instance* of the type class by providing an implementation of the required functions.

```

1  typeclass Shape(a) {
2    draw      : a -> unit;
3    move      : a -> int -> int -> a;
4    dimensions : a -> [ min_x : int, min_y : int, max_x : int, max_y : int ];
5  };
6
7  type point = Point(int,int);
8
9  instance Shape(point) {
10   draw(Point(x,y))      { draw_point(x,y)      };
11   move(Point(x,y), dx, dy) { Point(x + dx, y + dy) };
12   dimensions(Point(x,y)) { [ min_x = x, min_y = y, max_x = x, max_y = y ] };
13 };
14
15 type circle = Circle(int,int,int);
16
17 instance Shape(circle) {
18   draw(Circle(x,y,r))      { draw_circle(x,y,r)      };
19   move(Circle(x,y,r), dx, dy) { Circle(x + dx, y + dy, r) };
20   dimensions(Circle(x,y,r)) { [ min_x = x-r, min_y = y-r,

```

7 Objects

```
21         max_x = x+r, max_y = y+r ] };
22     };
```

Comparison with variant types There is an alternative solution based on variant types instead of objects. We could have defined

```
1  type shape =
2  | Point(int,int)
3  | Circle(int,int,int)
4  | Rectangle(int,int,int,int)
5  | Group(list(shape));
6
7  let draw(sh) {
8  case sh
9  | Point(x,y)      => draw_point(x,y)
10 | Circle(x,y,r)   => draw_circle(x,y,r)
11 | Rectangle(x,y,w,h) => draw_rectangle(x,y,w,h)
12 | Group(shapes)  => iter(draw, shapes)
13 };
14 let move(sh, dx, dy) {
15 case sh
16 | Point(x,y)      => Point(x + dx, y + dy)
17 | Circle(x,y,r)   => Circle(x + dx, y + dy, r)
18 | Rectangle(x,y,w,h) => Rectangle(x + dx, y + dy, w, h)
19 | Group(shapes)  => Group(map(fun (s) { move(s, dx, dy) }, shapes))
20 };
21 let dimensions(sh) {
22 case sh
23 | Point(x,y)      => ...
24 | Circle(x,y,r)   => ...
25 | Rectangle(x,y,w,h) => ...
26 | Group(shapes)  => ...
27 };
```

The only difference is the way we have grouped the code. In the object-based solution we collect all code pertaining to a given shape in one place, whereas when using variant types we collect all code pertaining to a given operation on shapes in one place. The difference becomes noticeable if we want to extend the program. If we add a new shape, say a triangle, the object-based approach is more convenient, we only need to define a new class. In the variant-type-based solution we have to modify every operation to add a new case. If, on the other hand, we add a new operation, like rotation, then the solution using variant types is more convenient. In the object-based approach we have to modify every class definition.

7.2 Subtyping

A type s is a *subtype* of a type t if values of type s can be used everywhere a value of type t is expected. This means that s is more specialised than t , or t more general than s . We write $s \leq t$ to denote this fact. As with type equivalence there are two different approaches to implement subtyping: *structural* and *by name*. In languages like Java where subtyping is defined by name, the programmer has to explicitly declare if one object type is to be a subtype of another. In languages with structural subtyping on the other hand, a type s is automatically a subtype of all types that are more general than s . This means that, if s and t both are object types, then s will be a subtype of t if s supports all the methods of t . For instance, if we have defined a class of shapes with methods `draw`, `move`, and `box` and a subclass of rectangles with an additional method `area`, then the rectangle class is a subtype of the shape class.

Programming languages have a certain choice in how exactly to define the subtyping relation. Let us discuss the possibilities for some of the usual types. It does make a difference whether we have *immutable* or *mutable* values. We start with the case where all values are immutable. It is possible to already define subtyping relations between basic types. For instance, we could have

$$\text{int16} \leq \text{int32} \leq \text{int64} \quad \text{or} \quad \text{uint16} \leq \text{int32}$$

What about

$$\text{uint32} \leq \text{int32} \quad \text{or} \quad \text{int32} \leq \text{float32} ?$$

For records we have

$$[l_1 : s_1, \dots, l_m : s_m, k_1 : t_1, \dots, k_n : t_n] \leq [l_1 : u_1, \dots, l_m : u_m]$$

if $s_i \leq u_i$ for all i , that is, if every label appearing in the second record is also present in the first one with a subtype of the corresponding type on the right-hand side.

Example

```

1  type shape      = [ x : int, y : int ];
2  type circle    = [ x : int, y : int, r : int ];
3  type rectangle = [ x : int, y : int, w : int, h : int ];
4
5  circle <: shape  and  rectangle <: shape
```

Exercise What is the subtype ordering for variant types?

What about functions? Suppose we have a function of type $f : a \rightarrow b$. When can we use it at a place where a function of type $c \rightarrow d$ is expected? f will get passed a value of type c (so $c <: a$) and it will return a value of type b where one of type d is expected (so $b <: d$).

```

1  let g(f : c -> d) = {
2    ...
3    let x : c = ...;
```

7 Objects

```
4   let y : d = f(x);
5   ...
6   };
```

This means that

$$a \rightarrow b <: c \rightarrow d \quad \text{iff} \quad c <: a \quad \text{and} \quad b <: d.$$

Note that the orders for the parameter and the return value are different. We say that functions types are *contravariant* (the order is reversed) in the parameter position and *covariant* (the order is the same) in the result position. In general a type constructor is contravariant in all types used as inputs and covariant in all types used as outputs. If an argument type is used both as input and output, the constructor is *invariant*.

Example

```
1  type shape      = [ x : int, y : int ];
2  type circle     = [ x : int, y : int, r : int ];
3  type rectangle  = [ x : int, y : int, w : int, h : int ];
4
5  shape -> circle <: rectangle -> circle <: rectangle -> shape
```

Example The most important example of invariant constructors is the case of mutable data structures.

```
1  type box(a) = [ data : a ];
2
3  let get(box : box(a)) : a {
4    box.data
5  };
6  let set(box : box(a), x : a) : unit {
7    box.data := x
8  };
```

When is $\text{box}(a) <: \text{box}(b)$? Suppose that $\text{box}(a) <: \text{box}(b)$. Then applying $\text{get} : \text{box}(b) \rightarrow b$ to a value of type $\text{box}(a)$, we need to get a value of some subtype of b . Hence, we must have $a <: b$. Furthermore, if we call $\text{set} : \text{box}(b) \rightarrow b \rightarrow \text{unit}$ with a box of type $\text{box}(a)$ and an element of type b , and then apply $\text{get} : \text{box}(a) \rightarrow a$ to that box, we need to get an element of type a . Hence, we also must have $b <: a$.

Subtyping for objects Many languages define simpler subtyping relations than the most general one we have described above. In particular, when determining whether some class is a subclass of another one, object-oriented languages frequently require the types of methods to match exactly instead of one being just a subtype of the other one. This makes type checking faster and, more importantly, the type system less complex.

In fact, this form of subtyping is simple enough that it can be emulated by a certain variant of parametric polymorphism called *row polymorphism*. The idea is to allow parameters in record (and object) types of the form

$$[l_0 : t_0, \dots, l_{n-1} : t_{n-1}, \alpha]$$

which can be instantiated with further label declarations. For instance, instantiating the α in the above record type with the value $k_0 : s_0, k_1 : s_1, \beta$ yields the record type

$$[l_0 : t_0, \dots, l_{n-1} : t_{n-1}, k_0 : s_0, k_1 : s_1, \beta]$$

Then we have a subtyping relation

$$[l_0 : t_0, \dots, \alpha] <: [k_0 : s_0, \dots, \beta]$$

between two such types if we can obtain the later by a suitable instantiation of the parameter α in the former. Hence, we can simulate object types with subtyping by identifying an object type

object $m_0 : t_0, \dots, m_{n-1} : t_{n-1}$ **end**

with the record type

$$[m_0 : t_0, \dots, m_{n-1} : t_{n-1}, \alpha]$$

In this context of subtyping for objects let us also mention the language Eiffel, where the definition allows subtypes when comparing methods. But the designer of Eiffel consciously chose to define subtyping for functions to be covariant in *both* types. This leads to an unsound type system since the programmer is allowed to pass arguments of unsupported types to a function (in which case Eiffel generates a runtime-exception). The reason for this decision was that it was felt that contravariance was too confusing for the programmer. But it is questionable whether this solution is any less confusion.

Let us conclude this section with an example showing the advantages of subtyping. One area where it can be superior to parametric polymorphism is one wants to use *heterogeneous data structures*. For instance, using subtyping it is possible to have a list containing both circles and rectangles, whereas when using parametric polymorphism we have to decide which of the two kinds of objects we want to put into the list.

7.3 Encapsulated state

We have shown above how to implement purely functional objects. Now it is time to add mutable state. We can do so by simply combining dynamic dispatch with side-effects.

```

1  type account = object {
2      deposit  : int -> unit;
3      withdraw : int -> unit;
4  };

```

7 Objects

```
5  let new_account(balance) {
6    object {
7      deposit(amount) { balance := balance + amount };
8      withdraw(amount) { balance := balance - amount };
9    }
10 };
```

As are more involved example let us give a version of the shape class with internal state.

```
1  type shape = object {
2    draw      : unit -> unit;
3    move      : int -> int -> unit;
4    dimensions : unit -> [ min_x : int, min_y : int, max_x : int, max_y : int ];
5  };
6
7  let new_point(x : int, y : int) : shape {
8    object {
9      draw()      { draw_point(x,y)      };
10     move(dx, dy) { x := x + dx; y := y + dy; };
11     dimensions() { [ min_x = x, min_y = y, max_x = x, max_y = y ] };
12   }
13 };
14
15 let new_circle(x : int, y : int, r : int) : shape {
16   object {
17     draw()      { draw_circle(x,y,r)      };
18     move(dx, dy) { x := x + dx; y := y + dy; };
19     dimensions() { [ min_x = x - r, min_y = y - r,
20                     max_x = x + r, max_y = y + r ] };
21   }
22 };
23
24 let new_rectangle(x : int, y : int, w : int, h : int) : shape {
25   object {
26     draw()      { draw_rectangle(x,y,w,h) };
27     move(dx, dy) { x := x + dx; y := y + dy; };
28     dimensions() { [ min_x = x,      min_y = y,
29                     max_x = x + w, max_y = y + h ] };
30   }
31 };
32
33 let new_group(shapes : list(shape)) {
34   object {
35     draw()      { iter(fun (s) { s.draw()      }, shapes) };
36     move(dx, dy) { iter(fun (s) { s.move(dx, dy) }, shapes) };
37   }
38 };
```



```

37     dimensions() { ... };
38   };
39 };

```

7.4 Inheritance

With the object framework introduced so far, we have to write every class from scratch. It would be desirable to share common code between classes. Besides requiring less typing, this also increases code maintainability as changes in the in question code do not have to be repeated for every class. On the negative side, one has to note that such sharing reduced code locality, as the implementation of a class becomes distributed over several parts of the program. We will call the mechanism for code sharing within a class framework *inheritance*, although strictly speaking this term only refers to using code from a parent class in a subclass. There are several ways to support inheritance, some more problematic than others.

Delegates Suppose we want to add classes for coloured shapes to the class hierarchy defined above. A coloured shape has two more methods to access the colour of the shape.

```

1  type coloured_shape = object {
2    draw      : unit -> unit;
3    move      : int -> int -> shape;
4    dimensions : unit -> [ min_x : int, min_y : int, max_x : int, max_y : int ];
5    colour    : colour;
6    set_colour : colour -> unit;
7  };

```

One easy way to create objects for such a class is to use an object of the parent class as is and implement the methods of the new class using those of the old one. An object used as part of another one in this way is called a *delegate*.

```

1  let new_coloured_point(x : int, y : int, c : colour) : coloured_shape =
2    let p = new_point(x,y);
3    object {
4      draw()      { p.draw() };
5      move()      { p.move() };
6      dimensions() { p.dimensions() };
7      colour()    { c };
8      set_colour(col) { c := col };
9    };

```

Adding methods In the above example we directly called the methods of the delegate without any changes. In this case we can simplify the code slightly as follows.

```

1  let new_coloured_point(x : int, y : int, c : colour) : coloured_shape =
2    let p = new_point(x,y);

```

7 Objects

```
3   [ draw      = p.draw,
4     move      = p.move,
5     dimensions = p.dimensions,
6     colour    = fun () { c },
7     set_colour = fun (col) { c := col } ];
```

Adding syntactic sugar we can neaten up the code further and obtain something like the following.

```
1   let new_coloured_point(x : int, y : int, c : colour) : coloured_shape =
2     let p = new_point(x,y);
3     object {
4       include p;
5       colour()      { c },
6       set_colour(col) { c := col };
7     }
```

Replacing methods Just adding new methods is not always enough. Sometimes we also want to change existing ones. Let us first see how to replace an old methods with a completely new one.

```
1   let new_coloured_point(x : int, y : int, c : colour) : coloured_shape =
2     let p = new_point(x,y);
3     [ draw      = fun () { set_drawing_colour(c); draw_point(x,y); },
4       move      = p.move,
5       dimensions = p.dimensions,
6       colour    = fun () { c },
7       set_colour = fun (col) { c := col } ];
```

In the following example, some methods of the superclass are mere stubs that are intended to be overwritten by each subclass. This is a common idiom in languages like C++ that use name equivalence for subtyping and that do not support object types (interfaces in Java's terminology). In such a language we can emulate object types in the following way via inheritance and subtyping.

```
1   class shape {
2     draw()          { () };
3     move(dx : int, dy : int) { () };
4     dimensions()    { [ min_x = 0, min_y = 0, max_x = 0, max_y = 0 ] };
5   };
6
7   class point(x : int, y : int) extends shape {
8     draw()          { draw_point(x,y) };
9     move(dx, dy) { x := x + dx; y := y + dy; };
10    dimensions() { [ min_x = x, min_y = y, max_x = x, max_y = y ] };
11  };
12
13  class circle(x : int, y : int, r : int) extends shape {
14    draw()          { draw_circle(x,y,r) };
15  };
```

```

15   move(dx, dy) { x := x + dx; y := y + dy; };
16   dimensions() { [ min_x = x - r, min_y = y - r,
17                   max_x = x + r, max_y = y + r ] };
18 };
19
20 class rectangle(x : int, y : int, w : int, h : int) extends shape {
21   ...
22 };
23
24 class group(shapes : list(shape)) extends shape {
25   ...
26 };

```

Modifying methods In the implementation of coloured points above we did repeat the code of the old methods in the definition of the new one. We can increase the amount of code reuse by using the old methods instead.

```

1  let new_coloured_point(x : int, y : int, c : colour) : coloured_shape =
2    let super = new_point(x,y);
3    [ draw      = fun () { set_drawing_colour(c); super.draw(); },
4      move      = super.move,
5      dimensions = super.dimensions,
6      colour    = fun () { c },
7      set_colour = fun (col) { c := col } ];

```

One question one has to address when designing an inheritance mechanism for a programming language is who is in command, the subclass or the superclass? That is, when invoking a method of a subclass, do we execute the function given in the subclass definition (which then may or may not call the function of the superclass), or do we execute the function of the superclass (which then can call the function of the subclass)? For instance, suppose we use a class hierarchy to model widgets in a graphical user interface. We might define a class for a general kind of text field and several subclasses for more special versions. Consider the method that gets called when the user presses a key. Do we want the superclass to first process this key press and then pass the keys it is not interested in to the subclass, or do we want it to be the other way round? The following examples illustrate the differences between these two approaches.

Example Let us discuss the various choices on how to use inheritance in the example of a class for buttons in a user interface. In most object-oriented GUI frameworks they are implemented using inheritance with modification.

```

1  type button = object {
2    button_down : unit -> unit;
3    button_up   : unit -> unit;
4    ...
5  };

```

7 Objects

```
6  let basic_button() {
7    object {
8      button_down(self) { ... draw button ... };
9      button_up(self)   { ... draw button ... };
10     ...
11   };
12 };
13 let my_button() {
14   let super = basic_button();
15   object {
16     include super;
17     button_down(self) {
18       super.button_down(self);
19       ... do something ...
20     };
21     ...
22   }
23 };
```

Note that in this solution, it is not obvious how a subclass is intended to call the button superclass. When should it call the `button_down` method of the superclass? At the beginning of its own method, at the end, somewhere in between? Should it call it at all? Here we see why it is sometimes better to be able to call subclass methods via `outer` instead of superclass methods via `super`.

We can clean this design up, by splitting the `button_down` method into two parts. One part to be overwritten by the superclass and one to be left alone.

```
1  type button = object {
2    button_down  : unit -> unit;
3    button_up    : unit -> unit;
4    button_pressed : unit -> unit;
5    ...
6  };
7  let basic_button() {
8    object {
9      button_down(self) {
10       ... draw button ...
11       self.button_pressed();
12     };
13     button_up(self) { ... draw button ... };
14     button_pressed(self) { () };
15     ...
16   }
17 };
18 let my_button() {
19   let super = basic_button();
```

```

20  object {
21    include super;
22    button_pressed(self) {
23      ... do something ...
24    };
25    ...
26  }
27 };

```

Finally, we can simplify our implementation further, by using a first-class function instead of inheritance.

```

1  type button = object {
2    button_down    : unit -> unit;
3    button_up      : unit -> unit;
4    ...
5  };
6  let basic_button(pressed : unit -> unit) {
7    object {
8      button_down(self) {
9        ... draw button ...
10       pressed();
11     },
12     button_up(self) {
13       ... draw button ...
14     };
15     ...
16   }
17 };

```

In this case we do not need to define new classes at all. We can simply use the base class as is.

Type classes Type classes also offer two of the forms of inheritance discussed above. Firstly, we can extend a given type class with new functions.

```

1  typeclass Eq(a) {
2    equal      : a -> a -> bool;
3    not_equal  : a -> a -> bool;
4  };
5  typeclass Eq(a) => Ord(a) {
6    type cmp = | LT | EQ | GT;
7    compare  : a -> a -> cmp;
8  };
9
10 instance Eq(int) {
11   equal(x,y) { prim_equal_int(x,y) };

```

7 Objects

```
12   not_equal(x,y) { not(equal(x,y)) };
13 };
14   instance Ord(int) {
15     compare(x,y) { if x < y then LT else if x > y then GT else EQ };
16 };
```

Secondly, a type class can offer a default implementation that may be overwritten by the instance.

```
1   typeclass Eq(a) {
2     equal      : a -> a -> bool;
3     not_equal  : a -> a -> bool;
4     not_equal(x,y) { not(equal(x,y)) };
5   };
6
7   instance Eq(int) {
8     equal(x,y) { prim_equal_int(x,y) };
9   };
```

Multiple inheritance Inheritance is mainly a mechanism to reuse code from existing objects. Sometimes one would like to use code of several objects at once. Therefore some languages (most notably C++) allow to define classes that extend several superclasses at the same time. This is called *multiple inheritance*. While adding more power to the language, multiple inheritance makes the object system also considerably more complicated. For instance, what happens if several of the superclasses have methods with the same name? Does this result in an error, or do we simply pick one of the methods for the subclass? Another problematic situation is the following one. Suppose we have two classes B and C that both inherit from some class A . What happens if we inherit a class D from both B and C ? Do we get two copies of the class A or only one? For these reasons, many modern languages do not support multiple inheritance and try to provide alternative, cleaner ways to achieve the same effects. For instance, Java does only support single inheritance, but it allows classes to implement multiple interfaces.

Mixins The main reason why multiple inheritance is problematic, is the fact that in most languages inheritance is the only mechanism to define class hierarchies. A declaration like

```
class B extends A { ... }
```

both declares B as a subclass of A and lets B inherit the methods of A . If we provide separate mechanisms for inheritance and the declaration of subtyping relationships, the object system becomes much simpler and cleaner.

How could an inheritance mechanism look like that is decoupled from subtyping? One example of such a mechanism is called *mixins*. A mixin is a function $F : I \rightarrow J$ that takes a class A of a specified object type I and produces a new class $F(A)$ of type J . Hence, a mixin is very similar to the parametrised modules we have described in Section 4.4. As an example let us consider a mixin that turns shapes into coloured shapes.

```
1   type coloured_shape = object {
```

```

2   draw      : unit -> unit,
3   move      : int -> int -> shape,
4   dimensions : unit -> [ min_x : int, min_y : int, max_x : int, max_y : int ],
5   colour    : colour,
6   set_colour : colour -> unit
7 };
8
9   let make_coloured(s : shape, c : colour) : coloured_shape =
10  [ draw      = s.draw,
11    move      = s.move,
12    dimensions = s.dimensions,
13    colour    = fun () { c },
14    set_colour = fun (col) { c := col } ];

```

So, instead of using inheritance to extend a class A to a subclass B , we can use a mixin F such that $B = F(A)$. In some cases, we can use mixins also to simulate multiple inheritance. Suppose that A is a common superclass of both B and C , and D inherits from B and C . If we can write $B = F(A)$ and $C = G(A)$ with mixins F and G , then we can try to express $A = H(G(F(A)))$ as an extension of $G(F(A))$ via a third mixin H .

7.5 Discussion

The problem with many object-oriented languages is that they offer a single mechanism (class definitions) that combines all the object features. This makes the language very complex and has led to much confusion about object-oriented design. A much cleaner and simpler design is to provide separate mechanisms for subtyping, dynamic dispatch, encapsulated state, and inheritance.

For instance, there is an old debate on whether subclasses should represent an ‘is-a’ or a ‘has-a’ relationship. Separating the aspects of object-oriented design, we see that an ‘is-a’ relationship is precisely modelled by the subtyping relation, whereas a ‘has-a’ relationship is more suitably modelled by some form of inheritance or the use of delegates.

Separation also allows one to only use those features necessary for a particular solution. For instance, if stateless objects are sufficient for the task at hand, we can avoid the added complexities involved with side-effects.

Let me conclude this chapter with a word of advice: while subtyping and object-oriented programming as a whole are quite powerful, but they are also quite complex. They are not always the best way to solve a problem. Only use them if they make the resulting program simpler. If a purely functional solution, or a plain imperative one, works as well, there is no need to resort to objects.

Also one can easily get carried away with designing elaborate class hierarchies, instead of writing code that actually does something. For instance, if you are about to define several helper classes to perform a single task, you should ask yourself whether you really need that many classes or whether a different approach would not offer a simpler solution. For instance, if it is possible to use them, higher-order functions and parametric polymorphism are usually the better approach.

8 Concurrency

So far in our language the evaluation order is completely *deterministic*. If we run a program several times, we observe the same ordering every time. In this chapter we study language features that cause the evaluation order to be *non-deterministic*. In this case we say that the program is *concurrent*. The most common case of concurrency is when the program consists of several threads of processes that are executed in *parallel*. But there are also examples of concurrency without parallelism.

Of course, concurrency increases the complexity of programs and makes them harder to reason about, in particular, when combined with side-effects, which as we have seen care about the evaluation order of expressions. Purely functional, concurrent programs are much better behaved.

On the language level concurrency manifests in having several independent ‘paths of execution’, that is, instead of having a single code location identifying the expression that is to be executed next, there can be several such locations. Each of the expressions pointed to will be executed eventually, but the ordering in which this happens is left unspecified. Such a path of execution is called a *fibre* of the program. If fibres are executed in parallel, they are also called *threads* or *processes*. Thus, a fibre is a part of the program with its own control flow. Within each fibre execution is linear, but the program execution can jump between fibres at certain points. Even without parallelism, fibres have the advantage that the program is no longer restricted to a single syntactic nesting and stack discipline. It can use several of them in parallel.

The main problem of concurrent programming is to organise the communication between different fibres. This is called *synchronisation*. There are two fundamentally different methods for synchronisation: *message passing* and *shared memory*. We will consider each one in turn.

8.1 Fibres

Before turning to the synchronisation problem let us first see how to implement fibres in our language. As real parallelism requires support from the operating system, our implementation will be non-parallel and based on a form of *cooperative multi-threading*. We start with the following functions.

```
1  make_ready : (unit -> unit) -> unit;  
2  schedule   : unit -> unit;  
3  spawn      : (unit -> unit) -> unit;  
4  yield      : unit -> unit;
```

The two low-level functions `make_ready` and `schedule` respectively add a fibre to the list of running fibres and execute the next fibre in the list. The high-level functions `spawn` and `yield` respectively create a new fibre and mark a place where we can switch from one fibre to another one. We can implement them as follows.

8 Concurrency

```
1  let ready_fibres = Queue.make ();
2
3  type terminate = | Terminated;
4
5  let make_ready(f) {
6    Queue.push(ready_fibres, f);
7  };
8
9  let schedule() {
10   if Queue.is_empty(ready_fibres) then
11     throw Terminated
12   else {
13     let f = Queue.pop(ready_fibres);
14     f();
15     schedule()
16   }
17 };
18
19 let start_scheduler() {
20   try
21     schedule()
22   catch e => case e
23     | Terminated => ()
24     | else          => print "uncaught exception" e
25 };
26
27 let spawn(f) {
28   letcc k => {
29     make_ready(k);
30     make_ready(f);
31     schedule()
32   }
33 };
34
35 let yield() {
36   letcc k => { make_ready(k); schedule() }
37 };
```

The module `Queue` contains a simple queue implementation where we can add elements at one end and remove them at the other one.

```
1  let f1() {
2    for i = 1 .. 10 {
3      print "fibre1:" i;
4      yield();
```

```

5   }
6   };
7   let f2() {
8     for i = 1 .. 10 {
9       print "fibre2:" i;
10      yield();
11    }
12  };
13  spawn(f1);
14  spawn(f2);
15  start_scheduler()

```

With only these two operation fibres of not of much use. We also need some operations for synchronisation of and communication between fibres. We start by defining a few primitive operations that can be then used to implement more complex communication mechanisms. These operations are based on the notion of a *condition*. A fibre can wait on such a condition and other fibres can wake them up again.

```

1  type condition(a);
2  new_condition : unit -> condition(a);
3  wait          : condition(a) -> a;
4  wait_multi   : list(condition(a)) -> a;
5  resume       : condition(a) -> a -> unit;

```

`new_condition` creates a new condition, `wait(c)` sends the current fibre to sleep, waiting on the condition `c`, and `resume(c)` wakes up all fibres waiting on `c`.

```

1  type trigger(a) = [ cont : a -> void, triggered : bool ];
2  type condition(a) = [ waiting : list(a -> void) ];
3
4  let make_trigger(k) {
5    [ cont = k, triggered = False ]
6  };
7
8  let resume_trigger(t,v) {
9    if t.triggered then
10     ()
11   else {
12     t.triggered := True;
13     make_ready(fun () { t.cont(v) });
14   }
15 };
16
17 let new_condition() {
18   [ waiting = [] ]
19 };

```

8 Concurrency

```
20
21 let resume(c,v) {
22     let waiting = c.waiting;
23     c.waiting := [];
24     List.iter(fun (k) { resume_trigger(k,v) },
25             waiting);
26 };
27
28 let wait(c) {
29     letcc k => {
30         let t = make_trigger(k);
31         c.waiting := [t | c.waiting];
32         yield();
33     }
34 };
35
36 let wait_multi(cs) {
37     letcc k => {
38         let t = make_trigger(k);
39         List.iter(fun (c) { c.waiting := [t | c.waiting] },
40                 cs);
41         yield();
42     }
43 };

1 let c1 = new_condition();
2 let c2 = new_condition();
3
4 let f1() {
5     for i = 1 .. 10 {
6         print "fibre1:" i;
7         resume(c2);
8         wait(c1);
9     };
10    resume(c2);
11 };
12
13 let f2() {
14     for i = 1 .. 10 {
15         print "fibre2:" i;
16         resume(c1);
17         wait(c2);
18     };
19    resume(c1);
```

```

20 };
21
22 spawn(f1);
23 spawn(f2);
24 start_scheduler()

```

8.2 Ramifications

When adding concurrency to a language many new phenomena and problems arise. Let us discuss a few of them.

Mutual exclusion When two regions of code in different fibres want to modify the same data structure, it is usually required that the control flows of the two fibres do not enter these regions simultaneously. This is called the *mutual exclusion problem* and most synchronisation problems can be reduced to it. Many of the synchronisation mechanisms we will introduce below have specifically been invented to ensure mutual exclusion.

Deadlocks A *deadlock* is a situation where several fibres each waiting on an action that can only be performed by one of the other waiting fibres.

Race conditions A *race condition* is a bug that is caused by the particular choices and timing of the scheduler.

Starvation If a fibre is ready but it is never executed because there is always another fibre that goes first, we say the fibre is *starving*.

Lifeness Lifeness is the opposite of starvation: every ready fibre is executed eventually.

Fairness When several fibres compete for a certain resource, we ideally want them to get access to the resource in equal amounts. This is called *fairness*.

8.3 Message passing

Having a basic implementation of fibres we can turn to communication mechanisms. We start with message passing. The central concept of message passing are objects called *channels* or *ports*. A channel is a line of communication between processes that supports two operations: one process can write data, a *message*, to the channel and the other one can read it. Channels come in many variants. They might be

- *synchronous*: a writer waits until the other process reads the message, a reader waits until the other process has supplied a message;

8 Concurrency

- *asynchronous and buffered*: a writer can continue immediately after sending the message, a reader must still wait until a message is available; there can be several messages waiting for the reader;
- *asynchronous and unbuffered*: a writer can continue immediately after sending the message, a reader must still wait until a message is available; there can be at most one message waiting for the reader; if the writer wants to send a second message, he blocks until the reader has read the first one;
- *one-directional*: a channel is split into two parts: one for reading and one for writing, if a process has access to only one of the parts, it can perform only the corresponding operation:
- *bidirectional*: each end of a channel can be used both for reading and for writing.

Before giving an example, let us describe the library functions we need to implement.

```
new_channel : unit -> channel(a)
send        : channel(a) * a -> unit
receive     : channel(a) -> a
```

```
new_channel()  creates a new channel
send(ch,x)     sends the value x over the channel ch
receive(ch)    reads a value from the channel ch
```

In our implementation, channels are synchronous and bidirectional.

```
1  type channel_state(a) = | Free | Reading | Written(a);
2  type channel(a)       = [ state   : channel_state(a),
3                          readers  : condition(a),
4                          writers  : condition(unit) ];
5
6  let new_channel() {
7    [ state   = Free,
8      readers = new_condition(),
9      writers = new_condition() ]
10 };
11
12 let receive(ch) {
13   case ch.state
14   | Free      => { ch.state := Reading; wait(ch.readers) }
15   | Written(v) => { ch.state := Free;   resume(ch.writers, ()); v }
16   | Reading   => error
17 };
18
19 let send(ch,v) {
20   case ch.state
21   | Free      => { ch.state := Written(v); wait(ch.writers); }
22   | Written(v) => error
```

```

23   | Reading    => { ch.state := Free; resume(ch.readers,v) }
24 };
25
26 let merge(ch1,ch2) {
27   let merge_fibre(ch1,ch2,c) {
28     while True {
29       case ch1.state
30       | Written(v) => send(c,receive(ch1))
31       | else      => case ch2.state
32         | Written(v) => send(c,receive(ch2))
33         | else      => { ch1.state := Reading;
34                       ch2.state := Reading;
35                       wait_multi([ch1, ch2]); }
36     }
37   };
38   let c = new_channel();
39   spawn(fun () { merge_fibre(ch1,ch2,c) });
40   c
41 };

```

Example As an example of how to use message passing, let us take a look at the well-known producer–consumer problem.

```

1  let produce(ch) {
2    while True {
3      let x = get_next_item();
4      send(ch,x);
5    };
6  };
7
8  let ch = new_channel();
9  spawn(fun () { produce(ch) });
10 spawn(fun () { consume(ch) });
11 start_scheduler()

```

```

let consume(ch) {
  while True {
    let x = receive(ch);
    process_item(x);
  };
};

```

Example Suppose we have a user interface where we want to implement drag-and-drop. The usual GUI frameworks have an event-loop where the program can register call-backs for various events like mouse clicks. When using such a framework, we face the problem of how to remember the program state between user inputs. The typical solution looks as follows.

```

1  type state = | Idle | Dragging(obj);
2
3  let state = Idle;
4

```

8 Concurrency

```
5  let mouse_down(x,y) {
6    case state
7    | Idle          => case object_under_pointer(x,y)
8                      | None          => Nothing
9                      | Some(obj) => state := Dragging(obj)
10   | Dragging(obj) => Nothing
11 };
12 let mouse_up(x,y) {
13   case state
14   | Dragging(obj) => { move_object_to(obj,x,y); state := Idle; }
15   | Idle          => Nothing
16 };
17 let mouse_move(x,y) {
18   case state
19   | Dragging(obj) => move_object_to(obj,x,y);
20   | Idle          => Nothing
21 };
22
23 register_call_back_mouse_down(mouse_down);
24 register_call_back_mouse_up(mouse_up);
25 register_call_back_mouse_move(mouse_move);
```

When having more states than 'idle' and 'dragging', this quickly become tedious. Using fibres we can avoid having to manage the program state explicitly.

```
1  type event = | Start(object) | Move(int,int) | Drop;
2
3  let drag_and_drop(ch) {
4    case receive(ch)
5    | Start obj =>
6      while True {
7        case receive(ch)
8          | Move(x,y) => move_object_to(obj,x,y);
9          | Drop      => break
10       };
11 };
12
13 let mouse_down(ch,x,y) {
14   case object_under_pointer(x,y)
15   | None      => Nothing
16   | Some(obj) => send(ch, Start(obj))
17 };
18 let mouse_up(ch,x,y) {
19   send(ch, Drop(x,y));
20 };
```



```

21  let mouse_move(ch,x,y) {
22    send(ch, Move(x,y));
23  };
24
25  let ch = make_channel();
26  spawn(drag_and_drop(ch));
27  register_callback_mouse_down(mouse_down(ch));
28  register_callback_mouse_up(mouse_up(ch));
29  register_callback_mouse_move(mouse_move(ch));

```

Inversion of control If we do not use fibres, communication between program units is asymmetric. One unit is in control and calls the functions provided by the other unit. In the above example, the event loop was in control and invokes the callbacks provided by the main program. Alternatively, we could have the main program be in control and call a library function to get the next event. But we must choose between these two options. Going from one to the other is called *inversion of control*. The choice between these two versions requires careful consideration, as it has a big influence on the structure of the whole program.

The big advantage of using fibres is that we do not need to choose: both parts can be in control at the same time and communicate via channels. Thus communication is symmetric.

When compared to shared-memory communication, which we will describe next, message passing has some overhead as messages must be constructed and passed to another process. But it scales really well to any number of processes and it works equally well on a single computer or on a distributed system. Furthermore, it is conceptually really simple and easy to use for the programmer. In my opinion it is therefore clearly superior to approaches relying on shared-memory.

Futures Futures are a simple synchronisation mechanism where we can evaluate a given expression in parallel and wait for the result. They work like a channel that can only be used a single time. The implementation is straightforward using the tools we already have. For instance, we can use channels.

```

1  let future(e) {
2    let ch = new_channel();
3    spawn(fun () { send(ch,e) });
4    ch
5  };
6
7  let get(f) { receive(f) };

```

We can also use single-assignment variables (if we use the convention that reading from an uninitialised single-assignment variable blocks until some other fibre assigns a value to it.)

```

1  let future(e) {
2    let x;
3    spawn(fun () { x := e });
4    x

```

```

5   };
6
7   let get(f) { let y = f; y };

```

8.4 Shared-memory

When several fibres or threads run on the same processor they can use the shared memory to communicate. Of course this relies on side-effects and, as the evaluation order matters with side-effects, the non-determinism of this order inherent in concurrency makes such program even harder to understand. Over the years people have invented several mechanisms and constructs to make shared memory communication easier to use and less error prone.

Atomic operations Before presenting the various synchronisation mechanisms let us introduce the primitive operations we need to implement them. These are operations that are *atomic* which means that, when executing such an operation we cannot observe any intermediate state. Either we see the state before the operation or we see the resulting state, but we can never find the operation to be halfway executed.

Usually processors provide a few special atomic instructions to build concurrency mechanisms with. Typical examples are a *test-and-set* and a *compare-and-swap* operation.

```

1   test_and_set    : location -> a -> a;
2   compare_and_swap : location -> a -> a -> bool;

```

The operation `test_and_set(x, a)` stores the value `a` in the variable `x` and returns the old value of `x`. `compare_and_swap(x, a, b)` compares the value in the memory location `x` with the value `a`. If they are the same, it sets the value of `x` to `b`. Otherwise, it leaves `x` unchanged. The return value indicates whether a change occurred.

Locks/Mutexes A *lock* (also called a *mutex*) is a data structure with two states. It can either be *locked* by a certain fibre (we say that the fibre *holds* the lock, or that it has *acquired* the lock) or it is *open*. There are two operations.

```

1   type lock;
2   lock   : lock -> unit;
3   unlock : lock -> unit;

```

When a fibre calls `lock` on an open lock, the lock changes state to *locked* and it is now held by the fibre. When called with a lock that is hold by another fibre, the operation blocks until that fibre unlocks it. What happens when a fibre calls `lock` on a lock that is held by itself depends on the particular implementation. Some implementations just block, which results in a deadlock (as the fibre waits on itself to release the lock). In this case, we call the locks *non-reentrant*. The more sensible solution is to allow a fibre to acquire a lock several times (of course, in this case it has to release the lock the same number of times before the lock is open again). Such locks are called *reentrant*.

For simplicities sake, we will implement non-reentrant locks.

```

1  type lock = [ locked : bool; waiting : condition() ];
2
3  let new_lock() {
4    [ locked = False, waiting = new_condition() ]
5  };
6
7  let lock(l) {
8    while test_and_set(l.locked, True) {
9      wait(l.waiting)
10   }
11 };
12
13 let unlock(l) {
14   l.locked := False;
15   resume(l.waiting);
16 };

```

Using locks manually is quite error prone. If several locks are involved it is important to lock and unlock them in the correct order. Also it is easy to forget some of the unlock calls. Some of these errors can be avoided if the language has built in support for locks, like the following

```

1  lock name                lock(name);
2  ...                      ⇒    ...
3  end                      unlock(name);

```

Condition variables A condition variable is a condition that is associated with a lock. Waiting on the condition also unlocks the lock and when it is woken up again it automatically acquires the lock again.

```

1  type cvar = [ cond : condition, lock : lock ];
2
3  let new_cvar(l) {
4    [ cond = new_condition(), lock = l ];
5  };
6
7  let wait_cvar(c) {
8    unlock(c.lock);
9    wait(c.cond);
10   lock(c.lock);
11 };
12
13 let resume_cvar(c) {
14   resume(c.cond);
15 };

```

Semaphores A semaphore is a generalisation of a lock. It takes the form of an integer counter that cannot go below zero. If a fibre tries to decrease the counter when it already is zero, it blocks instead until another fibre increases the counter again. So, we have two operations: one to increment the counter and one to decrement it.

```

1  type semaphore;
2  increment : semaphore -> unit;
3  decrement : semaphore -> unit;

```

The implementation is as follows.

```

1  type semaphore = [
2    count   : int,
3    lock    : lock,
4    waiting : cvar
5  ];
6
7  let new_semaphore() {
8    let l = new_lock();
9    let cv = new_cvar(1);
10   [ count   = 0,
11     lock    = l,
12     waiting = cv ]
13 };
14
15 let increment(sem) {
16   lock(sem.lock);
17   sem.count := sem.count + 1;
18   resume_cvar(sem.waiting);    // this automatically unlocks the lock
19 };
20
21 let decrement(sem) {
22   lock(sem.lock);
23   while sem.count == 0 {
24     wait_cvar(sem.waiting);
25   }
26   sem.count := sem.count - 1;
27   unlock(sem.lock);
28 };

```

Example The following producer–consumer implementation uses a buffer whose size is bounded by some constant n .

```

1  let lock   = new_semaphore();
2  let full   = new_semaphore();
3  let empty  = new_semaphore();

```

```

4  let n      = 10;
5  let buffer = new_buffer(n);
6
7  for i = 1 .. n {           // the initial value of empty \rmfamily should be n
8    increment(empty)
9  };
10
11 let producer() {
12   for i = 0 .. 10 {
13     let x = ... generate a value ...;
14     decrement(empty);
15     decrement(lock);
16     put(buffer, x);
17     increment(lock);
18     increment(full);
19   }
20 };
21
22 let consumer() {
23   for i = 0 .. 10 {
24     decrement(full);
25     decrement(lock);
26     let x = get(buffer);
27     increment(lock);
28     increment(empty);
29     ... process the value x ...
30   }
31 };

```

Monitors A *monitor* is an abstract data type that is protected by a lock. Each operation on the type first acquires the lock and releases it again upon return. This makes the operations atomic. For instance, a monitor for a queue implementation could look as follows.

```

1  type node(a) = [ ... ];
2
3  type queue(a) = [
4    lock : lock,
5    front : node(a),
6    back  : node(a)
7  ];
8
9  let make() {
10   let node = [ ... ];
11   [ lock = new_lock(), front = node, back = node ]

```

```

12  };
13
14  let pop(q) {
15    lock(q.lock);
16    ... remove the first node from the list ...
17    unlock(q.lock);
18    ... return the data stored in the removed node ...
19  };
20
21  let push(q,x) {
22    lock(q.lock);
23    ... add a new node at the end of the list ...
24    unlock(q.lock);
25  };

```

Transactional memory Transactional memory is a general mechanism to make arbitrary operations atomic. A *transaction* is a piece of code that can either succeed or fail with its task. When it fails it has no effect on the program, it is as if the transaction was never executed. Thus, transactional memory can be seen as a form of backtracking.

We add the following constructs to our language.

$$\langle expr \rangle ::= \dots \mid \mathbf{atomic}\{ \langle expr \rangle \} \mid \mathbf{abort} \mid \mathbf{retry}$$

An expression of the form `atomic { e }` evaluates the expression `e` as if it were atomic. That means that no other fibres can see intermediate states of the execution of `e`. They either see the state before its execution or the one after it.

In addition, the expression `e` can contain `abort` and `retry` statements to indicate, respectively, that the transaction has failed, or that it should be restarted from the beginning.

For a programmers perspective, transactional memory is by far the easiest to use mechanism for shared memory concurrency. It automatically avoids the typical errors associated with this form of concurrency, like deadlocks and race conditions. Furthermore, the resulting code is much more composable and modular than code written with other primitives.

Of course this convenience comes at a significant cost. Transactional memory is very hard to implement and it comes with a considerable overhead. The runtime cost is increased further by the fact that we sometimes need to execute a transaction several times for it to succeed. Finally, since transactional memory relies on backtracking, some operations cannot be used inside a transaction. In particular IO operations are not supported.

Discussion Shared-memory communication is very efficient, but it requires all processors to share memory. This becomes quickly impractical as their number increases. From a programmers perspective the main problem with shared-memory communication is that it is very error prone. The reason is that mechanisms for shared-memory communication require side-effects and we have seen that, when programming with side-effects, the order of execution becomes important. In a concurrent setting, this order is non-deterministic and it is therefore much harder

to reason about. This added complexity makes it almost impossible to write correct code in this setting.