

Why Functional Programming Matters

J. HUGHES

Department of Computing Science, University of Glasgow, Glasgow G12 8QQ

As software becomes more and more complex, it is more and more important to structure it well. Well-structured software is easy to write, easy to debug, and provides a collection of modules that can be re-used to reduce future programming costs. Conventional languages place conceptual limits on the way problems can be modularised. Functional languages push those limits back. In this paper we show that two features of functional languages in particular, higher-order functions and lazy evaluation, can contribute greatly to modularity. As examples, we manipulate lists and trees, program several numerical algorithms, and implement the alpha-beta heuristic (an Artificial Intelligence algorithm used in game-playing programs). Since modularity is the key to successful programming, functional languages are vitally important to the real world.

Received November 1988

1. INTRODUCTION

This paper is an attempt to convince the ‘real world’ that functional programming is vitally important, and also to help functional programmers exploit its advantages to the full by clarifying what those advantages are.

Functional programming is so called because a program consists entirely of functions. The main program itself is written as a function which receives the program’s input as its argument and delivers the program’s output as its result. Typically the main function is defined in terms of other functions, which in turn are defined in terms of still more functions, until at the bottom level the functions are language primitives. These functions are much like ordinary mathematical functions, and in this paper will be defined by ordinary equations. Our notation follows Turner’s language *Miranda*^{TM,1} but should be readable with no prior knowledge of functional languages. (*Miranda* is a trademark of Research Software Ltd.)

The special characteristics and advantages of functional programming are often summed up more or less as follows. Functional programs contain no assignment statements, so variables, once given a value, never change. More generally, functional programs contain no side-effects of any kind. A function call can have no effect other than to compute its result. This eliminates a major source of bugs, and also makes the order of execution irrelevant – since no side-effect can change the value of an expression, it can be evaluated at any time. This relieves the programmer of the burden of prescribing the flow of control. Since expressions can be evaluated at any time, one can freely replace variables by their values and vice versa – that is, programs are ‘referentially transparent’. This freedom helps make functional programs more tractable mathematically than their conventional counterparts.

Such a catalogue of ‘advantages’ is all very well, but one must not be surprised if outsiders don’t take it too seriously. It says a lot about what functional programming is *not* (it has no assignment, no side-effects, no flow of control) but not much about what it is. The functional programmer sounds rather like a medieval monk, denying himself the pleasures of life in the hope

that it will make him virtuous. To those more interested in material gains, these ‘advantages’ are not convincing.

Functional programmers argue that there *are* great material benefits – that a functional programmer is an order of magnitude more productive than his conventional counterpart, because functional programs are an order of magnitude shorter. Yet why should this be? The only faintly plausible reason one can suggest on the basis of these ‘advantages’ is that conventional programs consist of 90% assignment statements, and in functional programs these can be omitted! This is plainly ridiculous. If omitting assignment statements brought such enormous benefits then FORTRAN programmers would have been doing it for twenty years. It is a logical impossibility to make a language more powerful by omitting features, no matter how bad they may be.

Even a functional programmer should be dissatisfied with these so-called advantages, because they give him no help in exploiting the power of functional languages. One cannot write a program which is particularly lacking in assignment statements, or especially referentially transparent. There is no yardstick of program quality here, and therefore no ideal to aim at.

Clearly this characterisation of functional programming is inadequate. We must find something to put in its place – something which not only explains the power of functional programming, but also gives a clear indication of what the functional programmer should strive towards.

2. AN ANALOGY WITH STRUCTURED PROGRAMMING

It is helpful to draw an analogy between functional and structured programming. In the past, the characteristics and advantages of structured programming have been summed up more or less as follows. Structured programs contain no **goto** statements. Blocks in a structured program do not have multiple entries or exits. Structured programs are more tractable mathematically than their unstructured counterparts. These ‘advantages’ of structured programming are very similar in spirit to the ‘advantages’ of functional programming we discussed earlier. They are essentially negative statements, and

have led to much fruitless argument about ‘essential gotos’ and so on.

With the benefit of hindsight, it is clear that these properties of structured programs, although helpful, do not go to the heart of the matter. The most important difference between structured and unstructured programs is that the former are designed in a modular way. Modular design brings great productivity improvements. First of all, small modules can be coded quickly and easily. Secondly, general purpose modules can be reused, leading to faster development of subsequent programs. Thirdly, the modules of a program can be tested independently, helping to reduce the time spent debugging.

The absence of **gotos** and so on, has very little to do with this. It helps with ‘programming in the small’, whereas modular design helps with ‘programming in the large’. Thus one can enjoy the benefits of structured programming in FORTRAN or assembly language, even if it is a little more work.

It is now generally accepted that modular design is the key to successful programming, and recent languages such as Modula-II,² Ada,³ and Standard ML⁴ include features specifically designed to help improve modularity. However, there is a very important point that is often missed. When writing a modular program to solve a problem, one first divides the problem into sub-problems, then solves the sub-problems and combines the solutions. The ways in which one can divide up the original problem depend directly on the ways in which one can glue solutions together. Therefore, to increase one’s ability to modularise a problem conceptually, one must provide new kinds of glue in the programming language. Complicated scope rules and provision of separate compilation only help with clerical details; they offer no new conceptual tools for decomposing problems.

One can appreciate the importance of glue by an analogy with carpentry. A chair can be made quite easily by making the parts – seat, legs, back and so on – and sticking them together in the right way. But this depends on the ability to make joints and wood-glue. Lacking that ability, the only way to make a chair is to carve it in one piece out of a solid block of wood, a much harder task. This example demonstrates both the enormous power of modularisation and the importance of having the right glue.

Now let us return to functional programming. We shall argue in the remainder of this paper that functional languages provide two new, very important kinds of glue. We shall give many examples of programs that can be modularised in new ways, and thereby greatly simplified. This is the key to functional programming’s power – it allows greatly improved modularisation. It is also the goal for which functional programmers must strive – smaller and simpler and more general modules, glued together with the new glues we shall describe.

3. GLUEING FUNCTIONS TOGETHER

The first of the two new kinds of glue enables simple functions to be glued together to make more complex ones. It can be illustrated with a simple list-processing problem – adding up the elements of a list. We define lists by

$\text{listof } \alpha ::= [] | \alpha : (\text{listof } \alpha)$

which means that a list of α s is either $[]$, representing a list with no elements, or it is a pair of an α and another list of α s, built with the constructor ‘ $:$ ’ (pronounced ‘cons’). α here may stand for any type – for example, if α is integer then the definition says that a list of integers is either empty, or a cons of an integer and another list of integers. Following normal practice, we will write down lists simply by enclosing their elements in square brackets, rather than by writing ‘ $:$ ’ and $[]$ explicitly. This is simply a shorthand for notational convenience. For example,

$[1] \quad \text{means } 1 : []$
 $[1, 2, 3] \quad \text{means } 1 : (2 : (3 : []))$

The elements of a list can be added up by a recursive function **sum**. The function **sum** must be defined for two kinds of argument: an empty list ($[]$), and a cons. Since the sum of no numbers is zero, we define

$\text{sum } [] = 0$

and since the sum of a cons can be calculated by adding the first element of the list to the sum of the others, we can define

$\text{sum } (\text{num} : \text{list}) = \text{num} + \text{sum list}$

Examining this definition, we see that only the boxed parts below are specific to computing a sum.

$\text{sum } [] = 0$
 $\text{sum } (\text{num} : \text{list}) = \text{num} + \text{sum list}$

This means that the computation of a sum can be modularised by glueing together a general recursive pattern and the boxed parts. This recursive pattern is conventionally called **reduce** and so **sum** can be expressed as

$\text{sum} = \text{reduce } (+) \ 0$

The operator $+$ is enclosed in parentheses to make it clear that it is a parameter to **reduce** – otherwise this expression would appear to add 0 to a function!

The definition of **reduce** can be derived just by parameterising the definition of **sum**, giving

$(\text{reduce } (+) x) [] = x$
 $(\text{reduce } (+) x) (el : \text{list}) = el + ((\text{reduce } (+) x) \text{ list})$

where $+$ is a parameter which is a binary operator. Here we have written parentheses around $(\text{reduce } (+) x)$ to make it clear that it replaces **sum**. Conventionally the parentheses are omitted, and so $((\text{reduce } (+) x) \text{ list})$ is written $(\text{reduce } (+) x \text{ list})$. A function of three arguments such as **reduce**, applied to only two is taken to be a function of the one remaining argument, and in general, a function of n arguments applied to only $m (< n)$ is taken to be a function of the $n - m$ remaining ones. (This is called *partial application* or *currying*.) We will follow this convention in future.

Having modularised **sum** in this way, we can reap benefits by re-using the parts. The most interesting part is **reduce**, which can be used to write down a function for multiplying together the elements of a list with no further programming:

$\text{product} = \text{reduce } (*) \ 1$

It can also be used to test whether any of a list of booleans is true

`anytrue = reduce (∨) False`

or whether they are all true

`alltrue = reduce (∧) True`

One way to understand `(reduce (⊕) a)` is as a function that replaces all occurrences of `(:)` in a list by `(⊕)`, and all occurrences of `[]` by `a`. Taking the list `[1, 2, 3]` as an example, since this means

`1:(2:(3:[]))`

then `(reduce (+) 0)` converts it into

$$\begin{aligned} 1 + (2 + (3 + 0)) \\ = 6 \end{aligned}$$

and `(reduce (*) 1)` converts it into

$$\begin{aligned} 1 * (2 * (3 * 1)) \\ = 6 \end{aligned}$$

(An operator in parentheses stands for a two-argument function wherever it appears.)

Now it is obvious that `(reduce (:)[])` just copies a list. Since one list can be appended to another by consing its elements onto the front, we find

`a++b = reduce (: b a)`

As an example,

$$\begin{aligned} [1, 2] ++ [3, 4] \\ = \text{reduce} (: [3, 4] [1, 2]) \\ = (\text{reduce} (: [3, 4])) (1:(2:[])) \\ = 1:(2:[3, 4]) \\ = [1, 2, 3, 4] \end{aligned}$$

We can count the number of elements in a list using the function `length`, defined by

`length = reduce (⊕) 0`
where `el ⊕ len = 1 + len`

because `⊕` increments 0 as many times as there are conses.

Now suppose we want a function which doubles every element of a list. Having defined

`double x = 2 * x`

we want `doubleall` to satisfy

`doubleall (a:b:...) = (double a : double b :...)`

But this can also be programmed using `reduce!` We just replace `(:)` by `(⊕)`, where

`a ⊕ 1 = double a : 1`

So

`doubleall = reduce (⊕) []`

Now,

$$\begin{aligned} \text{doubleall} (1:(2:[])) \\ = 1 ⊕ (2 ⊕ []) \\ = \text{double } 1 : (\text{double } 2 : []) \\ = 2:(4:[]) \end{aligned}$$

as expected.

With one further modularisation we arrive at

`doubleall = map double`
`map f = reduce (⊕) Nil`
`where a ⊕ l = f a : l`

where `map` applies any function `f` to all the elements of a list. This is another generally useful function.

We can even write down a function to add up all the elements of a matrix, represented as a list of lists. It is

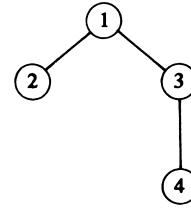
`summatrix = sum o map sum`

The `map sum` uses `sum` to add up all the rows, and then the left-most `sum` adds up the row totals to get the sum of the whole matrix.

These examples should be enough to convince the reader that a little modularisation can go a long way. By modularising a simple function (`sum`) as a combination of a ‘higher-order function’ and some simple arguments, we have arrived at a part (`reduce`) that can be used to write down many other functions on lists with no more programming effort. We do not need to stop with functions on lists. As another example, consider the datatype of ordered labelled trees, defined by

`treesof α ::= α @ (listof (treeof α))`

This definition says that a tree of `α`s is a pair built with the constructor ‘`@`’, with a label which is a `α`, and a list of subtrees which are also trees of `α`s. We refer to pairs built with ‘`@`’ as *nodes*. For example, the tree



would be represented by

`1 @ [2 @ [], 3 @ [4 @ []]]`

Instead of considering an example and abstracting a higher-order function from it, we will go straight to a function `redtree` analogous to `reduce`. Recall that `reduce` took two arguments: something to replace `(:)` with, and something to replace `[]` with. Since trees are built using `(@)`, `(:)` and `[]`, `redtree` must take three arguments – something to replace each of these with. Therefore we define

$$\begin{aligned} \text{redtree} (⊕) (⊗) a (\text{label} @ \text{subtrees}) \\ = \text{label} ⊕ \text{redtree}' (⊕) (⊗) a \text{subtrees} \\ \text{redtree}' (⊕) (⊗) a (\text{subtree} : \text{rest}) \\ = \text{redtree} (⊕) (⊗) a \text{subtree} \\ \quad ⊗ \\ \quad \text{redtree}' (⊕) (⊗) a \text{rest} \\ \text{redtree}' (⊕) (⊗) a [] \\ = a \end{aligned}$$

(Since trees contain values of two different types – nodes and lists – we have to implement `redtree` by two functions. In the definition above, `redtree` is applied to nodes while `redtree'` is applied to lists.)

Many interesting functions can be defined by glueing `redtree` and other functions together. For example, all

the labels in a tree of numbers can be added together using

```
sumtree = redtree (+) (+) 0
```

Taking the tree we wrote down earlier as an example, sumtree gives

```
sumtree (1 @ [2 @ [ ], 3 @ [4 @ [ ]]])  
= sumtree (1 @ ((2 @ [ ]):(3 @ ((4 @ [ ]):[ ])):[ ]))  
= 1 + ((2+0)+(3+((4+0)+0))+0)  
= 10
```

A list of all the labels in a tree can be computed using

```
labels = redtree (:) (+) []
```

The same example gives

```
labels (1 @ [2 @ [ ], 3 @ [4 @ [ ]]])  
= labels (1 @ ((2 @ [ ]):(3 @ ((4 @ [ ]):[ ])):[ ]))  
= 1:((2:[ ])+(3:((4:[ ])+[ ]))+[ ]))  
= [1, 2, 3, 4]
```

Finally, one can define a function analogous to map which applies a function f to all the labels in a tree:

```
maptree f = redtree (⊕) (:) []  
where label ⊕ subtrees = f label @ subtrees
```

All this can be achieved because functional languages allow functions which are indivisible in conventional programming languages to be expressed as a combination of parts – a general higher-order function and some particular specialising functions. Once defined, such higher-order functions allow many operations to be programmed very easily. Whenever a new datatype is defined higher-order functions should be written for processing it. This makes manipulating the datatype easy, and also localises knowledge about the details of its representation. The best analogy with conventional programming is with extensible languages – it is as though the programming language can be extended with new control structures whenever desired.

4. GLUING PROGRAMS TOGETHER

The other new kind of glue that functional languages provide enables whole programs to be glued together. Recall that a complete functional program is just a function from its input to its output. If f and g are such programs, then $(g \circ f)$, the composition of g and f, is a program which, when applied to its input, computes

$g(f(\text{input}))$

The program f computes its output which is used as input to program g. This might be implemented conventionally by storing the output from f in a temporary file. The problem with this is that the temporary file might occupy so much memory that it is impractical to glue the programs together in this way. Functional languages provide a solution to this problem. The two programs f and g are run together in strict synchronisation. The function f is only started once g tries to read some input, and only runs long enough to deliver the output g is trying to read. Then f is suspended and g is run until it tries to read another input. As an added bonus, if g terminates without reading all of f's output then f is aborted. The function f can even be a non-terminating program, producing an infinite amount

of output, since it will be terminated forcibly as soon as g is finished. This allows termination conditions to be separated from loop bodies – a powerful modularisation.

Since this method of evaluation runs f as little as possible, it is called ‘lazy evaluation’. It makes it practical to modularise a program as a generator which constructs a large number of possible answers, and a selector which chooses the appropriate one. While some other systems allow programs to be run together in this manner, only functional languages use lazy evaluation uniformly for every function call, allowing any part of a program to be modularised in this way. Lazy evaluation is perhaps the most powerful tool for modularisation in the functional programmer’s repertoire.

4.1. Newton–Raphson square roots

We will illustrate the power of lazy evaluation by programming some numerical algorithms. First of all, consider the Newton–Raphson algorithm for finding square roots. This algorithm computes the square root of a number N by starting from an initial approximation a_0 and computing better and better ones using the rule

$$a_{n+1} = (a_n + N/a_n)/2,$$

If the approximations converge to a limit a, then

$$a = (a + N/a)/2,$$

so

$$a = \sqrt{N}.$$

In fact the approximations converge rapidly to a limit. Square root programs take a tolerance ϵ and stop when two successive approximations differ by less than ϵ .

The algorithm is usually programmed more or less as follows:

```
C  N IS CALLED ZN HERE SO THAT IT HAS THE  
C  RIGHT TYPE  
X = A0  
Y = A0 + 2.*EPS  
C  THE VALUE OF Y DOESN'T MATTER SO  
C  LONG AS ABS(X-Y).GT.EPS  
100 IF (ABS(X-Y).LE.EPS) GOTO 200  
Y = X  
X = (X+ZN/X)/2.  
GOTO 100  
200 CONTINUE  
C  THE SQUARE ROOT OF ZN IS NOW IN X
```

This program is indivisible in conventional languages. We will express it in a more modular form using lazy evaluation, and then show some other uses to which the parts may be put.

Since the Newton–Raphson algorithm computes a sequence of approximations it is natural to represent this explicitly in the program by a list of approximations. Each approximation is derived from the previous one by the function

$$\text{next } n \ x = (x + n/x)/2$$

so $(\text{next } n)$ is the function mapping one approximation onto the next. Calling this function f, the sequence of approximations is

$$[a_0, f(a_0), f(f(a_0)), f(f(f(a_0))), \dots]$$

We can define a function to compute this:

```
repeat f a = a:(repeat f (f a))
```

so that the list of approximations can be computed by

```
repeat (next n) a0
```

`repeat` is an example of a function with an ‘infinite’ output – but it does not matter, because no more approximations will actually be computed than the rest of the program requires. The infinity is only potential: all it means is that any number of approximations can be computed if required, `repeat` itself places no limit.

The remainder of a square root finder is a function `within`, that takes a tolerance and a list of approximations and looks down the list for two successive approximations that differ by no more than a given tolerance. It can be defined by

```
within eps (a:b:rest)
```

```
= b,           if abs (a - b) ≤ eps  
= within eps (b:rest), otherwise
```

Putting the parts together,

```
sqrt a0 eps n = within eps (repeat (next n) a0)
```

Now that we have the parts of a square root finder, we can try combining them in different ways. One modification we might wish to make is to wait for the ratio between two successive approximations to approach one, rather than for the difference to approach zero. This is more appropriate for very small numbers (when the difference between successive approximations is small to start with) and for very large ones (when rounding error could be much larger than the tolerance). It is only necessary to define a replacement for `within`:

```
relative eps (a:b:rest)
```

```
= b,           if abs (a/b - 1) ≤ eps  
= relative eps (b:rest), otherwise
```

(The two functions `within` and `relative` could have been defined as instances of a still more general limit-finding function, parameterised by a test on two successive approximations. We chose not to do this, because it’s hard to see that the more general function would be useful for anything other than defining these two specific ones.)

Now a new version of `sqrt` can be defined by

```
relativesqrt a0 eps n = relative eps (repeat (next n) a0)
```

It is not necessary to rewrite the part that generates approximations.

4.2 Numerical differentiation

We have re-used the sequence of approximations to a square root. Of course, it is also possible to re-use `within` and `relative` with any numerical algorithm that generates a sequence of approximations. We will do so in a numerical differentiation algorithm.

The result of differentiating a function at a point is the slope of the function’s graph at that point. It can be estimated quite easily by evaluating the function at the given point and at another point nearby and computing the slope of a straight line between the two points. This assumes that, if the two points are close enough together

then the graph of the function will not curve much in between. This gives the definition

```
easydiff f x h = (f (x + h) - f x)/h
```

In order to get a good approximation the value of `h` should be very small. Unfortunately, if `h` is too small then the values `f(x+h)` and `f x` are very close together, and so the rounding error in the subtraction may swamp the result. How can the right value of `h` be chosen? One solution to this dilemma is to compute a sequence of approximations with smaller and smaller values of `h`, starting with a reasonably large one. Such a sequence should converge to the value of the derivative, but will become hopelessly inaccurate eventually due to rounding error. If (`within eps`) is used to select the first approximation that is accurate enough then the risk of rounding error affecting the result can be much reduced. We need a function to compute the sequence:

```
differentiate h0 f x
```

```
= map (easydiff f x) (repeat halve h0)
```

```
halve x = x/2
```

Here `h0` is the initial value of `h`, and successive values are obtained by repeated halving. Given this function, the derivative at any point can be computed by

```
within eps (differentiate h0 f x)
```

Even this solution is not very satisfactory because the sequence of approximations converges fairly slowly. A little simple mathematics can help here. The elements of the sequence can be expressed as

the right answer + an error term involving `h`

and it can be shown theoretically that the error term is roughly proportional to a power of `h`, so that it gets smaller as `h` gets smaller. Let the right answer be `A`, and the error term be `B * h^n`: Since each approximation is computed using a value of `h` twice that used for the next one, any two successive approximations can be expressed as

$$a_n = A + B * 2^n * h^n$$

$$a_{n+1} = A + B * h^n$$

Now the error term can be eliminated. We conclude

$$A = \frac{a_{n+1} * 2^n - a_n}{2^n - 1}$$

Of course, since the error term is only roughly a power of `h` this conclusion is also approximate, but it is a much better approximation. This improvement can be applied to all successive pairs of approximations using the function

```
elimerror n (a:b:rest)  
= ((b * 2^n - a)/(2^n - 1)):  
  (elimerror n (b:rest))
```

Eliminating error terms from a sequence of approximations yields another sequence which converges much more rapidly.

Now the derivative of a function `f` can be computed more efficiently as follows:

```
within eps (elimerror 1 (differentiate h0 f x))
```

`elimerror` only works on sequences of approximations which are computed using a parameter `h`, which is halved

between each approximation. However, if it is applied to such a sequence its result is also such a sequence! This means that a sequence of approximations can be improved more than once. A different error term is eliminated each time, and the resulting sequences converge faster and faster. So, one could compute a derivative very efficiently using

```
within eps (elimerror 3
  (elimerror 2
    (elimerror 1 (differentiate h0 fx)))
```

In numerical analysts terms, this is a fourth order method, and gives an accurate result very quickly. One could even define

```
elimall n s = s:(elimall (n+1) (elimerror n s))
```

which eliminates more and more error terms to give a sequence of sequences of approximations, each one converging faster than the last. Now

```
super s = map hd (elimall 1 s)
where hd (a:b) = a
```

selects the first approximation from each of these sequences, to form a very rapidly converging sequence indeed. This is really a very sophisticated algorithm – it uses a better and better numerical method as more and more approximations are computed. One could compute derivatives very efficiently indeed with the program

```
within eps (super (differentiate h0 fx))
```

This is probably a case of using a sledge-hammer to crack a nut, but the point is that even an algorithm as sophisticated as *super* is easily expressed when modularised using lazy evaluation.

4.3 Numerical integration

The last example we will discuss in this section is numerical integration. The problem may be stated very simply: given a real valued function *f* of one real argument, and two end-points *a* and *b*, estimate the area under the curve *f* describes between the end-points. The easiest way to estimate the area is to assume that *f* is nearly a straight line, in which case the area would be

```
easyintegrate f a b = (fa + fb) * (b - a)/2
```

Unfortunately this estimate is likely to be very inaccurate unless *a* and *b* are very close together. A better estimate can be made by dividing the interval from *a* to *b* into two, estimating the area on each half, and adding the results together. We can define a sequence of better and better approximations to the value of the integral by using the formula above for the first approximation, and then adding together better and better approximations to the integrals on each half to calculate the others. This sequence is computed by the function

```
integrate f a b
= (easyintegrate f a b):
  (map addpair (zip (integrate f a mid)
    (integrate f mid b)))
where mid = (a+b)/2
```

The function *zip* is another standard list-processing function. It takes two lists and returns a list of pairs, each

pair consisting of corresponding elements of the two lists. Thus the first pair consists of the first element of the first list and the first element of the second, and so on. The function *zip* can be defined by

```
zip (a:s) (b:t) = (a,b):(zip st)
```

In *integrate*, *zip* computes a list of pairs of corresponding approximations to the integrals on the two sub-intervals, and *map addpair* adds the elements of the pairs together to give a list of approximations to the original integral. The function *addpair* can be defined by

```
addpair (a,b) = a + b
```

Actually, this version of *integrate* is rather inefficient because it continually recomputes values of *f*. As it is written, *easyintegrate* evaluates *f* at *a* and at *b*, and then the recursive calls of *integrate* re-evaluate each of these. Also, (*f mid*) is evaluated in each recursive call. It is therefore preferable to use the following version which never recomputes a value of *f*.

```
integrate f a b = integ f a b (fa) (fb)
integ f a b fa fb
= ((fa + fb) * (b - a)/2):
  (map addpair (zip (integ f a m fa fm)
    (integ f m b fm fb)))
where m = (a+b)/2
      fm = fm
```

integrate computes an infinite list of better and better approximations to the integral, just as *differentiate* did in the section above. We can therefore just write down integration routines that *integrate* to any desired accuracy, as in

```
within eps (integrate f a b)
relative eps (integrate f a b)
```

This integration algorithm suffers from the same disadvantage as the first differentiation algorithm in the preceding sub-section – it converges rather slowly. Once again, it can be improved. The first approximation in the sequence is computed (by *easyintegrate*) using only two points, with a separation of *b-a*. The second approximation also uses the mid-point, so that the separation between neighbouring points is only $(b-a)/2$. The third approximation uses this method on each half-interval, so the separation between neighbouring points is only $(b-a)/4$. Clearly the separation between neighbouring points is halved between each approximation and the next. Taking this separation as ‘*h*’, the sequence is a candidate for improvement using the *elimerror* function defined in the preceding section. Therefore we can now write down quickly converging sequences of approximations to integrals, for example

```
super (integrate sin 0 4)
```

or

```
elimerror 4 (integrate f 0 1) where f x = 1/(1+x2)
```

(This latter sequence is an eighth-order method for computing $\pi/4$. The second approximation, which requires only five evaluations of *f* to compute, is correct to five decimal places.)

In this section we have taken a number of numerical algorithms and programmed them functionally, using lazy evaluation as glue to stick their parts together.

Thanks to this, we were able to modularise them in new ways, into generally useful functions such as `within`, `relative`, and `elimerror`. By combining these parts in various ways we programmed some quite good numerical algorithms very simply and easily.

5. AN EXAMPLE FROM ARTIFICIAL INTELLIGENCE

We have argued that functional languages are powerful primarily because they provide two new kinds of glue: higher-order functions and lazy evaluation. In this section we take a larger example from Artificial Intelligence and show how it can be programmed quite simply using these two kinds of glue.

The example we choose is the alpha-beta ‘heuristic’, an algorithm for estimating how good a position a game-player is in. The algorithm works by looking ahead to see how the game might develop, but avoids pursuing unprofitable lines.

Let game-positions be represented by objects of the type `position`. This type will vary from game to game, and we assume nothing about it. There must be some way of knowing what moves can be made from a position: assume that there is a function

`moves: position → listof position`

that takes a game-position as its argument and returns the list of all positions that can be reached from it in one move. Taking noughts and crosses (tic-tac-toe) as an example,

$$\begin{aligned} \text{moves } \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \end{array} &= \left[\begin{array}{|c|c|c|}\hline & & \\ \hline & X & \\ \hline & & \end{array}, \begin{array}{|c|c|c|}\hline & & \\ \hline & & X \\ \hline & & \end{array}, \begin{array}{|c|c|c|}\hline & & X \\ \hline & & \end{array} \right] \\ \text{moves } \begin{array}{|c|c|c|}\hline & & \\ \hline & X & \\ \hline & & \end{array} &= \left[\begin{array}{|c|c|c|}\hline & O & \\ \hline & & X \\ \hline & & \end{array}, \begin{array}{|c|c|c|}\hline & O & \\ \hline & & X \\ \hline & & \end{array} \right] \end{aligned}$$

This assumes that it is always possible to tell which player’s turn it is from a position. In noughts and crosses this can be done by counting the noughts and crosses. In a game like chess one would have to include the information explicitly in the type `position`.

Given the function `moves`, the first step is to build a gametree. This is a tree in which the nodes are labelled by positions, such that the children of a node are labelled with the positions that can be reached in one move from that node. That is, if a node is labelled with position `p`, then its children are labelled with the positions in `(moves p)`. A gametree may very well be infinite, if it is possible for a game to go on for ever with neither side winning. Gametrees are exactly like the trees we discussed in Section 2 – each node has a label (the position it represents) and a list of subnodes. We can therefore use the same datatype to represent them.

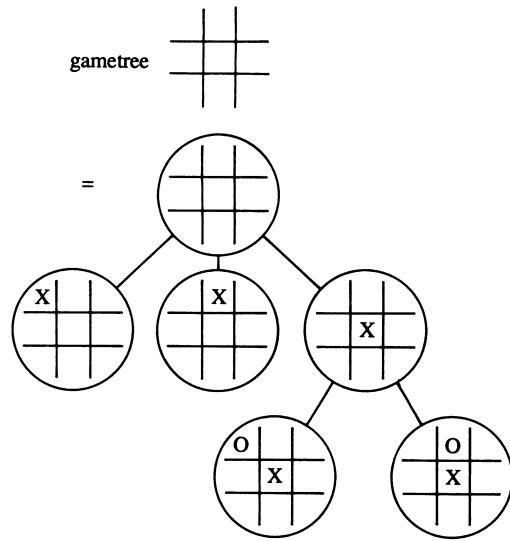
A gametree is built by repeated applications of `moves`. Starting from the root position, `moves` is used to generate the labels for the subtrees of the root. `moves` is then used again to generate the subtrees of the subtrees and so on. This pattern of recursion can be expressed as a higher-order function,

`reptree f a = a @ map (reptree f) (f a)`

Using this function another can be defined which constructs a gametree from a particular position

`gametree p = reptree moves p`

As an example,



(This diagram shows only a part of the gametree, of course.) The higher-order function used here (`reptree`) is analogous to the function `repeat` used to construct infinite lists in the preceding section.

The alpha-beta algorithm looks ahead from a given position to see whether the game will develop favourably or unfavourably, but in order to do so it must be able to make a rough estimate of the value of a position without looking ahead. This ‘static evaluation’ must be used at the limit of the look-ahead, and may be used to guide the algorithm earlier. The result of the static evaluation is a measure of the promise of a position from the computer’s point of view (assuming that the computer is playing the game against a human opponent). The larger the result, the better the position for the computer. The smaller the result, the worse the position. The simplest such function would return (say) +1 for positions where the computer has already won, -1 for positions where the computer has already lost, and 0 otherwise. In reality, the static evaluation function measures various things that make a position ‘look good’, for example material advantage and control of the centre in chess. Assume that we have such a function,

`static: position → number`

Since a gametree is a (treeof `position`), it can be converted into a (treeof `number`) by the function (`maptree static`), which statically evaluates all the positions in the tree (which may be infinitely many). This uses the function `maptree` defined in Section 2.

Given such a tree of static evaluations, what is the true value of the positions in it? In particular, what value should be ascribed to the root position? Not its static value, since this is only a rough guess. The value ascribed to a node must be determined from the true values of its subnodes. This can be done by assuming that each player makes the best moves he can. Remembering that a high value means a good position for the computer, it is clear that when it is the computer’s move from any position, it will choose the move leading to the sub-node with the

maximum true value. Similarly the opponent will choose the move leading to the sub-node with the minimum true value. Assuming that the computer and its opponent alternate turns, the true value of a node is computed by the function `maximise` if it is the computer's turn, and `minimise` if it is not:

```
maximise (n @ sub) = max (map minimise sub)
minimise (n @ sub) = min (map maximise sub)
```

Here `max` and `min` are functions on lists of numbers that return the maximum and minimum of the list respectively. These definitions are not complete because they recurse for ever – there is no base case. We must define the value of a node with no successors, and we take it to be the static evaluation of the node (its label). Therefore the static evaluation is used when either player has already won, or at the limit of look-ahead. The complete definitions of `maximise` and `minimise` are

```
maximise (n @ [ ]) = n
maximise (n @ sub)
  = max (map minimise sub), if sub ≠ [ ]
```

```
minimise (n @ [ ]) = n
minimise (n @ sub)
  = min (map maximise sub), if sub ≠ [ ]
```

One could almost write down a function at this stage that would take a position and return its true value. This would be:

```
evaluate = maximise o maptree static o gametree
```

There are two problems with this definition. First of all, it doesn't work for infinite trees. `Maximise` keeps on recursing until it finds a node with no subtrees – an end to the tree. If there is no end then `maximise` will return no result. The second problem is related – even finite gametrees (like the one for noughts and crosses) can be very large indeed. It is unrealistic to try to evaluate the whole of the gametree – the search must be limited to the next few moves. This can be done by pruning the tree to a fixed depth,

```
prune 0 (a @ x) = a @ [ ]
prune n (a @ x)
  = a @ map (prune (n - 1)) x, if n > 0
```

(`prune n`) takes a tree and 'cuts off' all nodes further than `n` from the root. If a gametree is pruned it forces `maximise` to use the static evaluation for nodes at depth `n`, instead of recursing further. The function `evaluate` can therefore be defined by, say,

```
evaluate = maximise o maptree static o
  prune 5 o gametree
```

which looks 5 moves ahead.

Already in this development we have used higher-order functions and lazy evaluation. Higher-order functions `reptree` and `maptree` allow us to construct and manipulate gametrees with ease. More importantly, lazy evaluation permits us to modularise `evaluate` in this way. Since `gametree` has a potentially infinite result, this program would never terminate without lazy evaluation. Instead of writing

```
prune 5 o gametree
```

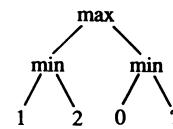
we would have to fold these two functions together into one which only constructed the first five levels of the tree.

Worse, even the first five levels may be too large to be held in memory at one time. In the program as we have written it, the function

```
maptree static o prune 5 o gametree
```

only constructs parts of the tree as `maximise` requires them. Since each part can be thrown away (reclaimed by the garbage collector) as soon as `maximise` has finished with it, the whole tree is never resident in memory. Only a small part of the tree is stored at a time. The lazy program is therefore efficient. Since this efficiency depends on an interaction between `maximise` (the last function in the chain of compositions) and `gametree` (the first), it could only be achieved without lazy evaluation by folding all the functions in the chain together into one big one. This is a drastic reduction in modularity, but is what is usually done. We can make improvements to this evaluation algorithm by tinkering with each part: this is relatively easy. A conventional programmer must modify the entire program as a unit, which is much harder.

So far we have only described simple minimaxing. The heart of the alpha-beta algorithm is the observation that one can often compute the value of `maximise` or `minimise` without looking at the whole tree. Consider the tree:



Strangely enough, it is unnecessary to know the value of the question mark in order to evaluate the tree. The left minimum evaluates to 1, but the right minimum clearly evaluates to something less than or equal to 0. Therefore the maximum of the two minima must be 1. This observation can be generalised and built into `maximise` and `minimise`.

We will do so by changing the way we represent the values of positions. Instead of a single number, we will use *ordered lists* of numbers, representing closer and closer approximations to the actual numeric result. The function `maximise'` will construct an increasing list of numbers, and `minimise'` will construct a decreasing list. So if

```
maximise' t = [1, 3, 4]
```

we can interpret successive elements of the result as giving us more and more information about the value of the position; namely 'I know it's at least one', 'I know it's at least three', 'I know it's at least four', and (at the end of the list) 'I know it's exactly four'. We hope that earlier elements of the list will be cheaper to compute than later ones, and that it will sometimes be unnecessary to look all the way down such a list to the exact value at the end.

We can recover a single numeric value from such a list by computing

```
max (maximise' t)
```

or

```
min (minimise' t)
```

To define `maximise'` and `minimise'` we just reprogram `maximise` and `minimise`, replacing numbers by ordered lists of approximations throughout. So instead of

```
maximise (n @ [ ]) = n
maximise (n @ sub)
= max (map minimise sub), if sub ≠ [ ]
we define
```

```
maximise' (n @ [ ]) = [n]
maximise' (n @ sub)
= max' (map minimise' sub), if sub ≠ [ ]
```

and `minimise'` is similar. Note that, given a leaf node we construct a list containing one approximation (the right answer), and that we need a function `max'` which corresponds to `max` for lists of approximations.

It is into this function `max'` that we can build the optimisation we discussed above. We expect

```
max' [[2,1], [0,?]]
```

to be [1] whatever the value of the question mark – so `max'` will not in general need to read all of its input to produce its result.

How can we define `max'`? It should produce an *increasing* list of approximations to the maximum of the minima of the lists in its argument. A good first approximation is the minimum of the first list in its argument – the maximum of the minima is certain to be at least that. So we will define

```
max' (first:rest) = min first:...
```

Now, if the second list in `max'`'s argument has a *larger* minimum than the first, then we can produce a better approximation to the maximum of the minima. If the second list's minimum is smaller than the first then it gives us no more information and we can go on to the third. So we will define a function `betterthan` so that

`betterthan` approx lists

constructs a list of approximations to the maximum of the lists' minima *better than the given one*. We can then define

```
max' (first:rest)
= min first:betterthan (min first) rest
```

The function `betterthan` could simply be defined by

```
betterthan a (next:rest)
= min next:betterthan (min next) rest,
  if min next > a
= betterthan a rest,
  if min next ≤ a
betterthan a [ ] = [ ]
```

But notice that we can test whether the minimum of a list is less than or equal to a given value without necessarily looking at all of the list. In our original example, we can see that

```
min [0,?] ≤ 1
```

without knowing the value of the question mark. So we can define a function `minleq` which performs this test, but sometimes doesn't look at all of its input. It is defined as follows:

```
minleq [ ] pot = False
minleq (num:rest) pot
```

```
= True, if num ≤ pot
= minleq rest pot, otherwise
```

Now we can use `minleq` in `betterthan`, to complete our optimised minimaxer. It is simple to write a new evaluator:

```
evaluate
= max o maximise' o maptree static o prune 8 o
  gametree
```

Thanks to lazy evaluation, the fact that `maximise'` looks at less of the tree means that the whole program runs more efficiently, just as the fact that `prune` looks at only part of an infinite tree enables the program to terminate. The optimisations in `maximise'`, although fairly simple, can have a dramatic effect on the speed of evaluation, and so can allow the evaluator to look further ahead.

Other optimisations can be made to the evaluator. For example, the alpha-beta algorithm just described works best if the best moves are considered first, since if one has found a very good move then there is no need to consider worse moves, other than to demonstrate that the opponent has at least one good reply to them. One might therefore wish to sort the subtrees at each node, putting those with the highest values first when it is the computer's move, and those with the lowest values first when it is not. This can be done with the function

```
highfirst (n @ sub)
= n @ (sort higher (map lowfirst sub))
lowfirst (n @ sub)
= n @ (sort (not o higher) (map highfirst sub))
higher (Node n1 sub1) (Node n2 sub 2)
= n1 > n2
```

where `sort` is a general purpose sorting function. The evaluator would now be defined by

```
evaluate
= max o maximise' o highfirst
  o maptree static o prune 8 o gametree
```

One might regard it as sufficient to consider the three best moves for the computer or the opponent, in order to restrict the search. To program this, it is only necessary to replace `highfirst` by `(taketree 3 o highfirst)`, where

```
taketree n = redtree (nodett n) (:)
nodett n label sub = label @ (take n sub)
```

`taketree` replaces all the nodes in a tree with nodes with at most `n` subnodes, using the function `(take n)` which returns the first `n` elements of a list (or fewer if the list is shorter than `n`).

Another improvement is to refine the pruning. The program above looks ahead a fixed depth even if the position is very dynamic – it may decide to look no further than a position in which the queen is threatened in chess, for example. It is usual to define certain 'dynamic' positions and not to allow look-ahead to stop in one of these. Assuming a function `dynamic` that recognises such positions, we need only add one equation to `prune` to do this:

```
prune 0 (pos @ sub)
= pos @ (map (prune 0) sub), if dynamic pos
```

Making such changes is easy in a program as modular as this one. As we remarked above, since the program

depends crucially for its efficiency on an interaction between maximise, the last function in the chain, and gametree, the first, it can only be written as a monolithic program without lazy evaluation. Such a program is hard to write, hard to modify, and very hard to understand.

6. CONCLUSION

In this paper, we have argued that modularity is the key to successful programming. Languages which aim to improve productivity must support modular programming well. But new scope rules and mechanisms for separate compilation are not enough – modularity means more than modules. Our ability to decompose a problem into parts depends directly on our ability to glue solutions together. To assist modular programming, a language must provide good glue. Functional programming languages provide two new kinds of glue – higher-order functions and lazy evaluation. Using these glues one can modularise programs in new and exciting ways, and we have shown many examples of this. Smaller and more general modules can be re-used more widely, easing subsequent programming. This explains why functional programs are so much smaller and easier to write than conventional ones. It also provides a target for functional programmers to aim at. If any part of a program is messy or complicated, the programmer should attempt to modularise it and to generalise the parts. He or she should expect to use higher-order functions and lazy evaluation as the tools for doing this.

REFERENCES

1. D. A. Turner, Miranda: a non-strict functional language with polymorphic types. In *Proceedings of the IFIP International Conference on Functional Programming and Computer Architecture*. LNCS 201, Springer Verlag (1985).
2. N. Wirth, *Programming in Modula-II*. Springer Verlag (1982).
3. United States Department of Defense. *The Programming Language Ada Reference Manual*. Springer Verlag (1980).
4. D. McQueen, *Modules for Standard ML*. AT & T Bell Laboratories (May, 1985).
5. D. A. Turner, The Semantic Elegance of applicative Languages. In *ACM Symposium on Functional Languages and Computer Architecture*. Wentworth (1981).
6. H. Abelson and G. J. Sussman, *The Structure and Interpretation of Computer Programs*. MIT Press (1984).
7. P. Henderson, Purely Functional Operating Systems. In *Functional Programming and its Applications*. Cambridge University Press (1982).

Of course, we are not the first to point out the power and elegance of higher-order functions and lazy evaluation. For example, Turner shows how both can be used to great advantage in a program for generating chemical structures.⁵ Abelson and Sussman stress that streams (lazy lists) are a powerful tool for structuring programs.⁶ Henderson has used streams to structure functional operating systems.⁷ However, perhaps we place more emphasis on *modularity* as the key to functional programming's power than do earlier papers.

Some believe that functional languages should be lazy, others believe they should not. Some compromise, and provide only lazy lists, with a special syntax for constructing them (as, for example, in SCHEME).⁶ This paper provides further evidence that lazy evaluation is too important to be relegated to second-class citizenship. It is perhaps the most powerful glue functional programmers possess. One should not obstruct access to such a vital tool.

Acknowledgements

This paper owes much to many conversations with Phil Wadler and Richard Bird in the Programming Research Group at Oxford. Magnus Bondesson at Chalmers University, Göteborg pointed out a serious error in an earlier version of one of the numerical algorithms, and thereby prompted development of many of the others. This work was carried out with the support of a Research Fellowship from the UK Science and Engineering Research Council.