Use these laws, together with the definition of *unzip*, to prove the required result.

## 4.5  The fold functions

We have seen in the case of the datatype *Nat* that many recursive definitions can be expressed very succinctly using a suitable fold operator. Exactly the same is true of lists. Consider the following definition of a function *h*:

$$
\begin{aligned}
h\,[\,] &= e \\
h\,(x : xs) &= x \oplus h\,xs
\end{aligned}
$$

The function *h* works by taking a list, replacing [ ] by *e* and (:) by $\oplus$, and evaluating the result. For example, *h* converts the list

$$x_1 : (x_2 : (x_3 : (x_4 : [\,])))$$

to the value

$$x_1 \oplus (x_2 \oplus (x_3 \oplus (x_4 \oplus e)))$$

Since (:) associates to the right, there is no need to put in parentheses in the first expression. However, we do have to put in parentheses in the second expression because we do not assume that $\oplus$ associates to the right.

The pattern of definition given by *h* is captured in a function *foldr* (pronounced 'fold right') defined as follows:

$$
\begin{aligned}
foldr &\;::\; (\alpha \to \beta \to \beta) \to \beta \to [\alpha] \to \beta \\
foldr\,f\,e\,[\,] &= e \\
foldr\,f\,e\,(x : xs) &= f\,x\,(foldr\,f\,e\,xs)
\end{aligned}
$$

We can now write $h = foldr\,(\oplus)\,e$. The first argument of *foldr* is a binary operator that takes an $\alpha$-value on its left and a $\beta$-value on its right, and delivers a $\beta$-value. The second argument of *foldr* is a $\beta$-value. The third argument is of type $[\alpha]$, and the result is a value of type $\beta$. In many cases, $\alpha$ and $\beta$ will be instantiated to the same type, for instance when $\oplus$ denotes an associative operation.

The single function *foldr* can be used to define almost every function on lists that we have met so far. Here are just some examples:

$$
\begin{aligned}
concat &\;::\; [[\alpha]] \to [\alpha] \\
concat &= foldr\,(+\!\!+)\,[\,]
\end{aligned}
$$

$$reverse \quad :: \quad [\alpha] \rightarrow [\alpha]$$
$$reverse \quad = \quad foldr\ snoc\ [\,]$$
$$\textbf{where}\ snoc\ x\ xs = xs \mathbin{+\!\!+} [x]$$

$$length \quad :: \quad [\alpha] \rightarrow Int$$
$$length \quad = \quad foldr\ oneplus\ 0$$
$$\textbf{where}\ oneplus\ x\ n = 1 + n$$

$$sum \quad :: \quad Num\ \alpha \Rightarrow [\alpha] \rightarrow \alpha$$
$$sum \quad = \quad foldr\ (+)\ 0$$

$$and \quad :: \quad [Bool] \rightarrow Bool$$
$$and \quad = \quad foldr\ (\wedge)\ True$$

$$map \quad :: \quad (\alpha \rightarrow \beta) \rightarrow [\alpha] \rightarrow [\beta]$$
$$map\ f \quad = \quad foldr\ (cons \cdot f)\ [\,]$$
$$\textbf{where}\ cons\ x\ xs = x : xs$$

$$unzip \quad :: \quad [(\alpha, \beta)] \rightarrow ([\alpha], [\beta])$$
$$unzip \quad = \quad foldr\ conss\ ([\,], [\,])$$
$$\textbf{where}\ conss\ (x, y)\ (xs, ys) = (x : xs, y : ys)$$

In some of these definitions the operator $\oplus$ that replaces $(:)$ is associative, and the constant $c$ that replaces $[\,]$ is the unit of $\oplus$. For example, $(+\!\!+, [\,])$, $(+, 0)$, $(\wedge, True)$, and so on. In other definitions, $\oplus$ is not associative; indeed, $\oplus$ may not even have a type appropriate for an associative operator, namely $\alpha \rightarrow \alpha \rightarrow \alpha$.

Not every function on lists can be defined as an instance of *foldr*. For example, *zip* cannot be so defined. Even for those that can, an alternative definition may be more efficient. To illustrate, suppose we want a function *decimal* that takes a list of digits and returns the corresponding decimal number; thus

$$decimal\,[x_0, x_1, \ldots, x_n] = \sum_{k=0}^{n} x_k 10^{(n-k)}$$

It is assumed that the most significant digit comes first in the list. One way to compute *decimal* efficiently is by a process of multiplying each digit by ten and adding in the following digit. For example,

$$decimal\,[x_0, x_1, x_2] \quad = \quad 10 \times (10 \times (10 \times 0 + x_0) + x_1) + x_2$$

This decomposition of a sum of powers is known as *Horner's rule*. Suppose we define $\oplus$ by $n \oplus x = 10 \times n + x$. Then we can rephrase the above equation as

$$decimal\,[x_0, x_1, x_2]\quad=\quad((0 \oplus x_0) \oplus x_1) \oplus x_2$$

This is almost like an instance of *foldr*, except that the grouping is the other way round, and the starting value appears on the left, not on the right. In fact, the computation is dual: instead of processing from right to left, the computation processes from left to right.

This example motivates the introduction of a second fold operator, called *foldl* (pronounced 'fold left'). Informally,

$$foldl\,(\oplus)\,e\,[x_0, x_1, \ldots, x_{n-1}]\quad=\quad(\cdots((e \oplus x_0) \oplus x_1)\cdots) \oplus x_{n-1}$$

The parentheses group from the left, which is the reason for the name. The full definition of *foldl* is

$$
\begin{array}{lll}
foldl & :: & (\beta \to \alpha \to \beta) \to \beta \to [\alpha] \to \beta \\
foldl\,f\,e\,[\,] & = & e \\
foldl\,f\,e\,(x:xs) & = & foldl\,f\,(f\,e\,x)\,xs
\end{array}
$$

For example,

$$
\begin{array}{ll}
& foldl\,(\oplus)\,e\,[x_0, x_1, x_2] \\
= & foldl\,(\oplus)\,(e \oplus x_0)\,[x_1, x_2] \\
= & foldl\,(\oplus)\,((e \oplus x_0) \oplus x_1)\,[x_2] \\
= & foldl\,(\oplus)\,(((e \oplus x_0) \oplus x_1) \oplus x_2)\,[\,] \\
= & ((e \oplus x_0) \oplus x_1) \oplus x_2
\end{array}
$$

If $\oplus$ is associative with unit $e$, then *foldr* $(\oplus)$ $e$ and *foldl* $(\oplus)$ $e$ define the same function on finite lists, as we will see in the following section.

As another example of the use of *foldl*, consider the following definition:

$$
\begin{array}{lll}
reverse' & :: & [\alpha] \to [\alpha] \\
reverse' & = & foldl\,cons\,[\,] \quad \textbf{where}\ cons\,xs\,x = x : xs
\end{array}
$$

Note the order of the arguments to *cons*; we have *cons* = *flip* (:), where the standard function *flip* is defined by *flip* $f\,x\,y = f\,y\,x$. The function *reverse'* reverses a finite list. For example,

$$
\begin{array}{ll}
& reverse'\,[x_0, x_1, x_2] \\
= & cons\,(cons\,(cons\,[\,]\,x_0)\,x_1)\,x_2 \\
= & cons\,(cons\,[x_0]\,x_1)\,x_2 \\
= & cons\,[x_1, x_0]\,x_2 \\
= & [x_2, x_1, x_0]
\end{array}
$$

One can prove that *reverse'* = *reverse* by induction, or as an instance of a more general result to be described in the following section. Of greater importance than the mere fact that *reverse* can be defined in a different way, is that *reverse'* gives a much more efficient program: *reverse'* takes time proportional to $n$ on a list of length $n$, while *reverse* takes time proportional to $n^2$.

There are a number of important laws concerning *foldr* and its relationship to *foldl*, but we will postpone discussion to the following section.

### 4.5.1  Fold over nonempty lists

Say we wish to find the maximum element in a finite list of elements drawn from an ordered type. We would like to do this by defining

$$maxlist \quad :: \quad (Ord\ \alpha) \Rightarrow [\alpha] \to \alpha$$
$$maxlist \quad = \quad foldr\ (\textbf{max})\ e$$

where ($x$ **max** $y$) returns the greater of $x$ and $y$. But what should we choose as the value of $e$? In other words, what is *maxlist* [ ]? The choice is not arbitrary because we naturally want *maxlist* [$x$] = $x$, and the definition gives ($x$ **max** $e$) instead. So $e$ has to be smaller than any other value in the type $\alpha$. Such a value may not exist. In Haskell there is a type class *Bounded* whose instances are types possessing maximum and minimum elements:

**class** *Bounded* $\alpha$ **where**
  *minBound, maxBound*  ::  $\alpha$

Hence one solution is to redefine *maxlist* to

$$maxlist \quad :: \quad (Ord\ \alpha, Bounded\ \alpha) \Rightarrow [\alpha] \to \alpha$$
$$maxlist \quad = \quad foldr\ (\textbf{max})\ minBound$$

An alternative solution that does not make use of *Bounded* is to introduce a new fold function, one that works on nonempty lists. In fact, there are two such functions, *foldr1* and *foldl1*. The definitions are

$$foldr1 \qquad\qquad :: \quad (\alpha \to \alpha \to \alpha) \to [\alpha] \to \alpha$$
$$foldr1\ f\ (x : xs) \quad = \quad \textbf{if}\ null\ xs\ \textbf{then}\ x\ \textbf{else}\ f\ x\ (foldr1\ f\ xs)$$

$$foldl1 \qquad\qquad :: \quad (\alpha \to \alpha \to \alpha) \to [\alpha] \to \alpha$$
$$foldl1\ f\ (x : xs) \quad = \quad foldl\ f\ x\ xs$$

For example,

$$foldr1\ (\oplus)\ [x_0, x_1, x_2, x_3] \quad = \quad x_0 \oplus (x_1 \oplus (x_2 \oplus x_3))$$
$$foldl1\ (\oplus)\ [x_0, x_1, x_2, x_3] \quad = \quad ((x_0 \oplus x_1) \oplus x_2) \oplus x_3$$

Now we can solve the problem by defining

$$maxlist \;=\; foldr1\,(\mathbf{max})$$

Of course, since **max** is associative, we could equally well have defined *maxlist* in terms of *foldl1*.

### 4.5.2   Scan left

Sometimes it is convenient to apply a *foldl* operation to every initial segment of a list. This is done by a function *scanl*, pronounced 'scan left'. For example,

$$scanl\,(\oplus)\,e\,[x_0,x_1,x_2] \;=\; [e,\,e\oplus x_0,\,(e\oplus x_0)\oplus x_1,\,((e\oplus x_0)\oplus x_1)\oplus x_2]$$

In particular, *scanl* $(+)$ 0 computes the list of accumulated sums of a list of numbers, and $scanl\,(\times)\,1\,[1\mathbin{..}n]$ computes a list of the first $n$ factorial numbers. More substantial examples of the use of *scanl* are given in Sections 5.2 and 5.3. We will give two programs for *scanl*; the first is the clearest, while the second is more efficient.

For the first program we will need the function *inits* that returns the list of all initial segments of a list. For example,

$$inits\,[x_0,x_1,x_2] \;=\; [[\,],\,[x_0],\,[x_0,x_1],\,[x_0,x_1,x_2]]$$

The empty list has only one segment, namely the empty list itself; a list $(x:xs)$ has the empty list as its shortest initial segment, and all the other initial segments begin with $x$ and are followed by an initial segment of $xs$. Hence

$$
\begin{aligned}
inits &\;::\; [\alpha] \to [[\alpha]] \\
inits\,[\,] &\;=\; [[\,]] \\
inits\,(x:xs) &\;=\; [\,] : map\,(x:)\,(inits\,xs)
\end{aligned}
$$

The function *inits* can be defined more succinctly as an instance of *foldr*:

$$inits \;=\; foldr\,f\,[[\,]] \quad \textbf{where } f\,x\,xss = [\,] : map\,(x:)\,xss$$

Now we define

$$
\begin{aligned}
scanl &\;::\; (\alpha \to \beta \to \beta) \to \beta \to [\alpha] \to [\alpha] \\
scanl\,f\,e &\;=\; map\,(foldl\,f\,e) \cdot inits
\end{aligned}
$$

This is the clearest definition of *scanl* but it leads to an inefficient program. The function $f$ is applied $k$ times in the evaluation of *foldl* $f\,e$ on a list of length $k$ and, since the initial segments of a list of length $n$ are lists with lengths $0,1,\ldots,n$, the function $f$ is applied about $n^2/2$ times in total.

Let us now synthesise a more efficient program. The synthesis is by an induction argument on *xs* so we lay out the calculation in the same way.

**Case** ([ ]). Here we get *scanl f e* [ ] = [*e*]. The details are routine and are left as an exercise.

**Case** (*x* : *xs*). We calculate

$$scanl\ f\ e\ (x:xs)$$

= {definition of *scanl*}

$$map\ (foldl\ f\ e)\ (inits\ (x:xs))$$

= {definition of *inits*}

$$map\ (foldl\ f\ e)\ ([\ ]:map\ (x:)\ (inits\ xs))$$

= {second equation for *map*}

$$foldl\ f\ e\ [\ ]:map\ (foldl\ f\ e)\ (map\ (x:)\ (inits\ xs))$$

= {first equation for *foldl*}

$$e:map\ (foldl\ f\ e)\ (map\ (x:)\ (inits\ xs))$$

= {since *map* is a functor}

$$e:map\ (foldl\ f\ e \cdot (x:))\ (inits\ xs)$$

= {claim: *foldl f e* · (*x* :) = *foldl f* (*f e x*)}

$$e:map\ (foldl\ f\ (f\ e\ x))\ (inits\ xs)$$

= {definition of *scanl*}

$$e:scanl\ f\ (f\ e\ x)\ xs$$

The claim is left as a short exercise. In summary, we have derived

$$\begin{array}{lll} scanl\ f\ e\ [\ ] & = & [e] \\ scanl\ f\ e\ (x:xs) & = & e:scanl\ f\ (f\ e\ x)\ xs \end{array}$$

This program is more efficient in that the function *f* is applied exactly *n* times on a list of length *n*.

### 4.5.3  Scan right

The dual computation is given by *scanr*, defined by

$$scanr\ f\ a\ =\ map\ (foldr\ f\ a) \cdot tails$$

The function *tails* returns the tail segments of a list. For example,

$$tails\ [x_0, x_1, x_2]\ =\ [[x_0, x_1, x_2], [x_1, x_2], [x_2], [\ ]]$$

Note that while *inits* produces a list of initial segments in increasing order of length, *tails* produces the tail segments in decreasing order of length. We can define *tails* by

$$
\begin{aligned}
tails\ [\ ] \quad &=\quad [[\ ]] \\
tails\ (x : xs) \quad &=\quad (x : xs) : tails\ xs
\end{aligned}
$$

The corresponding efficient program for *scanr* is given by

$$
\begin{aligned}
scanr\ f\ a\ [\ ] \quad &=\quad [a] \\
scanr\ f\ a\ (x : xs) \quad &=\quad f\ x\ (head\ ys) : ys \\
&\qquad \textbf{where}\ ys = scanr\ f\ a\ xs
\end{aligned}
$$

We will leave the synthesis of the more efficient version as an exercise.

### 4.5.4   Scans over nonempty lists

We can also define two further scan functions, *scanl1* and *scanr1*. We have

$$
\begin{aligned}
scanl1\ f \quad &=\quad map\ (foldl1\ f) \cdot inits1 \\
scanr1\ f \quad &=\quad map\ (foldr1\ f) \cdot tails1
\end{aligned}
$$

where *inits1* and *tails1* return the nonempty initial and final segments of a nonempty list. For example,

$$inits1\ (x : xs)\ =\ map\ (x :)\ (inits\ xs)$$

Using *foldl1* $f \cdot (x :) = foldl\ f\ x$, we obtain

$$
\begin{aligned}
scanl1 \quad &::\quad (\alpha \to \alpha \to \alpha) \to [\alpha] \to [\alpha] \\
scanl1\ f\ (x : xs) \quad &=\quad scanl\ f\ x\ xs
\end{aligned}
$$

The efficient version of *scanr1* is left as an exercise.

### Exercises

**4.5.1** Define *filter p* as an instance of *foldr*.

**4.5.2** Consider the two functions *takeWhile* and *dropWhile*, which were discussed in Exercise 4.3.7. Define *takeWhile* as an instance of *foldr*. Can *dropWhile* be defined as an instance of a fold function?

**4.5.3** Which, if either, of the following equations is true?

$$foldl \ (-) \ x \ xs \quad = \quad x - sum \ xs$$
$$foldr \ (-) \ x \ xs \quad = \quad x - sum \ xs$$

**4.5.4** Consider the following definition of a function *insert*:

$$insert \ x \ xs \quad = \quad takeWhile \ (\leq x) \ xs \ \text{\tt ++} \ [x] \ \text{\tt ++} \ dropWhile \ (\leq x) \ xs$$

Show that if *xs* is in nondecreasing order, then so is *insert x xs* for all *x*. Use *insert* to define a function *isort* for sorting a list into nondecreasing order. Use *foldr*.

**4.5.5** The function *remdups* removes adjacent duplicates from a list. For example, *remdups* $[1, 2, 2, 3, 3, 3, 1, 1] = [1, 2, 3, 1]$. Define *remdups* using either *foldl* or *foldr*.

**4.5.6** Given a list $xs = [x_1, x_2, \ldots, x_n]$ of numbers, the sequence of successive maxima *ssm xs* is the longest subsequence $[x_{j_1}, x_{j_2}, \ldots, x_{j_m}]$ such that $j_1 = 1$ and $x_j < x_{j_k}$ for $j < j_k$. For example, the sequence of successive maxima of $[3, 1, 3, 4, 9, 2, 10, 7]$ is $[3, 4, 9, 10]$. Define *ssm* in terms of *foldl*.

**4.5.7** Prove that *scanl f e* $[\ ] = [e]$.

**4.5.8** Justify the efficient program for *scanr*.

**4.5.9** What list does *scanl* $(/) \ [1 .. n]$ produce?

**4.5.10** The mathematical constant *e* is defined by $e = \sum_{n=0}^{\infty} 1/n!$. Write down an expression that can be used to evaluate *e* to some reasonable measure of accuracy.

**4.5.11** Consider the datatype *Liste α* defined in Exercise 4.1.4. The fold function for this datatype is defined by

$$folde \qquad\qquad :: \quad (\beta \to \alpha \to \beta) \to \beta \to Liste \ \alpha \to \beta$$
$$folde \ f \ e \ Nil \qquad = \quad e$$
$$folde \ f \ e \ (Snoc \ xs \ x) \quad = \quad f \ (folde \ f \ e \ xs) \ x$$

The type assigned to *folde* is very nearly the same type as that assigned to *foldl*, the only difference being that $[\alpha]$ is replaced by *Liste α*. In fact,

$$folde \ f \ e \quad = \quad foldl \ f \ e \cdot convert$$

where *convert* $:: [\alpha] \to Liste \ \alpha$. Define the function *convert* and prove by induction that the above equation holds.

**4.5.12** Derive an efficient program for *scanr1*.

## 4.6  Laws of fold

There are a number of important laws concerning *foldr* and its relationship with *foldl*. As we saw in Section 3.3, instead of having to prove a property of a recursively defined function over a recursive datatype by writing down an explicit induction proof, one can often phrase the property as an instance of one of the laws of the fold operator for the datatype.

### 4.6.1  Duality theorems

The first three laws are called *duality theorems* and concern the relationship between *foldr* and *foldl*. The *first duality theorem* is as follows:

**First duality theorem.** Suppose $\oplus$ is associative with unit $e$. Then

$$foldr\ (\oplus)\ e\ xs\ \ =\ \ foldl\ (\oplus)\ e\ xs$$

for all finite lists $xs$.

For example, we could have defined

$$
\begin{array}{lcl}
sum & = & foldl\ (+)\ 0 \\
and & = & foldl\ (\wedge)\ True \\
concat & = & foldl\ (+\!\!+)\ [\,] \\
\end{array}
$$

However, as we will elaborate in Chapter 7, it is sometimes more efficient to implement a function using *foldl*, and sometimes more efficient to use *foldr*.

The *second duality theorem* is a generalisation of the first.

**Second duality theorem.** Suppose $\oplus$, $\otimes$, and $e$ are such that for all $x$, $y$, and $z$ we have

$$
\begin{array}{lcl}
x \oplus (y \otimes z) & = & (x \oplus y) \otimes z \\
x \oplus e & = & e \otimes x \\
\end{array}
$$

In other words, $\oplus$ and $\otimes$ associate with each other, and $e$ on the right of $\oplus$ is equivalent to $e$ on the left of $\otimes$. Then

$$foldr\ (\oplus)\ e\ xs\ \ =\ \ foldl\ (\otimes)\ e\ xs$$

for all finite lists $xs$.

**Proof.** The proof of the second duality theorem is by induction on $xs$.

**Case** ($[\,]$). Both sides simplify to $e$.

**Case** (*x* : *xs*). For the left-hand side we reason

$$foldr \ (\oplus) \ e \ (x : xs)$$

$$= \quad \{\text{definition of } foldr\}$$

$$x \oplus foldr \ (\oplus) \ e \ xs$$

$$= \quad \{\text{induction hypothesis}\}$$

$$x \oplus foldl \ (\otimes) \ e \ xs$$

For the right-hand side we reason

$$foldl \ (\otimes) \ e \ (x : xs)$$

$$= \quad \{\text{definition of } foldl\}$$

$$foldl \ (\otimes) \ (e \otimes x) \ xs$$

$$= \quad \{\text{assumption}\}$$

$$foldl \ (\otimes) \ (x \oplus e) \ xs$$

The two sides have simplified to two different expressions. To show that they are equal we have to resort to a second induction proof. However, if one takes as the second induction hypothesis the obvious assertion that

$$x \oplus foldl \ (\otimes) \ e \ xs \quad = \quad foldl \ (\otimes) \ (x \oplus e) \ xs$$

for all *x* and finite lists *xs*, then the induction doesn't work. The only way to appreciate this fact is to try it, a task we leave as an exercise. What is needed is a *stronger* hypothesis, namely

$$x \oplus foldl \ (\otimes) \ y \ xs \quad = \quad foldl \ (\otimes) \ (x \oplus y) \ xs$$

for all *x*, *y*, and finite lists *xs*. As we saw in Section 4.2 in the discussion of *reverse*, making an induction hypothesis stronger can often simplify a proof.

**Case** ([ ]). Both sides simplify to $x \oplus a$.

**Case** (*z* : *xs*). (Note the use of *z* to avoid confusion with the variable *x* already present in the assertion.) For the left-hand side we reason

$$x \oplus foldl \ (\otimes) \ y \ (z; \ xs)$$

$=$     {definition}

$\quad x \oplus foldl \; (\otimes) \; (y \otimes z) \; xs$

$=$     {induction hypothesis}

$\quad foldl \; (\otimes) \; (x \oplus (y \otimes z)) \; xs$

$=$     {assumption}

$\quad foldl \; (\otimes) \; ((x \oplus y) \otimes z) \; xs$

For the right-hand side we reason

$\quad foldl \; (\otimes) \; (x \oplus y) \; (z : xs)$

$=$     {definition}

$\quad foldl \; (\otimes) \; ((x \oplus y) \otimes z) \; xs$

<div align="right">□</div>

The second duality theorem has the first duality theorem as a special case, namely when $(\oplus) = (\otimes)$.

To illustrate the second duality theorem, consider the following definitions of *length* and *reverse*:

| | | | |
|---|---|---|---|
| *length* | $=$ | *foldr oneplus* 0, | where *oneplus* $x \; n = 1 + n$ |
| *length* | $=$ | *foldl plusone* 0, | where *plusone* $n \; x = n + 1$ |
| | | | |
| *reverse* | $=$ | *foldr snoc* [ ], | where *snoc* $x \; xs = xs \mathbin{+\!\!+} [x]$ |
| *reverse* | $=$ | *foldl cons* [ ], | where *cons* $xs \; x = x : xs$ |

The functions *oneplus*, *plusone*, and 0 meet the conditions of the second duality theorem, as do *snoc*, *cons*, and [ ]. We leave the verification as an exercise. Hence the two definitions of *length* and *reverse* are equivalent on all finite lists. It is not obvious whether there is any practical difference between the two definitions of *length*, but the second program for *reverse* is the more efficient of the two,

The proof of the *third duality theorem* is left as an exercise.

**Third duality theorem.** For all finite lists $xs$,

$\quad foldr \; f \; e \; xs \quad = \quad foldl \; (flip \; f) \; e \; (reverse \; xs)$

where $flip \; f \; x \; y = f \; y \; x$.

To illustrate the third duality theorem, consider

$\quad foldr \; (:) \; [\;] \; xs \quad = \quad foldl \; (flip \; (:)) \; [\;] \; (reverse \; xs)$

Since *foldr* (:) [ ] = *id* and *foldl* (*flip* (:)) [ ] = *reverse*, we obtain

$$xs \quad = \quad reverse\,(reverse\,xs)$$

for all finite lists *xs*, a result we have already proved directly.

### 4.6.2  Fusion

The duality theorems deal with the relationship between the two kinds of fold operator on lists; the fusion theorems deal with one kind of fold operator only. The fusion theorems for *foldr* and *foldl* are as follows.

**Fusion theorem for** *foldr*. If $f$ is strict, $f\,a = b$, and $f\,(g\,x\,y) = h\,x\,(f\,y)$ for all $x$ and $y$, then

$$f \cdot foldr\;g\;a \quad = \quad foldr\;h\;b$$

**Fusion theorem for** *foldl*. If $f$ is strict, $f\,a = b$, and $f\,(g\,x\,y) = h\,(f\,x)\,y$ for all $x$ and $y$, then

$$f \cdot foldl\;g\;a \quad = \quad foldl\;h\;b$$

The proof of the fusion theorem for *foldr* is very similar to the one in the previous chapter for *foldn*, and is left as an exercise. The proof of the fusion theorem for *foldl* is trickier, and is explored in Exercise 4.6.3.

The fusion law is very general, and there are useful special cases. Four of these are considered next.

### 4.6.3  Fold-map fusion

Recall that *map* $g$ = *foldr* (*cons* $\cdot$ $g$) [ ], where *cons* $x\,xs = x : xs$. Consider what conditions ensure that

$$foldr\;f\;a \cdot map\;g \quad = \quad foldr\;h\;b$$

Looking at the conditions of the fusion theorem for *foldr*, we have that *foldr f a* is a strict function, and the second condition is satisfied by taking $a = b$. The third condition expands to

$$foldr\;f\;a\;((cons \cdot g)\;x\;xs) = h\,x\,(foldr\;f\;a\;xs)$$

$$\equiv \qquad \{\text{definition of } cons\}$$

$$foldr\;f\;a\;(g\,x : xs) = h\,x\,(foldr\;f\;a\;xs)$$

$$\equiv \qquad \{\text{definition of } foldr\}$$

$$f\ (g\ x)\ (foldr\ f\ a\ xs) = h\ x\ (foldr\ f\ a\ xs)$$

$$\Leftarrow \quad \{generalisation\}$$

$$f\ (g\ x)\ y = h\ x\ y$$

The symbol $\Leftarrow$ is pronounced 'follows from'. Since $h\ x\ y = f\ (g\ x)\ y$ if and only if $h\ x = f\ (g\ x)$ if and only if $h = f \cdot g$, we get the result that

$$foldr\ f\ a \cdot map\ g \quad = \quad foldr\ (f \cdot g)\ a$$

In words, a fold after a map can always be expressed as a single fold. We will call this law the *fold-map fusion* law.

### 4.6.4   Fold-concat fusion

As a second example, recall that $concat = foldr\ (+\!\!+)\ [\ ]$. Consider what conditions ensure that

$$foldr\ f\ a \cdot concat \quad = \quad foldr\ h\ b$$

Again, looking at the conditions for fusion, we have that $foldr\ f\ a$ is strict, and the second condition is satisfied by taking $a = b$. The third condition expands to

$$foldr\ f\ a\ (xs +\!\!+ ys) = h\ xs\ (foldr\ f\ a\ ys)$$

$$\equiv \quad \{claim:\ foldr\ f\ a\ (xs +\!\!+ ys) = foldr\ f\ (foldr\ f\ a\ ys)\ xs\}$$

$$foldr\ f\ (foldr\ f\ a\ ys)\ xs = h\ xs\ (foldr\ f\ a\ ys)$$

Hence we can take $h\ xs\ y = foldr\ f\ y\ xs$ or, more shortly, $h = flip\ (foldr\ f)$. The proof of the claim is left as an exercise. Hence we have shown that

$$foldr\ f\ a \cdot concat \quad = \quad foldr\ (flip\ (foldr\ f))\ a$$

We will call this the *fold-concat fusion* law.

### 4.6.5   Bookkeeping law

As a third example, consider what conditions on $f$ and $a$ ensure that

$$foldr\ f\ a \cdot concat \quad = \quad foldr\ f\ a \cdot map\ (foldr\ f\ a)$$

We will call this particular equation the *bookkeeping* law. The special instance $sum \cdot concat = sum \cdot map\ sum$ is exploited by every bookkeeper who sums each individual day's accounts, and then gets the year's total by summing the sums.

We can use the fold-concat and fold-map fusion laws to rewrite the book-keeping law in the form

$$foldr \ (flip \ (foldr \ f)) \ a \ = \ foldr \ (f \cdot foldr \ f \ a) \ a$$

This equation holds if $flip \ (foldr \ f) = f \cdot foldr \ f \ a$. Applying both sides to $xs$ and then $y$, we get

$$foldr \ f \ y \ xs \ = \ f \ (foldr \ f \ a \ xs) \ y$$

This equation holds if $f$ is associative with unit $a$, a proof we will leave as an exercise.

In summary, the bookkeeping law is valid whenever the arguments to *foldr* are an associative operation and its unit. For example, since $sum = foldr \ (+) \ 0$, the bookkeeper's accounting method is justified. As another example,

$$concat \cdot concat \ = \ concat \cdot map \ concat$$

This equation says that flattening a list of lists of lists into one long list can be done by concatenating outside-in or inside-out.

### 4.6.6  Fold-scan fusion

Consider the equation

$$1 + x_0 + x_0 \times x_1 + x_0 \times x_1 \times x_2 \ = \ 1 + x_0 \times (1 + x_1 \times (1 + x_2))$$

The equation generalises to an arbitrary list of numbers and can be expressed in the form

$$foldr1 \ (+) \cdot scanl \ (\times) \ 1 \ = \ foldr \ (\odot) \ 1$$

where $x \odot y = 1 + x \times y$. This leads to a more general question: when does the law

$$foldr1 \ (\oplus) \cdot scanl \ (\otimes) \ e \ = \ foldr \ (\odot) \ e$$

hold, where $x \odot y = e \oplus (x \otimes y)$?

Under the assumption that $\otimes$ is associative and $e$ is the unit of $\otimes$, we can express *scanl* as an instance of *foldr*:

$$scanl \ (\otimes) \ e \ = \ foldr \ g \ [e] \quad \text{where } g \ x \ xs = e : map \ (x\otimes) \ xs$$

The proof is left as an exercise. Having rewritten *scanl* as an instance of *foldr*, we can see that the required law is yet another instance of the fusion law for

*foldr*. The first two conditions are immediate, and the third expands to

$$foldr1 \ (\oplus) \ (g \ x \ xs) = x \odot (foldr1 \ (\oplus) \ xs)$$

$\equiv$     {definition of $g$}

$$foldr1 \ (\oplus) \ (e : map \ (x\otimes) \ xs) = x \odot (foldr1 \ (\oplus) \ xs)$$

$\equiv$     {definition of *foldr1* and $\odot$}

$$e \oplus foldr1 \ (\oplus) \ (map \ (x\otimes) \ xs) = e \oplus (x \otimes foldr1 \ (\oplus) \ xs)$$

$\Leftarrow$     {substitution}

$$foldr1 \ (\oplus) \cdot map \ (x\otimes) = (x\otimes) \cdot foldr1 \ (\oplus)$$

The last equation holds if $\otimes$ distributes over $\oplus$, that is, if

$$x \otimes (y \oplus z) \quad = \quad (x \otimes y) \oplus (x \otimes z)$$

for all $x$, $y$, and $z$. The proof is left as an exercise.

In summary, if $\otimes$ is associative with unit $e$, and $\otimes$ distributes over $\oplus$, then $foldr1 \ (\oplus) \cdot scanl \ (\otimes) \ e = foldr \ (\odot) \ e$ where $x \odot y = e \oplus (x \otimes y)$. We will call this the *fold-scan* fusion law.

In the remainder of this section we will give an example of how these laws can be used in the synthesis of efficient programs.

### 4.6.7   Example: the maximum segment sum

The maximum segment sum problem is a famous one and its history is described in J. Bentley's *Programming Pearls* (1987). Given a sequence of numbers it is required to compute the maximum of the sums of all segments in the sequence. For example, the sequence

$$[-1, 2, -3, 5, -2, 1, 3, -2, -2, -3, 6]$$

has maximum sum 7, the sum of the segment $[5, -2, 1, 3]$. On the other hand, the sequence $[-1, -2, -3]$ has a maximum segment sum of zero, since the empty sequence $[\ ]$ is a segment of every list and its sum is zero. It follows that the maximum segment sum is always nonnegative. (The exercises deal with a variant problem in which the empty sequence is excluded, so the maximum segment sum can be negative.)

The problem can be specified formally by

$$mss \quad :: \quad [Int] \to Int$$
$$mss \quad = \quad maxlist \cdot map \ sum \cdot segs$$

where *segs* returns a list of all segments of a list. This function can be defined in a number of ways, including

$$segs \quad = \quad concat \cdot map\ inits \cdot tails$$

This definition corresponds to the process of taking all the initial segments of all the tail segments. For example,

$$segs\ \text{``abc''} \quad = \quad [\text{``''}, \text{``a''}, \text{``ab''}, \text{``abc''}, \text{``''}, \text{``b''}, \text{``bc''}, \text{``''}, \text{``c''}, \text{``''}]$$

The empty sequence appears four times in this list, once for every tail segment of "abc".

Direct evaluation of *mss* will take a number of steps proportional to $n^3$ on a list of length $n$. There are about $n^2$ segments, and summing them all takes about $n^3$ steps. It is not obvious that we can do better than cubic time for this problem.

However, let us see what we can achieve with some calculation. Installing the definition of *segs* in *mss*, we have to compute

$$maxlist \cdot map\ sum \cdot concat \cdot map\ inits \cdot tails$$

The only thing we can do with this expression (apart from going backwards) is to use the law $map\ f \cdot concat = concat \cdot map\ (map\ f)$ to replace the second and third terms, obtaining

$$maxlist \cdot concat \cdot map\ (map\ sum) \cdot map\ inits \cdot tails$$

There are two things we can do now: we can either apply the bookkeeping law to the first two terms, or apply the law $map\ (f \cdot g) = map\ f \cdot map\ g$ to the third and fourth terms. Of course, we can also do both these steps at once, in which case the result will be

$$maxlist \cdot map\ maxlist \cdot map\ (map\ sum \cdot inits) \cdot tails$$

For the next step we can apply the *map* law once more and obtain

$$maxlist \cdot map\ (maxlist \cdot map\ sum \cdot inits) \cdot tails$$

The term $map\ sum \cdot inits$ can be replaced by a *scanl*, a substitution that could have been done one step earlier:

$$maxlist \cdot map\ (maxlist \cdot scanl\ (+)\ 0) \cdot tails$$

Since $maxlist = foldr1\ (\textbf{max})$ and $(+)$ distributes over $(\textbf{max})$, the fold-scan fusion law now gives

$$maxlist \cdot map\ (foldr\ (\odot)\ 0) \cdot tails$$

where $x \odot y = 0 \max (x + y)$. Finally, the last two terms can be replaced by a *scanr*, and we obtain

$$mss \quad = \quad maxlist \cdot scanr (\odot) 0$$

The result is a linear-time program for the maximum segment sum.

**Exercises**

**4.6.1** Prove the third duality theorem.

**4.6.2** Prove the fusion theorem for *foldr*.

**4.6.3** The fusion theorem for *foldl* states that $f \cdot foldl\ g\ a = foldl\ h\ b$ provided $f$ is strict, $f\ a = b$, and $f\ (g\ x\ y) = h\ (f\ x)\ y$ for all $x$ and $y$.
In the induction proof the case $(x : xs)$ requires us to show that

$$f\ (foldl\ g\ (g\ a\ x)\ xs) \quad = \quad foldl\ h\ (h\ (f\ a)\ x)\ xs$$

To do this we need a second induction hypothesis, namely that

$$f\ (foldl\ g\ (g\ x\ y)\ zs) \quad = \quad foldl\ h\ (h\ (f\ x)\ y)\ zs$$

for all $x$ and $y$ and for all lists $zs$. Prove this equation and complete the proof of the fusion law.

**4.6.4** Prove the equations

$$foldr\ f\ a\ (xs + ys) \quad = \quad foldr\ f\ (foldr\ f\ a\ ys)\ xs$$
$$foldl\ f\ a\ (xs + ys) \quad = \quad foldl\ f\ (foldl\ f\ a\ xs)\ ys$$

What restrictions on $xs$ and $ys$ are necessary, if any?

**4.6.5** Prove that *scanl* $(\otimes)\ e$ can be written as a *foldr* if $\otimes$ is associative with unit $e$. Can *scanl* $(\otimes)\ e$ be written as a *foldl* without any assumptions on $\otimes$ and $e$?

**4.6.6** Prove that if $\otimes$ distributes over $\oplus$, then

$$foldr1 (\oplus) \cdot map (x\otimes) \quad = \quad (x\otimes) \cdot foldr1 (\oplus)$$

**4.6.7** It is tempting to try to push the calculation of the maximum segment sum one step further and use the fold-scan fusion law to replace *maxlist* $\cdot$ *scanr* $(\odot) 0$ by a *foldl*. Why doesn't this work? (See Exercise 4.6.10 for a version of fold-scan fusion that does work.)

**4.6.8** Give a fusion law for *foldr1* and use it to give conditions on $f$, $g$, and $h$ so that the law

$$foldr1\ f \cdot scanl1\ g\ =\ foldr1\ h$$

is valid.

**4.6.9** Define *segs1* by *segs1* = *concat* · *map inits1* · *tails1*, where *inits1* returns the nonempty initial segments of a list, and *tails1* the nonempty tail segments. Use this version to state and solve a version of the maximum segment sum problem that leaves out empty segments.

**4.6.10** There is a version of the fold-scan fusion law that does not depend on any properties of the operators involved. The law is

$$foldl1\ (\oplus) \cdot scanl\ (\otimes)\ e\ =\ fst \cdot foldl\ (\odot)\ (a, a)$$

with $(x, y) \odot z\ =\ (x \oplus t, t)$, where $t\ =\ y \otimes z$. Prove this version of fold-scan fusion.

## 4.7 Chapter notes

Haskell provides many more standard functions on lists than we have described in this chapter. For full details, consult the Haskell standard prelude and the library of list functions.

A great deal of work has been done on automatic or machine-aided generation of induction proofs; see, for example, Boyer and Moore (1979), Gordon, Milner, and Wadsworth (1979), Paulson (1983), Martin and Nipkow (1990).

The relationship between polymorphic functions and naturality conditions is explored in Wadler (1989); see also Bird and de Moor (1997).

The fusion laws for fold operators were systematised in Malcolm (1990), Fokkinga (1992), Jeuring (1993), and Meijer (1992).

The maximum segment sum problem is described in Bentley (1987). The derivation recorded in this chapter was first given in Bird (1989).