

Lecture 1: First-order data types and pattern-matching

Noam Zeilberger

version: September 6, 2023

A **data type** is a type defined by a collection of **constructors**, each of which can take any number of arguments of different types. Informally, the concept of data type is relevant in any situation where you have something that can be built out of different kinds of things – of course there are many such situations, which is why the concept is so important! In this lecture, we will learn how data types are defined, how to write functions over them by pattern-matching, and how to perform some simple reasoning about the behavior of these functions. We will restrict to data types that are *first-order*, in the sense that the constructors we consider do not themselves take functions as arguments – although they can take other data types as arguments, and they can moreover be polymorphic. Such types are also sometimes called **algebraic** data types, since (as we will see) they obey laws similar to the algebraic laws for sums and products.

1 First example: the booleans

Let's begin by considering one of the simplest examples of a data type, the type of booleans, which in the Haskell Prelude is defined as follows:

```
data Bool = False | True
```

The way to read this is that *Bool* is a data type with two constructors taking no arguments, that is:

```
False :: Bool  
True  :: Bool
```

The definition also guarantees that these are the *only* two ways to build a value of boolean type, and it is this guarantee that enables us to define functions over the booleans by pattern-matching against these two cases.

For example, boolean negation can be written like so:

```
not :: Bool → Bool  
not False = True  
not True  = False
```

Here we have declared *not* to be a function from booleans to booleans, and defined it by its action on the two possible input values. One easy consequence of this definition is that *not* is an involution, in the sense that

$$\text{not} (\text{not } x) = x$$

for all values $x :: \text{bool}$. To prove this, it suffices to consider the two values $x = \text{false}$ and $x = \text{true}$ and apply the definition of *not*:

$$\begin{aligned} \text{not } (\text{not } \text{False}) &= \text{not } \text{True} = \text{False} \\ \text{not } (\text{not } \text{True}) &= \text{not } \text{False} = \text{True} \end{aligned}$$

Although this “theorem” may not look that impressive, it provides us a first example of how to use reasoning-by-cases combined with purely equational reasoning to prove something about the behavior of a functional program – eventually we will see many more interesting examples!

Next let’s define boolean conjunction:

$$\begin{aligned} \text{both} &:: \text{Bool} \rightarrow \text{Bool} \rightarrow \text{Bool} \\ \text{both } \text{False } \text{False} &= \text{False} \\ \text{both } \text{False } \text{True} &= \text{False} \\ \text{both } \text{True } \text{False} &= \text{False} \\ \text{both } \text{True } \text{True} &= \text{True} \end{aligned}$$

The first thing to pay attention to here is the type of *both*. Arrow associates to the right by default, so the type declaration that we wrote above is the same as writing $\text{both} :: \text{Bool} \rightarrow (\text{Bool} \rightarrow \text{Bool})$. Literally, this says that *both* is a function from booleans to functions from booleans to booleans, which logically is essentially the same thing as a function that takes a pair of booleans to a boolean. In Haskell, functions of multiple arguments are usually expressed like this as having an iterated function type, in so-called “curried” form.¹ Otherwise, we can notice that the above definition by pattern-matching essentially expresses the truth table of boolean conjunction.

The definition of boolean conjunction in the Haskell Prelude, written $(\&\&)$, is actually slightly different:

$$\begin{aligned} (\&\&) &:: \text{Bool} \rightarrow \text{Bool} \rightarrow \text{Bool} \\ \text{False } \&\& _ &= \text{False} \\ \text{True } \&\& b &= b \end{aligned}$$

One inessential difference with the definition of *both* is that $(\&\&)$ is defined as an infix operator, but this is just “syntactic sugar”. A more significant, albeit subtle, difference between the two definitions of conjunction is that $(\&\&)$ implements *short-circuit evaluation*. That’s because Haskell is a so-called “call-by-need” language, and arguments to functions are only evaluated when they are needed to compute the result. In this case, once the conjunction function sees that its first argument is *False*, it doesn’t even need to examine its second argument (here, indicated by matching against the “wildcard” pattern $_$) to determine that the result should be *False*.

Although subtle, this example helps to illustrate the important difference in functional programming between *values* and *expressions*. A value of a given data type is, by definition, built using one of its constructors. An expression of that type, on the other hand, can be defined using all of the mechanisms of the language, and in general denotes a *computation* that is meant to reduce to

¹After the logician Haskell Curry (1900–1982), who employed this technique in his work on combinatory logic, although the technique was actually invented earlier by Moses Schönfinkel (1888–1942).

a value. For example, *not False* is an expression of type *Bool*, which evaluates to the value *True*. In Haskell, as in most functional languages, it is possible to write down expressions that ultimately never reduce to a value. For instance, we can give a circular definition of a boolean that just loops back on itself:

```
loop :: Bool
loop = loop
```

Although *loop* is not a value, it is a valid expression of type *Bool*, and we can even use it to observe the difference between the two versions of conjunction. Consider the following two expressions:

both False loop vs. *False && loop*

Since *loop* never reduces to either *True* or *False*, none of the four clauses in the definition of *both* can be invoked, and so the first expression gets stuck in an infinite loop. On the other hand, the second expression immediately reduces to the value *False*, by the first clause in the definition of (*&&*).

It is worth verifying this for yourself! You can do so by loading the file `DataCode.hs` associated to these notes (available on the course webpage) into the Haskell interpreter, and typing in the above expressions. You should see something roughly like the following interaction (since the first expression doesn't terminate, we type Control-C to interrupt the computation):

```
$ ghci DataCode
GHCi, version 8.8.4: https://www.haskell.org/ghc/  :? for help
[1 of 1] Compiling Main                ( DataCode.hs, interpreted )
Ok, one module loaded.
*Main> let loop = loop
*Main> both False loop
^CInterrupted.
*Main> False && loop
False
*Main>
```

Exercise 1.1. Define the boolean disjunction function by pattern-matching, both with and without short-circuit evaluation, and write down an expression that allows one to observe the difference in behavior between the two definitions.

2 Sums and products

Besides defining particular types like *Bool*, data type declarations also provide a general way of combining one or more types to form a new type. Two important instances of this are called **sum types** and **product types**.

In Haskell, sum types are defined as follows:

```
data Either a b = Left a | Right b
```

Here, *Either* is called a **type constructor**, since it takes two types *a* and *b* and defines a new type *Either a b*. The data declaration automatically introduces two (value) constructors:

$Left :: a \rightarrow Either\ a\ b$
 $Right :: b \rightarrow Either\ a\ b$

which, by definition, are the only two ways of forming a value of the sum type. In set-theoretic terms, the set of values of a sum type may be considered as a disjoint union $\{Left\ x \mid x \in a\} \cup \{Right\ y \mid y \in b\}$ of two sets of “tagged” values.

The fact that the values of a sum type are tagged *Left* or *Right* is precisely what allows us to define functions out of them by pattern-matching, without risk of ambiguity. In general, if $f :: a \rightarrow c$ and $g :: b \rightarrow c$ are two functions taking a and b , respectively, to the same target type c , then we can define a single function taking their sum to the same target by

$h :: Either\ a\ b \rightarrow c$
 $h\ (Left\ x) = f\ x$
 $h\ (Right\ y) = g\ y$

As a simple example, here is a function that coerces either a boolean or an integer² to an integer:

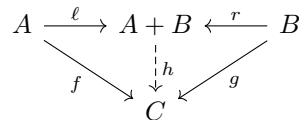
$asInt :: Either\ Bool\ Int \rightarrow Int$
 $asInt\ (Left\ b) = \mathbf{if\ } b\ \mathbf{then\ } 1\ \mathbf{else\ } 0$
 $asInt\ (Right\ n) = n$

Here the function $f :: Bool \rightarrow Int$ is given by $f\ b = \mathbf{if\ } b\ \mathbf{then\ } 1\ \mathbf{else\ } 0$, while $g :: Int \rightarrow Int$ is the identity function $g\ n = n$. As another example, we can write a function that checks whether a value of this same sum type is in fact a boolean:

$isBool :: Either\ Bool\ Int \rightarrow Bool$
 $isBool\ (Left\ b) = True$
 $isBool\ (Right\ n) = False$

In this case, $f :: Bool \rightarrow Int$ is the constant-true function $f = b \rightarrow True$ while $g :: Int \rightarrow Int$ is the constant-false function $g = n \rightarrow False$.

It is also worth mentioning that the definition of sum types in Haskell is formally quite closely related to (and under some interpretations an instance of) the definition of “coproducts” in category theory. Without getting into the completely formal definition, roughly speaking, a “category” is just a collection of objects and a collection of arrows between objects that can be composed (like paths in a graph) in an associative way.³ In any category, one can define the *coproduct* of two objects A and B as an object $A + B$ equipped with arrows $\ell : A \rightarrow A + B$ and $r : B \rightarrow A + B$, such that for any pair of arrows $f : A \rightarrow C$ and $g : B \rightarrow C$ there exists a unique $h : A + B \rightarrow C$ such that $f = h \circ \ell$ and $g = h \circ r$, as summarized in the following diagram:



²Haskell has several different types of integers, including arbitrary-length integers (or so-called “bignums”). Here we are using the type *Int* of machine-integers.

³Two examples are the category *Set* whose objects are sets and whose arrows are functions between sets, and the category *Vec* whose objects are vector spaces and whose arrows are linear transformations. Any partial order may also be considered as a category with a unique arrow $x \rightarrow y$ iff $x \leq y$.

Again without getting into the formal details, we can see that there is a close analogy between the arrows $\ell : A \rightarrow A + B$ and $r : B \rightarrow A + B$ in the definition of the coproduct and the constructors $Left :: a \rightarrow Either\ a\ b$ and $Right :: b \rightarrow Either\ a\ b$ of the sum type, and between the arrow $h : A + B \rightarrow C$ defined by the “universal property” of the coproduct and the canonical function $h :: Either\ a\ b \rightarrow c$ we defined above by pattern-matching.⁴

Whereas sum types describe values that can take multiple possible forms, product types describe values that contain multiple components. Haskell has built-in product types, written (a, b) where a and b are types. A value of type (a, b) is a pair (u, v) , where $u :: a$ and $v :: b$.⁵ The Haskell Prelude also provides built-in projection functions

$$\begin{aligned}fst &:: (a, b) \rightarrow a \\snd &:: (a, b) \rightarrow b\end{aligned}$$

which satisfy the equations $fst\ (u, v) = u$ and $snd\ (u, v) = v$.

If Haskell didn’t already have built-in product types, we could define them for ourselves as a data type, for example like so:

```
data Both a b = Pair a b
deriving (Show, Eq)
```

Of course, the names that we gave to the type and value constructors are arbitrary.⁶ The “deriving” line is just a bit of Haskell syntax to automatically derive pretty-printing and equality-testing routines for values of our newly-defined data type, which is useful for testing. (We will talk more about *Show*, *Eq*, and other type classes in the next lecture.) Once again, the way to read the data type declaration itself is that it introduces a single constructor taking two arguments

$$Pair :: a \rightarrow b \rightarrow Both\ a\ b$$

and asserts that this is the *only* way to build a value of the product type. In other words, a value of type *Both a b* necessarily contains both a value of type a and a value of type b .

Using pattern-matching, we can write projection functions for our own version of product types:

$$\begin{aligned}projFst &:: Both\ a\ b \rightarrow a \\projFst\ (Pair\ u\ v) &= u \\projSnd &:: Both\ a\ b \rightarrow b \\projSnd\ (Pair\ u\ v) &= v\end{aligned}$$

As a slightly more interesting example, we can write a pair of coercion functions witnessing the fact that having both A and either B or C is essentially the same thing as having either A and B or having A and C :

⁴Haskell is a bit idiosyncratic in using the double colon $::$ for the typing judgment, as opposed to the more standard single colon $:$ in other programming languages (such as OCaml) and in logic and mathematics. I will try to stick to the Haskell notation, although I may mix them from time to time, as in the preceding paragraph.

⁵ This kind of overloading of notation for type constructors and value constructors is common in Haskell. Although it can be confusing at first, it does have the advantage of reducing the number of notations you need to memorize!

⁶And if we were to follow typical Haskell naming conventions, it might be a bit more idiomatic to reuse the same name for both, see Footnote 5.

```

coerceTo :: Both a (Either b c) → Either (Both a b) (Both a c)
coerceTo (Pair x (Left y)) = Left (Pair x y)
coerceTo (Pair x (Right y)) = Right (Pair x y)
coerceFrom :: Either (Both a b) (Both a c) → Both a (Either b c)
coerceFrom (Left (Pair x y)) = Pair x (Left y)
coerceFrom (Right (Pair x y)) = Pair x (Right y)

```

It is not hard to check that the equations

$$\text{coerceFrom} (\text{coerceTo } u) = u \quad \text{and} \quad \text{coerceTo} (\text{coerceFrom } v) = v$$

hold for all $u :: \text{Both } a \text{ (Either } b \text{ c)}$ and $v :: \text{Either (Both } a \text{ b) (Both } a \text{ c)}$. This means that not only can we convert values of one type into values of the other, and vice versa, but moreover these transformations are reversible. Formally, we say that the pair of coercion functions realize a **type isomorphism** in the sense, that there are a pair of functions

$$f :: A \rightarrow B \quad \text{and} \quad g :: B \rightarrow A$$

such that $g (f x) = x$ for all $x :: A$, and $f (g y) = y$ for all $y :: B$, i.e., f and g are inverses. In this case, we can prove that

$$\text{Both } a \text{ (Either } b \text{ c)} \cong \text{Either (Both } a \text{ b) (Both } a \text{ c)} \quad (1)$$

which is a more mathematically precise way of saying that the two types are “essentially the same”. This isomorphism (1) expressing the *distributivity* of product types over sum types is an analogue of the algebraic identity

$$a(b + c) = ab + ac \quad (2)$$

familiar from high school algebra, which (as you will see in the exercises below) is an instance of a more general phenomenon.

Finally, both sum types and product types can be generalized from these binary versions to combine any number of types, including no types at all. The nullary product type is called the **unit type** and is written `()` in Haskell, with a single value also written `()`. Of course we can define our own isomorphic version, as a data type with a single constructor taking no arguments:

```

data Unit = U
deriving (Show, Eq)

```

The nullary sum type is called the **zero type** (or *void type*), and corresponds to a data type with *no constructors*. We can express this in Haskell (since the 2010 revision of the language) using the following syntax:

```

data Zero

```

Exercise 2.1. For each of the following valid type isomorphisms, write a pair of coercion functions realizing said isomorphism:

$$\text{Either } a \text{ (Either } b \text{ c)} \cong \text{Either (Either } a \text{ b) } c \quad (1)$$

$$\text{Either } a \text{ b} \cong \text{Either } b \text{ a} \quad (2)$$

$$\text{Both } a \text{ (Both } b \text{ c)} \cong \text{Both (Both } a \text{ b) } c \quad (3)$$

$$\text{Both } a \text{ b} \cong \text{Both } b \text{ a} \quad (4)$$

$$\text{Both Unit } a \cong a \quad (5)$$

$$\text{Either Zero } a \cong a \quad (6)$$

Exercise 2.2. Write a pair of functions of the following types:

$$\begin{aligned} f &:: \text{Either } a \text{ (Both } b \text{ } c) \rightarrow \text{Both (Either } a \text{ } b) \text{ (Either } a \text{ } c) \\ g &:: \text{Both (Either } a \text{ } b) \text{ (Either } a \text{ } c) \rightarrow \text{Either } a \text{ (Both } b \text{ } c) \end{aligned}$$

Does this show that $\text{Either } a \text{ (Both } b \text{ } c) \cong \text{Both (Either } a \text{ } b) \text{ (Either } a \text{ } c)$? Explain why or why not.

3 Lists and structural induction

Lists are ubiquitous in functional programming, appearing frequently both in the interfaces to functions as well as in the intermediate results of computations. In Haskell, the list type is notated $[a]$, with the type a of the underlying elements of the list surrounded in brackets. Modulo syntax, lists are defined by the following data type declaration:

$$\mathbf{data} [a] = [] \mid a : [a]$$

That is, a list of as is either the empty list, written $[]$ and pronounced “nil”, or constructed by prepending an element of type a to a list of as , with this binary operation written $x : xs$ and pronounced (for historical reasons) “ x cons xs ”. Note that lists are an example of a *recursive* data type, since they are defined in terms of themselves.

The above is not quite valid Haskell syntax, since Haskell does not allow operators like $[_]$ to be defined by data declarations. If we want, we can define our own isomorphic version of lists via the following perfectly legal Haskell data type declaration:

$$\begin{aligned} \mathbf{data} List\ a &= Nil \mid Cons\ a\ (List\ a) \\ \mathbf{deriving} &\ (Show, Eq) \end{aligned}$$

The following pair of functions realize an isomorphism $[a] \cong List\ a$:

$$\begin{aligned} toList &:: [a] \rightarrow List\ a \\ toList [] &= Nil \\ toList (x : xs) &= Cons\ x\ (toList\ xs) \\ fromList &:: List\ a \rightarrow [a] \\ fromList Nil &= [] \\ fromList (Cons\ x\ xs) &= x : fromList\ xs \end{aligned}$$

However, let’s stick to Haskell’s built-in syntax for lists in the rest of these notes.

Two basic operations that one can perform on a list are to extract its head and its tail, defined as follows:

$$\begin{aligned} head &:: [a] \rightarrow a \\ head (x : xs) &= x \\ tail &:: [a] \rightarrow [a] \\ tail (x : xs) &= xs \end{aligned}$$

Observe that these definitions are only partial, in the sense that they do not specify a behavior when the input is nil $[]$. By default, Haskell does not require

pattern-matching definitions to cover all possible patterns, and in this case it is clear that there is no good choice of values to return (since the empty list does not have a head or tail). It is generally considered good practice to cover all possible patterns for the input type, though, and to explicitly raise an exception in the case of unexpected input – indeed, that’s what’s done for the definitions of *head* and *tail* that appear in the Haskell Prelude. In any case, the Haskell compiler will generate a pattern-matching exception by default in the case that the input does not match any one of the provided patterns, so that the behavior of the function on any input is always well-defined.

Another basic operation on lists is concatenation. The concatenation of two lists xs and ys is notated $xs \ ++ \ ys$ in Haskell, and defined as follows:

$$\begin{aligned} (+) &:: [a] \rightarrow [a] \rightarrow [a] \\ [] \ ++ \ ys &= ys \\ (x : xs) \ ++ \ ys &= x : (xs \ ++ \ ys) \end{aligned}$$

Observe that this is a *recursive* definition: the concatenation of $x : xs$ with ys is defined as the list obtained by cons’ing x to the front of the concatenation of xs with ys . In general, recursive definitions may lead to non-termination (as we saw with the example of *loop* above), but in this case the circularity is not “vicious” since the first argument to $(+)$ gets smaller. Indeed, we can prove that $xs \ ++ \ ys$ is a well-defined list by induction on the length of xs , or alternatively by *structural induction* on the list xs itself.

The principle of **structural induction over lists** says that for any property $P(xs)$ of a list xs , if we can show both

1. $P([])$ holds; and
2. for any element x and list xs , $P(xs)$ implies $P(x : xs)$

then $P(xs)$ holds for all lists xs . To be a bit more precise, we should specify the type of the elements of the list. If $P(xs)$ is a property of lists of elements of type a , and we can show (1) above as well as that

2. for any element $x :: a$ and list $xs :: [a]$, $P(xs)$ implies $P(x : xs)$

then $P(xs)$ holds for all lists $xs :: [a]$.

Now let’s take the time to spell out how to prove that $(+)$ is well-defined by structural induction. (Although it may seem like overkill for this example, it will serve as a useful template for proving more interesting properties.) To that end, we take $P(xs :: [a])$ to be the property that “for every list $ys :: [a]$, there is a well-defined list $zs :: [a]$ such that $xs \ ++ \ ys = zs$ ”. We need to show two things:

1. $P([])$: by definition $[] \ ++ \ ys = ys$, so we take $zs = ys$.
2. $P(xs)$ implies $P(x : xs)$: by the inductive assumption, there is a list zs' such that $xs \ ++ \ ys = zs'$. But by definition $(x : xs) \ ++ \ ys = x : (xs \ ++ \ ys) = x : zs'$, so we take $zs = x : zs'$.

We have thus established that $P(xs)$ holds for all lists $xs :: [a]$, i.e., that for all $xs :: [a]$ and $ys :: [a]$, there is a well-defined $zs :: [a]$ such that $xs \ ++ \ ys = zs$.

Exercise 3.1. Classically, a *monoid* is defined as a set equipped with an associative binary operation and an identity element. For example, the natural numbers form a monoid under addition with zero as the identity element, written $(\mathbb{N}, +, 0)$, and also a monoid under multiplication with one as the identity element, written $(\mathbb{N}, \times, 1)$. Show that lists form a monoid under concatenation with `nil` as the identity element, by proving that the equations:

$$[] \# xs = xs \tag{1}$$

$$xs \# [] = xs \tag{2}$$

$$xs \# (ys \# zs) = (xs \# ys) \# zs \tag{3}$$

hold for all lists xs , ys , and zs . *Hint:* Use structural induction for (2) and (3).

Exercise 3.2. A *monoid homomorphism* from $(X, *, e_X)$ to (Y, \bullet, e_Y) is defined as a function $\phi : X \rightarrow Y$ that is compatible with both monoid structures, in the sense that $\phi(x_1 * x_2) = \phi(x_1) \bullet \phi(x_2)$ for all $x_1, x_2 \in X$ and moreover $\phi(e_X) = e_Y$. Prove that the length function:

$$\begin{aligned} \text{length} &:: [a] \rightarrow \text{Int} \\ \text{length} [] &= 0 \\ \text{length} (x : xs) &= 1 + \text{length} xs \end{aligned}$$

defines a monoid homomorphism from lists under concatenation to the natural numbers under addition. *Hint:* Use structural induction, and some properties of addition.

Exercise 3.3. Let (b, \otimes, e) be an arbitrary monoid represented in Haskell, i.e., a type b , a binary operation $(\otimes) : b \rightarrow b \rightarrow b$, and a value $e :: b$ such that

$$e \otimes v = v = v \otimes e \quad u \otimes (v \otimes w) = (u \otimes v) \otimes w$$

for all values $u, v, w :: b$. Suppose moreover that you are given a function $f : a \rightarrow b$. Write a function $\phi : [a] \rightarrow b$ satisfying both of the following properties:

- ϕ behaves like f for singleton lists, i.e., $\phi [x] = f x$ for all $x :: a$, and
- ϕ is a monoid homomorphism.

Finally, explain how to recover the *length* function from Exercise 3.2 as ϕ for a particular choice of monoid (b, \otimes, e) and function $f : a \rightarrow b$.

Exercise 3.4. Prove that the function ϕ you defined in Exercise 3.3 is essentially unique, in the following sense: if $\psi :: [a] \rightarrow b$ is any other monoid homomorphism such that $\psi [x] = f x$ for all $x :: a$, then $\psi xs = \phi xs$ for all $xs :: [a]$. (As a result of all these properties, we can conclude that the type of lists $[a]$ equipped with concatenation forms the *free monoid* over a .)

4 Maybe types

It sometimes happens that we want a function to try to compute something that *may* yield a value of some type, but might also fail, and we want the function to tell us explicitly when it fails. In Haskell, this is achieved by having the function return a value of the *Maybe* type, which is defined as follows:

```
data Maybe a = Nothing | Just a
```

The type *Maybe a* is called a **maybe type**, or alternatively *option type* (which is the terminology in Standard ML and OCaml). Observe that $Maybe\ a \cong Either\ ()\ a$, meaning that maybe types can be isomorphically encoded using sum types and the unit type – although maybe types are so commonly used that they deserve their own name and syntax!

To give an example application of maybe types, the Prelude defines a function *lookup* which tries to find the value paired with a key in a list of key-value pairs:⁷

```
lookup :: Eq a => a -> [(a, b)] -> Maybe b
lookup k [] = Nothing
lookup k ((k', v) : kvs)
  | k == k' = Just v
  | otherwise = lookup k kvs
```

As another example, the Standard Library function *elemIndex* tries to find the index of an element within a list:

```
elemIndex :: Eq a => a -> [a] -> Maybe Int
elemIndex x [] = Nothing
elemIndex x (x' : xs)
  | x == x' = Just 0
  | otherwise = case elemIndex x xs of
    Nothing -> Nothing
    Just i -> Just (i + 1)
```

Notice that here we use pattern-matching on the result of the recursive call, via Haskell’s **case** construct.

Exercise 4.1. Prove the type isomorphism

$$Both\ (Maybe\ a)\ (Maybe\ b) \cong Maybe\ (Either\ (Either\ a\ b)\ (Both\ a\ b))$$

which is the analogue of the algebraic identity $(1 + a)(1 + b) = 1 + a + b + ab$.

5 Introducing accumulators

Although the connections between functional programming and mathematics are fascinating, eventually one has to contend with the fact that in programming, we care not only about mathematical correctness but also about *efficiency*. There may be different ways of implementing the same function that although equivalent in terms of their input-output behavior, differ wildly in terms of their consumption of time, space, or other resources – and understanding these costs is an important part of learning functional programming.

The use of *accumulators* is a general technique that allows to derive relatively efficient functional programs by in a sense mimicking the behavior of a corresponding imperative program. Consider the list reversal function. We can

⁷This definition makes use of Haskell’s *pattern guards*, which extend ordinary pattern-matching against constructors by allowing tests of arbitrary boolean conditions.

give a straightforward recursive implementation that reduces list reversal to list concatenation:

$$\begin{aligned} \text{reverse} &:: [a] \rightarrow [a] \\ \text{reverse} [] &= [] \\ \text{reverse} (x : xs) &= \text{reverse } xs \text{ ++ } [x] \end{aligned}$$

It is easy to verify that this implementation of list reversal is correct, and even to establish that it satisfies expected properties such as $\text{reverse } (xs \text{ ++ } ys) = \text{reverse } ys \text{ ++ } \text{reverse } xs$ and $\text{reverse } (\text{reverse } xs) = xs$. However, this implementation of *reverse* is inefficient. The issue is that computing the concatenation of two lists takes time proportional to the length of the first list. Given a list *xs* of length *n*, *reverse xs* makes *n* calls to concatenation of the form $xs' \text{ ++ } [x]$, where xs' is the result of calling *reverse* on a sublist of *xs* of length $0 \leq k \leq n$. Since the *reverse* function is length-preserving, the overall time complexity of this implementation of *reverse* is $\Theta(n^2)$ in the length of the input list.

On the other hand, there is a simple imperative algorithm for reversing a list in $\Theta(n)$ time, using an auxiliary stack:

1. Initialize the stack to be empty.
2. While the input list is non-empty, push its head onto the stack, and keep processing its tail.
3. Once the input list is empty, return the contents of the stack.

Perhaps surprisingly, we can turn this imperative solution into an elegant functional program! The idea is to define a helper function which, in addition to the input list *xs*, takes the auxiliary “stack” (which is in fact just another list) as an extra input argument *ys*:

$$\begin{aligned} \text{revacc} &:: [a] \rightarrow [a] \rightarrow [a] \\ \text{revacc} [] \text{ } ys &= ys \\ \text{revacc} (x : xs) \text{ } ys &= \text{revacc } xs \text{ } (x : ys) \end{aligned}$$

The first clause $\text{revacc} [] \text{ } ys = ys$ of the definition corresponds to step (3) of the imperative algorithm above, while the clause $\text{revacc} (x : xs) \text{ } ys = \text{revacc } xs \text{ } (x : ys)$ corresponds to step (2). Finally, step (1) is implemented by (re-)defining *reverse* in terms of *revacc*:

$$\text{reverse } xs = \text{revacc } xs []$$

We have thus obtained an efficient, purely functional implementation of list reversal. Here’s an example of it in action:

$$\begin{aligned} \text{reverse } [1, 2, 3, 4] & \\ &= \text{revacc } [1, 2, 3, 4] [] \\ &= \text{revacc } [2, 3, 4] [1] \\ &= \text{revacc } [3, 4] [2, 1] \\ &= \text{revacc } [4] [3, 2, 1] \\ &= \text{revacc } [] [4, 3, 2, 1] \\ &= [4, 3, 2, 1] \end{aligned}$$

We can see that the second argument of *revacc* really behaves like a stack, and that the reversal is indeed computed in linear time.

The second argument of *revacc* is said to be an “accumulator”, since it accumulates intermediate results on the way towards computing the final answer. That gives a very operational way of understanding the accumulator technique, but there is also a slightly more conceptual way of understanding it, related to the idea that to solve a particular problem, oftentimes it can be helpful to try to solve a *more general problem*. In this instance, the function *revacc* actually solves the following more general problem: given two lists, compute the reversal of the first list concatenated with the second list, that is, $revacc\ xs\ ys = reverse\ xs \ ++\ ys$. After writing this more general function, we can then reduce list reversal to the case $ys = []$.

Another classic example of the use of accumulators is computing Fibonacci numbers. Mathematically, the Fibonacci numbers F_n are defined by the recurrence $F_n = F_{n-1} + F_{n-2}$ together with initial values $F_0 = 0$, $F_1 = 1$, which translates directly to the following Haskell code:⁸

```
fib :: Integer -> Integer
fib n
  | n == 0 = 0
  | n == 1 = 1
  | n >= 2 = fib (n - 1) + fib (n - 2)
```

Although mathematically correct, this recursive computation of the Fibonacci numbers is horribly inefficient: each call to *fib* spawns two recursive calls to *fib* while only decreasing the input by either 1 or 2, and so *fib n* has complexity exponential in n . We can see this by running the program on a few inputs:

```
*Main> :set +s -- ask ghci to print time and space usage
*Main> fib 10
55
(0.02 secs, 123,512 bytes)
*Main> fib 20
6765
(0.08 secs, 6,423,944 bytes)
*Main> fib 30
832040
(2.38 secs, 781,578,344 bytes)
*Main> fib 31
1346269
(3.58 secs, 1,264,577,008 bytes)
*Main> fib 32
2178309
(6.05 secs, 2,046,084,072 bytes)
```

We observe that the behaviour is indeed exponential, with each increment of n multiplying both the time and space usage by roughly the golden ratio.

On the other hand, there is a much more efficient imperative algorithm for computing F_n in linear time, using a pair of auxiliary variables a and b :

⁸Here we use the type *Integer* of arbitrary-length integers, rather than machine integers.

- Initialize $a \leftarrow 0$ and $b \leftarrow 1$.
- While $n > 0$, simultaneously update $(a, b) \leftarrow (b, a + b)$, and decrement n .
- Once $n = 0$, return the value of a .

Again, this imperative solution can be transformed almost mechanically into a purely functional one. First we define a helper function *fibacc* taking two extra parameters:

```
fibacc :: Integer -> Integer -> Integer -> Integer
fibacc n a b
  | n == 0 = a
  | n >= 1 = fibacc (n - 1) b (a + b)
```

And then we redefine *fib* as an appropriate call to *fibacc*:

```
fib n = fibacc n 0 1
```

As we saw before with *revacc*, the function *fibacc* may be seen as solving a more general problem than computing Fibonacci numbers: *fibacc n a b* actually computes the *n*th entry of a *generalized Fibonacci sequence*, defined by the same recurrence but with initial values *a* and *b*. (For example, *fibacc n 2 1* is the *n*th “Lucas number”.)

We can observe that the new version indeed computes Fibonacci numbers much more efficiently:

```
*Main> fib n = fibacc n 0 1 -- redefine fib
*Main> fib 32
2178309
(0.00 secs, 82,288 bytes)
*Main> fib 100
354224848179261915075
(0.01 secs, 114,400 bytes)
*Main> fib 1000
43466557686937456435688527675040625802564660517371780402481729089 \
53655541794905189040387984007925516929592259308032263477520968962 \
32398733224711616429964409065331879382989696499285160037044761377 \
95166849228875
(0.01 secs, 637,736 bytes)
```

6 Trees

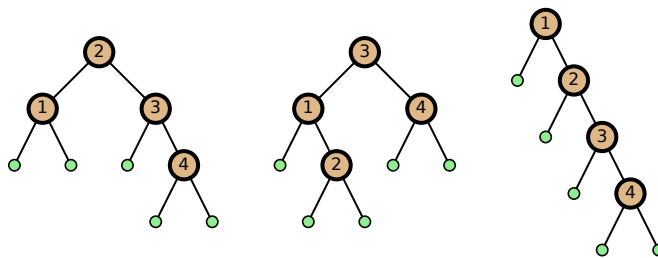
Trees provide another important example of a data type. In fact, there are many different kinds of “trees” in computer science and mathematics, but here for concreteness we will consider binary trees with unlabelled leaves and labelled binary nodes. In Haskell, these can be defined via the following recursive data type declaration:

```
data BinTree a = Leaf | Node a (BinTree a) (BinTree a)
  deriving (Show, Eq)
```

which, again to be clear, introduces the constructors *Leaf* and *Node* with the following types:

```
Leaf :: BinTree a
Node :: a -> BinTree a -> BinTree a -> BinTree a
```

So a binary tree is either a *Leaf*, containing no data, or a *Node*, which contains a value together with a left subtree and a right subtree. The binary trees



are represented as the following values, for example:

```
t1, t2, t3 :: BinTree Int
t1 = Node 2 (Node 1 Leaf Leaf) (Node 3 Leaf (Node 4 Leaf Leaf))
t2 = Node 3 (Node 1 Leaf (Node 2 Leaf Leaf)) (Node 4 Leaf Leaf)
t3 = Node 1 Leaf (Node 2 Leaf (Node 3 Leaf (Node 4 Leaf Leaf)))
```

Many different operations on binary trees can be naturally defined using pattern-matching and recursion:

(Compute the number of nodes in a binary tree.)

```
nodes :: BinTree a -> Int
nodes Leaf = 0
nodes (Node _ tL tR) = 1 + nodes tL + nodes tR
```

```
*Main> (nodes t1, nodes t2, nodes t3)
(4,4,4)
```

(Compute the number of leaves in a binary tree.)

```
leaves :: BinTree a -> Int
leaves Leaf = 1
leaves (Node _ tL tR) = leaves tL + leaves tR
```

```
*Main> (leaves t1, leaves t2, leaves t3)
(5,5,5)
```

(Compute the height of a binary tree.)

```
height :: BinTree a -> Int
height Leaf = 0
height (Node _ tL tR) = 1 + max (height tL) (height tR)
```

```
*Main> (height t1, height t2, height t3)
(3,3,4)
```

(Compute the mirror image of a binary tree.)

```
mirror :: BinTree a → BinTree a
mirror Leaf = Leaf
mirror (Node x tL tR) = Node x (mirror tR) (mirror tL)
```

```
*Main> mirror t2
Node 3 (Node 4 Leaf Leaf) (Node 1 (Node 2 Leaf Leaf) Leaf)
```

Just like lists, trees also admit a useful principle of structural induction. The principle of **structural induction over binary trees** says that for any property $P(t)$ of a binary tree t (with nodes labelled by values of type a), if we can show both

1. $P(\text{Leaf})$ holds; and
2. for any element x and pair of trees tL and tR , $P(tL)$ and $P(tR)$ implies $P(\text{Node } x \ tL \ tR)$

then $P(t)$ holds for all binary trees t . Again, to be a bit more precise we should state the principle of structural induction for a fixed type of values stored at the nodes of the binary tree. Then the inductive step becomes

2. for any element $x :: a$ and pair of trees $tL :: \text{BinTree } a$ and $tR :: \text{BinTree } a$, $P(tL)$ and $P(tR)$ implies $P(\text{Node } x \ tL \ tR)$

with conclusion that $P(t)$ holds for all binary trees $t :: \text{BinTree } a$.

Exercise 6.1. Define a data type of unary-binary trees with labelled leaves. (“Unary-binary” means that every internal node has 1 or 2 children.)

Exercise 6.2. Prove that $\text{height } (\text{mirror } t) = \text{height } t$ for all binary trees t .

Exercise 6.3. Write a function $\text{isBST} :: \text{Ord } a \Rightarrow \text{BinTree } a \rightarrow \text{Bool}$ which determines whether a binary tree is a binary search tree, by first writing a function $\text{inorder} :: \text{BinTree } a \rightarrow [a]$ that returns the list of node labels of a binary tree in an inorder traversal, as well as a function $\text{isSorted} :: \text{Ord } a \Rightarrow [a] \rightarrow \text{Bool}$ that determines whether a list is sorted, and composing them.