

FUNCTIONAL PEARL

Certified, total serialisers with an application to Huffman encoding

RALF HINZE 

RPTU Kaiserslautern-Landau, Germany
(e-mail: ralf.hinze@cs.rptu.de)

1 Introduction

The other day, I was assembling lecture material for a course on Agda. Pursuing an application-driven approach, I was looking for correctness proofs of popular algorithms. One of my all-time favourites is Huffman data compression (Huffman, 1952). Even though it is probably safe to assume that you are familiar with this algorithmic gem, a brief reminder of the essential idea may not be amiss.

A common representation of a textual document on a computer uses the ASCII character encoding scheme. The scheme uses seven or eight bits to represent a single character. The idea behind Huffman compression is to leverage the fact that some characters appear more frequently than others. Huffman encoding moves away from the fixed-length encoding of ASCII to a variable-length encoding, in which more frequently used characters have a shorter bit encoding than rarer ones. As an example, consider the text

`eight_in_the_evening` (1.1)

and observe that the letter e occurs four times, whereas v occurs only once. The code table shown below gives a possible encoding for the eight different letters occurring in the text.

<code>┌ 111</code>	<code>e 01</code>	<code>g 1101</code>	<code>h 000</code>	(1.2)
<code>i 100</code>	<code>n 101</code>	<code>t 001</code>	<code>v 1100</code>	

The more frequent a character, the shorter the code. Using the encoding above, the text (1.1) is compressed to the following bit string:

`01100110100000111110010111100100001111011100011011001011101`

The compressed text only occupies 59 bits, comparing favourably to the 160 bits of the standard ASCII encoding.

There is an impressive proof of correctness and optimality given in Coq (Théry, 2004). Unfortunately, the formalisation is slightly too elaborate to be presented in a single

lecture as it consists of more than 6K lines of code. The excellent text book on “Verified Functional Programming” (Stump, 2016) contains a case study on “Huffman Encoding and Decoding.” Alas, even though framed in Agda, the chapter treats neither correctness nor optimality.

So I decided to give it a go myself. Compared to the formalisation in Coq, my goals were more modest and more ambitious at the same time. More modest, as I decided to concentrate on correctness: decoding an encoded text gives back the original string. More ambitious, as I was aiming for a compositional design: Huffman compression should work seamlessly with standard bit serialisers.

Setting up a library for de/serialisers turned out to be a challenge in its own right. Indeed, we dedicate half of the pearl to its design and implementation. From an implementation perspective, a serialiser is like a pretty *printer*, except that its output, a bit string, is not particularly pretty. Likewise, a deserialiser can be seen as a simple recursive-descent *parser*. A parser is a partial function: it fails with an error message if it does not recognise the input; a recursive-descent parser may also fail to terminate. The library addresses the two sources of partiality in different ways. The first type of failure is unavoidable: if a deserialiser is applied to faulty data, it may unexpectedly exhaust the input. However, a deserialiser applied to the output of its mate should return the original value, like a Huffman decoder. To establish correctness, the library offers a set of *proof combinators* that allow the user to link parsers to printers. The second type of failure, non-termination, should ideally be avoided altogether: we aim for *total parsing combinators*, parsers that come with a termination guarantee. It will be interesting to see how this goal can be achieved in the spirit of combinator libraries, hiding the nitty-gritty details from the user. Now, without further ado, let’s see the core library in action.

The complete Agda code can be found in the accompanying material. Appendix 1 lists some basic types and functions.

2 Sneak preview

To illustrate the trinity of *printer*, *parser*, and *proof combinators*, we use the type of bushes, defined below, as a running example. (A similar datatype is required later for Huffman encoding, so our efforts are not wasted.)

```
data Bush : Set where
  Leaf : Bush
  Fork : Bush → Bush → Bush
```

For each type A whose elements we wish to serialise, we define a module called *Codec-A*.

```
module Codec-Bush where
```

Printing combinators. A bit string or code is generated using ϵ , $_ _ _$, 0, and I. The encoder for bushes employs all four combinators.

```

encode : Bush → Code
encode (Leaf)      = 0 · ε
encode (Fork t u) = I · encode t · encode u · ε

```

An encoder for elements of type A is typically defined by *structural* recursion on A , *encode* is no exception. Each constructor is encoded by a suitable number of bits followed by the encodings of its arguments, if any. For the example at hand, a single bit suffices. In general, we require $\lceil \log_2 n \rceil$ bits if the datatype comprises n constructors. (Of course, we may also apply Huffman’s idea to the encoding of constructors using variable-length instead of fixed-length bit codes.) Concatenation $_ \cdot _$ is associative with ϵ as its neutral element, so the trailing ϵ s seem a bit odd—odd, but convenient as we shall see shortly.

Parsing combinators. Decoders can be most conveniently defined using *do*-notation as the type of decoders has the structure of a monad.

```

{-# TERMINATING #-}
decode' = 0: return Leaf
          I: do t ← decode'; u ← decode'; return (Fork t u)

```

The decoder performs a case analysis on the first bit of the implicit input bit string, and then, it decodes the arguments, if any. In the best, or, worst tradition of Haskell, *decode'* is defined recursively. Unfortunately, Agda is not able to establish its termination. All is safe, however: the case analysis $0: \dots I: \dots$ consumes a bit, so both recursive calls process strictly smaller bit strings. The pragma preceding the definition communicates our findings to Agda.

We have come across a well-known problem in parsing: a recursive-descent parser may fail to terminate if the underlying grammar is left-recursive. Consequently, when decoders are defined recursively, we need to make sure that each recursive call is guarded by at least one case combinator. Are we happy to accept the potential threat of non-termination? No! We should be able to do better than that. And indeed, the parser library features a grown-up version of case analysis that allows for safe recursive calls.

```

decode : Decoder Bush
decode = 0-rec: (λ bush → return Leaf)
          I-rec: (λ bush → do t ← bush; u ← bush; return (Fork t u))

```

The combinator $0\text{-}rec\text{:}_I\text{-}rec\text{:}_$ can be seen as a recursion principle for bit strings; it passes the ability to issue recursive calls to each branch. All is well now: Agda happily accepts the definition; *decode* terminates on all inputs.

Proof combinators. We have created two independent artefacts: an encoder and a decoder. It remains to relate the two: we require that decoding an encoded value gives back the original value. To capture this property, we make use of a ternary relation, relating bit strings, values, and decoders: $s \langle a \rangle p$ holds iff p “applied” to s returns a . The correctness criterion then reads

$$\forall (a : A) \rightarrow \text{encode } a \langle a \rangle \text{decode.} \quad (2.1)$$

The library offers six proof combinators for discharging the proof obligation, presented as proof rules below. (Each rule *is* actually the type of the proof combinator—“propositions as types” at work.)

$$\frac{}{\epsilon \langle a \rangle \text{ return } a} \epsilon\text{-axiom} \quad \frac{s \langle a \rangle p \quad t \langle b \rangle q \quad a}{s \cdot t \langle b \rangle p \gg\approx q} \cdot\text{-rule}$$

The axiom relates the unit of the monoid to the “unit” of the monad; the rule relates the multiplication of the monoid to the “multiplication” of the monad.

$$\frac{s \langle a \rangle p}{\text{O} \cdot s \langle a \rangle (\text{O}; p \text{ I}; q)} \text{O-rule} \quad \frac{t \langle a \rangle q}{\text{I} \cdot t \langle a \rangle (\text{O}; p \text{ I}; q)} \text{I-rule}$$

The case parser consumes a bit and then acts as one of the branches. The recursive variant additionally passes “itself” to the branch.

$$\frac{s \langle a \rangle P \quad (\text{O-rec}: P \text{ I-rec}: Q)}{\text{O} \cdot s \langle a \rangle (\text{O-rec}: P \text{ I-rec}: Q)} \text{O-Rule} \quad \frac{t \langle a \rangle Q \quad (\text{O-rec}: P \text{ I-rec}: Q)}{\text{I} \cdot s \langle a \rangle (\text{O-rec}: P \text{ I-rec}: Q)} \text{I-Rule}$$

The proof of correctness then typically proceeds by structural recursion on the type of the to-be-encoded values, here the type of bushes.

$$\begin{aligned} \text{correct}' (\text{Leaf}) &= \text{O-Rule } \epsilon\text{-axiom} \\ \text{correct}' (\text{Fork } t \ u) &= \text{I-Rule } (\cdot\text{-rule } (\text{correct}' \ t) \ (\cdot\text{-rule } (\text{correct}' \ u) \ \epsilon\text{-axiom})) \end{aligned}$$

Two remarks are in order.

The definition of *encode* features trailing ϵ s so that all three artefacts, *printers*, *parsers*, and *proofs*, exhibit *exactly* the same structure. Any deviation incurs additional proof effort as monoid and monad laws do not hold on the nose, for example, $s \cdot \epsilon$ is not definitionally equal to s , only propositionally.

Each proof combinator takes a couple of implicit arguments—the variables appearing in the proof rules. Agda happily infers these arguments from the data provided, with one exception: the $\cdot\text{-rule}$ needs p and q as *explicit* arguments. Consequently, the proof is actually more verbose.

$$\begin{aligned} \text{correct} &: \forall (t : \text{Bush}) \rightarrow \text{encode } t \langle t \rangle \text{ decode} \\ \text{correct} (\text{Leaf}) &= \text{O-Rule } \epsilon\text{-axiom} \\ \text{correct} (\text{Fork } t \ u) &= \\ &\text{I-Rule } (\cdot\text{-rule } \text{decode } (\lambda \ t' \rightarrow \text{do } u' \leftarrow \text{decode}; \text{return } (\text{Fork } t' \ u')) \ (\text{correct } t) \\ &\quad (\cdot\text{-rule } \text{decode } (\lambda \ u' \rightarrow \text{return } (\text{Fork } t \ u')) \ (\text{correct } u) \\ &\quad (\epsilon\text{-axiom})) \end{aligned}$$

On the positive side, the additional arguments show very clearly how the parsing process proceeds, creating a *Fork* in two steps.

Codecs. For convenience, the three entities are wrapped up in a record.

$$\begin{aligned} \text{record } \text{Codec } (A : \text{Set}) &: \text{Set } \text{where} \\ \text{field } \text{encode} &: A \rightarrow \text{Code} \\ \text{field } \text{decode} &: \text{Decoder } A \\ \text{field } \text{correct} &: \forall (a : A) \rightarrow \text{encode } a \langle a \rangle \text{ decode} \end{aligned}$$

```

data Natural : Set where
  zero : Natural
  succ : Natural → Natural
module Codec-Natural where
  encode : Natural → Code
  encode zero = 0 · ε
  encode (succ n) = I · encode n · ε
  decode : Decoder Natural
  decode = 0-rec: (λ natural → return zero)
           I-rec: (λ natural → do n ← natural; return (succ n))
  correct : ∀ (n : Natural) → encode n ⟨ n ⟩ decode
  correct zero = 0-Rule ε-axiom
  correct (succ n) = I-Rule (·-rule decode (λ n' → return (succ n'))) (correct n)
                    (ε-axiom))

Natural : Codec Natural
Natural = record {Codec-Natural}

```

Fig. 1. A codec for natural numbers.

For each serialisable type $A : \text{Set}$ we define a codec $\mathbb{A} : \text{Codec } A$. You may want to view *Codec* as a property of a type, or as a type class and *Codec A* as an instance dictionary.

```

Bush : Codec Bush
Bush = record {Codec-Bush}

```

The definition makes use of a nifty Agda feature: the syntax `record {M} : R` constructs a record of type R , using suitable components from the module M .

For reference, Figure 1 shows a codec for a second example, the naturals. This completes the description of the interface. Next, we turn to the implementation, starting with parsing combinators as their design dictates everything else.

3 Parsing combinators

The goal is clear: our parsing combinators should come with a termination guarantee as Agda is not able to establish termination of recursively defined decoders. To understand the cure, let us delve a bit deeper into the problem.

First, we introduce the type of bit strings. A bit string is either empty, written `[]`, or a binary digit, `0` or `1`, followed by a bit string.

```

data Bits : Set where
  [] : Bits
  0 1 : Bits → Bits

```

The problem. A decoder is a deterministic parser, combining state, bit strings of type *Bits*, and partiality, the *Maybe* monad (see Appendix 1). Not using any special combinators, a decoder for bushes can be implemented as follows.

```

{-# TERMINATING #-}
decode : Bits → Maybe (Bush × Bits)
decode []      = undefined
decode (0 bs) = return (Leaf , bs)
decode (1 bs) = do (t , bs') ← decode bs
                  (u , bs'') ← decode bs'
                  return (Fork t u , bs'')

```

Observe how the state, the bit string, is threaded through the program. Putting the termination spectacles on, the first recursive call is benign as bs is an immediate substring of $1\ bs$. The second recursive call is, however, problematic as $decode$ is applied to the string returned from the first call. Agda has no clue that bs' is actually smaller than bs .

The recursion pattern is an instance of *course-of-values recursion* or *strong recursion* where a function recurses on all smaller values, for example, where $f(n + 1)$ is computed from $f(0), \dots, f(n)$. Perhaps a short detour is worthwhile to remind ourselves of this induction and recursion principle.

Detour: strong versus structural recursion. Say, $P : \mathbb{N} \rightarrow Set$ is a property of the naturals, established by strong induction. Strong induction can be reduced to structural induction by proving a stronger property: $Q\ m = P\ 0 \wedge \dots \wedge P\ m$, or, avoiding ellipsis $Q\ m = \forall n \rightarrow m \succ n \rightarrow P\ n$.

Say, $f : \mathbb{N} \rightarrow A$ is a function from the naturals, defined by strong recursion. Strong recursion can be reduced to structural recursion by defining a stronger function: $g\ m : \forall n \rightarrow m \succ n \rightarrow A$. Consider as an example,

```

f : (n : ℕ) → ℕ
f (zero)    = 1
f (succ n) = sf n where
  sf : (i : ℕ) → ℕ
  sf (zero)  = f (zero)
  sf (succ i) = f (succ i) + sf i

```

This is the simplest, non-trivial example I could think of: $f(n + 1)$ is given by the sum of all smaller values: $f(0) + \dots + f(n)$. (Do you see what f computes?) Its stronger counterpart is given by

```

g : (m : ℕ) → ∀ n → m ≻ n → ℕ
g m (zero) (zero) = 1
g (succ m) (succ n) (succ m ≻ n) = sg n m ≻ n where
  sg : (i : ℕ) → m ≻ i → ℕ
  sg (zero) below = g m (zero) below
  sg (succ i) below = g m (succ i) below + sg i (≻-transitive below succ-n ≻ n)

```

If we remove the first and the third argument of g and the second argument of the helper function sg , we get back the original definition: $g\ m\ n\ m \succ n \equiv f\ n$ for all natural numbers m and n with $m \succ n$: $m \succ n$. In particular, f can be reduced to its counterpart: $f\ n = g\ n\ n\ \succ$ -reflexive. Observe that g is defined by *structural* recursion on the first argument; the helper function maintains the invariant that the “actual” argument i is at most m .

The solution. Returning to our application, we would like to define decoders by strong induction over the length of bit strings. As a preparatory step, we replace the original definition of *Bits* by an *indexed type* that records the length.

```
data Bits : ℕ → Set where
  [] : Bits (zero)
  0 1 : Bits n → Bits (succ n)
```

Compared to the number-theoretic example above, there is one further complication: we need to relate the length of the residual bit string *returned* by a parser to the length of its input. Exact numbers are not called for; it suffices to know that the length of the output is at most the length of the input. To this end, we introduce $Bits_{\preceq} n$, the type of bit strings of length below a given upper bound n .

```
record Bits_{\preceq} (n : ℕ) : Set where
  constructor _such-that_
  field {size} : ℕ
       bits : Bits size
       below : n \succeq size
```

The record features three fields, the first of which is hidden, indicated by curly braces. If the decoder transmogrifies, say, $bs : Bits\ i$ to $bs' : Bits\ j$ with $i \succeq j$, then the residual output of type $Bits_{\preceq}\ i$ is given by $bs'\ \text{such-that}\ i \succeq j$.

Finally, we have all the gadgets in place, to define a decoder that is approved by Agda's termination checker. As to be expected, its type is more involved. In a sense, *decode* says: I can decode a bit string of size i , but I am happy to accept any bit string of smaller size; I promise to return a bush and a string, the length of which is at most the size of the string given to me.

```
decode : (i : ℕ) → ∀ {j} → i \succeq j → Bits\ j → Maybe (Bush × Bits_{\preceq}\ j)
decode i (zero) [] = undefined
decode (succ i) (succ i \succeq j) (0 bs) = return (Leaf , bs such-that succ-n \succeq n)
decode (succ i) (succ i \succeq j) (1 bs) =
  do (t , bs' such-that j \succeq k) ← decode i i \succeq j bs
     (u , bs'' such-that k \succeq m) ← decode i (\succeq-transitive i \succeq j j \succeq k) bs'
     return (Fork t u , bs'' such-that \succeq-transitive succ-n \succeq n (\succeq-transitive j \succeq k k \succeq m))
```

Voilà, *decode* is defined by *structural recursion* on the natural number i . To reduce clutter, the actual size of the bit string is passed implicitly. In the recursive case, we have $bs : Bits\ j$, $bs' : Bits\ k$, and $bs'' : Bits\ m$ with $j \succeq k$ and $k \succeq m$. All that remains to be done is to hide the plumbing of state and proofs.

Hiding the plumbing. A decoder is a function from bit strings to an optional pair of things and bit strings, promising not to increase the length of the string.

```
IDecoder : Set → (ℕ → Set)
IDecoder A i = ∀ {j} → i \succeq j → Bits\ j → Maybe (A × Bits_{\preceq}\ j)
```

Both *IDecoder A* and $Bits_{\preceq}$ are indexed by natural numbers, and they are functions of type $\mathbb{N} \rightarrow Set$. But there is more to them, they are actually functors: *IDecoder A* is a

covariant functor from the preorder $(\mathbb{N}, \succcurlyeq)$, viewed as a category, to the category of sets and total functions. Its action on arrows is defined

$$\begin{aligned} \text{lower} &: (i \succcurlyeq j) \rightarrow (\text{IDecoder } A \ i \rightarrow \text{IDecoder } A \ j) \\ \text{lower } i \succcurlyeq j &= \lambda j \succcurlyeq k \ bs \rightarrow p \ (\succcurlyeq\text{-transitive } i \succcurlyeq j \ j \succcurlyeq k) \ bs \end{aligned}$$

The function *lower* makes precise the idea that a decoder can be applied to any shorter bit string. By contrast, *Bits* \preccurlyeq is a contravariant functor: we may relax the upper bound.

$$\begin{aligned} \text{raise} &: (i \succcurlyeq j) \rightarrow (\text{Bits}\preccurlyeq \ j \rightarrow \text{Bits}\preccurlyeq \ i) \\ \text{raise } i \succcurlyeq j \ (bs \ \text{such-that } j \succcurlyeq k) &= bs \ \text{such-that } \succcurlyeq\text{-transitive } i \succcurlyeq j \ j \succcurlyeq k \end{aligned}$$

Turning to the implementation of the monad operations, the definition of *return* is standard, except for the proof that the length is unchanged.

$$\begin{aligned} \text{return} &: A \rightarrow \text{IDecoder } A \ i \\ \text{return } a &= \lambda i \succcurlyeq j \ bs \rightarrow \text{Maybe.return } (a \ , \ bs \ \text{such-that } \succcurlyeq\text{-reflexive}) \end{aligned}$$

The implementation of monadic bind is more interesting.

$$\begin{aligned} \text{relax} &: i \succcurlyeq j \rightarrow \text{Maybe } (A \times \text{Bits}\preccurlyeq \ j) \rightarrow \text{Maybe } (A \times \text{Bits}\preccurlyeq \ i) \\ \text{relax } i \succcurlyeq j &= \text{Maybe.map } (\lambda (a \ , \ bs) \rightarrow (a \ , \ \text{raise } i \succcurlyeq j \ bs)) \\ _ \bowtie _ &: \text{IDecoder } A \ i \rightarrow (A \rightarrow \text{IDecoder } B \ i) \rightarrow \text{IDecoder } B \ i \\ p \bowtie q &= \lambda i \succcurlyeq j \ bs \rightarrow \\ & \quad p \ i \succcurlyeq j \ bs \ \text{Maybe.}\bowtie \ \lambda (a \ , \ bs' \ \text{such-that } j \succcurlyeq k) \rightarrow \\ & \quad \text{relax } j \succcurlyeq k \ (\text{lower } i \succcurlyeq j \ (q \ a) \ j \succcurlyeq k \ bs') \end{aligned}$$

The value *a* returned by *p* is passed to the continuation *q*. Since we wish to apply *q a* to the residual input *bs'*, we need to lower its index. This move has a knock-on effect: the upper bound of *q*'s result is consequently too low and needs to be relaxed in a final step.

The case combinator dispatches on the first bit of the input bit string. It fails if the input is exhausted. This is *the* single source of undefinedness—the correctness proofs guarantee that this failure never happens for “legal” inputs.

$$\begin{aligned} \text{O_I_} &: (p \ q : \text{IDecoder } A \ i) \rightarrow \text{IDecoder } A \ i \\ (\text{O} : p \ \text{I} : q) \ (\text{zero}) \ \square &= \text{undefined} \\ (\text{O} : p \ \text{I} : q) \ (\text{succ } i \succcurlyeq j) \ (\mathbf{0} \ bs) &= \text{relax } \text{succ-}n \succcurlyeq n \ ((\text{lower } \text{succ-}n \succcurlyeq n \ p) \ i \succcurlyeq j \ bs) \\ (\text{O} : p \ \text{I} : q) \ (\text{succ } i \succcurlyeq j) \ (\mathbf{1} \ bs) &= \text{relax } \text{succ-}n \succcurlyeq n \ ((\text{lower } \text{succ-}n \succcurlyeq n \ q) \ i \succcurlyeq j \ bs) \end{aligned}$$

Again, we first lower the index and then adjust the upper bound of the result. The proof *succ- $n \succcurlyeq n$* : *succ* *n* \succcurlyeq *n* records the fact that *bs* is one element shorter compared to $\mathbf{0} \ bs$ or $\mathbf{1} \ bs$.

It is time to bring in the harvest! The implementation of recursive case analysis merits careful study.

$$\begin{aligned} \text{O-rec_I-rec_} &: (P \ Q : \forall \{i\} \rightarrow \text{IDecoder } A \ i \rightarrow \text{IDecoder } A \ i) \\ & \quad \rightarrow (\forall \{i\} \rightarrow \text{IDecoder } A \ i) \\ \text{O-rec} : P \ \text{I-rec} : Q &= \text{recurse } _ \ \text{where} \\ \text{recurse} &: \forall i \rightarrow \text{IDecoder } _ \ i \\ \text{recurse } i & \quad (\text{zero}) \ \square = \text{undefined} \end{aligned}$$

$recurse (succ\ i) (succ\ i \succ j) (\mathbf{0}\ bs) = relax\ succ\text{-}n \succ n (P (recurse\ i) i \succ j\ bs)$
 $recurse (succ\ i) (succ\ i \succ j) (\mathbf{1}\ bs) = relax\ succ\text{-}n \succ n (Q (recurse\ i) i \succ j\ bs)$

The recursor is defined by structural induction on its first argument. The entire set-up—the use of indexed types, keeping track of bounds—serves the sole purpose of enabling this definition. Observe that the type of branches involves a local quantifier: the branches work for any index, in particular, they happily accept the recursive call at index $i - 1$.

Ultimately, the four combinators enable the user to construct decoders that work for any index.

$Decoder : Set \rightarrow Set$
 $Decoder\ A = \forall \{i\} \rightarrow IDecoder\ A\ i$

4 Pretty-printing combinators

For reasons of efficiency, encoders use John Hughes’ representation of lists (Hughes, 1986), specialised to bit strings.

$Code : Set$
 $Code = \forall \{i\} \rightarrow Bits\ i \rightarrow Bits \succ i$

Encoders are dual to decoders. While decoders promise not to increase the size of the bit string, encoders guarantee not to decrease its size. To capture this invariant, we need the dual of $Bits \preceq$, bit strings of size *above* a given bound:

record $Bits \succ (n : \mathbb{N}) : Set$ **where**
constructor $_such\text{-}that_$
field $\{size\} : \mathbb{N}$
 $bits : Bits\ size$
 $above : n \preceq size$

Composition of functions *and* proofs then serves as concatenation with the identity as its neutral element.

$\epsilon : Code$
 $\epsilon = \lambda bs \rightarrow bs\ such\text{-}that\ \preceq\text{-}reflexive$
 $_._ : Code \rightarrow Code \rightarrow Code$
 $f \cdot g = \lambda bs \rightarrow let\ bs' \ such\text{-}that\ i \preceq j = g\ bs\ in$
 $let\ bs'' \ such\text{-}that\ j \preceq k = f\ bs'\ in$
 $bs'' \ such\text{-}that\ \preceq\text{-}transitive\ i \preceq j \preceq k$

The combinators $\mathbf{0}$ and \mathbf{I} increase the size by one.

$\mathbf{0}\ \mathbf{I} : Code$
 $\mathbf{0} = \lambda bs \rightarrow \mathbf{0}\ bs\ such\text{-}that\ n \preceq succ\text{-}n$
 $\mathbf{I} = \lambda bs \rightarrow \mathbf{1}\ bs\ such\text{-}that\ n \preceq succ\text{-}n$

As an amusing aside, the proof that $Code$ is a monoid makes use of the fact that (\mathbb{N}, \preceq) is a category: $\preceq\text{-}transitive$ is associative with $\preceq\text{-}reflexive$ as its neutral element.

5 Proof combinators

Decoders are partial functions. Given an arbitrary bit string, a decoder might fail, unexpectedly exhausting the input. The correctness criterion (2.1) guarantees that this does not happen if a decoder is applied to an encoding. This guarantee is, in a sense, the sole purpose of the whole exercise. The ternary relation $_(_) _$ underlying (2.1) relates a code, an element, and a decoder.

$$\begin{aligned} _(_) _ &: \text{Code} \rightarrow A \rightarrow \text{Decoder } A \rightarrow \text{Set} \\ \text{code } \langle a \rangle \text{ decode} &= \\ \forall k \rightarrow (bs : \text{Bits } k) \rightarrow \forall i \rightarrow (i \succcurlyeq j : i \succcurlyeq \text{size } (\text{code } bs)) \rightarrow \\ &\text{decode } i \succcurlyeq j (\text{bits } (\text{code } bs)) \equiv \text{Just } (a, bs \text{ such-that above } (\text{code } bs)) \end{aligned}$$

The function *code* prepends some bits to *bs*; the decoder removes this prefix returning the value *a* and *bs*. This relation must hold for any upper bound *i*.

The implementation of the proof combinators is fairly straightforward. The code is, however, not too instructive, so we content ourselves with two illustrative examples. (The definition of the remaining combinators can be found in the accompanying material.)

The axiom *ε-axiom* holds definitionally: Agda can show its correctness solely by unfolding definitions, witnessed by the use of $\equiv\text{-reflexive} : a \equiv a$.

$$\begin{aligned} \epsilon\text{-axiom} &: \epsilon \langle a \rangle \text{ return } a \\ \epsilon\text{-axiom } i \text{ bs } j \succcurlyeq k &= \equiv\text{-reflexive} \end{aligned}$$

To establish the *0-rule*,

$$\begin{aligned} \text{0-rule} &: (s \langle a \rangle p) \rightarrow (0 \cdot s \langle a \rangle (0 : p \text{ I} : q)) \\ \text{0-rule premise } i \text{ bs } (\text{succ } j) (\text{succ } j \succcurlyeq k) &= \\ \equiv\text{-congruent } (\text{relax succ-}n \succcurlyeq n) (\text{premise } i \text{ bs } (\text{succ } j) (\succcurlyeq\text{transitive succ-}n \succcurlyeq n \text{ } j \succcurlyeq k)) & \end{aligned}$$

we apply the premise to the argument of *relax succ- $n \succcurlyeq n$* .

Now that the core library is in place, we can finally deal with the motivating example, Huffman Encoding. Its implementation turns out to be a nice exercise in datatype-generic programming.

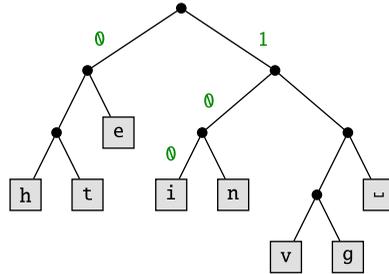
6 Application: Huffman encoding

Before diving into the mechanics of Huffman compression, let us state the guiding principles of the design: (1) we would like to seamlessly combine statistical compression with standard serialisers; (2) we aim to avoid irrelevant detail, such as: Is the code unambiguous? Does the alphabet contain at least two characters? Huffman *encoding* typically relies on a dictionary, mapping characters to codes; Huffman *decoding* makes use of a code tree. Does the dictionary contain a code for each to-be-encoded character? Is the code tree actually related to the dictionary? And so forth, and so forth ... In other words, we are aiming for a clean, lean interface.

We proceed in three steps, defining a codec for a single character, for a sequence of characters, and for Huffman code trees.

6.1 Encoding and decoding characters

The central property of variable-length encodings is that no code is a prefix of any other code. This property is guaranteed by construction if a code tree is used instead of a code table. The code tree corresponding to Table (1.2) is depicted below.



The character ‘i’, for example, is encoded by $1 (0 (0 []))$. So, as a first requirement, we assume that a code tree is given to us. A code tree is either a singleton, written $[c]$, or a binary node $t \wedge u$, where t and u are code trees (the bushes of Section 2 are code trees without codes).

data *Tree* (*Alphabet* : *Set*) : *Set* **where**
 $[c]$: *Alphabet* \rightarrow *Tree Alphabet*
 $t \wedge u$: *Tree Alphabet* \rightarrow *Tree Alphabet* \rightarrow *Tree Alphabet*

The code tree above is represented by

$code = ((['h'] \wedge ['t']) \wedge ['e']) \wedge ((['i'] \wedge ['n']) \wedge ((['v'] \wedge ['g']) \wedge ['␣']))$

Using this encoding, we can compress the text (1.1), but not **genius**, as neither u nor s have an encoding. This motivates the second requirement: we need evidence that the to-be-compressed character is contained in the code tree.

data $_{\in}$ (c : *Alphabet*) : *Tree Alphabet* \rightarrow *Set* **where**
here : $c \in [c]$
left : $c \in t \rightarrow c \in (t \wedge u)$
right : $c \in u \rightarrow c \in (t \wedge u)$

The evidence $c \in t$ can be seen as a path into the tree. For example,

i-am-in : ‘i’ \in *code*
i-am-in = *right* (*left* (*left here*))

certifies that ‘i’ is in the code tree *code*: starting at the root, it can be found going to the right and then twice to the left.

Perhaps surprisingly, no further requirements are necessary: our Huffman compressor takes a pair such as (‘i’, *i-am-in*), consisting of a single character and a proof that this character is contained in the code tree. Observe that the *type* of the second component depends on the *value* of the first component (in the “propositions as types” paradigm a proposition such as $c \in t$ is a type). In other words, the data are a *dependent pair*, an element of a Σ type (*Alphabet* is an implicit type argument, indicated by curly braces).

$$\begin{aligned} In &: \{Alphabet : Set\} \rightarrow Tree\ Alphabet \rightarrow Set \\ In\ \{Alphabet\}\ t &= \Sigma\ Alphabet\ (\lambda\ c \rightarrow c \in t) \end{aligned}$$

Elements of the type $In\ t$ can be seen as *self-certifying data*: data that can be readily serialised.

module *Codec-In where*

$$\begin{aligned} encode &: \forall \{t : Tree\ Alphabet\} \rightarrow In\ t \rightarrow Code \\ encode\ (c\ ,\ here) &= \epsilon \\ encode\ (c\ ,\ left\ p) &= 0 \cdot encode\ (c\ ,\ p) \cdot \epsilon \\ encode\ (c\ ,\ right\ p) &= 1 \cdot encode\ (c\ ,\ p) \cdot \epsilon \end{aligned}$$

The encoder turns the evidence, the path into the code tree, into a bit string, ignoring both the code tree and the character!

The decoder has to resurrect the path. For this to work out, the code tree is required.

$$\begin{aligned} decode &: \forall (t : Tree\ Alphabet) \rightarrow Decoder\ (In\ t) \\ decode\ [c] &= return\ (c\ ,\ here) \\ decode\ (t\ \wedge\ u) &= 0: (\mathbf{do}\ (c\ ,\ p) \leftarrow decode\ t; return\ (c\ ,\ left\ p)) \\ &\quad 1: (\mathbf{do}\ (c\ ,\ p) \leftarrow decode\ u; return\ (c\ ,\ right\ p)) \end{aligned}$$

The bits steer the search, the decoder stops if a leaf is reached. Note that we do *not* use the recursive case combinator; *decode* is recursively defined by induction on the structure of the code tree.

The proof of correctness replicates the structure of encoder and decoder.

$$\begin{aligned} correct &: \forall \{t : Tree\ Alphabet\} \rightarrow \forall (c : In\ t) \rightarrow encode\ c\ \langle\ c\ \rangle\ decode\ t \\ correct\ (c\ ,\ here) &= \epsilon\text{-axiom} \\ correct\ (c\ ,\ left\ p) &= 0\text{-rule}\ (\text{-rule}\ (decode\ _) (\lambda\ (c'\ ,\ p') \rightarrow \\ &\quad return\ (c'\ ,\ left\ p'))) (correct\ (c\ ,\ p))\ \epsilon\text{-axiom} \\ correct\ (c\ ,\ right\ p) &= 1\text{-rule}\ (\text{-rule}\ (decode\ _) (\lambda\ (c'\ ,\ p') \rightarrow \\ &\quad return\ (c'\ ,\ right\ p'))) (correct\ (c\ ,\ p))\ \epsilon\text{-axiom} \end{aligned}$$

For brevity, we henceforth omit the correctness proof—in all the examples, it can be mechanically derived from the encoder.

Provided with a code tree, the Huffman serialiser happily deals with a single self-certified character.

$$\llbracket _ \rrbracket : \forall (t : Tree\ Alphabet) \rightarrow Codec\ (In\ t)$$

Before we tackle the encoding of strings, it is worthwhile recording what we do *not* need to assume. We do not require that the character has a unique code; it may appear multiple times in the code tree. (Perhaps, bit strings are used to exchange secret messages. Using different codes for the same letter may slightly improve security.) Neither do we require that each element of *Alphabet* appears in the code tree. (If the underlying alphabet is large—the current Unicode standard defines 144,697 characters—then a text will rarely include each character.) Finally, the compressor works for arbitrary code trees, including trees that consist of a single leaf. These extreme cases feature a fantastic compression rate as no bits are required. (Usually some care must be exercised to avoid sending the decoder

into an infinite loop. The decoder presented in Stump (2016) exercises this care but fails with an error message instead!)

6.2 Encoding and decoding sequences

Let us be slightly more ambitious by not only considering strings, sequences of characters, but sequences of arbitrary encodable values. There are at least two options for defining a suitable container type; we discuss each in turn.

The obvious choice is to use standard lists, given by the definition below.

```
data List (Elem : Set) : Set where
  [] : List Elem
  _::_ : Elem → List Elem → List Elem
```

The type features two data constructors; to distinguish between empty and non-empty lists we need one bit.

```
module Codec-List {A : Set} (A : Codec A) where
  encode : List A → Code
  encode [] = 0 · ε
  encode (a :: as) = I · encode A a · encode as · ε
  decode : Decoder (List A)
  decode = 0-rec: (λ list → return [])
           I-rec: (λ list → do a ← decode A; as ← list; return (a :: as))
```

The codec for *List* takes a codec for *A* to a codec for *List A*.

```
ℓList : {A : Set} → Codec A → Codec (List A)
```

You may recognise the type pattern if you are familiar with Haskell's type classes or with datatype-generic programming.

Unfortunately, if we compose the list encoder with our Huffman encoder, then the compression factor suffers, as the code for each character is effectively prolonged by one bit. By contrast, the standard Huffman compressor, which operates on strings, simply concatenates the character codes; the decompressor keeps decoding until the input is exhausted. Alas, this approach is in conflict with our modularity requirements: such a codec cannot be used as a component of other codecs. (Imagine compressing tree-structured data such as HTML, using the statistical compressor for leaves, paragraphs of text.)

We can save the additional bit per list item if we keep track of the total number of items. In other words, we may want to use a vector, a length-indexed list, instead of a standard list.

```
data Vector (Elem : Set) : Natural → Set where
  ⟨⟩ : Vector Elem zero
  _→_ : Elem → Vector Elem n → Vector Elem (succ n)
```

The constructor names are chosen to reflect that a vector is really a nested pair. Like the type of lists, *Vector* features two data constructors. However, only one of these is available

for each given length—there is only one shape of vector. So the good news is that *no bits* are required to encode the constructors.

```

module Codec-Vector {A : Set} (A : Codec A) where
  encode : {i : Natural} → Vector A i → Code
  encode ⟨⟩ = ε
  encode (a , as) = encode A a · encode as · ε
  decode : (i : Natural) → Decoder (Vector A i)
  decode (zero) = return ⟨⟩
  decode (succ i) = do a ← decode A; as ← decode i; return (a , as)

```

The encoder and the decoder are driven by the length argument, which is passed implicitly to *encode* and explicitly to *decode*.

```

Vector : {A : Set} → Codec A → (i : Natural) → Codec (Vector A i)

```

Of course, for our application at hand we need to store the length alongside the vector, another case for dependent pairs:

```

List' : Set → Set
List' A = Σ Natural (Vector A)

```

Inspecting the type on the right-hand side, we note that we already have a codec for naturals and vectors, so we are left with defining a codec for Σ -types.

Encoding dependent pairs is not too hard, the main challenge is to get the type signature right. A dependent pair $(a , b) : \Sigma A B$ consists of a field $a : A$ and a field $b : B a$. This dependency is reflected in the types of the codecs:

```

module Codec-Σ {A : Set} (A : Codec A)
  {B : A → Set} (B : (a : A) → Codec (B a)) where
  encode : Σ A B → Code
  encode (a , b) = encode A a · encode (B a) b · ε
  decode : Decoder (Σ A B)
  decode = do a ← decode A; b ← decode (B a); return (a , b)

```

The implementation is identical to a codec for standard pairs, except that the first field is passed to the codec for the second field.

```

Sigma : {A : Set} → Codec A →
  {B : A → Set} → ((a : A) → Codec (B a)) →
  Codec (Σ A B)

```

Assembling the various bits and pieces,

```

List' A = Sigma Natural (Vector A)

```

is a space-efficient codec for sequences. Well, not quite. A moment's reflection reveals that we have not gained anything: $\text{List } A$ and $\text{List}' A$ produce bit strings of exactly the same length. (This tells you something about the relation between lists, naturals, and vectors.) We have chosen to represent naturals in unary. Consequently, the vector length n

is encoded by a bit string of length $n + 1$. That said, a cure is readily at hand: we simply switch to a binary or ternary representation.

We seize the opportunity and define a general codec for representation changers. Given an isomorphism $A \cong B$, we can turn a codec for A into a codec for B . The functions witnessing the isomorphism are called *to* and *fro*.

```

module Codec-Iso (iso : A ≅ B) (A : Codec A) where
  encode : B → Code
  encode b = encode A (fro iso b)

  decode : Decoder B
  decode = do a ← decode A; return (to iso a)

  Map : (A ≅ B) → (Codec A → Codec B)

```

We chose to replace unary numbers by ternary¹ numbers, obtaining the final implementation of the encoder for sequences:

```

List' : {A : Set} → Codec A → Codec (List' A)
List' A = Sigma (Map Ternary≅Natural Ternary) (Vector A)

```

Our library is taking shape. s

We can, for instance, confirm that the introductory example (1.1) is compressed to 59 bits—we assume below that *text* : *Vector (In code) 20* is provided from somewhere.

```

_ : size (encode (Vector (In code) 20) text []) ≡ 59
_ = reflexive

_ : size (encode (List' (In code)) (20 , text) []) ≡ 67
_ = reflexive

```

Additional 8 bits are needed to store the length.

6.3 Encoding and decoding code trees

All that remains to be done is to define a codec for Huffman trees. The type of code trees is a standard container type, so this is a routine exercise by now—the codec for bushes, see Section 2, serves nicely as a blueprint.

```

module Codec-Code-Tree {A : Set} (A : Codec A) where
  encode : Tree A → Code
  encode [ a ] = 0 · encode A a · ε
  encode (t λ u) = 1 · encode t · encode u · ε

  decode : Decoder (Tree A)
  decode = 0-rec: (λ tree → do a ← decode A; return [ a ])
           1-rec: (λ tree → do t ← tree; u ← tree; return (t λ u))

  Tree : {A : Set} → Codec A → Codec (Tree A)

```

¹ Why? Well, the details are actually irrelevant for the application at hand, but since you asked: We deal with variable-width, not fixed-width numbers. Hence, the datatype *Ternary* features four constructors, one for each ternary digit plus an end-of-list constructor, see Appendix 1. Four constructors can be conveniently encoded using two bits, not wasting precious bit space. For example, the ternary representation of 20, the length of (1.1), is 0t 1 3 2, which occupies $2 \cdot 4 = 8$ bits.

A Huffman compressor stores the code tree alongside the encoded text. Since the encoding depends on the tree this is another example of a dependent pair—actually, quite a nice one. (Once you start looking closely, you notice that dependent pairs are everywhere.)

$$\text{Huffman} : \text{Set} \rightarrow \text{Set}$$

$$\text{Huffman Alphabet} = \Sigma (\text{Tree Alphabet}) (\lambda t \rightarrow \text{List}' (\text{In } t))$$

An element of the type is a dependent pair consisting of a code tree and a length-encoded list of self-certifying characters. The corresponding codec is obtained simply by changing the font.

$$\text{Huffman} : \{\text{Alphabet} : \text{Set}\} \rightarrow \text{Codec Alphabet} \rightarrow \text{Codec} (\text{Huffman Alphabet})$$

$$\text{Huffman Alphabet} = \text{Sigma} (\text{Tree Alphabet}) (\lambda t \rightarrow \text{List}' (\text{In } t))$$

Voilà, again. The compositional approach shows its strengths. In particular, it is easy to change aspects. For example, to compress a list of lists of characters, we simply replace the inner codec by $\text{List}' (\text{List}' (\text{Compress } t))$.

7 Conclusion

Overall, this was an enjoyable exercise in interface design.

The six proof rules of Section 2 nicely summarise the interface; the user of the library needs to know little more. The implementation of the recursive case combinator posed the greatest challenge: how to hide the details of the termination proof behind the abstraction barrier of the combinator library? Our solution draws inspiration from Guillaume Allais' impressive library of total parsing combinators (Allais, 2018). Unfortunately, his approach was not immediately applicable as it is based on the fundamental assumption that a successful parse consumes at least one character, too strong an assumption for our application.

The combination of statistical and structural compression was a breeze, once the idea of self-certifying data was in place. The definition of codecs is an instance of datatype-generic programming (Jansson & Jeuring, 2002). It was pleasing to see that the technique carries over effortlessly to dependent types and proofs.

The modular implementation of the Huffman compressor provides some additional insight into the mechanics of Huffman encoding. The usual treatment of sequences is inherently non-modular and imposes an additional constraint: singleton code trees are not admissible. Our modular compressor needs roughly $2 \cdot \log_3 n$ additional bits to store the length n of the encoded text—a small price one is, perhaps, willing to pay for modularity.

Last but not least, Huffman-encoded data make a nice example for dependent pairs, an observation that is perhaps obvious, but which I have not seen spelled out before.

Conflicts of Interest

None.

Supplementary materials

For supplementary material for this article, please visit <http://doi.org/10.1017/S095679682200017X>

References

- Allais, G. (2018). *agdarsec* — Total parser combinators. *Journées Francophones des Langages Applicatifs*. Available at: <http://gallais.github.io/index.html>.
- Huffman, D. A. (1952). A method for the construction of minimum-redundancy codes. *Proc. IRE.* **40**(9), 1098–1101.
- Hughes, R. J. M. (1986). A novel representation of lists and its application to the function “reverse”. *Inf. Process. Lett.* **22**(3), 141–144.
- Jansson, P. & Jeuring, J. (2002). Polytypic data conversion programs. *Sci. Comput. Program.* **43**(1), 35–75.
- Stump, A. (2016). *Verified Functional Programming in Agda*. Association for Computing Machinery and Morgan & Claypool.
- Théry, L. (2004). Formalising Huffman’s algorithm. Research report. Università degli Studi dell’Aquila.

1 Appendix

Natural numbers. The standard ordering on the naturals, m succeeds n , is defined

```
data _>_ : ℕ → ℕ → Set where
  zero : n > zero
  succ : n > m → succ n > succ m
```

The datatype overloads the constructors *zero* and *succ*: the proof *zero* states that *zero* (the natural) is the least element; *succ* captures that *succ* (the successor function) is monotone. It is straightforward to show that \succ is reflexive and transitive and that $\text{succ } n \succ n$.

```
>reflexive : n > n
>transitive : n > m → m > k → n > k
succ-n>n : succ n > n
```

The inverse relation, m precedes n , is given by $m \preccurlyeq n = n \succ m$.

Partial functions. In Agda, a partial function from A to B can be modelled by a total function of type $A \rightarrow \text{Maybe } B$.

```
data Maybe (Elem : Set) : Set where
  Nothing : Maybe Elem
  Just    : Elem → Maybe Elem

undefined : Maybe A
undefined = Nothing
```

The type constructor *Maybe* forms a monad: *return* is the identity partial function and *bind*, “ \succ ,” denotes postfix application of a partial function.

```
return : A → Maybe A
return a = Just a
```

$$\begin{aligned}
\gg\ &: \text{Maybe } A \rightarrow (A \rightarrow \text{Maybe } B) \rightarrow \text{Maybe } B \\
\text{Nothing} \gg k &= \text{Nothing} \\
\text{Just } a \gg k &= k\ a
\end{aligned}$$

It is straightforward to verify the monad laws.

Ternary numbers. For completeness, here is the type of ternary numbers.

```

data Ternary : Set where
  0t : Ternary
  _1 _2 _3 : Ternary → Ternary
Eau-de-Cologne = 0t 1 2 3 3 3 1 1 1

```