

Arithmetic coding with folds and unfolds

*Richard Bird and Jeremy Gibbons
with Barney Stratford
University of Oxford
AFP, August 2002*

1. Introduction

- Arithmetic coding (AC) is a method for data compression
- AC can be more *effective* than Huffman coding and other rival schemes
- AC is well-suited to *adaptive* encoding
- Substantial effort has gone into making AC as *efficient* as competing methods
- AC was popularized in 1987 by Witten et al, who provided an 'accessible implementation' of 300 lines of C

Goals

- Our aim is to develop a 30-line Haskell program for the basic encoding and decoding algorithms
- The development makes use of the algebra of folds and unfolds
- We start with a naive executable program for encoding, then refine it
- Along the way, we will present some of the theory of folds and unfolds
- Be warned: there is a fair amount of arithmetic in arithmetic coding!

1.1. Arithmetic coding, naively

- Break source message into *symbols*, where a symbol is either a character or some logical grouping of characters
- Associate each distinct symbol with a semi-open *interval* of the unit interval $[0..1)$
- Successively *narrow* the unit interval by the intervals associated with each symbol of the message in turn
- Represent the final interval by choosing some *fraction* within it

1.2. Intervals

type *Fraction* = *Ratio Integer*
type *Interval* = (*Fraction*, *Fraction*)
unit :: *Interval*
unit = (0, 1)
within :: *Fraction* → *Interval* → *Bool*
within *x* (*l*, *r*) = $l \leq x \wedge x < r$
pick :: *Interval* → *Fraction*
pick (*l*, *r*) = $(l + r) / 2$

We will assume that $0 \leq l < r \leq 1$ for every $(l, r) :: \textit{Interval}$.

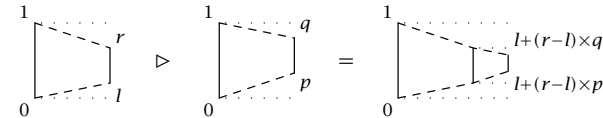
All we require of *pick* is that

pick *int* *within* *int*

Narrowing

(\triangleright) :: *Interval* → *Interval* → *Interval*
 $(l, r) \triangleright (p, q) = (l + (r - l) \times p, l + (r - l) \times q)$

Diagrammatically:



1.3. Models

Model is an abstract type representing a finite mapping from *Symbols* to *Intervals*, with associated functions

encodeSym :: *Model* → *Symbol* → *Interval*
decodeSym :: *Model* → *Fraction* → *Symbol*

We assume that for any $m :: \textit{Model}$ and $x :: \textit{Fraction}$ such that x *within* *unit*, there is a unique $s :: \textit{Symbol}$ such that

x *within* *encodeSym* *m* *s*

(so intervals in the model partition *unit*). Consequently,

$s = \textit{decodeSym} *m* $x \equiv x$ *within* *encodeSym* *m* *s*$

Adaptive models

Rather than a single fixed model for a whole message, the model can *adapt* as the message is read.

Therefore we suppose also

newModel :: *Model* → *Symbol* → *Model*

As long as the decoder performs the same adaptations as the message is reconstructed, the message can be retrieved.

Crucially, *there is no need to transmit the model with the message*.

1.4. Encoding

$$\begin{aligned}
 \text{encode}_0 &:: \text{Model} \rightarrow [\text{Symbol}] \rightarrow \text{Fraction} \\
 \text{encode}_0 m &= \text{pick} \cdot \text{foldl } (\triangleright) \text{ unit} \cdot \text{encodeSyms } m \\
 \text{encodeSyms} &:: \text{Model} \rightarrow [\text{Symbol}] \rightarrow [\text{Interval}] \\
 \text{encodeSyms } m \text{ ss} &= \text{unfoldr nextInt } (m, \text{ss}) \\
 \text{nextInt} &:: (\text{Model}, [\text{Symbol}]) \rightarrow \\
 &\quad \text{Maybe } (\text{Interval}, (\text{Model}, [\text{Symbol}])) \\
 \text{nextInt } (m, []) &= \text{Nothing} \\
 \text{nextInt } (m, s : \text{ss}) &= \text{Just } (\text{encodeSym } m \ s, (\text{newModel } m \ s, \text{ss}))
 \end{aligned}$$

1.5. Folds

foldl iterates over a list, from left to right:

$$\begin{aligned}
 \text{foldl} &:: (\beta \rightarrow \alpha \rightarrow \beta) \rightarrow \beta \rightarrow [\alpha] \rightarrow \beta \\
 \text{foldl } f \ e \ [] &= e \\
 \text{foldl } f \ e \ (x : \text{xs}) &= \text{foldl } f \ (f \ e \ x) \ \text{xs}
 \end{aligned}$$

Thus,

$$\text{foldl } (\oplus) \ e \ [x, y, z] = ((e \oplus x) \oplus y) \oplus z$$

Another fold: *foldr*

foldr iterates over a list, from right to left:

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

Thus,

$$\text{foldr } (\oplus) \ e \ [x, y, z] = x \oplus (y \oplus (z \oplus e))$$

Crucial fact: universal property — for strict h ,

$$h = \text{foldr } f \ e \equiv h \ [] = e \wedge h \ (x : \text{xs}) = f \ x \ (h \ \text{xs})$$

A relationship between *foldl* and *foldr*

Theorem (First Duality Theorem). If \oplus is associative with unit e , then

$$\text{foldl } f \ e \ \text{xs} = \text{foldr } f \ e \ \text{xs}$$

for finite lists xs .

Conversely,

Theorem Third Homomorphism Theorem. If $h = \text{foldl } f_1 \ e = \text{foldr } f_2 \ e$, then there is an associative f with unit e such that $h = \text{foldr } f \ e$

foldl vs foldr

From duality, we have

$$\text{foldl } (\triangleright) \text{ unit} = \text{foldr } (\triangleright) \text{ unit}$$

So why not use the (more familiar) *foldr*?

Lemma.

$$\text{foldl } (\triangleright) \text{ unit} \cdot \text{encodeSyms } m = \text{snd} \cdot \text{foldl } \text{step } (m, \text{unit})$$

where

$$\text{step } (m, \text{int}) s = (\text{newModel } m \ s, \text{int} \triangleright \text{encodeSym } m \ s)$$

...but there is no corresponding result for *foldr*.

1.6. unfoldr and Maybe

Recall the Haskell standard type *Maybe*:

$$\mathbf{data} \text{ Maybe } \alpha = \text{Just } \alpha \mid \text{Nothing}$$

and the function *unfoldr*:

$$\begin{aligned} \text{unfoldr} &:: (\beta \rightarrow \text{Maybe } (\alpha, \beta)) \rightarrow \beta \rightarrow [\alpha] \\ \text{unfoldr } f \ b &= \mathbf{case} \ f \ b \ \mathbf{of} \\ &\quad \text{Just } (a, b') \rightarrow a : \text{unfoldr } f \ b' \\ &\quad \text{Nothing} \rightarrow [] \end{aligned}$$

A relationship between unfoldr and foldr

The Haskell Library Report states:

The *unfoldr* function undoes a *foldr* operation. . . :

$$\text{unfoldr } f' (\text{foldr } f \ z \ xs) = xs$$

if the following holds:

$$f' (f \ x \ y) = \text{Just } (x, y)$$

$$f' \ z = \text{Nothing}$$

(and that's essentially all it says on unfolds!).

1.7. Decoding

$$\text{decode}_0 :: \text{Model} \rightarrow \text{Fraction} \rightarrow [\text{Symbol}]$$

specified by

$$ss \ \underline{\text{begins}} \ \text{decode}_0 \ m \ (\text{encode}_0 \ m \ ss)$$

for all *ss* (where *xs* begins *ys* if *ys* = *xs* ++ *xs'* for some *xs'*).

So if

$$\text{compress } m \ ss = (\text{length } ss, \text{encode}_0 \ m \ ss)$$

we can define

$$\text{decompress } m \ (n, x) = \text{take } n \ (\text{decode}_0 \ m \ x)$$

2.2. Defining *fromBits*

```

fromBits :: [Bit] → Fraction
fromBits = foldr pack (1/2)
pack      :: Bit → Fraction → Fraction
pack b x  = (b + x)/2

```

Proposition.

fromBits (toBits int) within int

2.3. Hylomorphisms

Consider *hylomorphisms*, functions h of the form

$$h = \text{foldr } f \ e \cdot \text{unfoldr } g$$

Recall that

$$\begin{aligned} \text{unfoldr } g \ z &= \text{case } g \ z \ \text{of} \\ &\quad \text{Nothing} \rightarrow [] \\ &\quad \text{Just } (x, z') \rightarrow x : \text{unfoldr } g \ z' \end{aligned}$$

$$\begin{aligned} \text{foldr } f \ e \ [] &= e \\ \text{foldr } f \ e \ (x : xs) &= f \ x \ (\text{foldr } f \ e \ xs) \end{aligned}$$

The two loops can be *fused*, and the intermediate list *deforested*:

$$\begin{aligned} h \ z &= \text{case } g \ z \ \text{of} \\ &\quad \text{Nothing} \rightarrow e \\ &\quad \text{Just } (x, z') \rightarrow f \ x \ (h \ z') \end{aligned}$$

Two hyls

For $\text{pick} = \text{fromBits} \cdot \text{toBits}$, we get

$$\begin{aligned} \text{pick } (l, r) \\ \mid r \leq 1/2 &= \text{pick } (2 \times l, 2 \times r) / 2 \\ \mid 1/2 \leq l &= (1 + \text{pick } (2 \times l - 1, 2 \times r - 1)) / 2 \\ \mid l < 1/2 < r &= 1/2 \end{aligned}$$

For $\text{size} = \text{foldr } \text{count } 0 \cdot \text{toBits}$ where $\text{count } a \ n = 1 + n$, we get

$$\begin{aligned} \text{size } (l, r) \\ \mid r \leq 1/2 &= 1 + \text{size } (2 \times l, 2 \times r) \\ \mid 1/2 \leq l &= 1 + \text{size } (2 \times l - 1, 2 \times r - 1) \\ \mid l < 1/2 < r &= 0 \end{aligned}$$

Proof principle for hyls

For hyls

$$h = \text{foldr } f \ e \cdot \text{unfoldr } g$$

to show $P(x, h \ x)$ for every x for which $h \ x \neq \perp$, it suffices to perform a case analysis on $g \ x$ and show that

Case $g \ x = \text{Nothing}$: $P(x, e)$

Case $g \ x = \text{Just } (a, y)$: $P(y, z) \Rightarrow P(x, f \ a \ z)$

provided that f is strict in its second argument.

2.4. Summary of first refinement

We define

$$\begin{aligned} \text{encode}_1 &:: \text{Model} \rightarrow [\text{Symbol}] \rightarrow [\text{Bit}] \\ \text{encode}_1 m &= \text{toBits} \cdot \text{foldl} (\triangleright) \text{unit} \cdot \text{encodeSyms} m \end{aligned}$$

Informally, execution of encode_1 consumes all its input before delivering any output.

Formally,

$$\text{encode}_1 m ss = \perp$$

for all partial or infinite lists ss .

Can we do better?

2.5. Streaming

The function *stream* alternates between production and consumption:

$$\begin{aligned} \text{stream} &:: (\beta \rightarrow ([\alpha], \beta)) \rightarrow (\beta \rightarrow \gamma \rightarrow \beta) \rightarrow \beta \rightarrow [\gamma] \rightarrow [\alpha] \\ \text{stream } f \ g \ z \ xs &= \gamma s ++ \text{case } xs \text{ of} \\ &\quad [] \rightarrow [] \\ &\quad (x : xs) \rightarrow \text{stream } f \ g \ (g \ z' \ x) \ xs \\ &\quad \text{where } (\gamma s, z') = f \ z \end{aligned}$$

(Informally, γ is the input type, β the state, and α the output type.)

2.6. evolve: a variation on *unfoldr*

$$\begin{aligned} \text{evolve} &:: (\beta \rightarrow \text{Maybe} (\alpha, \beta)) \rightarrow \beta \rightarrow ([\alpha], \beta) \\ \text{evolve } f \ z &= \text{case } f \ z \text{ of} \\ &\quad \text{Nothing} \rightarrow ([], z) \\ &\quad \text{Just } (y, z') \rightarrow (y : \gamma s, z'') \\ &\quad \text{where } (\gamma s, z'') = \text{evolve } f \ z' \end{aligned}$$

Claim: *evolve* yields the unfold and the final state.

$$\text{evolve } f \ z = (\text{unfoldr } f \ z, \text{loop } (\text{fmap } \text{snd} \cdot f) \ z)$$

where the *fmap* here is the map for *Maybe*, and

$$\begin{aligned} \text{loop} &:: (\beta \rightarrow \text{Maybe } \beta) \rightarrow \beta \rightarrow \beta \\ \text{loop } f \ z &= \text{case } f \ z \text{ of} \\ &\quad \text{Nothing} \rightarrow z \\ &\quad \text{Just } z' \rightarrow \text{loop } f \ z' \end{aligned}$$

2.7. The Streaming Theorem

Definition. The *streaming condition* for f and g is:

$$f \ z = \text{Just } (y, z') \Rightarrow f \ (g \ z \ x) = \text{Just } (y, g \ z' \ x)$$

for all z, y, z' and x .

Lemma. If the streaming condition holds for f and g , then

$$f \ z = \text{Just } (y, z') \Rightarrow f \ (\text{foldl } g \ z \ xs) = \text{Just } (y, \text{foldl } g \ z' \ xs)$$

for all z, y, z' and finite lists xs .

Theorem. If the streaming condition holds for f and g , then

$$\text{unfoldr } f \ (\text{foldl } g \ z \ xs) = \text{stream } (\text{evolve } f) \ g \ z \ xs$$

for all *finite* lists xs .

Proof of Lemma

By induction on xs :

Case []: Immediate.

Case $x : xs$: Assume $f z = \text{Just } (y, z')$. So by streaming condition,
 $f (g z x) = \text{Just } (y, g z x')$. Now we have

$$\begin{aligned}
 & f (\text{foldl } g z (x : xs)) \\
 = & \quad \{\text{definition of foldl}\} \\
 & f (\text{foldl } g (g z x) xs) \\
 = & \quad \{\text{induction}\} \\
 & \text{Just } (y, \text{foldl } g (g z' x) xs) \\
 = & \quad \{\text{definition of foldl}\} \\
 & \text{Just } (y, \text{foldl } g z' (x : xs))
 \end{aligned}$$

Proof of Streaming Theorem

A *double induction* on xs and n to show that

$$\begin{aligned}
 & \text{approx } n (\text{unfoldr } f (\text{foldl } g z xs)) \\
 = & \quad \text{approx } n (\text{stream } (\text{evolve } f) g z xs)
 \end{aligned}$$

if the streaming condition holds for f and g .

Case $xs = []$: Then

$$\begin{aligned}
 \text{lhs} &= \text{approx } n (\text{unfoldr } f z) \\
 \text{rhs} &= \text{approx } n (\text{fst } (\text{evolve } f z))
 \end{aligned}$$

Case $xs = x : xs'$: We perform a case analysis on $f z$.

Subcase $f z = \text{Nothing}$: Then

$$\begin{aligned}
 \text{lhs} &= \text{approx } n (\text{unfoldr } f (\text{foldl } g (g z x) xs')) \\
 &= \quad \{\text{induction}\} \\
 & \quad \text{approx } n (\text{stream } (\text{evolve } f) g (g z x) xs') = \text{rhs}
 \end{aligned}$$

Approx lemma

Recall *approx lemma* (IFPH §9.3):

Lemma. Let

$$\begin{aligned}
 & \text{approx} :: \text{Integer} \rightarrow [\alpha] \rightarrow [\alpha] \\
 & \text{approx } (n + 1) [] = [] \\
 & \text{approx } (n + 1) (x : xs) = x : \text{approx } n xs
 \end{aligned}$$

Then two lists xs and ys are equal iff

$$(\forall n : 0 \leq n : \text{approx } n xs = \text{approx } n ys)$$

Subcase $f z = \text{Just } (y, z')$: We perform an induction on n .

Subsubcase $n = 0$: Trivial, since $\text{approx } 0 xs = \perp$ for any xs .

Subsubcase $n = n' + 1$: Then

$$\begin{aligned}
 \text{lhs} &= \text{approx } (n' + 1) (\text{unfoldr } f (\text{foldl } g z (x : xs'))) \\
 &= \quad \{\text{by lemma}\} \\
 & \quad \text{approx } (n' + 1) (y : \text{unfoldr } f (\text{foldl } g z' (x : xs'))) \\
 &= \quad \{\text{approx}\} \\
 & \quad y : \text{approx } n' (\text{unfoldr } f (\text{foldl } g z' (x : xs'))) \\
 &= \quad \{\text{induction}\} \\
 & \quad y : \text{approx } n' (\text{stream } (\text{evolve } f) g z' (x : xs')) \\
 &= \quad \{\text{approx}\} \\
 & \quad \text{approx } (n' + 1) (y : \text{stream } (\text{evolve } f) g z' (x : xs')) \\
 &= \quad \{\text{stream, evolve}\} \\
 & \quad \text{approx } (n' + 1) (\text{stream } (\text{evolve } f) g z (x : xs')) = \text{rhs}
 \end{aligned}$$

2.8. Streaming condition for *nextBit* and \triangleright

Recall

$$\begin{aligned} \text{nextBit} &:: \text{Interval} \rightarrow \text{Maybe} (\text{Bit}, \text{Interval}) \\ \text{nextBit } (l, r) & \\ \quad | r \leq 1/2 &= \text{Just } (0, (0, 2) \triangleright (l, r)) \\ \quad | 1/2 \leq l &= \text{Just } (1, (-1, 1) \triangleright (l, r)) \\ \quad | l < 1/2 < r &= \text{Nothing} \end{aligned}$$

We leave it as an exercise to show that the streaming condition follows from associativity of \triangleright , and the fact that $\text{int}_1 \triangleright \text{int}_2$ is contained within int_1 .

2.9. Summary of second refinement

Previously,

$$\begin{aligned} \text{encode}_1 &:: \text{Model} \rightarrow [\text{Symbol}] \rightarrow [\text{Bit}] \\ \text{encode}_1 m &= \text{unfoldr } \text{nextBit} \cdot \text{foldl } (\triangleright) \text{unit} \cdot \text{encodeSyms } m \end{aligned}$$

Now,

$$\begin{aligned} \text{encode}_2 &:: \text{Model} \rightarrow [\text{Symbol}] \rightarrow [\text{Bit}] \\ \text{encode}_2 m &= \text{stream } (\text{evolve } \text{nextBit}) (\triangleright) \text{unit} \cdot \text{encodeSyms } m \end{aligned}$$

Note that $\text{encode}_1 m \neq \text{encode}_2 m$, but they are equal on all finite symbol sequences.

3. Decoding and interval expansion

The function $\text{decode}_2 :: \text{Model} \rightarrow [\text{Bit}] \rightarrow [\text{Symbol}]$ is specified by

$$\text{ss } \underline{\text{begins}} \text{ decode}_2 m (\text{encode}_2 m \text{ ss})$$

for all finite sequences of symbols ss , where

$$\begin{aligned} \text{encode}_2 &:: \text{Model} \rightarrow [\text{Symbol}] \rightarrow [\text{Bit}] \\ \text{encode}_2 m &= \text{stream } (\text{evolve } \text{nextBit}) (\triangleright) \text{unit} \cdot \text{encodeSyms } m \end{aligned}$$

How can we implement decode_2 ?

3.1. The Stream Inversion Theorem

Theorem. Let $\text{process } z = \text{foldr } (\oplus) e \cdot \text{stream } f g z$.

Suppose \ominus and extract can be found so that

1. $(y \oplus w) \ominus y = w$, for all y and w .
2. $w = \text{process } z (x : \text{xs}) \Rightarrow \text{extract } (w, z) = x$, for all w, z, x and xs .

Then $\text{xs } \underline{\text{begins}} \text{ unfoldr } \text{nextx } (\text{process } z \text{ xs}, z)$, where

$$\begin{aligned} \text{nextx } (w, z) &= \text{Just } (x, (\text{foldl } (\ominus) w \text{ ys}, g z' x)) \\ &\quad \text{where } x &= \text{extract } (w, z) \\ &\quad \text{ys}, z' &= f z \end{aligned}$$

3.2. Proof of the Inversion Theorem

The proof is by induction on xs .

Case []: Observe firstly that $unfoldr\ nextx\ (process\ z\ xs,\ z)$ evaluates to an infinite list (of possibly undefined elements). Since [] begins every finite or infinite list, the case is established.

Case $x : xs$: Let $w = process\ z\ (x : xs)$, so $x = extract\ (w,\ z)$ and

$$unfoldr\ nextx\ (w,\ z) = x : unfoldr\ nextx\ (foldl\ (\ominus)\ w\ ys,\ g\ z'\ x)$$

where $(ys,\ z') = f\ z$. We therefore have to show

$$xs\ \underline{begins}\ unfoldr\ nextx\ (foldl\ (\ominus)\ w\ ys,\ g\ z'\ x)$$

Proof, continued

By induction, $xs\ \underline{begins}\ unfoldr\ nextx\ (foldl\ (\ominus)\ w\ ys,\ g\ z'\ x)$ if

$$foldl\ (\ominus)\ w\ ys = process\ (g\ z'\ x)\ xs$$

By definition of *stream* and $f\ z = (ys,\ z')$ we have

$$w = foldr\ (\oplus)\ e\ (ys\ ++\ stream\ f\ g\ (g\ z'\ x)\ xs)$$

Furthermore,

$$foldl\ (\ominus)\ (foldr\ (\oplus)\ e\ (ys\ ++\ zs))\ ys = foldr\ (\oplus)\ e\ zs$$

for all ys and zs . The proof is left as an exercise.

Proof, continued

Hence

$$\begin{aligned} & foldl\ (\ominus)\ w\ ys \\ = & \quad \{\text{above expression for } w\} \\ & foldl\ (\ominus)\ (foldr\ (\oplus)\ e\ (ys\ ++\ stream\ f\ g\ (g\ z'\ x)\ xs)) \\ = & \quad \{\text{above identity}\} \\ & foldr\ (\oplus)\ e\ (stream\ f\ g\ (g\ z'\ x)\ xs) \\ = & \quad \{\text{definition of } process\} \\ & process\ (g\ z'\ x)\ xs \end{aligned}$$

Summarising, $foldl\ (\ominus)\ w\ ys = process\ (g\ z'\ x)\ xs$, establishing the case and the theorem.

3.3. Applying the Inversion Theorem: a preliminary step

Recall that in the first lecture we showed

$$foldl\ (\triangleright)\ unit \cdot encodeSyms\ m = snd \cdot foldl\ step\ (m,\ unit)$$

where

$$step\ (m,\ int)\ s = (newModel\ m\ s,\ int\ \triangleright\ encodeSym\ m\ s)$$

This identity allows us to fuse *encodeSyms* into the narrowing process:

$$encode_1\ m = unfoldr\ nextBitM \cdot foldl\ step\ (m,\ unit)$$

where *nextBitM* is identical to *nextBit* except that it propagates the model as an additional argument. The Streaming Theorem is again applicable and we obtain

$$encode_2\ m = stream\ (evolve\ nextBitM)\ step\ (m,\ unit)$$

3.4. Applying the Inversion Theorem: \ominus and *extract*

Let $process(m, int) = fromBits \cdot stream(evolve nextBitM) step(m, int)$.

Recall that $fromBits = foldr(\oplus)(1/2)$ where $b \oplus x = b + x/2$. Defining \ominus by $x \ominus b = 2 \times x - b$, we therefore have $(b \oplus x) \ominus b = x$.

Let $x = process(m, int)(s : ss)$, so x within $(int \triangleright encodeSym m s)$. If $int = (l, r)$ and $encodeSym m s = (p, q)$, then

$$\begin{aligned} & l + (r - l) \times p \leq x < l + (r - l) \times q \\ \equiv & \quad \{\text{arithmetic}\} \\ & p \leq x - l / r - l < q \\ \equiv & \quad \{\text{definition of } decodeSym\} \\ & s = decodeSym m (x - l / r - l) \end{aligned}$$

Hence $extract(x, (m, (l, r))) = decodeSym m (x - l / r - l)$.

3.5. Definition of $decode_2$

We therefore obtain

$$\begin{aligned} decode_2 m bs &= unfoldr nextSym (fromBits bs, (m, unit)) \\ nextSym(x, (m, int)) &= Just(s, (x', (newModel m s, \\ & \quad \quad \quad int' \triangleright (encodeSym m s)))) \\ \text{where } s &= decodeSym m (resize x int) \\ x' &= foldl(\lambda b x \rightarrow 2 \times x - b) x bs \\ (bs, int') &= evolve nextBit int \\ \\ resize x (l, r) &= x - l / r - l \end{aligned}$$

Note, however, that $decode_2$ inspects all of bs before producing any output, so is not an incremental algorithm.

3.6. Interval expansion

In the final lecture we are going to replace fractional arithmetic by arithmetic with limited-precision integers. Before doing so we need a preparatory step: interval expansion.

Quoting from Howard and Vitter:

The idea is to prevent the current interval from narrowing too much when the endpoints are close to $1/2$ but straddle $1/2$. In that case we do not yet know the next output bit, but we do know that whatever it is, the *following* bit will have the opposite value; we merely keep track of that fact, and expand the current interval about $1/2$. This follow-on procedure may be repeated any number of times, so the current interval is always strictly longer than $1/4$.

Interval expansion, formally

Formally, interval expansion is a data refinement in which an interval (l, r) is represented by a triple of the form $(n, (l', r'))$ satisfying

$$l' = scale(n, l) \quad \text{and} \quad r' = scale(n, r)$$

where

$$scale(n, x) = 2^n \times (x - 1/2) + 1/2$$

and subject to $0 \leq l' < r' \leq 1$. In particular, $(0, (l, r))$ is one possible representation of (l, r) .

A *fully-expanded* interval is a triple $(n, (l', r'))$ in which n is as large as possible. Intervals will be fully-expanded immediately before narrowing.

The remainder of this lecture is devoted to installing this data refinement.

Definition of *expand*

Since

$$\begin{aligned} 0 \leq 2 \times (l - 1/2) + 1/2 &\equiv 1/4 \leq l \\ 2 \times (r - 1/2) + 1/2 \leq 1 &\equiv r \leq 3/4 \end{aligned}$$

we can further expand $(n, (l, r))$ if $1/4 \leq l$ and $r \leq 3/4$. Hence we define

$$\begin{aligned} \mathit{expand} (n, (l, r)) \\ | 1/4 \leq l \wedge r \leq 3/4 &= \mathit{expand} (n+1, (2 \times l - 1/2, 2 \times r - 1/2)) \\ | \mathbf{otherwise} &= (n, (l, r)) \end{aligned}$$

In particular, $\mathit{expand} (0, (l, r))$ returns a fully-expanded interval for (l, r) .

3.7. Interval contraction

The converse of *expand* is given by

$$\mathit{contract} (n, (l, r)) = (\mathit{rescale} (n, l), \mathit{rescale} (n, r))$$

where $\mathit{rescale} (n, x) = (x - 1/2) / 2^n + 1/2$.

We leave it as an exercise to verify that

$$\begin{aligned} \mathit{contract} \cdot \mathit{expand} &= \mathit{contract} \\ \mathit{contract} (n, \mathit{int1} \triangleright \mathit{int2}) &= \mathit{contract} (n, \mathit{int1}) \triangleright \mathit{int2} \end{aligned}$$

Consequently, defining *enarrow* (“expand and narrow”) by

$$\begin{aligned} \mathit{enarrow} \mathit{ei} \mathit{int2} &= (n, \mathit{int1} \triangleright \mathit{int2}) \\ &\quad \mathbf{where} (n, \mathit{int1}) = \mathit{expand} \mathit{ei} \end{aligned}$$

we have $\mathit{contract} (\mathit{enarrow} \mathit{ei} \mathit{int}) = \mathit{contract} \mathit{ei} \triangleright \mathit{int}$.

3.8. Aim

Our aim is to show

$$\begin{aligned} &\mathit{stream} (\mathit{evolve} \mathit{nextBit}) (\triangleright) (\mathit{contract} \mathit{ei}) \\ &= \mathit{concat} \cdot \mathit{stream} (\mathit{evolve} \mathit{nextBits}) \mathit{enarrow} \mathit{ei} \end{aligned}$$

for the following definition of *nextBits*:

$$\begin{aligned} \mathit{nextBits} (n, (l, r)) \\ | r \leq 1/2 &= \mathit{Just} (\mathit{bits} \ n \ 0, (0, (2 \times l, 2 \times r))) \\ | 1/2 \leq l &= \mathit{Just} (\mathit{bits} \ n \ 1, (0, (2 \times l - 1, 2 \times r - 1))) \\ | \mathbf{otherwise} &= \mathit{Nothing} \end{aligned}$$

where $\mathit{bits} \ n \ b = b : \mathit{replicate} \ n \ (1 - b)$ returns a b followed by a sequence of n copies of $1 - b$.

3.9. Stream fusion

Say that $k :: [\alpha] \rightarrow [\beta]$ is *distributive* if

$$\begin{aligned} k (xs ++ ys) &= k \ xs ++ k \ ys \\ k [] &= [] \end{aligned}$$

Examples of distributive functions are *id*, *map f*, and *concat*.

Theorem. Suppose k is distributive. Then

$$k \cdot \mathit{stream} \ f \ g \ z = \mathit{stream} \ f' \ g' \ (h \ z)$$

if the following two conditions hold for all z, z', ys and x :

$$\begin{aligned} f \ z = (ys, z') &\Rightarrow f' (h \ z) = (k \ ys, h \ z') \\ h (g \ z \ x) &= g' (h \ z) \ x \end{aligned}$$

We leave the proof as an exercise in fixpoint induction.

Applying stream fusion

We have that *concat* is distributive, and that

$$\text{contract } (\text{enarrow } ei \text{ int}) = (\text{contract } ei) \triangleright \text{int}$$

It remains to verify the remaining fusion condition:

$$\begin{aligned} & \text{evolve nextBits } ei1 = (bss, ei2) \\ \Rightarrow & \text{evolve nextBit } (\text{contract } ei1) = (\text{concat } bss, \text{contract } ei2) \end{aligned}$$

Again, we relegate the proof to the exercises.

3.10. Conclusion

We have derived

$$\text{encode}_3 m = \text{concat} \cdot \text{stream nextBits enarrow } (0, \text{unit}) \cdot \text{encodeSyms } m$$

In the new version we expand intervals (l, r) satisfying $l < 1/2 < r$ immediately before narrowing. Consequently, narrowing is applied only when $l < 1/4$ and $1/2 < r$, or $l < 1/2$ and $3/4 < r$.

In either case, $1/4 < r - l$.

This fact will be crucial in the final lecture.

4. From fractions to integers

We now want to replace fractional arithmetic by arithmetic with limited-precision integers.

In the final version intervals take the form (l, r) , where l and r are integers in the range $0 \leq l < r \leq w$ and w is a *fixed* power of two. This pair represents the interval $(l/w, r/w)$.

Intervals in each model m take the form (p, q, d) , where p and q are integers in the range $0 \leq p < q \leq d$ and d is an integer which is fixed for m and called the *denominator* for m . This triple represents the interval $(p/d, q/d)$.

4.1. Integer narrowing

The narrowing function is redefined as follows:

$$(l, r) \triangleright (p, q, d) = (l + [(r-l) \times p/d], l + [(r-l) \times q/d])$$

Equivalently

$$(l, r) \triangleright (p, q, d) = (l + ((r-l) \times p) \text{ div } d, l + ((r-l) \times q) \text{ div } d)$$

There are various problems with this idea, and a number of obstacles that have to be overcome.

Problems

- The revised definition of narrowing completely changes the specification: encoding will now produce different outputs than before. In general, the effectiveness of compression will be reduced.
- Worse, \blacktriangleright is not associative and none of the foregoing development applies.
- Unless we take steps to avoid it, intervals can *collapse* to the empty interval when $\lfloor (r-l) \times p/d \rfloor = \lfloor (r-l) \times q/d \rfloor$.

4.2. Change of specification

Fortunately, we can recover all of the previous development. Observe that

$$(l, r) \blacktriangleright (p, q, d) = (l/w, r/w) \triangleright (p'/d, q'/d)$$

where

$$\begin{aligned} p' &= d/r-l \times \lfloor (r-l) \times p/d \rfloor \\ q' &= d/r-l \times \lfloor (r-l) \times q/d \rfloor \end{aligned}$$

Hence, provided $p' < q'$, integer narrowing can be viewed as fractional narrowing applied in an *adjusted* model. The adjustment to the model depends on the current interval (l, r) but that doesn't matter. What does follow is that our development remains valid because we have allowed for changing models in the specification.

4.3. When intervals collapse

It is left as an exercise to show that

$$(\forall p, q: 0 \leq p < q \leq d: \lfloor (r-l) \times p/d \rfloor < \lfloor (r-l) \times q/d \rfloor)$$

if and only if $d \leq r - l$.

Hence we have to ensure that the width of each interval is at least d before narrowing. But interval expansion guarantees that the width of each (expanded) interval is greater than $w/4$ before narrowing, so interval collapse is avoided if $w/4 \geq d$.

Since $w \times d \leq w \times w/4 = 2^{2 \times e - 2}$ if $w = 2^e$, we have to ensure that our limited-precision arithmetic is accurate to $2 \times e - 2$ bits.

4.4. Final version of encode

$$\text{encode } m = \text{concat} \cdot \text{stream} (\text{evolve nextBits}) \text{enarrow unit} \cdot \text{encodeSyms } m$$

$$\text{unit} = (0, (0, w))$$

$$\text{enarrow } ei \text{ int2} = (n, \text{int1} \blacktriangleright \text{int2})$$

$$\text{where } (n, \text{int1}) = \text{expand } ei$$

$$\text{expand } (n, (l, r))$$

$$| w/4 \leq l \wedge r \leq 3 \times w/4 = \text{expand } (n+1, (2 \times l - w/2, 2 \times r - w/2))$$

$$| \text{otherwise} = (n, (l, r))$$

$$\text{nextBits } (n, (l, r))$$

$$| r \leq w/2 = \text{Just } (\text{bits } n \ 0, (0, (2 \times l, 2 \times r)))$$

$$| w/2 \leq l = \text{Just } (\text{bits } n \ 1, (0, (2 \times l - w, 2 \times r - w)))$$

$$| \text{otherwise} = \text{Nothing}$$

4.5. A final problem

Notwithstanding everything that has gone before, encoding is *not* guaranteed to work with any form of limited-precision arithmetic!!

Can you see why?

4.6. Decoding in the integer version

To see how to perform decoding with limited-precision arithmetic, we need some preliminary investigation.

Suppose bs is the output of *encode* when started off with an initial interval ei on a symbol sequence $s : ss$.

Setting $x = w \times \text{fromBits } bs$, we know that x is a fractional number satisfying

$$x \text{ within } \text{contract } (\text{enarrow } ei \text{ (encodeSym } m \text{ s)})$$

How can we calculate s given x , m and ei ? We need to do this in order to define *extract*.

BUG!

4.7. Another Stream Inversion Theorem

Theorem. Let $\text{process } z = \text{foldr } (\oplus) e \cdot \text{stream } (\text{evolve } f) g \text{ z}$.

Suppose \ominus and *extract* can be found so that

1. $(y \oplus w) \ominus y = w$, for all y and w .
2. $w = \text{process } z \text{ (} x : xs \text{)} \wedge f \text{ z} = \text{Nothing} \Rightarrow \text{extract } (w, z) = x$, for all w , z , x and xs .

Then xs *begins* *unfoldr nextx* ($\text{process } z \text{ xs, z}$), where

$$\begin{aligned} \text{nextx } (w, z) &= \text{Just } (x, (w', z')) \\ \text{where } x &= \text{extract } (w, z) \\ (ys, z') &= \text{evolve } f \text{ (} g \text{ z } x \text{)} \\ w' &= \text{foldl } (\ominus) w \text{ ys} \end{aligned}$$

Floor calculations

The following property characterizes floors: for all integers n and fractions x ,

$$n \leq \lfloor x \rfloor \equiv n \leq x$$

Ignorance of this simple rule has marred every published paper on arithmetic coding that we have read.

The calculation

$$\begin{aligned}
 & x \text{ *within* } \text{contract } (enarrow\ ei\ (encodeSym\ m\ s)) \\
 \equiv & \quad \{setting\ (n,\ (l,\ r)) =\ expand\ ei\ \} \\
 & x \text{ *within* } (\text{contract } (n,\ (l,\ r)\ \blacktriangleright\ encodeSym\ m\ s)) \\
 \equiv & \quad \{setting\ y =\ scale\ (n,\ x)\ \} \\
 & y \text{ *within* } ((l,\ r)\ \blacktriangleright\ encodeSym\ m\ s) \\
 \equiv & \quad \{setting\ (p,\ q,\ d) =\ encodeSym\ m\ s\ \} \\
 & l + \lfloor (r - l) \times p/d \rfloor \leq y < l + \lfloor (r - l) \times q/d \rfloor \\
 \equiv & \quad \{arithmetic\ \} \\
 & \lfloor (r - l) \times p/d \rfloor \leq y - l < \lfloor (r - l) \times q/d \rfloor \\
 \equiv & \quad \{rule\ of\ floors,\ setting\ k = \lfloor y \rfloor\ \} \\
 & \lfloor (r - l) \times p/d \rfloor \leq k - l < \lfloor (r - l) \times q/d \rfloor
 \end{aligned}$$

The calculation, continued

$$\begin{aligned}
 & \lfloor (r - l) \times p/d \rfloor \leq k - l < \lfloor (r - l) \times q/d \rfloor \\
 \equiv & \quad \{arithmetic\ \} \\
 & \lfloor (r - l) \times p/d \rfloor < k - l + 1 \leq \lfloor (r - l) \times q/d \rfloor \\
 \equiv & \quad \{rule\ of\ floors\ \} \\
 & (r - l) \times p/d < k - l + 1 \leq (r - l) \times q/d \\
 \equiv & \quad \{arithmetic\ \} \\
 & p \leq ((k - l + 1) \times d - 1) / (r - l) < q \\
 \equiv & \quad \{rule\ of\ floors\ \} \\
 & p \leq \lfloor ((k - l + 1) \times d - 1) / (r - l) \rfloor < q
 \end{aligned}$$

Defining *extract*

Hence, redefining *decodeSym* to have type $Model \rightarrow Int \rightarrow Symbol$, we have

$$\begin{aligned}
 extract\ (bs,\ (m,\ ei)) &= decodeSym\ m\ t \\
 \text{where } t &= ((k - l + 1) \times d - 1) \underline{div}\ (r - l) \\
 d &= denom\ m \\
 k &= \lfloor scale\ (n,\ w \times fromBits\ bs) \rfloor \\
 (n,\ (l,\ r)) &= expand\ ei
 \end{aligned}$$

But we do not need to inspect all of *bs* to compute *k*; we only need the first $n + e$ bits, where $w = 2^e$.

This is the clue to making *decode* incremental.

4.8. A final data refinement

Suppose we represent bs by a pair (x, ds) where $x :: Int$ and $ds :: [Bit]$ are defined by

$$\begin{aligned} x &= \text{foldl } (\lambda x b \rightarrow 2 \times x + b) 0 \text{ cs} \\ (cs, ds) &= \text{splitAt } e \text{ (bs ++ 1 : replicate (e - 1) 0)} \end{aligned}$$

Recalling that $w = 2^e$ and

$$\text{fromBits } bs = \text{toFrac } 0 \text{ (bs ++ 1 : replicate } n \text{ 0)}$$

for any n , we have

$$w \times \text{fromBits } bs = x + \text{toFrac } 0 \text{ ds}$$

In particular, $x = \lfloor w \times \text{fromBits } bs \rfloor$.

Redefinition of $extract$

Recalling that $scale(n, x) = 2^n \times (x - w/2) + w/2$, we can compute $\lfloor scale(n, w \times \text{fromBits } bs) \rfloor$ by computing $fscale(n, (x, ds))$, where

$$\begin{aligned} fscale(n, (x, ds)) &= \\ &\text{foldl } (\lambda x b \rightarrow 2 \times x + b - w/2) x \text{ (take } n \text{ ds)} \end{aligned}$$

Hence

$$\begin{aligned} extract((x, ds), (m, ei)) &= \text{decodeSym } m \ t \\ \text{where } t &= ((k - l + 1) \times d - 1) \text{ div } (r - l) \\ d &= \text{denom } m \\ k &= fscale(n, (x, ds)) \\ (n, (l, r)) &= \text{expand } ei \end{aligned}$$

4.9. An alternative

The representation of bs by (x, ds) can be computed another way, though it involves inspecting all the bs :

$$\begin{aligned} (x, ds) &= \text{foldr } \text{step } (w/2, []) \text{ bs} \\ \text{step } b \ (x, ds) &= (y, d : ds) \\ &\quad \text{where } (y, d) = (x + b \times w) \text{ divMod } 2 \end{aligned}$$

The point of this alternative is that we have

$$\text{foldr } \text{step } (w/2, []) \cdot \text{concat} = \text{foldr } (\oplus) (w/2, [])$$

where $bs \oplus (x, ds) = \text{foldr } \text{step } (x, ds) \text{ bs}$.

We leave the proof as an exercise.

Preparing for inversion

Moreover, we can invert \oplus . Define \ominus by

$$\begin{aligned} (x, ds) \ominus bs &= \text{foldl } \text{unstep } (x, ds) \text{ bs} \\ \text{unstep } (x, ds) \ b &= (2 \times x - b \times w + \text{head } ds, \text{tail } ds) \end{aligned}$$

Then we have

$$\begin{aligned} & (bs \oplus (x, ds)) \ominus bs \\ &= \{ \text{definitions of } \oplus \text{ and } \ominus \} \\ & \text{foldl } \text{unstep } (\text{foldr } \text{step } (x, ds) \text{ bs}) \text{ bs} \\ &= \{ \text{since } \text{unstep } (\text{step } b \ (x, ds)) \ b = (x, ds) \} \\ & (x, ds) \end{aligned}$$

Having defined $extract$ and \ominus we can now define $decode$.

4.10. Defining *decode*

```

decode m bs =
  unfoldr nextSym ((x, ds), (m, unit))
  where x      = foldl ( $\lambda x b \rightarrow 2 \times x + b$ ) 0 cs
        (cs, ds) = splitAt e (bs # 1 : replicate (e - 1) 0)

nextSym ((x, ds), (m, ei)) =
  Just (s, ((x', ds'), (newModel m s, ei')))
  where s      = extract ((x, ds), (m, ei))
        (x', ds') = foldl ( $\ominus$ ) (x, ds) bss
        (bss, ei') = evolve nextBits (enarrow ei (encodeSym m s))

```

4.11. Final conclusions

- There is a lot of arithmetic in Arithmetic Coding — not only the arithmetic of fractions and integers, but also the arithmetic of folds and unfolds.
- To the best of our knowledge, no previous description of Arithmetic Coding has ever tackled the formal basis for why the method works.
- Functional programming provides a convenient vehicle for expressing and reasoning about the various algorithms.
- Structured recursion operators, such as *foldl*, *unfoldr* and *stream*, are ubiquitous.

Selected references

- P. G. Howard and J. S. Vitter. Arithmetic coding for data compression. *Proc. IEEE*, Vol 82, 6, 857-865, 1994.
- J. Jiang. Novel design of arithmetic coding for data compression. *IEE Proc. Comput. Dig. Tech.*, Vol 142, 6 (Nov) 419-424, 1995.
- A. Moffat, R. M. Neal, and I. H. Witten. Arithmetic coding revisited. *ACM Trans. on Inf. Systems* Vol 16, No 3, 256-294, July 1998.
- I. H. Witten, R. M. Neal, and J. G. Cleary. Arithmetic coding for data compression. *C. ACM*, Vol 30, No 6, 520-540, June 1987.