

A Combinator Library for Prefix-Free Codes

Emin Karayel

April 19, 2022

Abstract

This entry contains a set of binary encodings for primitive data types, such as natural numbers, integers, floating-point numbers as well as combinators to construct encodings for products, lists, sets or functions of/between such types.

For natural numbers and integers, the entry contains various encodings, such as Elias-Gamma-Codes and exponential Golomb Codes, which are efficient variable-length codes in use by current compression formats.

A use-case for this library is measuring the persisted size of a complex data structure without having to hand-craft a dedicated encoding for it, independent of Isabelle's internal representation.

1 Introduction

theory *Prefix-Free-Code-Combinators*

imports

HOL-Library.Extended-Real

HOL-Library.Float

HOL-Library.FuncSet

HOL-Library.List-Lexorder

HOL-Library.Log-Nat

HOL-Library.Sublist

begin

The encoders are represented as partial prefix-free functions. The advantage of prefix free codes is that they can be easily combined by concatenation. The approach of using prefix free codes (on the byte-level) for the representation of complex data structures is common in many industry encoding libraries (cf. [2]).

The reason for representing encoders using partial functions, stems from some use-cases where the objects to be encoded may be in a much smaller sets, as their type may suggest. For example a natural number may be known to have a given range, or a function may be encodable because it has a finite domain.

Note: Prefix-free codes can also be automatically derived using Huffmans' algorithm, which was formalized by Blanchette [1]. This is especially useful if it is possible to transmit a dictionary before the data. On the other hand these standard codes are useful, when the above is impractical and/or the distribution of the input is unknown or expected to be close to the one's implied by standard codes.

The following section contains general definitions and results, followed by Section 3 to 10 where encoders for primitive types and combinators are defined. Each construct is accompanied by lemmas verifying that they form prefix free codes as well as bounds on the bit count to encode the data. Section 11 concludes with a few examples.

2 Encodings

fun *opt-prefix* **where**
opt-prefix (Some x) (Some y) = *prefix* x y |
opt-prefix - - = *False*

definition *opt-comp* x y = (*opt-prefix* x y \vee *opt-prefix* y x)

fun *opt-append* :: *bool list option* \Rightarrow *bool list option* \Rightarrow *bool list option*
where
opt-append (Some x) (Some y) = *Some* (x@y) |
opt-append - - = *None*

lemma *opt-comp-sym*: *opt-comp* x y = *opt-comp* y x
 \langle *proof* \rangle

lemma *opt-comp-append*:
assumes *opt-comp* (opt-append x y) z
shows *opt-comp* x z
 \langle *proof* \rangle

lemma *opt-comp-append-2*:
assumes *opt-comp* x (opt-append y z)
shows *opt-comp* x y
 \langle *proof* \rangle

lemma *opt-comp-append-3*:
assumes *opt-comp* (opt-append x y) (opt-append x z)
shows *opt-comp* y z
 \langle *proof* \rangle

type-synonym 'a *encoding* = 'a \rightarrow *bool list*

definition *is-encoding* :: 'a encoding \Rightarrow bool
where *is-encoding* f = ($\forall x y. \text{opt-prefix } (f x) (f y) \longrightarrow x = y$)

An encoding function is represented as partial functions into lists of booleans, where each list element represents a bit. Such a function is defined to be an encoding, if it is prefix-free on its domain. This is similar to the formalization by Hibon and Paulson [4] except for the use of partial functions for the practical reasons described in Section 1.

lemma *is-encodingI*:
assumes $\bigwedge x x' y y'. e x = \text{Some } x' \Longrightarrow e y = \text{Some } y' \Longrightarrow$
 $\text{prefix } x' y' \Longrightarrow x = y$
shows *is-encoding* e
 $\langle \text{proof} \rangle$

lemma *is-encodingI-2*:
assumes $\bigwedge x y. \text{opt-comp } (e x) (e y) \Longrightarrow x = y$
shows *is-encoding* e
 $\langle \text{proof} \rangle$

lemma *encoding-triv*: *is-encoding* Map.empty
 $\langle \text{proof} \rangle$

lemma *is-encodingD*:
assumes *is-encoding* e
assumes *opt-comp* (e x) (e y)
shows $x = y$
 $\langle \text{proof} \rangle$

lemma *encoding-imp-inj*:
assumes *is-encoding* f
shows *inj-on* f (dom f)
 $\langle \text{proof} \rangle$

fun *bit-count* :: bool list option \Rightarrow ereal **where**
bit-count None = ∞ |
bit-count (Some x) = ereal (length x)

lemma *bit-count-finite-imp-dom*:
bit-count (f x) < $\infty \Longrightarrow x \in \text{dom } f$
 $\langle \text{proof} \rangle$

lemma *bit-count-append*:
bit-count (opt-append x y) = *bit-count* x + *bit-count* y
 $\langle \text{proof} \rangle$

3 (Dependent) Products

definition *encode-dependent-prod* ::
'a encoding \Rightarrow ('a \Rightarrow 'b encoding) \Rightarrow ('a \times 'b) encoding
(**infixr** \bowtie_e 65)
where
 encode-dependent-prod e f x =
 opt-append (e (fst x)) (f (fst x) (snd x))

lemma *dependent-encoding*:
 assumes *is-encoding* e1
 assumes $\bigwedge x. x \in \text{dom } e1 \implies \text{is-encoding } (e2\ x)$
 shows *is-encoding* (e1 \bowtie_e e2)
<proof>

lemma *dependent-bit-count*:
 bit-count ((e1 \bowtie_e e2) (x1,x2)) =
 bit-count (e1 x1) + *bit-count* (e2 x1 x2)
<proof>

lemma *dependent-bit-count-2*:
 bit-count ((e1 \bowtie_e e2) x) =
 bit-count (e1 (fst x)) + *bit-count* (e2 (fst x) (snd x))
<proof>

This abbreviation is for non-dependent products.

abbreviation *encode-prod* ::
'a encoding \Rightarrow 'b encoding \Rightarrow ('a \times 'b) encoding
(**infixr** \times_e 65)
where
 encode-prod e1 e2 \equiv e1 \bowtie_e ($\lambda\cdot$. e2)

4 Composition

lemma *encoding-compose*:
 assumes *is-encoding* f
 assumes *inj-on* g {x. p x}
 shows *is-encoding* ($\lambda x. \text{if } p\ x \text{ then } f\ (g\ x) \text{ else } \text{None}$)
<proof>

lemma *encoding-compose-2*:
 assumes *is-encoding* f
 assumes *inj* g
 shows *is-encoding* ($\lambda x. f\ (g\ x)$)
<proof>

5 Natural Numbers

fun *encode-bounded-nat* :: *nat* \Rightarrow *nat* \Rightarrow *bool list* **where**
encode-bounded-nat (*Suc l*) *n* =
 (let *r* = *n* \geq (2^l) in *r* # *encode-bounded-nat l* (*n* - of-bool *r* * 2^l)) |
encode-bounded-nat 0 - = []

lemma *encode-bounded-nat-prefix-free*:
fixes *u v l* :: *nat*
assumes *u* < 2^l
assumes *v* < 2^l
assumes *prefix* (*encode-bounded-nat l u*) (*encode-bounded-nat l v*)
shows *u* = *v*
 <proof>

definition *Nb_e* :: *nat* \Rightarrow *nat* *encoding*
where *Nb_e l n* = (
 if *n* < *l*
 then *Some* (*encode-bounded-nat* (*floorlog 2 (l-1)*) *n*)
 else *None*)

Nb_e l is encoding for natural numbers strictly smaller than *l* using a fixed length encoding.

lemma *bounded-nat-bit-count*:
bit-count (*Nb_e l y*) = (if *y* < *l* then *floorlog 2 (l-1)* else ∞)
 <proof>

lemma *bounded-nat-bit-count-2*:
assumes *y* < *l*
shows *bit-count* (*Nb_e l y*) = *floorlog 2 (l-1)*
 <proof>

lemma *dom* (*Nb_e l*) = {..*l*}
 <proof>

lemma *bounded-nat-encoding: is-encoding* (*Nb_e l*)
 <proof>

fun *encode-unary-nat* :: *nat* \Rightarrow *bool list* **where**
encode-unary-nat (*Suc l*) = *False* # (*encode-unary-nat l*) |
encode-unary-nat 0 = [*True*]

lemma *encode-unary-nat-prefix-free*:
fixes *u v* :: *nat*
assumes *prefix* (*encode-unary-nat u*) (*encode-unary-nat v*)
shows *u* = *v*
 <proof>

definition *Nu_e* :: *nat* *encoding*

where $Nu_e\ n = \text{Some } (\text{encode-unary-nat } n)$

Nu_e is encoding for natural numbers using unary encoding. It is inefficient except for special cases, where the probability of large numbers decreases exponentially with its magnitude.

lemma *unary-nat-bit-count*:

$\text{bit-count } (Nu_e\ n) = \text{Suc } n$
 $\langle \text{proof} \rangle$

lemma *unary-encoding: is-encoding* Nu_e

$\langle \text{proof} \rangle$

Encoding for positive numbers using Elias-Gamma code.

definition $Ng_e :: \text{nat encoding where}$

$Ng_e\ n =$
 $(\text{if } n > 0$
 $\text{then } (Nu_e\ \times_e\ (\lambda r. Nb_e\ (2^{\wedge}r)))$
 $\text{(let } r = \text{floorlog } 2\ n - 1 \text{ in } (r, n - 2^{\wedge}r))$
 $\text{else None})$

Ng_e is an encoding for positive numbers using Elias-Gamma encoding[3].

lemma *elias-gamma-bit-count*:

$\text{bit-count } (Ng_e\ n) = (\text{if } n > 0 \text{ then } 2 * \lfloor \log 2\ n \rfloor + 1 \text{ else } (\infty :: \text{ereal}))$
 $\langle \text{proof} \rangle$

lemma *elias-gamma-encoding: is-encoding* Ng_e

$\langle \text{proof} \rangle$

definition $N_e :: \text{nat encoding where } N_e\ x = Ng_e\ (x+1)$

N_e is an encoding for all natural numbers using exponential Golomb encoding [6]. Exponential Golomb codes are also used in video compression applications [5].

lemma *exp-golomb-encoding: is-encoding* N_e

$\langle \text{proof} \rangle$

lemma *exp-golomb-bit-count-exact*:

$\text{bit-count } (N_e\ n) = 2 * \lfloor \log 2\ (n+1) \rfloor + 1$
 $\langle \text{proof} \rangle$

lemma *exp-golomb-bit-count*:

$\text{bit-count } (N_e\ n) \leq (2 * \log 2\ (\text{real } n+1) + 1)$
 $\langle \text{proof} \rangle$

lemma *exp-golomb-bit-count-est*:

assumes $n \leq m$

shows $\text{bit-count } (N_e\ n) \leq (2 * \log 2\ (\text{real } m+1) + 1)$

$\langle \text{proof} \rangle$

6 Integers

definition $I_e :: \text{int encoding where}$

$$I_e x = N_e (\text{nat } (\text{if } x \leq 0 \text{ then } (-2 * x) \text{ else } (2*x-1)))$$

I_e is an encoding for integers using exponential Golomb codes by embedding the integers into the natural numbers, specifically the positive numbers are embedded into the odd-numbers and the negative numbers are embedded into the even numbers. The embedding has the benefit, that the bit count for an integer only depends on its absolute value.

lemma *int-encoding: is-encoding* I_e

$\langle \text{proof} \rangle$

lemma *int-bit-count: bit-count* $(I_e n) = 2 * \lfloor \log 2 (2*|n|+1) \rfloor + 1$

$\langle \text{proof} \rangle$

lemma *int-bit-count-1:*

assumes $\text{abs } n > 0$

shows $\text{bit-count } (I_e n) = 2 * \lfloor \log 2 |n| \rfloor + 3$

$\langle \text{proof} \rangle$

lemma *int-bit-count-est-1:*

assumes $|n| \leq r$

shows $\text{bit-count } (I_e n) \leq 2 * \log 2 (r+1) + 3$

$\langle \text{proof} \rangle$

lemma *int-bit-count-est:*

assumes $|n| \leq r$

shows $\text{bit-count } (I_e n) \leq 2 * \log 2 (2*r+1) + 1$

$\langle \text{proof} \rangle$

7 Lists

definition Lf_e **where**

$Lf_e e n xs =$

$(\text{if length } xs = n$

$\text{then fold } (\lambda x y. \text{opt-append } y (e x)) xs (\text{Some } [])$

$\text{else None})$

$Lf_e e n$ is an encoding for lists of length n , where the elements are encoding using the encoder e .

lemma *fixed-list-encoding:*

assumes *is-encoding* e

shows *is-encoding* $(Lf_e e n)$

$\langle \text{proof} \rangle$

lemma *fixed-list-bit-count:*

$bit\text{-}count (Lf_e e n xs) =$
 (if length $xs = n$ then $(\sum x \leftarrow xs. bit\text{-}count (e x))$ else ∞)
 $\langle proof \rangle$

definition L_e

where $L_e e xs = (Nu_e \bowtie_e (\lambda n. Lf_e e n)) (length\ xs, xs)$

$L_e e$ is an encoding for arbitrary length lists, where the elements are encoding using the encoder e .

lemma *list-encoding*:

assumes *is-encoding* e

shows *is-encoding* $(L_e e)$

$\langle proof \rangle$

lemma *sum-list-triv-ereal*:

fixes $a :: \text{ereal}$

shows *sum-list* $(map (\lambda-. a) xs) = length\ xs * a$

$\langle proof \rangle$

lemma *list-bit-count*:

$bit\text{-}count (L_e e xs) = (\sum x \leftarrow xs. bit\text{-}count (e x) + 1) + 1$

$\langle proof \rangle$

8 Functions

definition *encode-fun* $:: 'a\ list \Rightarrow 'b\ encoding \Rightarrow ('a \Rightarrow 'b)\ encoding$

(**infixr** \rightarrow_e 65) **where**

encode-fun $xs\ e\ f =$

(if $f \in \text{extensional}\ (set\ xs)$

then $(Lf_e e (length\ xs) (map\ f\ xs))$

else $None$)

$xs \rightarrow_e e$ is an encoding for functions whose domain is *set* xs , where the values are encoding using the encoder e .

lemma *fun-encoding*:

assumes *is-encoding* e

shows *is-encoding* $(xs \rightarrow_e e)$

$\langle proof \rangle$

lemma *fun-bit-count*:

$bit\text{-}count ((xs \rightarrow_e e) f) =$

(if $f \in \text{extensional}\ (set\ xs)$ then $(\sum x \leftarrow xs. bit\text{-}count (e (f x)))$
 else ∞)

$\langle proof \rangle$

lemma *fun-bit-count-est*:

assumes $f \in \text{extensional}\ (set\ xs)$

assumes $\bigwedge x. x \in set\ xs \implies bit\text{-}count (e (f x)) \leq a$

shows $\text{bit-count } ((xs \rightarrow_e e) f) \leq \text{ereal } (\text{real } (\text{length } xs)) * a$
 <proof>

9 Finite Sets

definition $S_e :: 'a \text{ encoding} \Rightarrow 'a \text{ set encoding}$ **where**

$S_e e S =$
 (if finite $S \wedge S \subseteq \text{dom } e$
 then $(L_e e (\text{linorder.sorted-key-list-of-set } (\leq) (the \circ e) S))$
 else None)

$S_e e$ is an encoding for finite sets whose elements are encoded using the encoder e .

lemma *set-encoding:*

assumes *is-encoding* e
shows *is-encoding* $(S_e e)$

<proof>

lemma *set-bit-count:*

assumes *is-encoding* e
shows $\text{bit-count } (S_e e S) = (\text{if finite } S \text{ then } (\sum x \in S. \text{bit-count } (e x) + 1) + 1 \text{ else } \infty)$

<proof>

lemma *sum-triv-ereal:*

fixes $a :: \text{ereal}$
assumes *finite* S
shows $(\sum - \in S. a) = \text{card } S * a$

<proof>

lemma *set-bit-count-est:*

assumes *is-encoding* f
assumes *finite* S
assumes $\text{card } S \leq m$
assumes $0 \leq a$
assumes $\bigwedge x. x \in S \implies \text{bit-count } (f x) \leq a$
shows $\text{bit-count } (S_e f S) \leq \text{ereal } (\text{real } m) * (a + 1) + 1$

<proof>

10 Floating point numbers

definition F_e **where** $F_e f = (I_e \times_e I_e)$ (*mantissa* f , *exponent* f)

lemma *float-encoding:*

is-encoding F_e

<proof>

lemma *suc-n-le-2-pow-n:*

```

fixes  $n :: nat$ 
shows  $n + 1 \leq 2 \wedge n$ 
<proof>

```

```

lemma float-bit-count-1:
   $bit\_count (F_e f) \leq 6 + 2 * (\log 2 (|mantissa f| + 1) +$ 
     $\log 2 (|exponent f| + 1))$  (is ?lhs  $\leq$  ?rhs)
<proof>

```

The following establishes an estimate for the bit count of a floating point number in non-normalized representation:

```

lemma float-bit-count-2:
  fixes  $m :: int$ 
  fixes  $e :: int$ 
  defines  $f \equiv float\_of (m * 2^{powr e})$ 
  shows  $bit\_count (F_e f) \leq$ 
     $6 + 2 * (\log 2 (|m| + 2) + \log 2 (|e| + 1))$ 
<proof>

```

```

lemma float-bit-count-zero:
   $bit\_count (F_e (float\_of 0)) = 2$ 
<proof>

```

end

11 Examples

```

theory Examples
imports Prefix-Free-Code-Combinators
begin

```

The following introduces a few examples for encoders:

```

notepad
begin
  define example1 where  $example1 = N_e \times_e N_e$ 

```

This is an encoder for a pair of natural numbers using exponential Golomb codes.

Given a pair it is possible to estimate the number of bits necessary to encode it using the *bit-count* lemmas.

```

have  $bit\_count (example1 (0,1)) = 4$ 
by (simp add:example1-def dependent-bit-count exp-golomb-bit-count-exact)

```

Note that a finite bit count automatically implies that the encoded element is in the domain of the encoding function. This

means usually it is possible to establish a bound on the size of the datastructure and verify that the value is encodable simultaneously.

```

hence (0,1) ∈ dom example1
by (intro bit-count-finite-imp-dom, simp)

```

```

define example2
where example2 = [0..<42] →e Nbe 314

```

The second example illustrates the use of the combinator (\rightarrow_e), which allows encoding functions with a known finite encodable domain, here we assume the values are smaller than $314::'a$ on the domain $\{..<42::'a\}$.

```

have bit-count (example2 f) = 42*9 (is ?lhs = ?rhs)
if a:f ∈ {0..<42} →E {0..<314} for f
proof -
have ?lhs = (∑ x←[0..<42]. bit-count (Nbe 314 (f x)))
using a by (simp add:example2-def fun-bit-count PiE-def)
also have ... = (∑ x←[0..<42]. ereal (floorlog 2 313))
using a Pi-def PiE-def bounded-nat-bit-count
by (intro arg-cong[where f=sum-list] map-cong, auto)
also have ... = ?rhs
by (simp add: compute-floorlog sum-list-triv)
finally show ?thesis by simp
qed

```

```

define example3
where example3 = Ne ⋈e (λn. [0..<42] →e Nbe n)

```

The third example is more complex and illustrates the use of dependent encoders, consider a function with domain $\{..<42\}$ whose values are natural numbers in the interval $\{... Let us assume the bound is not known in advance and needs to be encoded as well. This can be done using a dependent product encoding, where the first component encodes the bound and the second component is an encoder parameterized by that value.$

end

end

References

- [1] J. C. Blanchette. The textbook proof of huffman's algorithm. *Archive of Formal Proofs*, Oct. 2008. <https://isa-afp.org/entries/Huffman.html>, Formal proof development.

- [2] C. Bormann and P. E. Hoffman. Concise Binary Object Representation (CBOR). RFC 8949, Dec. 2020.
- [3] P. Elias. Universal codeword sets and representations of the integers. *IEEE Transactions on Information Theory*, 21(2):194–203, 1975.
- [4] Q. Hibon and L. C. Paulson. Source coding theorem. *Archive of Formal Proofs*, Oct. 2016. https://isa-afp.org/entries/Source_Coding_Theorem.html, Formal proof development.
- [5] I. E. Richardson. *H.264 Transform and Coding*, chapter 7, pages 179–221. John Wiley & Sons, Ltd, 2010.
- [6] J. Teuhola. A compression method for clustered bit-vectors. *Information Processing Letters*, 7(6):308–311, 1978.