

CRAM format specification (version 4.0)

samtools-devel@lists.sourceforge.net

4 Apr 2019

The master version of this document can be found at <https://github.com/samtools/hts-specs>.
This printing is version 924537d-dirty from that repository, last modified on the date shown above.

license: Apache 2.0

Range coding

The range coder is a byte-wise arithmetic coder that operates by repeatedly reducing a probability range (for example 0.0 to 1.0) one symbol (byte) at a time with the complete compressed data can be represented by any value within the final range.

This is easiest demonstrated with a worked example, so let us imagine we have an alphabet of 4 symbols, ‘t’, ‘c’, ‘g’, and ‘a’ with probabilities 0.2, 0.3, 0.3 and 0.2 respectively. We can construct a cumulative distribution table and apply probability ranges to each of the symbols:

Symbol	Probability	Range low	Range high
t	0.2	0.0	0.2
c	0.3	0.2	0.5
g	0.3	0.5	0.8
a	0.2	0.8	1.0

As a *conceptual example* (note: this is not how it is implemented in practice, see below) using arbitrary precision floating point mathematics this could operate as follows.

If we wish to encode a message, such as “cat” then we will encode one symbol at a time (‘c’, ‘a’, ‘t’) successively reducing the initial range of 0.0 to 1.0 by the cumulative distribution for that symbol. At each point the new range is adjusted to be the proportion of the previous range covered by the cumulative symbol range. See the table footnotes below for the worked mathematics.

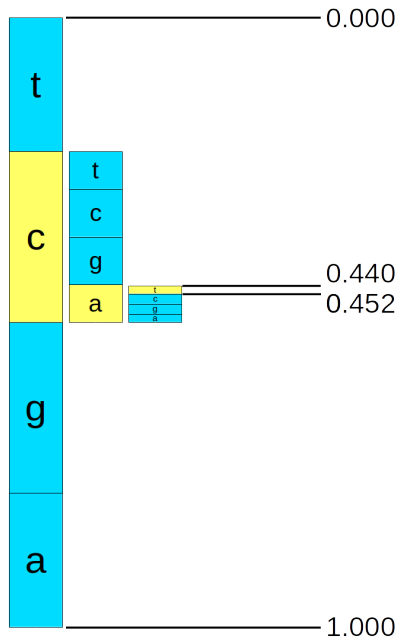
Range low	Range high	Symbol	Symbol low	Symbol high
0.000	1.000	c	0.2	0.5
0.200	0.500	a	0.8	1.0
0.440 ^a	0.500 ^a	t	0.0	0.2
0.440	0.452	<end>		

a. Old range 0.2 to 0.5 plus symbol range 0.8 to 1.0 gives an updated range of 0.44 to 0.5:

$$0.2 + 0.8 \times (0.5 - 0.2) = 0.44$$

$$0.2 + 1.0 \times (0.5 - 0.2) = 0.50$$

Our final range is 0.44 to 0.452 with any value in that range representing “cat”, thus 0.45 would suffice. A pictorial example of this process is below.



Decoding is simply the reverse of this. In the above picture we can see that 0.45 would read off ‘c’, ‘a’ and ‘t’ by repeatedly comparing the symbol ranges to the current range and using those to identify the symbol and produce a new range.

Range low	Range high	Fraction into range	Symbol
0.000	1.000	0.450	c
0.200	0.500	0.833 ^a	a
0.440 ^b	0.500	0.167	t

a. 0.45 into range 0.2 to 0.5: $(0.45 - 0.2)/(0.5 - 0.2) = 0.833$.

This falls within the 0.8 to 1.0 symbol range for ‘a’.

b. ‘a’ symbol range 0.8 to 1.0 applied to range 0.2 to 0.5: $0.2 + 0.8 \times (0.5 - 0.2) = 0.44$ and $0.2 + 1.0 \times (0.5 - 0.2) = 0.5$.

Note: The above example not how the actual implementation works¹. For efficiency, we use integer values having a starting range of 0 to $2^{32} - 1$. We write out the top 8-bits of the range when low and high become the same value. Special care needs to be taken to handle small values are are numerically close but straddling a top byte value, such as 0x37ffba20 to 0x38000034. The decoder does not need to do anything special here, but the encoder must track the number of 0xff or 0x00 values to emit in order to avoid needing arbitrary precision integers.

Pseudocode for the range codec decoding follows. This implementation uses code (next few bytes in the current bit-stream) and range instead of low and high, both 32-bit unsigned integers. This specification focuses on decoding, but given the additional complexity of the precision overflows in encoder we describe this implementation too.

RANGECREATE initialises the range coder, reading the first bytes of the compressed data stream.

```

1: procedure RANGECREATE
2:   range  $\leftarrow 2^{32} - 1$                                 ▷ Maximum 32-bit unsigned value
3:   code  $\leftarrow 0$                                          ▷ 32-bit unsigned
4:   for i  $\leftarrow 0$  to 5 do
5:     code  $\leftarrow (code \ll 8) + \text{READBYTE}$ 
6:   end for
7: end procedure

```

Decoding each symbol is in two parts; getting the current frequency and updating the range.

```

1: function RANGEGETFREQ(tot_freq)
2:   range  $\leftarrow range / tot\_freq$ 

```

¹This implementation was designed by Eugene Shelwein, based on Michael Schindler's earlier work.

```

3:   return code/range
4: end function
1: procedure RANGEDECODE(sym_low, sym_freq, tot_freq)
2:   code  $\leftarrow$  code - sym_low  $\times$  range
3:   range  $\leftarrow$  range  $\times$  sym_freq
4:   while range <  $2^{24}$  do                                      $\triangleright$  Renormalise
5:     range  $\leftarrow$  range << 8
6:     code  $\leftarrow$  (code << 8) + READBYTE
7:   end while
8: end procedure

```

As mentioned above, the encoder is more complex as it cannot shift out the top byte until it has determined the value. This can take a considerable while if our current low / high (*low* + *range*) are very close but span a byte boundary, such as 0x37ffba20 to 0x38000034, where ultimately we will later emit either 0x37 or 0x38. To handle this case, when the range gets too small but the top bytes still differ, the encoder caches the top byte of low (0x37) and keeps track of how many 0xff or 0x00 values will need to be written out once we finally observe which value the range has shrunk to.

The RANGEENCODE function is a straight forward reversal of the RANGEDECODE, with the exception of the special code for shifting the top byte out of the *low* variable.

```

1: procedure RANGEENCODE(sym_low, sym_freq, tot_freq)
2:   old_low  $\leftarrow$  low
3:   range  $\leftarrow$  range / tot_freq
4:   low  $\leftarrow$  low + sym_low  $\times$  range
5:   range  $\leftarrow$  range  $\times$  sym_freq
6:   if low < old_low then
7:     carry  $\leftarrow$  1                                              $\triangleright$  overflow
8:   end if
9:   while range <  $2^{24}$  do                                      $\triangleright$  Renormalise
10:    range  $\leftarrow$  range << 8
11:    RANGESHIFTLOW
12:  end while
13: end procedure

```

RANGESHIFTLOW is the main heart of the encoder renormalisation. It tracks the total number of extra bytes to emit and *carry* indicates whether they are a string of 0xFF or 0x00 values.

```

1: procedure RANGESHIFTLOW
2:   if low < 0xff000000 or carry  $\neq$  0 then
3:     if carry = 0 then
4:       WRITEBYTE(cache)                                          $\triangleright$  top byte cache plus FFs
5:       while FFnum > 0 do
6:         WRITEBYTE(0xff)
7:         FFnum  $\leftarrow$  FFnum - 1
8:       end while
9:     else
10:      WRITEBYTE(cache + 1)                                        $\triangleright$  top byte cache + 1 plus 00s
11:      while FFnum > 0 do
12:        WRITEBYTE(0)
13:        FFnum  $\leftarrow$  FFnum - 1
14:      end while
15:    end if
16:    cache  $\leftarrow$  low >> 24                                      $\triangleright$  Copy of top byte ready for next flush
17:    carry  $\leftarrow$  0
18:  else
19:    FFnum  $\leftarrow$  FFnum + 1
20:  end if
21:  low  $\leftarrow$  low << 8

```

22: **end procedure**

For completeness, the Encoder initialisation and finish functions are below.

```
1: procedure RANGEENCODESTART
2:   low       $\leftarrow 0$ 
3:   range     $\leftarrow 2^{32} - 1$ 
4:   FFnum  $\leftarrow 0$ 
5:   carry     $\leftarrow 0$ 
6:   cache     $\leftarrow 0$ 
7: end procedure
1: procedure RANGEENCODEEND
2:   for i  $\leftarrow 0$  to 5 do                                 $\triangleright$  Flush any residual state in low
3:     RANGE SHIFTLOW
4:   end for
5: end procedure
```

Statistical Modelling

The probabilities passed to the range coder may be fixed for all scenarios (as we had in the “cat” example), or they may be adaptive and context aware. For example the letter ‘u’ occurs around 3% of time in English text, but if the previous letter was ‘q’ it is close to 100% and if the previous letter was ‘u’ it is close to 0%. Using the previous letter is known as an Order-1 entropy encoder, but the context can be anything. We can also adaptively adjust our probabilities as we encode or decode, learning the likelihoods and thus avoiding needing to store frequency tables in the data stream covering all possible contexts.

To do this we use a statistical model, containing an array of symbols S and their frequencies F . The sum of these frequencies must be less than $2^{16} - 32$. When they get too high, they are renormalised by approximately halving the frequencies (ensuring none drop to zero).

Typically an array of models are used where the array index represents the current context.

To encode any symbol the entropy encoder needs to know the frequency of the symbol to encode, the cumulative frequencies of all symbols prior to this symbol, and the total of all frequencies. For decoding a cumulative frequency is obtained given the frequency total and the appropriate symbol is found matching this frequency. Symbol frequencies are updated after each encode or decode call and the symbols are kept in order of most-frequent symbol first in order to reduce the overhead of scanning through the cumulative frequencies.

MODELCREATE initialises a model by setting every symbol to have a frequency of 1. (At no point do we permit any symbol to have zero frequency.)

```
1: procedure MODELCREATE(num_sym)
2:   total_freq  $\leftarrow$  num_sym
3:   max_sym  $\leftarrow$  num_sym - 1
4:   for i  $\leftarrow 0$  to max_sym do
5:      $S_i \leftarrow i$ 
6:      $F_i \leftarrow 1$ 
7:   end for
8: end procedure
```

MODELDECODE is called once for each decoded symbol. It returns the next symbol and updates the model frequencies automatically.

```
1: function MODELDECODE(rc)
2:   freq  $\leftarrow$  rc.RANGEGETFREQUENCY(total_freq)
3:   x  $\leftarrow 0$ 
4:   acc  $\leftarrow 0$ 
5:   while acc +  $F_x \leq$  freq do
6:     acc  $\leftarrow$  acc +  $F_x$ 
7:     x  $\leftarrow$  x + 1
8:   end while
9:   rc.RANGEDECODE(acc,  $F_x$ , total_freq)
10:   $F_x \leftarrow F_x + 8$                                  $\triangleright$  Update model frequencies
```

```

11:   $total\_freq \leftarrow total\_freq + 8$ 
12:  if  $total\_freq > 2^{16} - 32$  then
13:      MODELRENORMALISE
14:  end if
15:   $sym \leftarrow S_x$ 
16:  if  $x > 0$  and  $F_x > F_{x-1}$  then
17:      SWAPELEMENT( $F, x, x - 1$ ) ▷ swap  $F_x$  with  $F_{x-1}$ 
18:      SWAPELEMENT( $S, x, x - 1$ ) ▷ swap  $S_x$  with  $S_{x-1}$ 
19:  end if
20:  return  $sym$ 
21: end function

1: procedure SWAPELEMENT( $A, i, j$ )
2:    $tmp \leftarrow A_i$ 
3:    $A_i \leftarrow A_j$ 
4:    $A_j \leftarrow tmp$ 
5: end procedure

```

MODELRENORMALISE is called whenever the total frequencies get too high. The frequencies are halved, taking sure to avoid any zero frequencies being created.

```

1: procedure MODELRENORMALISE
2:    $total\_freq \leftarrow 0$ 
3:   for  $i \leftarrow 0$  to  $max\_sym$  do
4:        $F_i \leftarrow F_i - (F_i \text{ div } 2)$ 
5:        $total\_freq \leftarrow total\_freq + F_i$ 
6:   end for
7: end procedure

```

Order-0 and Order-1 Encoding

We can combine the model defined above and the range coder to provide a simple function to perform Order-0 entropy decoder.

```

1: function DECODEORDER0( $len$ )
2:    $max\_sym \leftarrow \text{READBYTE}$ 
3:    $model\_lit \leftarrow \text{MODELCREATE}(max\_sym)$ 

4:   for  $i \leftarrow 0$  to  $len - 1$  do
5:        $out_i \leftarrow model\_lit.\text{MODELDECODE}(rc)$ 
6:   end for
7:   return  $out$ 
8: end function

```

The Order-1 variant simply uses an array of models and selects the appropriate model based on the previous value encoded or decoded. This array index is our “context”.

```

1: function DECODEORDER1( $len$ )
2:    $max\_sym \leftarrow \text{READBYTE}$ 
3:   for  $i \leftarrow 0$  to  $max\_sym - 1$  do
4:        $model\_lit_i \leftarrow \text{MODELCREATE}(max\_sym)$ 
5:   end for

6:    $last \leftarrow 0$ 
7:   for  $i \leftarrow 0$  to  $len - 1$  do
8:        $out_i \leftarrow model\_lit_{last}.\text{MODELDECODE}(rc)$ 
9:        $last \leftarrow out_i$ 
10:  end for
11:  return  $out$ 
12: end function

```

RLE with Order-0 and Order-1 Encoding

The DECODEORDER0 and DECODEORDER1 codecs can be expanded to include a count of how many runs of each symbol should be decoded. Both order 0 and order 1 variants are possible.

After the symbol is decoded, the run length must be decoded to indicate how many *extra* copies of this symbol occur. Long runs are broken into a series of lengths of no more than 3. If length 3 is decoded it indicates we must decode an additional length and add to the current one. The context used for the run length model is the symbol itself for the initial run, 256 for the first continuation run (if ≥ 4) and 257 for any further continuation runs. Thus encoding 10 ‘A’ characters would first store symbol ‘A’ followed by run length 3 (with context ‘A’), length 3 (context 256), length 3 (context 257), and length 1 (context 258).

For example, if we have the string “RRRRUNN” we will decode symbol ‘R’ run 3, symbol ‘U’ run 0, symbol ‘N’ run 1.

```
1: function DECODERLE0(len)
2:   max_sym  $\leftarrow$  READBYTE
3:   if max_sym = 0 then
4:     max_sym  $\leftarrow$  256
5:   end if
6:   model_lit  $\leftarrow$  MODELCREATE(max_sym)
7:   for i  $\leftarrow$  0 to 257 do
8:     model_runi  $\leftarrow$  MODELCREATE(4)
9:   end for

10:  i  $\leftarrow$  0
11:  while i < len do
12:    outi  $\leftarrow$  model_lit.MODELDECODE(rc)
13:    part  $\leftarrow$  model_runouti.MODELDECODE(rc)
14:    run  $\leftarrow$  part
15:    rctx  $\leftarrow$  256
16:    while part = 3 do
17:      part  $\leftarrow$  model_runrctx.MODELDECODE(rc)
18:      rctx  $\leftarrow$  257
19:      run  $\leftarrow$  run + part
20:    end while
21:    for j  $\leftarrow$  1 to run do
22:      outi+j  $\leftarrow$  outi
23:    end for
24:    i  $\leftarrow$  run + 1
25:  end while
26:  return out
27: end function
```

The order-1 run length variant is identical to order-0 except the previous symbol is used as the context for the next literal. The context for the run length does not change.

```
1: function DECODERLE1(len)
2:   max_sym  $\leftarrow$  READBYTE
3:   for i  $\leftarrow$  0 to max_sym - 1 do
4:     model_liti  $\leftarrow$  MODELCREATE(max_sym)
5:   end for
6:   for i  $\leftarrow$  0 to 257 do
7:     model_runi  $\leftarrow$  MODELCREATE(4)
8:   end for

9:  last  $\leftarrow$  0
10:  i  $\leftarrow$  0
11:  while i < len do
12:    outi  $\leftarrow$  model_litlast.MODELDECODE(rc)
13:    last  $\leftarrow$  outi
```

```

14:   $part \leftarrow model\_run_{last}.MODELDECODE(rc)$ 
15:   $run \leftarrow part$ 
16:   $rctx \leftarrow 256$ 
17:  while  $part = 3$  do
18:     $part \leftarrow model\_run_{rctx}.MODELDECODE(rc)$ 
19:     $rctx \leftarrow 257$ 
20:     $run \leftarrow run + part$ 
21:  end while
22:  for  $j \leftarrow 1$  to  $run$  do
23:     $out_{i+j} \leftarrow last$ 
24:  end for
25:   $i \leftarrow run + 1$ 
26: end while
27: return  $out$ 
28: end function

```

General Purpose Entropy Encoder

We wrap up the Order-0 and 1 entropy encoder, both with and without run length encoding, into a data stream that specifies the type of encoded data and also permits a number of additional transformations to be applied. These transformations support bit packing (for example a data block with only 4 distinct values can be packed with 4 values per byte), no-op for tiny data blocks where entropy encoding would grow the data, dictionary lookups for common 4-byte values and 4-way interleaving of the 8-bit components of a 32-bit value.

Bits	Type	Name	Description
8	uint8	<i>flag</i>	Data format bit field
<i>Unless</i> NO SIZE <i>flag is set:</i>			
?	uint7	<i>ulen</i>	Uncompressed length
<i>If</i> X4 <i>flag is set:</i>			
?	uint7	<i>clen1</i>	Compressed sub-block length 1
?	uint7	<i>clen2</i>	Compressed sub-block length 2
?	uint7	<i>clen3</i>	Compressed sub-block length 3
?	uint7	<i>clen4</i>	Compressed sub-block length 4
?	uint8[]	<i>cdata1</i>	Compressed data sub-block 1 (recurse)
?	uint8[]	<i>cdata2</i>	Compressed data sub-block 2 (recurse)
?	uint8[]	<i>cdata3</i>	Compressed data sub-block 3 (recurse)
?	uint8[]	<i>cdata4</i>	Compressed data sub-block 4 (recurse)
<i>If</i> CAT <i>flag is set (and X4 flag is unset):</i>			
?	uint8[]	<i>udata</i>	Uncompressed data stream
<i>If</i> DICTIONARY <i>flag is set (and neither X4 or CAT flags are set):</i>			
?	uint8[]	<i>dict_meta</i>	Dictionary lookup table
<i>If</i> PACK <i>flag is set (and neither X4 or CAT flags are set):</i>			
?	uint8[]	<i>pack_meta</i>	Pack lookup table
<i>If neither X4 or CAT flags are set:</i>			
?	uint8[]	<i>cdata</i>	Entropy encoded data stream (see ORDER / RLE flags)

The first byte of our generalised data stream is a bit-flag detailing the type of transformations and entropy encoders to be combined, followed by optional meta-data, followed by the actual compressed data stream. The bit-flags are defined below, but note not all combinations are permitted.

Bit AND value	Code	Description
1	ORDER	Order-0 or Order-1 entropy coding.
2	reserved	Reserved (for possible order-2/3)
4	DICT	map 32-bit values to 8-bit
8	X4	4-way interleaving of byte streams.
16	NoSIZE	original size is not recorded (used by X4)
32	CAT	Data is uncompressed
64	RLE	Run length encoding, with runs and literals encoded separately
128	PACK	Pack 2, 4, 8 or infinite symbols per byte.

Of these X4 is the most complex. The uncompressed data must be a multiple of four bytes long. Each 4th byte is sent to its own stream producing 4 interleaved streams, so the 1st stream will hold data from byte 0, 4, 8, etc while the 2nd stream will hold data from byte 1, 5, 9, etc. Each of those four streams is then itself compressed using this compression format. For example an input block of small unsigned 32-bit little-endian numbers may use RLE for the first three streams as they are mostly zero, and a non-RLE Order-0 entropy encoder of the last stream. Normally our data format will include the decoded size, but with X4 we can omit this from the internal four compressed streams as we know they will each be a quarter of the outer stream.

The data layout differs for each of these bit types, as described below in the ARITHDECODE function. Some of these can be used in combination, so the order needs to be observed. The Dict and Pack formats have meta data. This is decoded first (dict then pack as appropriate), before entropy decoding and finally expanding any specified pack and dict transformations in that order. For example value 193 indicates a byte stream should be decoded with an RLE aware order-1 entropy encoder and then unpacked.

```

1: function ARITHDECODE(len)
2:   flags ← READBYTE
3:   if flags AND NoSIZE ≠ 0 then
4:     out_len ← READINT7
5:   end if
6:   if flags AND X4 then
7:     out ← DECODEX4(len) return out
8:   end if
9:   if flags AND CAT then
10:    data ← DECODECAT
11:  end if
12:  if flags AND DICT then
13:    D ← DECODEDICTMETA
14:  end if
15:  if flags AND PACK then
16:    P ← DECODEPACKMETA
17:  end if
18:  if flags AND RLE then
19:    if flags AND ORDER then
20:      data ← DECODERLE1
21:    else
22:      data ← DECODERLE0
23:    end if
24:  else
25:    if flags AND ORDER then
26:      data ← DECODEORDER1
27:    else
28:      data ← DECODEORDER0
29:    end if
30:  end if
31:  if flags AND PACK then
32:    data ← DECODEPACK(data, P)
33:  end if

```

▷ Decode meta-data

▷ Swap order of this with Dict for consistency?

▷ Entropy Decoding

▷ Apply data transformations


```

34:   if flags AND DICT then
35:       data ← DECODEDICT(data, D)
36:   end if
37:   out ← data
38: end function

```

The specifics of each sub-format are described below, in the order (minus meta-data specific shuffling) they are applied.

- **X4**: The byte stream consists of a 7-bit encoded uncompressed length, which must be a multiple of 4, followed by 4 compressed lengths also 7-bit encoded. Each of these compressed byte streams is then itself a valid *cdata* stream and will recurse again, each starting with their own format flag. The total uncompressed byte stream is then an interleaving of one byte in turn from each of the 4 substreams (in order of 1st to 4th). Thus an array of 32-bit unsigned integers could be unpacked using X4 to compress each of the 8-bit components together with their own algorithm.

```

1: function DECODEX4(len)
2:   plen ← out_len / 4
3:   for i ← 0 to 3 do                                     ▷ Fetch 4 compressed lengths
4:       cleni ← READINT7
5:   end for
6:   X0 ← ARITHDECODE(plen)                                  ▷ Decode 4 streams
7:   X1 ← ARITHDECODE(plen)
8:   X2 ← ARITHDECODE(plen)
9:   X3 ← ARITHDECODE(plen)

10:  i ← 0
11:  for j ← 0 to plen - 1 do                                ▷ Interleave
12:      outi+0 ← X0j
13:      outi+1 ← X1j
14:      outi+2 ← X2j
15:      outi+3 ← X3j
16:      i ← i + 4
17:  end for
18:  return out
19: end function

```

- **NoSIZE**: Do not store the size of the uncompressed data stream. This information is not required when the data stream is one of the four sub-streams in the X4 format.
- **CAT**: If present, all other bit flags should be nul, with the possible exception of NoSIZE.

The uncompressed data stream is the same as the compressed stream. This is useful for very short data where the overheads of compressing are too high.

```

1: function DECODECAT(len)
2:   for i ← 0 to len - 1 do
3:       outi ← READBYTE
4:   end for
5:   return out
6: end function

```

- **ORDER**: Bit field defining order-0 (unset) or order-1 (set) entropy encoding, as described above by the DECODEORDER0 and DECODEORDER1 functions.
- **RLE**: Bit field defining whether the Order-0 and Order-1 encoding should also use a run-length. When set, the DECODERLE0 and DECODERLE1 functions will be used instead of DECODEORDER0 and DECODEORDER1.

- **PACK**: Data containing only 1, 2, 4 or 16 distinct values can have multiple values packed into a single byte (infinite, 8, 4 or 2). The distinct symbol values do not need to be adjacent as a mapping table *M* converts quantised value *x* to original symbol *M_x* where *x* is the number of distinct possible symbols.

The packed format is split into uncompressed meta-data (below) and the compressed packed data.

Bytes	Type	Name	Description
1	byte	<i>nsym</i>	Number of distinct symbols
<i>nsym</i>	byte[]	<i>syms</i>	Symbol map

The first meta-dat byte holds *nsym*, the number of distinct values, followed by *nsym* bytes to construct the *M* map. If *nsym* = 1 then the byte stream is a stream of constant values and no bit-packing is done (we know every value already). If *nsym* = 2 then each symbol is 1 bit (8 per byte), if $2 < nsym \leq 4$ symbols are 2 bits each (4 per byte) and if $4 < nsym \leq 16$ symbols are 4 bits each (2 per byte). It is not permitted to have *nsym* > 16 as bit packing is not possible. Bits are unpacked from low to high.

Decoding this meta-data is implemented by the DECODEPACKMETA function.

```

1: function DECODEPACKMETA(M, nsym)
2:   nsym ← READBYTE
3:   for i ← 0 to nsym − 1 do
4:     Mi ← READBYTE
5:   end for
6: end function

```

The meta-data is unpacked at before entropy decoding. Once the main data block has been decoded, the byte stream is then expanded with the DECODEPACK function below.

```

1: function DECODEPACK(len)
2:   if nsym ≤ 1 then                                     ▷ Constant value
3:     for i ← 0 to len − 1 do
4:       outi ← M0
5:     end for

6:   else if nsym ≤ 2 then                                   ▷ 1 bit per value
7:     for i ← 0 to len − 1 do
8:       if i mod 8 = 0 then
9:         v ← READBYTE
10:      end if
11:      outi ← M(v AND 1)
12:      v = v >> 1
13:    end for

14:   else if nsym ≤ 4 then                                   ▷ 2 bits per value
15:     for i ← 0 to len − 1 do
16:       if i mod 4 = 0 then
17:         v ← READBYTE
18:       end if
19:       outi ← M(v AND 3)
20:       v = v >> 2
21:     end for

22:   else if nsym ≤ 16 then                                   ▷ 4 bits per value
23:     for i ← 0 to len − 1 do
24:       if i mod 2 = 0 then
25:         v ← READBYTE
26:       end if
27:       outi ← M(v AND 15)
28:       v = v >> 4
29:     end for

30:   else
31:     ERROR
32:   end if

33:   return out
34: end function

```

- **DICT:**

As with PACK, DICT has disjoint uncompressed meta-data and compressed data. The meta-data describes numeric dictionary.

If our data consists of a few distinct 32-bit values then compression can be gained by using a lookup table and converting these to 8-bit values. We use the term *stride* to indicate the size of elements in the input data (4-byte for 32-bit values).

For example, with values $V = \{0x1234 = 0, 0x1235 = 1, 0x1236 = 2, 0xbeef = 3\}$ we can replace this input data stream:

36	12	0	0	35	12	0	0	ef	be	0	0	36	12	0	0	34	12	0	0
----	----	---	---	----	----	---	---	----	----	---	---	----	----	---	---	----	----	---	---

with a copy of map V (see below for format) and a new data byte-stream:

2	1	3	2	0
---	---	---	---	---

Here we describe the reversal of this; the dictionary decoding method. In the example above our elements of V were all 2-bytes long (16-bit). Currently the DICT format is restricted to a *stride* of 4 and element size 2. The map format itself is described in a meta-data block:

Bytes	Type	Name	Description
$\frac{1}{2}$	nibble(high)	v	Size in bytes for each value in dictionary V (must be 2)
$\frac{1}{2}$	nibble(low)	s	Stride size in bytes for output (must be 4)
1	byte	n	Number of items in dictionary
?	byte[]	V	Dictionary of n v -byte values

The first byte indicates the multi-byte stride size s to map in the low 4 bits and the size v of the multi-byte values in the top 4 bits, where $s \geq v$. Currently these must be 4 and 2 respectively.

The next byte holds the number of values n in the map V . Following this is a run-length encoding of the map. The meta-data stream consists of a byte holding the number of literals in a row in the top 4 bits and the subsequent number of runs of consecutive literals in a row in the bottom 4 bits. For each literal, the next v bytes (little endian) represent the value.

During decode, these literals will be expressed as s -byte little endian values. For each run the value is numerically one higher than before and is not stored.

For example $V = \{1000, 1100, 1200, 1201, 1202, 1203, 1204, 2200\}$ is stored in the meta-data stream as:

24	08	43	e8	03	4c	04	b0	04	01	98	08
v 2	n	run 4	1000		1100		1200		run 0	2200	
s 4		lit 3							lit 1		

```

1: function DECODEDICTMETA( $len$ )
2:    $sv \leftarrow \text{READBYTE}$ 
3:    $stride \leftarrow sv \bmod 16$ 
4:    $vsiz \leftarrow sv \div 16$ 
5:    $n \leftarrow \text{READBYTE}$ 
6:    $i \leftarrow 0$ 
7:   while  $i < n$  do
8:      $x \leftarrow \text{READBYTE}$ 
9:      $lit \leftarrow x \bmod 16$ 
10:     $run \leftarrow x \div 16$ 
11:    for  $j \leftarrow 0$  to  $lit - 1$  do
12:       $D_{i+j} \leftarrow \text{READINT16}$ 
13:    end for
14:     $i \leftarrow i + lit$ 
15:    if  $run > 0$  then
16:      for  $j \leftarrow 0$  to  $run - 1$  do

```

▷ Currently must be 4
 ▷ Currently must be 2
 ▷ For $vsiz = 2$

```

17:          $D_{i+j} \leftarrow D_{i+j-1} + 1$ 
18:     end for
19:      $i \leftarrow i + run$ 
20: end if
21: end while
22: return  $D$ ,  $stride$ 
23: end function

```

The compressed data stream follows (either immediately after this or following the PACK meta-data), encoded according to the other format bit flags. Once the bytes have been uncompressed a new data stream is generated by replacing each byte value x with the unsigned v -byte value V_x , thus growing the byte stream by v times.

```

1: function DECODEDICT( $data$ ,  $out\_len$ ) ▷ Only for  $stride = 4$ 
2:    $i \leftarrow 0$ 
3:   for  $j \leftarrow 0$  to  $out\_len - 1$  do
4:      $v \leftarrow data_j$ 
5:      $out_{i+0} \leftarrow (D_v \gg 0) \text{ AND } 255$ 
6:      $out_{i+1} \leftarrow (D_v \gg 8) \text{ AND } 255$ 
7:      $out_{i+2} \leftarrow (D_v \gg 16) \text{ AND } 255$ 
8:      $out_{i+3} \leftarrow (D_v \gg 24) \text{ AND } 255$ 
9:      $i \leftarrow i + 4$ 
10:  end for
11:  return  $out$ 
12: end function

```

Name tokenisation codec

Sequence names (identifiers) typically follow a structured pattern and compression based on columns within those structures usually leads to smaller sizes.

As an example, take a series of names:

```

I17_08765:2:123:61541:01763#9
I17_08765:2:123:1636:08611#9
I17_08765:2:124:45613:16161#9

```

We may wish to tokenise each of these into 7 tokens, e.g. “I17_08765:2:”, “123”, “:”, “61541”, “:”, “01763” and “#9”. Some of these are multi-byte strings, some single characters, and some numeric, possibly with a leading zero. We also observe some regularly have values that match the previous line (the initial prefix string, colons, “#9”) while others are numerically very close to the value in the previous line (124 vs 123).

The name tokeniser compares each name against a previous name (which is not necessarily the one immediately prior) and encodes this name either as a series of differences to the previous name or marking it as an exact duplicate. A maximum of 128 tokens is permitted within any single read name.

The tokens and values are stored in an array $T_{pos,type}$ of byte streams, where pos 0 is reserved for name meta-data (whether a duplicate name) and pos 1 onwards is for the first, second and later tokens. $Type$ is one of the token types listed below, corresponding to the type of data being stored. Some token types may also have associated values. Type TYPE (0) holds the token type itself and that type is then used to retrieve the associated value(s) if appropriate. Thus multiple types at the same token position will have their values encoded in distinct data streams, e.g. is position 5 is of type either DIGITS or DDELTA then data streams will exist for $T_{5,TYPE}$, $T_{5,DIGITS}$ and $T_{5,DDELTA}$. Decoding per name continues until a token of type END is observed.

The sequence name (identifier) tokenisation relies heavily on the General Purpose Entropy Encoder described above.

A simplistic pseudocode for decoding the n^{th} individual name is:

```

1: function DECODENAME( $n$ )
2:    $type \leftarrow GET\_TYPE(0, TYPE)$ 
3:    $dist \leftarrow GET\_INT32(0, TYPE)$ 

```

```

4:  if  $type = \text{DUP}$  then
5:       $name_n \leftarrow name_{n-dist}$ 
6:      return  $name_n$ 
7:  end if

8:   $pos \leftarrow 0$ 
9:  repeat
10:      $pos \leftarrow pos + 1$ 
11:      $type \leftarrow \text{GET\_TYPE}(pos, \text{TYPE})$ 
12:     (Append to  $name_n$  based on  $type$ )
13:  until  $type = \text{END}$ 
14:  return  $name_n$ 
15: end function

```

Token IDs (types) are listed below.

ID	Type	Value	Description
0	TYPE	Type	Used to determine the type of token at a given position.
5	DUP	Integer (distance)	The entire name is a duplicate of an earlier one. Used in position 0 only.
6	DIFF	Integer (distance)	The entire name is differs to earlier ones. Used in position 0 only.
1	STRING	String	A string of characters
2	CHAR	Byte	A single character
7	DIGITS	$0 \leq \text{Int} < 2^{32}$	A numerical value, not containing a leadng zero
4	DIGITS0	$0 \leq \text{Int} < 2^{32}$	A numerical value possibly starting in leading zeros
?	DZLEN	Int length	Length of DIGITS0 token.
11	DDELTA	$0 \leq \text{Int} < 256$	A numeric value being stored as the difference to the value of this token on the previous name
12	DDELTA0	$0 \leq \text{Int} < 256$	As DDELTA, but for numeric values starting with leading zeros
13	MATCH	(none)	This token is identical type and value to the same position in the previous name
14	END	(none)	Marks end of name

More detail on the token types is given below.

- **TYPE**: This is the first token fetched at each token position. It holds the type of the token at this position, which in turn may then required retrieval from type-specific data streams at this position.
For position 0, the TYPE field indicates whether this record is an exact duplicate of a prior read name or has been encoded as a delta to an earlier one.
- **DUP, DIFF**: These types are fetched for position 0, at the start of each new identifier. The value is an integer value describing how many reads before this (with 1 being the immediately previous name) we are comparing against. When we refer to “previous name” below, we always mean the one indicated by the DIFF field and not the one immediately prior to the current name.
By convention the first record will have a DIFF of zero and no delta or match operations are permitted.
- **STRING**: We fetch one byte at a time from the value array, appending to the name buffer until the byte retrieved is zero. The zero byte is not stored in the name buffer. For purposes of token type MATCH, a match is defined as entirely matching the string including the length.
- **CHAR**: Fetch one single byte from the value array and append to the name buffer.
- **DIGITS**: Fetch 4 bytes from the value array and intrepret these as a little endian unsigned integer. This is appended to the name buffer as string of base-10 digits, most significant digit first. Larger values may be represented, but will require multiple DIGITS tokens. Negative values may be encoded by treating the minus sign as a CHAR or STRING and storing the absolute value.
- **DIGITS0, DZLEN**: This fetches the 4 bytes value from $T_{pos, \text{DIGITS0}}$ and a 1 byte length from $T_{pos, \text{DZLEN}}$. As per DIGITS, the value is intrepreted as a little endian unsigned integer. The length indicates the total size of the numeric value when displayed in base 10 which must be greater than $\log_{10}(\text{value})$ with any

remaining length indicating the number of leading zeros. For example if DIGITS0 value is 123 and DZLEN length is 5 the string “00123” must be appended to the name.

For purposes of the MATCH type, both primary and secondary lengths must match.

- **DDELTA**: Fetch a 1 byte value and add this to the DIGITS value from the previous name. It is invalid to have a DDELTA token when the previous name does not have a DIGITS token at this position.

For the purposes of a MATCH type, the DDELTA value is assumed to match meaning we have another increment by the same amount.

- **DDELTA0**: As per DDELTA, but the 1 byte value retrieved is added to the DIGITS0 value in the previous name. No DZLEN value is retrieved, with the length from the previous name being used instead.

For the purposes of a MATCH type, the DDELTA0 value is assumed to match meaning we have another increment by the same amount and display to the same length.

- **MATCH**: This token matches the token at the same position in the previous name. (The previous name is permitted to also have a MATCH token at this position, in which case it recurses to its previous name.)

The definition of MATCH is token type specific, as described above. No value is needed for MATCH tokens.

- **END**: Marks the end of the name. A nul byte is added to the name output buffer. No value is needed for END tokens.

Given a complex name and both position and type specific values, this can lead to many separate data streams. These are serialised into a single byte stream.

The packed data stream starts with two unsigned little endiand 32-bit integers holding the total size of uncompressed name buffer and the number of read names. This is followed the array elements themselves.

Token types, *ttype* holds one of the token ID values listed above in the list above, plus special values to indicate certain additional flags. Bit 6 (64) set indicates that this entire token data stream is a duplicate of one earlier. Bit 7 (128) set indicates the token is the first token at a new position.

The total size of the serialised stream needs to be already known, in order to determine when the token types finish.

Bytes	Type	Name	Description
4	uint32	<i>uncomp_length</i>	Length of uncompressed name buffer
4	uint32	<i>num_reads</i>	Number of read names
<i>For each token data stream</i>			
1	uint8	<i>ttype</i>	Token type code.
<i>If ttype AND 64 (duplicate)</i>			
1	uint8	<i>dup_pos</i>	Duplicate from this token position
1	uint8	<i>dup_type</i>	Duplicate from this token type ID
<i>else if not duplicate</i>			
?	i7	<i>clen</i>	compressed length (7-bit encoding)
<i>clen</i>	<i>cdata</i>	stream	compressed data stream

TODO: write simple pseudocode for this stage.

The *cdata* stream itself is as described in the General Purpose Entropy Encoder section above, with the ARITHDECODE function.

FQZComp quality codec

The FQZComp quality codec uses an adaptive statistical model to predict the next quality value in a given context (comprised of previous quality values, position along this sequence, whether the sequence is the second in a pair, and a running total of number of times the quality has changed in this sequence).

For each quality value, the models produce probabilities for all possible next quality values, which are passed into an arithmetic entropy encoder to encode or decode the actual next quality value. The models are then updated based on the actual next quality in order to learn the statistical properties of the quality data stream. This step wise update process is identical for both encoding and decoding.

The algorithm is a generalisation on the original fqzcomp program, described in *Compression of FASTQ and SAM Format Sequencing Data* by Bonfield JK, Mahoney MV (2013). PLoS ONE 8(3): e59190. <https://doi.org/10.1371/journal.pone.0059190>

FQZComp Models

The FQZComp process utilises knowledge of the read lengths, complement (qualities reversed) status, and a generic parameter selector, but in order to maintain a strict separation between CRAM data series this knowledge is stored (duplicated) within the quality data stream itself. Note the complement model is only needed in CRAM 3.1 as CRAM 4 natively stores the quality in the original orientation already. Both reversed and duplication models have no context and are boolean values.

The parameter selector model also has no context associated with it and encodes *max_sel* distinct values. The selector value may be quantised further using *stab* to reduce the selector to fewer sets of parameters. This is useful if we wish to use the selector bits directly in the context using the same parameters. The selector is arbitrary and may be used for distinguishing READ1 from READ2, as a precalculated “delta” instead of the running total, distinguishing perfect alignments from imperfect ones, or any other factor that is shown to improve quality predictability and increase compression ratio (average quality, number of mismatches, tile, swathe, proximity to tile edge, etc).

The quality model has a 16-bit context used to address an array of 2^{16} models, each model permitting *max_sym* distinct quality values. The context used is defined by the FQZcomp parameters, of which there may be multiple sets, selected using the selector model. There are 4 read length models each having *max_sym* of 256. Each model is used for the 4 successive bytes in a 32-bit length value.

The entropy encoder used is shared between all models, so the bit streams are multiplexed together.

The 16-bit quality value context is constructed by adding sub-contexts together consisting of previous quality values, position along the current record, a running count (per record) of how many times the quality value has differed to the previous one (delta), and an arbitrary stored selector value, each shifted to a defined location within the combined context value (*qloc*, *ploc*, *dloc*, *rloc* and *sloc* respectively). The qual, pos and delta sub-contexts are computed from the previous data while the selector, if used, is read directly from the compressed data stream. The selector may be used to switch parameter sets, or simply to group quality strings into arbitrary user-defined sub-sets. The numeric values for each of these components can be passed through lookup tables (*qtab* for quality, *ptab* for positions, *dtab* for running delta and *stab* for turning the selector *s* into a parameter index *x*). These all convert the monotonically increasing range $0 \rightarrow M$ to a (usually smaller) monotonically increasing $0 \rightarrow N$. For example if we wish to use the approximate position along a 100 byte string, we may uniformly map $0 \rightarrow 99$ to $0 \rightarrow 7$ to utilise 3 bits of our 16-bit combined context.

As some sequencing instruments produce binned qualities, e.g. 0, 10, 25, 35, these values are squashed to incremental values from 0 to *max_sym* - 1 where *max_sym* is the maximum number of distinct quality values observed. If this transform is required, the flag *have_qmap* will be set and a mapping table (*qmap*) will hold the original quality values. The decoded qualities will be the smaller mapped range.

The quality sub-context is constructed by shifting left the previous quality sub-context by *qshift* bits and adding the current quality after passing through the *qmap* squashing process and if defined through the *qtab* lookup table. The quality context is limited to *qbits* long and is added to the combined context starting at bit *qloc*. The quality sub-context is reset to zero at the start of each new record.²

The position context is simply the number of remaining quality values in this record, so is a value starting at record length (minus 1) and decrementing. As with the quality context it may be passed through a lookup table *ptab* before shifting left by *ploc* bits and adding to the combined context.

Delta is a count of the number of times the quality value has changed from one value to a different one. Thus a run of identical values will not increase delta. It gets reset to zero at the start of every record. It may be adjusted by the *dtab* lookup table and is shifted by *dloc* before adding to the combined context.

The selector value may also be used as a sub-context, if the *do_sel* paramter is set. The initial context value (reset per record) is defined within each parameter set, providing a more general purpose alternative to adding

²For example if we have 4 quality values in use - 0, 10, 25 and 35 - we will be encoding quality values 0, 1, 2 and 3. We may wish to define *qbits* to be 6 and *qshift* to be 2 such that the previous 3 quality values can be used as context, for the prediction of the next quality value. There will likely be little reason to use *qtab* in this scenario, but an encoder could define *qtab* to convert {0, 1, 2, 3} to {0, 0, 0, 1} and use *qshift* of 1 instead, giving us knowledge of which of the previous 6 values were maximum quality.

the selector value at a defined location (*sloc*) into the context.

Thus the full context can be updated after each decoded quality with the following pseudocode. Note for brevity this is assuming the various parameters referred are global and updateable.

```

1: function FQZUPDATECONTEXT(params, q)
2:   ctx  $\leftarrow$  params.context ▷ Also the initial value
3:   qctx  $\leftarrow$  (qctx  $\ll$  params.qshift) + qtabq
4:   ctx  $\leftarrow$  ctx + ((qctx AND ( $2^{\text{params.qbits}} - 1$ ))  $\ll$  params.qloc)
5:   if params.pflags AND 32 then ▷ have_ptab
6:     p  $\leftarrow$  MIN(pos, 1023)
7:     ctx  $\leftarrow$  ctx + (ptabp  $\ll$  params.ploc)
8:   end if
9:   if params.pflags AND 64 then ▷ have_dtab
10:    d  $\leftarrow$  MIN(delta, 255)
11:    ctx  $\leftarrow$  ctx + (dtabd  $\ll$  params.dloc)
12:    if prevq  $\neq$  q then
13:      delta  $\leftarrow$  delta + 1
14:    end if
15:    prevq  $\leftarrow$  q
16:  end if
17:  if params.pflags AND 8 then ▷ do_sel
18:    ctx  $\leftarrow$  ctx + (sel  $\ll$  params.sloc)
19:  end if
20:  return ctx AND ( $2^{16} - 1$ )
21: end function

```

In summary context is produced using the following models:

Model	Max symbol	Context size	Description
<i>model_qual</i>	<i>max_sym</i>	2^{16}	Primary model for quality values
<i>model_len</i>	256	4	Read length models with the context 0-3 being successive byte numbers (little endian order)
<i>model_rev</i>	2	none	Used if <i>pflags.do_rev</i> is defined. Indicating which strings to reverse.
<i>model_dup</i>	2	none	Used if <i>pflags.do_dup</i> is defined. Indicates if this whole string is a duplicate of the last one
<i>model_sel</i>	<i>max_sel</i>	none	Used if <i>hflags.multi_param</i> is defined and <i>nparam</i> > 1.

FQZComp Data Stream

The start of an FQZComp data stream consists of the parameters used by the decoder. The data layout is as follows.

DECODEFQZPARAMS below describes the pseudocode for reading the parameter block.

```
1: procedure DECODEFQZPARAMS
2:   vers ← READBYTE
3:   if vers ≠ 5 then
4:     ERROR
5:   end if
6:   gflags ← READBYTE
7:   if gflags AND 1 then                                     ▷ multi_param
8:     nparam ← READBYTE
9:     max_sel ← nparam
10:  else
11:    nparam ← 1
12:    max_sel ← 0
13:  end if
14:  if gflags AND 2 then                                     ▷ have_stab
15:    max_sel ← READBYTE
16:    stab ← READARRAY(256)
17:  end if
18:  max_sym ← 0
19:  for p ← 0 to nparam − 1 do
20:    paramp ← DECODEFQZSINGLEPARAM
21:    if max_sym < paramp.max_sym then
22:      max_sym ← paramp.max_sym                             ▷ Maximum across all param sets
23:    end if
24:  end for
25: end procedure

1: function DECODEFQZSINGLEPARAM
2:   p.context ← READUINT16
3:   p.flags ← READBYTE
4:   p.max_sym ← READBYTE
5:   p.first_len ← 1
6:   x ← READBYTE
7:   p.qbits ← x div 16
8:   p.qshift ← x mod 16
9:   x ← READBYTE
10:  p.qloc ← x div 16
11:  p.sloc ← x mod 16
12:  x ← READBYTE
13:  p.ploc ← x div 16
14:  p.dloc ← x mod 16
15:  if p.flags AND 1 then                                     ▷ Have qmap
16:    for i ← 0 to p.max_sym − 1 do
17:      p.qmapi ← READBYTE
18:    end for
19:  end if
20:  if p.flags AND 128 then                                   ▷ Have qtab
21:    p.qtab ← READARRAY(256)
22:  else
23:    for i ← 0 to 256 do
24:      p.qtabi ← i
25:    end for
26:  end if
27:  if p.flags AND 16 then                                   ▷ Have ptab
28:    p.ptab ← READARRAY(1024)
```

Bits	Type	Name	Description
8	uint8	<i>version</i>	FQZComp format version: must be 5
8	uint8	<i>gflags</i>	Global FQZcomp bit-flags. From lowest bit to highest: 1: <i>multi_param</i> : indicates more than one parameter block is present. Otherwise set <i>nparam</i> = 1 2: <i>have_stab</i> : indicates the parameter selector is mapped through <i>stab</i> . Otherwise set <i>stab_i</i> = <i>i</i> 4: <i>do_rev</i> : <i>model_revcomp</i> will be used. (CRAM v3.1)
<i>If multi_param gflag is set:</i>			
8	uint8	<i>nparam</i>	Number of parameter blocks (defaults to 1)
<i>If have_stab gflag is set:</i>			
8 variable	uint8 table	<i>max_sel</i> <i>stab</i>	Maximum parameter selector value Parameter selector table
<i>Parameter block: repeated nparam times: (selected via model_sel and stab)</i>			
16	uint16	<i>context</i>	Starting context value
8	uint8	<i>pflags</i>	Per-parameter block bit-flags. From lowest bit to highest: 1: Reserved 2: <i>do_dedup</i> : <i>model_dup</i> will be used 4: <i>do_len</i> : <i>model_len</i> will be used for every record. 8: <i>do_sel</i> : <i>model_sel</i> will be used. 16: <i>have_qmap</i> : indicates quality map is present 32: <i>have_ptab</i> : Load <i>ptab</i> , otherwise position contexts are unused 64: <i>have_dtab</i> : Load <i>dtab</i> , otherwise delta contexts are unused 128: <i>have_qtab</i> : Load <i>qtab</i> , otherwise set <i>qtab_i</i> = <i>i</i>
8	uint8	<i>max_sym</i>	Total number of distinct quality values
4	uint4 (high)	<i>qbits</i>	Total number of bits for quality context
4	uint4 (low)	<i>qshift</i>	Left bit shift per successive quality in quality context
4	uint4 (high)	<i>qloc</i>	Bit position of quality context
4	uint4 (low)	<i>sloc</i>	Bit position of selector context
4	uint4 (high)	<i>ploc</i>	Bit position of position context
4	uint4 (low)	<i>dloc</i>	Bit position of delta context
<i>If have_map pflag is set:</i>			
variable	uint8[<i>max_sym</i>]	<i>qmap</i>	Map for unbinning quality values.
<i>If have_qtab pflag is set:</i>			
variable	table	<i>qtab</i>	Quality context lookup table
<i>If have_tab pflag is set:</i>			
variable	table	<i>ptab</i>	Position context lookup table
<i>If have_dtab pflag is set:</i>			
variable	table	<i>dtab</i>	Delta context lookup table

```

29:   end if
30:   if  $p.flags$  AND 64 then ▷ Have dtab
31:      $p.qtab \leftarrow \text{READARRAY}(256)$ 
32:   end if
33:   return  $p$ 
34: end function

```

FQZCREATEMODELS creates the decoder models based on the above parameters and the shared range coder.

```

1: procedure FQZCREATEMODELS
2:    $rc \leftarrow \text{RANGECREATE}$ 
3:   for  $i \leftarrow 0$  to 3 do
4:      $model\_len_i \leftarrow \text{MODELCREATE}(256)$ 
5:   end for
6:   for  $i \leftarrow 0$  to  $2^{16} - 1$  do
7:      $model\_qual_i \leftarrow \text{MODELCREATE}(max\_sym)$ 
8:   end for
9:    $model\_dup \leftarrow \text{MODELCREATE}(2)$ 
10:   $model\_rev \leftarrow \text{MODELCREATE}(2)$ 
11:  if  $max\_sel > 0$  then
12:     $model\_sel \leftarrow \text{MODELCREATE}(max\_sel)$ 
13:  end if
14: end procedure

```

READARRAY reads an array A of size n which maps values 0 to $n - 1$ to a smaller range (0 to $m - 1$), both monotonically increasing. For efficiency this is done using a two-level run length encoding.

Assuming $m < n$ there will be runs of the same value. We measure run lengths for all values (even if they are zero). For example an array $A = \{0, 1, 3, 4, 5, 6, 7, 7, 7\}$ may be converted to run lengths $R = \{1, 1, 0, 1, 1, 1, 1, 4\}$. This array R is no longer monotonically increasing but still has repeated values, so is run-length encoded by storing the number of additional values whenever the last two lengths match. This converts R to $R2 = \{1, 1, +0, 0, 1, 1, +2, 4\}$ where the ‘+’ symbol is shown purely to indicate the values representing the additional run-length copy numbers.

Finally, for coping with runs of 255 or more any run value of 255 is assumed to be part of a larger run and the next value is read and added to run-length until it is no longer 255. For example run length 600 would be represented as 255 255 90.

The final array $R2$ is the stored data stream. The decoder process is the reverse of the above, starting by creating R and then A , for example using the following pseudocode.

```

1: function READARRAY( $n$ )
2:    $i, j, z \leftarrow 0$ 
3:    $last \leftarrow -1$ 
4:   while  $z < n$  do ▷ Convert  $R2$  to  $R$ 
5:      $run \leftarrow \text{READBYTE}$ 
6:      $R_j \leftarrow run$ 
7:      $j \leftarrow j + 1$ 
8:      $z \leftarrow z + run$ 
9:     if  $run = last$  then
10:       $copy \leftarrow \text{READBYTE}$ 
11:      for  $x \leftarrow 1$  to  $copy$  do
12:         $R_j \leftarrow run$ 
13:         $j \leftarrow j + 1$ 
14:      end for
15:       $z \leftarrow z + run \times copy$ 
16:     end if
17:   end while
18:    $i, j, z \leftarrow 0$ 

```

```

19:   while  $z < n$  do                                     ▷ Convert  $R$  to  $A$ 
20:        $run\_len \leftarrow 0$ 
21:       repeat
22:            $part \leftarrow R_j$ 
23:            $j \leftarrow j + 1$ 
24:            $run\_len \leftarrow run\_len + part$ 
25:       until  $part \neq 255$ 
26:       for  $x \leftarrow 1$  to  $run\_len$  do
27:            $A_z \leftarrow i$ 
28:            $z \leftarrow z + 1$ 
29:       end for
30:        $i \leftarrow i + 1$ 
31:   end while

32:   return  $A$ 
33: end function

```

The main loop decodes data in the following order per read: read length (if not fixed), the flag for whether this is read 2 (if needed), a bit flag to indicate if the quality is duplicated (if needed), followed by record length number of quality values using various data gathered since the start of this read as context.

The output of this function is an array of quality values in the variable *output*, indexed with the i^{th} value via *output_i*. The output buffer is a concatenation of all quality values for each record. The record lengths are recorded, but note this is the number of qualities encoded in CRAM for this sequence record and this does not necessarily have to match the number of base calls (for example where qualities are explicitly specified for SNP bases but not elsewhere).

```

1: function DECODEFQZNEWRECORD
2:    $sel \leftarrow 0$ 
3:    $x \leftarrow 0$ 
4:   if  $max\_sel > 0$  then                                     ▷ Find parameter selector
5:        $sel \leftarrow model\_sel.MODELDECODE(rc)$ 
6:       if  $have\_stab$  then
7:            $x \leftarrow stab_{sel}$ 
8:       end if
9:   end if
10:   $param \leftarrow params_x$ 

11:  if  $param.do\_len$  or  $param.first\_len$  then                     ▷ Decode read length
12:       $rec\_len \leftarrow DECODELENGTH(rc)$ 
13:       $param.last\_len \leftarrow rec\_len$ 
14:      if  $param.do\_len = 0$  then
15:           $param.first\_len = 0$ 
16:      end if
17:  else
18:       $rec\_len \leftarrow param.last\_len$ 
19:  end if
20:   $pos \leftarrow rec\_len$ 

21:  if  $param.do\_rev$  then                                       ▷ Check if needs reversal
22:       $rev_{rec} \leftarrow model\_rev.MODELDECODE(rc)$ 
23:       $len_{rec} \leftarrow rec\_len$ 
24:  end if
25:   $rec \leftarrow rec + 1$ 

26:   $is\_dup \leftarrow 0$ 
27:  if  $do\_dedup$  then                                           ▷ Duplicate last string if appropriate
28:      if  $model\_dup.MODELDECODE(rc) > 0$  then
29:           $is\_dup \leftarrow 1$ 
30:      end if

```

```

31:   end if

32:    $qctx \leftarrow 0$ 
33:    $\delta \leftarrow 0$ 
34:    $prevq \leftarrow 0$ 
35:   return  $x$  ▷ Tabulated parameter selector
36: end function

1: procedure DECODEFQZ
2:    $i \leftarrow 0$  ▷ Position in total quality block
3:    $pos \leftarrow 0$  ▷ Remaining base count current quality string
4: next_record:
5:   while  $i < buf\_len$  do
6:     if  $pos = 0$  then ▷ Reset state at start of each new record
7:        $x \leftarrow \text{DECODEFQZNEWRECORD}$ 
8:       if  $is\_dup = 1$  then
9:         for  $j \leftarrow 0$  to  $rec\_len - 1$  do
10:             $output_{i+j} \leftarrow output_{i+j-rec\_len}$ 
11:          end for
12:           $i \leftarrow i + rec\_len$ 
13:           $rec \leftarrow rec + 1$ 
14:          go to next_record
15:        end if

16:         $param \leftarrow params_x$ 
17:         $last \leftarrow param.context$ 
18:      end if

19:       $q \leftarrow model\_qual_{ctx}.\text{MODELDECODE}(rc)$  ▷ Decode a single quality value
20:      if  $param.have\_qmap$  then
21:         $output_i \leftarrow qmap_q$ 
22:      else
23:         $output_i \leftarrow q$ 
24:      end if

25:       $ctx \leftarrow \text{FQZUPDATECONTEXT}(param, q)$  ▷ Also updates qlast, prevq and delta

26:       $i \leftarrow i + 1$ 
27:       $pos \leftarrow pos - 1$ 
28:    end while
29:    if  $do\_rev$  then
30:      REVERSEQUALITIES( $output$ ,  $rev$ ,  $len$ )
31:    end if
32: end procedure

```

Where $\text{MIN}(a, b)$ returns the smallest integer value from a and b .

```

1: function MIN(a,b)
2:   if  $a < b$  then
3:     return  $a$ 
4:   else
5:     return  $a$ 
6:   end if
7: end function

```

Read lengths are encoded as 4 8-bit bytes, each having its own model.

```

1: function DECODELENGTH(rc)
2:    $rec\_len \leftarrow model\_len_0.\text{MODELDECODE}(rc)$ 
3:    $rec\_len \leftarrow rec\_len + (model\_len_1.\text{MODELDECODE}(rc) < 8)$ 
4:    $rec\_len \leftarrow rec\_len + (model\_len_2.\text{MODELDECODE}(rc) < 16)$ 
5:    $rec\_len \leftarrow rec\_len + (model\_len_3.\text{MODELDECODE}(rc) < 24)$ 

```

```

6:   return last_len
7: end function

```

For CRAMv4 quality values are stored in their original FASTQ orientation. For CRAMv3 they are stored in their alignment orientation and it may be beneficial for compression purposes to reverse them first. If so *do_rev* will be set and the REVERSEQUALITIES procedure called below after decoding.

```

1: procedure REVERSEQUALITIES(qual, qual_len, rev, len)
2:   rec  $\leftarrow$  0
3:   for i  $\leftarrow$  0 to qual_len - 1 do
4:     if rev_rec  $\neq$  0 then
5:       k  $\leftarrow$  len_rec - 1
6:       for j  $\leftarrow$  0 to len_rec/2 do
7:         tmp  $\leftarrow$  quali+j
8:         quali+j  $\leftarrow$  quali+k
9:         quali+k  $\leftarrow$  tmp
10:      end for
11:    end if
12:  end for
13: end procedure

```