

NPFCODE Implementation

Yasin OZTURK

December 2018

1 Introduction

NPFCODE encodes given data V by using different D values to output payload (PL) and disambiguation information (DI) in varying sizes. Increasing the D size reduces the size of DI but also increases the size of PL. When D size gets larger than 8, it becomes harder to compute DI and PL because of the byte-aligned memory architecture. NPFCODE splits incoming data by D -bits chunks and outputs each chunk's minimum binary representation (MBR) as payload and corresponding code length in bits as disambiguation information. MBR removes the most significant bit and returns the remaining value with its length. A pseudo-code is given below for calculating MBR value of a unsigned char value.

Algorithm 1: MBR_Char(value, mbrValue, mbrLength)

Input: $value \in [0, 255]$

Output: $mbrValue$: Minimum Binary Representation of $value$.

$mbrValue \in [0, 127]$

$mbrLength$: The number of bits after the most significant bit.

$mbrLength \in [1, 7]$

```
1  $mbrValue = 0, mbrLength = 0;$ 
2 if  $value == 0 \parallel value == 1$  then return ;
3  $orgValue = value;$ 
4 while  $value \neq 1$  do
5    $value = value \gg 1;$ 
6    $mbrLength = mbrLength + 1;$ 
7 end
8  $bitMask = (1 \ll mbrLength) - 1;$ 
9  $mbrValue = orgValue \& bitMask;$ 
```

NPFCODE iterates through V by chunks of D -bits long data and calculates responding MBR values and lengths. Since there can be only 2^D many distinct values to calculate MBR for, these MBR values of every D -bit long data can be pre-processed and saved. It should be noted that resulting MBR lengths can never be equal to D . In order to output PL, all MBR values must be concatenated. Since MBR values can have varying lengths, it is not computationally cheap to concatenate each MBR value. The reason behind that is the resulting concatenated value's length can be different than 8. To resolve the inefficiency caused by shifting concatenated value to fit in 8 bits of data, a Huffman table will be implemented to store the output values of the concatenation and the remaining bits.

Huffman table can be considered as a two-dimensional matrix where the rows of the matrix represent the reverse MBR of the remaining bits of concatenated value and the columns represent the next D-bits length of data. Since remaining bits of data can have at most 7 bits, reverse MBR of the remaining bits can have up to 8 bits of data. The reason behind to get reverse MBR of values is, same MBR values can have different lengths and could be derived from different values. In order to eliminate the ambiguity, it is required to use the original data as an index. It should be noted that 0 and 1 cannot be a valid index since reverse MBR cannot produce these values. When there is no remaining data left, 1 will be used as the index for columns of Huffman table.

$v_1 = 13$ (1101) \Rightarrow MBR(13) = value: 5 (101), length: 3

$v_2 = 37$ (100101) \Rightarrow MBR(37) = value: 5 (00101), length 5

Algorithm 2: ReverseMBR_Char(*mbrValue*, *mbrLength*)

Input: *mbrValue* $\in [0, 128]$

mbrLength $\in [1, 7]$

Output: Reversed MBR value of *mbrValue*, $\in [0, 255]$

1 *return* ($1 \ll mbrLength$) | *mbrValue*;

Each cell of the Huffman table holds MBR length of the corresponding column index (*mbrLength*), output value (*output*), a flag indicating whether output value has more or equal to 8 bits of data (*canOutput*) and reverse MBR value of the remaining bits (*nextTableId*). The pseudo-code for creating the MBR Huffman table for D=8 is as follows:

Algorithm 3: CreateMBRHuffmanTable_D8()

Input: $permAlph[v], v \in [0, 255]$: permutated alphabet

Output: $mT[t, r], t \in [1, 255], r \in [0, 255]$

```
1 for (t = 1; t ≤ 255; t++) do
2   inV = 0, inL = 0 ;          /* Initital value, initial length */
3   if t != 1 then MBR_Char(t, inV, inL) ;
4   for (r = 0; r ≤ 255; r++) do
5     tL = 0, rL = 0, rem = 0 ; /* Total length, remaining length,
6       remaining */
7     v = permAlph[r] ;          /* value */
8     if v ≤ 253 then
9       mV = 0, mL = 0 ;          /* MBR value, MBR length */
10      v = v + 2;
11      MBR_Char(v, mV, mL);
12      mT[t, v].mbrLength = mL;
13      tL = inL + mL;
14      if tL ≥ 8 then
15        mT[t, v].canOutput = 1;
16        rL = tL - 8;
17        mT[t, v].output = (inV << (8 - inL)) | mV >> rL;
18        /* Get the last rL bits of mV, these bits are
19           not outputted and will be used in the next
20           iteration */
21        rem = mV & ((1 << rL) - 1);
22      else
23        mT[t, v].canOutput = 0;
24        rL = tL;
25        rem = (inV << mL) | mV;
26      end
27      mT[t, v].nextTableId = ReverseMBR_Char(rem, rL);
28    else
29      mT[t, v].canOutput = 1;
30      rL = inL;
31      if value == 254 then
32        mT[t, v].output = inV << (8 - inL);
33      else if value == 255 then
34        mT[t, v].output = (inV << (8 - inL)) | (1 >> rL);
35        if rL > 0 then rem = 1;
36      end
37      mT[t, v].nextTableId = ReverseMBR_Char(rem, rL);
38      mT[t, v].mbrLength = 8;
39    end
40  end
41 end
```

Created MBR Huffman table has 256 rows representing the all possible remaining bits from the last D-bit length chunk of V, and has 256 columns representing the all possible values that current D-bit length chunk of V can have. Each cell of the table holds 4 bytes of data. The resulting table's total size can be calculated as 256 kiB for D=8.

MBR Huffman table helps us to construct PL, but in order to decode the PL, DI must be constructed as well. DI holds the corresponding MBR lengths of each D-bit length chunk of V. These length representations is shown in the table below:

MBR Length	Codeword
1	0000001
2	000001
3	00001
4	0001
5	001
6	01
7	1
8	0000000

Table 1: MBR Length representations for D=8 in DI

Creating DI also suffers from the memory architecture. Each MBR length codewords must be concatenated to create DI, since most of the codewords contain less than 8 bits of data, it is not computationally cheap to achieve on-the-go concatenation. However, these DI values can also be pre-processed and stored in a Huffman table. Similar to the previously implemented MBR Huffman table, this table also has 256 rows. Each row value corresponds to the reverse MBR value of the remaining values from previous MBR length codeword representations. Columns of the DI Huffman table represents every possible MBR lengths, ranging from 1 to 8. Each cell of the table holds a flag indicating whether the current concatenated codeword's length is bigger than or equal to 8 (*canOutput*), an 8-bit long output value (*output*), and reverse MBR value of the remaining bits to indicate the next row of the table (*nextTableId*). DI Huffman table takes 6 kiB of space for D=8. The pseudo-code for creating DI Huffman table for D=8 is given below:

Algorithm 4: CreateDisInfoHuffmanTable_D8()

Output: $dT[t, r], t \in [1, 255], r \in [0, 7]$

```
1 for ( $t = 1; t \leq 255; t++$ ) do
2    $inV = 0, inL = 0$  ;           /* Initital value, initial length */
3   if  $t \neq 1$  then MBR_Char( $t, inV, inL$ ) ;
4   for ( $r = 0; r < 8; r++$ ) do
5      $v = 1, l = 8 - (r + 1)$  ;           /* value, length */
6     if  $r == 7$  then  $v = 0, l = 7$  ;
7      $tL = inL + l, rL = 0, rem = 0$  ;    /* Total length, remaining
      length, remaining */
8     if  $tL \geq 8$  then
9        $rem = v$ ;
10       $rL = tL - 8$ ;
11       $dT[t, r].canOutput = 1$ ;
12       $dT[t, r].output = inV << (8 - inL)$ ;
13      if  $rL == 0$  then
14         $dT[t, r].output | = v$ ;
15         $rem = 0$ ;
16      end
17    else
18       $rL = tL$ ;
19       $dT[t, r].canOutput = 0$ ;
20       $rem = (inV << l) | v$ ;
21    end
22     $dT[t, r].nextTableId = Reverse\_MBR(rem, rL)$ ;
23  end
24 end
```

Encoding of a file is done by utilizing previously computed Huffman tables. Each D-bit long chunk of V is read and fed to the MBR Huffman table. If selected cell indicates an available output, the indicated value will be written to the PL. Then, cell's *mbrLength* value is fed to DI Huffman table. This cycle continues till file has been completely traversed.

Some information about the encoding process is placed as a header to the DI in the encoding phase. This header stores information such as remaining data, skip amount, D value and whether encoding is seeded or not. Remaining data signifies how many of the file's last bytes have not been encoded. When D size is 8, this information is irrelevant since all of the data can be traversed byte by byte. But when D is set to 20, there can be cases where the last 2 bytes of V might be left out. In such cases, remaining data is set to signify how many bytes worth of data is not included in the PL. These remaining bytes of V are added to the end of the header. In the decoding phase, these left out bytes of V are added to the end of the decoded file.

In some cases, when all the file has been traversed, DI Huffman table's and MBR Huffman table's current cells might not indicate an available output but there might still be some bits left that need to be outputted. After the traversing of V has been completed, these cells are checked. MBR values of these cells row value is written to corresponding files. It is known that these bytes last $8 - MBRlength$ many bits are 0. When decoding the last written DI byte, this many zeroes might result in incorrect number of codeword lengths. Let's say that last DI Huffman table's cell is in row 2. MBR value of 2 is 0 and MBR length is 1. Last 7 bits of last DI byte are not used and are zero. When D=8, 0000000 bit sequence signifies codeword length of 8. This last codeword length must be skipped when decoding the PL, since there is no corresponding codeword in PL. The number of skips that must be done when reading the last DI byte must be written to the header. Header byte structure is given in the table below.

D Value	Skip Amount	Remaining data	Seeded
000 → D4	00 → No Skip	00 → No Remaining	0 → Not Seeded
001 → D8	01 → Skip 1	01 → Remaining 1	1 → Seeded
010 → D12	10 → Skip 2	10 → Remaining 2	
011 → D16	11 → Skip 3	11 → Remaining 3	
100 → D20			

Table 2: Header byte structure

Algorithm 5: Encode_D8()

Input: *file* : Input file
mT : Pre-computed MBR table
dT : Pre-computed disambiguation information table
Output: *pL* : Payload file
dI : Disambiguation Information file

```
1 fL = length(file);
2 cP = 0, pLI = 0, dII = 4, dITI = 1, mTI = 1 ;    /* current position
   (cP), payload index(pLI), disambiguation information
   index(dII), disambiguation information table index(dITI),
   mbrtable index (mTI) */
3 while cP < fL do
4   v = file[cP];
5   cP = cP + 1;
6   pL[pLI] = mT[mTI, v].output;
7   pLI = pLI + mT[mTI, v].canOutput;
8   mTI = mT[mTI, v].nextTableId;
9   dI[dII] = dT[dITI, mT[mTI, v].mbrLength - 1].output;
10  dII = dII + dT[dITI, mT[mTI, v].mbrLength - 1].canOutput;
11  dITI = dT[dITI, mT[mTI, v].mbrLength - 1].nextTableId;
12 end
13 skip = 0, dType = 32, remCount = 0, seeded = 0 ;    /* skip amount, D
   size, remaining data count */
14 if dITI != 1 then
15   lastDIOut = 0, lastDILen = 0 ;    /* last disambiguation
   information output, last disambiguation information
   length */
16   MBR_Char(dITI, lastDIOut, lastDILen);
17   dI[dII] = lastDIOut << (8 - lastDILen);
18   if mTI == 1 then skip = 8;
19 end
20 if mTI != 1 then
21   lastPLOut = 0, lastPLLen = 0 ;    /* last payload output, last
   payload length */
22   MBR_Char(mTI, lastPLOut, lastPLLen);
23   pL[pLI] = lastPLOut << (8 - lastPLLen)
24 end
25 dI[0] = dType | skip | remCount | seeded;
26 dI[1] = 0, dI[2] = 0, dI[3] = 0;
```

Decoding phase reconstructs the original data from PL and DI. Decoding schema also utilizes pre-computed Huffman tables for efficiency and speed. In this implementation, 3 tables have been used, 2 of them being more significant and used in converting codeword lengths and codewords, the other one for calculating the current bit position in the read PL data.

The table made for parsing DI has only 6 rows, corresponding to the number of zero bits remained from the previous iteration. The remaining zero bit count cannot be 7, because 7 zero bits corresponds to codeword length of 1. Each row has 256 columns for the current DI byte. This table's cell holds the codeword lengths to split the PL. Each cell can hold at most 8 length. When all the bits of the current DI byte is 1, each bit signifies 7 bit long codeword length, which can be at most 8. PL data must be split according to these codeword lengths to achieve decoding. Reverse DI Huffman table takes 17.5 kiB of space for D=8. The pseudo-code for creating the table for D=8 is given below:

Algorithm 6: CreateReverseDisInfoHuffmanTable.D8()

```

Output:  $revDIT[t, r], t \in [0, 6], r \in [0, 255]$ 
1 for ( $t = 0; t \leq 6; t++$ ) do
2    $inL = t;$  /* initial length */
3    $tL = inL + 8;$  /* total length */
4    $bM = 32768;$  /* bit mask : 10000000 00000000 */
5   for ( $r = 0; r \leq 255; r++$ ) do
6      $revDIT[t, r].outputC = 0;$  /* output count */
7      $alignedV = r << (16 - tL);$  /* left aligned value */
8      $zeroC = 0;$  /* zero count */
9     for ( $i = 0; i < tL; i++$ ) do
10       $lmBit = (alignedV \& bM) == 0 ? 0 : 1;$  /* left most bit */
11       $alignedV = alignedV << 1;$ 
12      if  $lmBit == 0$  then
13         $zeroC++;$ 
14        if  $zeroC == 7$  then
15           $revDIT[t, r].outputs[revDIT[t, r].outputC] = 8;$ 
16           $revDIT[t, r].outputC++;$ 
17           $zeroC = 0;$ 
18        end
19      else
20         $revDIT[t, r].outputs[revDIT[t, r].outputC] = 7 - zeroC;$ 
21         $revDIT[t, r].outputC++;$ 
22         $zeroC = 0;$ 
23      end
24    end
25     $revDIT[t, r].nextTableId = zeroC;$ 
26  end
27 end

```

The PL data byte must be split with the codeword lengths obtained from the reverse DI Huffman table. For each byte, there are 8 positions to split the data. All possible combinations of this splitting can be calculated prior to decoding. This pre-computed Decode Huffman table has 256 rows for every possible PL data byte and 8 rows for each possible codeword length. The read byte is split according to the codeword length and remaining bits are stored to find the next cell in the table. Decode Huffman table takes 6 kiB of space for D=8. The pseudo-code for creating the table for D=8 is given below:

Algorithm 7: CreateDecodeHuffmanTable_D8()

Input: $invPermAlph[v], v \in [0, 255]$: inverse permuted alphabet
Output: $decT[t, r], t \in [0, 255], r \in [0, 7]$

```

1 for ( $t = 0; t \leq 255; t++$ ) do
2   for ( $r = 0; r < 8; r++$ ) do
3      $l = r + 1;$ 
4      $code = t >> (8 - l);$ 
5     if  $l \neq 8$  then
6        $codeVal = ReverseMBR\_Char(code, l);$ 
7        $decT[t, r].output = invPermAlph[codeVal - 2];$ 
8     else
9       if  $code == 0$  then  $decT[t, r].output = invPermAlph[254];$  ;
10      else if  $code == 1$  then
11         $decT[t, r].output = invPermAlph[255];$  ;
12      end
13       $decT[t, r].remaining = t << l;$ 
14       $decT[t, r].remainingLength = 8 - l;$ 
15    end
16  end
17 end

```

A trivial Huffman table can be used to determine whether the current PL data is completely decoded or not. This Length Huffman table raises a flag when all the bits of the current byte is traversed and helps to get the next byte of the PL. This table has 8 rows corresponding to the number of bits that are not used in the current PL byte, and 8 columns for the codeword lengths. When the sum of the remaining bits and codeword length passes 8, *getNextByte* value of the cell's is set to 1. This pre-computed table helps to remove if conditions from the main decoding loop. Length Huffman table takes 128 bytes of space for D=8. The pseudo-code for creating Length Huffman table for D=8 is given below:

Algorithm 8: CreateLengthHuffmanTable.D8()

Output: $lenT[t, r], t \in [0, 7], r \in [0, 7]$

```
1 for ( $t = 0; t \leq 7; t++$ ) do
2   for ( $r = 0; r \leq 7; r++$ ) do
3      $tL = t + r + 1$  ;                               /* total length */
4     if  $tL \geq 8$  then
5        $lenT[t, r].getNextByte = 1$ ;
6        $lenT[t, r].remainingLength = tL - 8$ ;
7     else
8        $lenT[t, r].getNextByte = 0$ ;
9        $lenT[t, r].remainingLength = tL$ ;
10    end
11  end
12 end
```

Decoding of PL and DI utilizes previously created Huffman tables for vectorizing the generated code and speeding up the decoding phase. After header bytes of the DI is read, each byte of DI and PL are traversed in a nested loop. First loop reads the current DI byte and determines the codeword lengths using Reverse DI Huffman table. According to the read cell from the table, all of the codeword lengths obtained from this cell is iterated with a for loop. Within the for loop, a payload short is constructed with the previous and current PL byte using Length Huffman table. This payload short is used to determine the next row value of the Decode Huffman table. This table's initial row value is set to the first byte of the PL. Current codeword length is used as the column value for Decode Huffman table. With obtained row and column values, corresponding cell is selected and its *output* value is written to *file*. This cycle continues until all DI values are iterated. The pseudo-code for the decoding schema is given below:

Algorithm 9: Decode.D8()

Input: pL : Payload file
 dI : Disambiguation information file
 $decT$: Pre-computed decode table
 $lenT$: Pre-computed length table
 $revDIT$: Pre-computed reverse disambiguation information table
Output: $file$: Decoded file

```
1  $pLL = \text{length}(pL), dIL = \text{length}(dI);$ 
2  $fileI = 0, pLI = 1;$  /* file index, payload index */
3  $curPLByte = pL[pLI], prevPLByte = pL[pLI - 1];$ 
4  $pLShort = prevPLByte << 8 | curPLByte;$ 
5  $pLShortBI = 0;$  /* payload short bit index */
6  $decTI = pL[0], revDITI = 0, lenTI = 0;$  /* decode table index,
   reverse dis. info. table index, length table index */
7  $hData = dI[0];$  /* header data */
8  $dType = hData >> 5, skip = (hData >> 3) \& 3;$ 
9  $remCount = (hData >> 1) \& 3, seeded = hData \& 3;$ 
10  $dII = 4;$  /* dis. info. index, skip header */
11 while  $dII < dIL$  do
12    $v = dI[dII];$ 
13    $dII = dII + 1;$ 
14    $row = revDIT[revDITI, v];$ 
15   for  $i = 0; i < row.outputC; i++$  do
16     if  $dII == dIL \&\& i == row.outputC - skip$  then break;
17      $l = row.outputs[i];$ 
18      $lenRow = lenT[lenTI, l - 1];$ 
19      $lenTI = lenRow.remainingLength;$ 
20      $pLI = pLI + lenRow.getNextByte;$ 
21      $prevPLByte = curPLByte;$ 
22      $curPLByte = pL[pLI];$ 
23      $pLShort = (prevPLByte << 8) | curPLByte;$ 
24      $decRow = decT[decTI, l - 1];$ 
25      $file[fileI] = decRow.output;$ 
26      $fileI = fileI + 1;$ 
27      $sV = (pLShort << pLShortBI) >> (8 + decRow.remLength);$ 
28      $decTI = decRow.remaining | sV;$ 
29      $pLShortBI = lenRow.remainingLength;$ 
30   end
31    $revDITI = row.nextTableId;$ 
32 end
```

For smaller D values such as 4 and 8, creating the Huffman tables to store pre-processed values work well. When D values gets larger, the space required to store all the tables increases. Using these relatively large tables causes page faults and slows down the encoding and decoding phases. For D values larger than 8, only disambiguation information Huffman table (dT) and reverse disambiguation information Huffman table ($revDIT$) are used. Other pre-processed values are calculated on-the-fly in encoding and decoding loops.