

# Method for Encoding and Decoding Arbitrary Computer Files in DNA Fragments

(Version 2.0: 21 August 2013)

Nick Goldman

EMBL-European Bioinformatics Institute, Hinxton, United Kingdom

## 1 Encoding

**1.1:** An arbitrary computer file is represented as a string  $S_\emptyset$  of bytes (often interpreted as a number between  $\emptyset$  and  $2^8 - 1$ , i.e. a value in the set  $\{\emptyset \dots 255\}$ )\*.

**1.2:**  $S_\emptyset$  is encoded using a given Huffman code, converting it to base-3. This generates the string  $S_1$  of characters in  $\{\emptyset, 1, 2\}$ , each such character called a ‘trit’.

**1.3:** Write  $len()$  for the function that computes the length (in characters) of a string, and define  $n = len(S_1)$ . Represent  $n$  in base-3 and prepend  $\emptyset$ s to generate a string  $S_2$  of trits such that  $len(S_2) = 25$ . Form the string concatenation  $S_4 = S_2 \cdot S_1 \cdot S_3$ , where  $S_3$  is a string of at most 24  $\emptyset$ s chosen so that  $len(S_4)$  is an integer multiple of 25.

**1.4:**  $S_4$  is converted to a DNA string  $S_5$  of characters in  $\{A, C, G, T\}$  with no repeated nucleotides (nt) using the scheme illustrated in Figure 1. (The first trit of  $S_4$  is coded using the ‘A’ row of the table. For each subsequent trit, characters are taken from the row defined by the previous character conversion.)

**1.5:** Define  $N = len(S_5)$ , and let  $ID$  be a 2-trit string identifying the original file and unique within a given experiment (permitting mixing of DNA from upto  $3^2 = 9$  different files  $S_\emptyset$  in one experiment). Split  $S_5$  into overlapping segments of length 100 nt, each offset from the previous by 25 nt. This means there will be  $\frac{N}{25} - 3$  segments, conveniently indexed  $i = \emptyset \dots \frac{N}{25} - 4$ ; segment  $i$  is denoted  $F_i$  and contains (DNA) characters  $25i \dots 25i + 99$  of  $S_5$ .

Each segment  $F_i$  is further processed as follows:

---

\* $\emptyset$  is used throughout to represent the number zero, to avoid confusion with letters o and O.

previous nt written	next trit to encode		
	$\emptyset$	1	2
A	C	G	T
C	G	T	A
G	T	A	C
T	A	C	G

**Figure 1: Base-3 to DNA encoding ensuring no repeated nucleotides.** For each trit  $t$  to be encoded, select the row labelled with the previous nt used and the column labelled  $t$  and encode using the nt in the corresponding table cell.

**1.6:** If  $i$  is odd, reverse complement  $F_i$ .

**1.7:** A keystream method is used to add randomness to the pattern of bases of each segment  $F_i$ . The  $100$  DNA characters of  $F_i$  define 99 single-trit values by converting characters A, C, G, T to numbers  $\emptyset, 1, 2, 3$ ; then forming the successive differences by subtracting (base-4) each such number from its successor ( $d_2 = n_2 - n_1, d_3 = n_3 - n_2 \dots d_{100} = n_{100} - n_{99}$ ); and then subtracting 1 from each of these numbers. (The fact there are no repeated nt in  $F_i$  guarantees that the results of this process are valid trits.) To each of these 99 trits, add (base-3) the corresponding trit from the  $j$ th following random keystream, where  $j = i \bmod 4$ :

**keystream  $\emptyset$ :**  $\emptyset\emptyset 2\emptyset\emptyset\emptyset\emptyset 1\emptyset 11\emptyset 1\emptyset 211111212221\emptyset\emptyset 11122221\emptyset 1\emptyset 1\emptyset 2121222\emptyset 2$

$2\emptyset\emptyset\emptyset 2212\emptyset 1\emptyset 2\emptyset 221\emptyset\emptyset 21211121\emptyset\emptyset\emptyset 212222\emptyset 21211121122221$

**keystream 1:**  $2\emptyset 2\emptyset 2\emptyset 12212121\emptyset 12\emptyset\emptyset\emptyset 12\emptyset\emptyset 21\emptyset 222112\emptyset 2\emptyset 222\emptyset 222222\emptyset 22$

$\emptyset\emptyset\emptyset 1221\emptyset 12111\emptyset 2212112\emptyset 2\emptyset 2\emptyset 222112211122\emptyset 2\emptyset\emptyset 2121\emptyset 22$

**keystream 2:**  $2212211\emptyset 12\emptyset\emptyset 22112\emptyset 22\emptyset\emptyset 11\emptyset\emptyset 22221\emptyset\emptyset\emptyset\emptyset\emptyset 2\emptyset 2\emptyset\emptyset\emptyset 21121\emptyset 2$

$1\emptyset 2\emptyset 1221\emptyset\emptyset\emptyset 212\emptyset 1\emptyset 1\emptyset 21\emptyset 2\emptyset 2\emptyset\emptyset 2\emptyset\emptyset\emptyset 1\emptyset 1\emptyset 2\emptyset\emptyset 221211\emptyset\emptyset 1\emptyset\emptyset$

**keystream 3:**  $1\emptyset\emptyset 1221\emptyset\emptyset\emptyset 111121\emptyset\emptyset 12\emptyset 21\emptyset\emptyset 2\emptyset\emptyset 111\emptyset 22\emptyset 11221221\emptyset\emptyset 1\emptyset\emptyset 12$

$\emptyset 122212\emptyset\emptyset\emptyset\emptyset 2122\emptyset\emptyset 22\emptyset 122\emptyset 22\emptyset 11\emptyset\emptyset\emptyset 1\emptyset 21222211\emptyset 222\emptyset 2\emptyset$

The 99 resulting trits give the ‘randomized’ version of  $F_i$  by the reverse of the procedure described at the start of this step: add 1 to each to generate base-4

differences ( $d'_2 \dots d'_{1000}$ ); then, starting with  $n_1$  (unaltered from above), add (base-4) successive differences to generate 1000 successive values in  $\{\emptyset, 1, 2, 3\}$  ( $n_1, n_1 + d'_2, n_1 + d'_2 + d'_3 \dots n_1 + d'_2 + d'_3 + \dots + d'_{1000}$ ); and then convert values  $\emptyset, 1, 2, 3$  to nt A, C, G, T.

**1.8:** Let  $i\mathcal{B}$  be the base-3 representation of  $i$ , appending enough leading  $\emptyset$ s so that  $len(i\mathcal{B}) = 12$ . Compute  $P$  as the sum (mod 3) of the odd-positioned trits in  $ID$  and  $i\mathcal{B}$ , i.e.  $ID_1 + i\mathcal{B}_1 + i\mathcal{B}_3 + i\mathcal{B}_5 + i\mathcal{B}_7 + i\mathcal{B}_9 + i\mathcal{B}_{11}$ . ( $P$  acts as a ‘parity trit’—analogous to a parity bit—to check for errors in the encoded information about  $ID$  and  $i$ .)

**1.9:** Form the indexing information string  $IX = ID . i\mathcal{B} . P$  (comprising  $2+12+1 = 15$  trits). Append the DNA-encoded version of  $IX$  to  $F_i$  using the same strategy as at step 1.4 above, starting with the code table row defined by the last character of  $F_i$ , to give indexed segment  $F'_i$ .

**1.10:** Form  $F''_i$  by prepending A or T and appending C or G to  $F'_i$ —choosing between A and T, and between C and G, randomly if possible but always such that there are no repeated nt. (This ensures that we can distinguish a DNA segment that has been reverse complemented during DNA sequencing from one that has not—the former will start with G|C and end with T|A; the latter will start A|T and end C|G.)

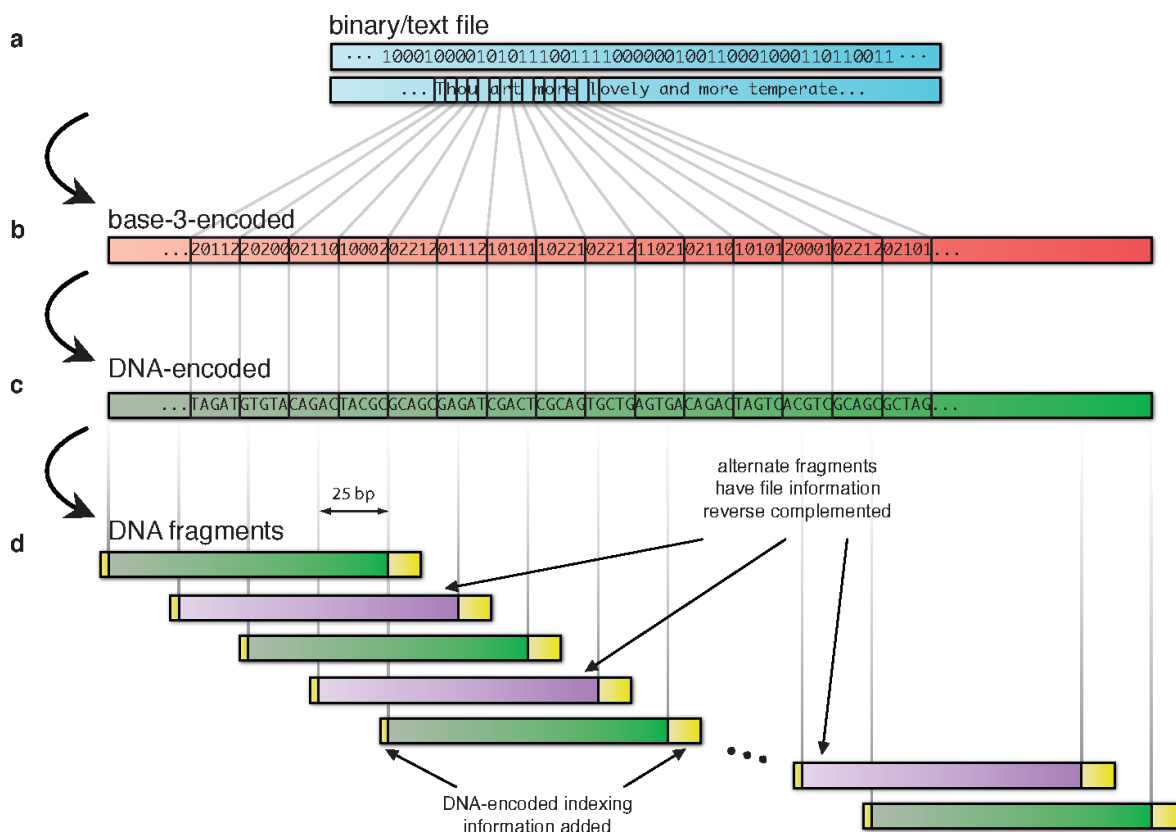
See Figure 2 for a schematic representation of the DNA-encoding of computer files.

**1.11:** The segments  $F''_i$  are synthesized as actual DNA oligonucleotides and stored, and may be supplied for sequencing.

## 2 Example

**2.1:** As it is difficult to represent all possible bytes in this document, we use a simple example of a file comprising just 18 bytes that happen to be easily represented via ASCII codes (for clarity, we show spaces as  $\sqcup$ ):

$S_\emptyset = \text{Birney}\sqcup\text{and}\sqcup\text{Goldman}$



**Figure 2: Schematic of DNA-fragment-encoding of computer files.** The computer file (in any format, e.g. binary or text) shown in blue in (a) (step 1.1 above) is encoded in base-3 (red in b; 1.2–1.3) and then as DNA with no repeated nt (green in c; 1.4). This representation of the complete file is then split into overlapping segments (1.5), alternate segments reverse complemented (mauve; 1.6), and segment indexing and direction-determining information added (yellow; see d; 1.8–1.10). This diagram does not explicitly depict the ‘keystream randomization’ step (1.7 above).

**2.2:** Using the Huffman code defined in the example file `View_huff3.cd`, we convert the bytes of  $S_0$  (shown above/below  $S_1$  for illustrative purposes) into base-3<sup>†</sup>:

<sup>†</sup>Throughout what follows, spaces are not part of base-3 or DNA strings but are included to assist in ‘reading’ how strings have been derived in each step.



$$F_1 = \text{ACTACACAGTCGACTACGCTGTACT GCAGAGTGCTGTCTCACGTGATGAC} \\ \text{GTGCTGCATGATATCTACAGTCATC GTCTATCGAGATACGCTACGTACGT}$$

**2.6:** Only  $i = 1$  is odd, so  $F_1$  is reverse complemented:

$$F_1 = \text{ACGTACGTAGCGTATCTCGATAGACGATGACTGTAGATATCATGCAGCAC} \\ \text{GTCATCACGTGAGACAGCACTCTGCAGTACAGCGTAGTCGACTGTGTAGT}$$

**2.7:** Fragments  $F_\emptyset$  and  $F_1$  are randomized using keystreams  $\emptyset$  and 1 (step 1.7), respectively. Full details of the example are omitted, but after keystream randomization these become:

$$F_\emptyset = \text{CGTGTACGACTCGACAGAGATGCAGAGACTAGACACTCTCTATACGCATC} \\ \text{TCAGACACTCGTAGTCTACTGACATCGCACGAGAGAGAGAGCAGCGCACTCA}$$

$$F_1 = \text{ATATATAGCGTGATAGCGTCAGTCGCATCGCACACGAGCTCTCTTAGCA} \\ \text{TAGCGTCTAGATGCGACGAGCGAGCTCATGTGATCTGCGTACTCTCTACA}$$

**2.8:** For  $i = \emptyset$ ,  $i\mathcal{B} = \emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset$  (length 12) and the sum (mod 3) of the odd-positioned trits of  $ID$  and  $i\mathcal{B}$  is  $P = 1 + \emptyset + \emptyset + \emptyset + \emptyset + \emptyset + \emptyset \pmod{3} = 1$ . For  $i = 1$ ,  $i\mathcal{B} = \emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset 1$  and  $P = 1 + \emptyset + \emptyset + \emptyset + \emptyset + \emptyset + \emptyset \pmod{3} = 1$ .

**2.9:** For  $i = \emptyset$ ,  $IX = ID \cdot i\mathcal{B} \cdot P = 12 \emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset 1$ ; for  $i = 1$ ,  $IX = 12 \emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset 1$ .

So:

$$F'_\emptyset = \text{CGTGTACGACTCGACAGAGATGCAGAGACTAGACACTCTCTATACGCATC} \\ \text{TCAGACACTCGTAGTCTACTGACATCGCACGAGAGAGAGCAGCGCACTCA} \\ \text{GC GTACGTACGTAC T (length } 1\emptyset\emptyset + 15 = 115)$$

$$F'_1 = \text{ATATATAGCGTGATAGCGTCAGTCGCATCGCACACGAGCTCTCTTAGCA} \\ \text{TAGCGTCTAGATGCGACGAGCGAGCTCATGTGATCTGCGTACTCTCTACA} \\ \text{GC GTACGTACGTAG A (length 115)}$$

**2.10:** Prepend A|T and append C|G (note that in this example we have three random choices, but no choice at the start of  $F_1''$ ):

$F_\emptyset'' =$  A CGTGTACGACTCGACAGAGATGCAGAGACTAGACACTCTCTATACGCATC  
 TCAGACACTCGTAGTCTACTGACATCGCACGAGAGAGAGCAGCGCACTCA  
 GCGTACGTACGTACT G (length  $1 + 115 + 1 = 117$ )

$F_1'' =$  T ATATATAGCGTGATAGCGTCAGTCGCATCGCACACGAGCTCTCTCTAGCA  
 TAGCGTCTAGATGCGACGAGCGAGCTCATGTGATCTGCGTACTCTCTACA  
 GCGTACGTACGTAGA C (length 117)

### 3 Decoding

Decoding is simply the reverse of encoding, starting with sequenced DNA fragments  $F_i''$  of length 117 nt. Reverse complementation during the DNA sequencing procedure (e.g. during PCR reactions) can be identified for subsequent reversal by observing whether fragments start with A|T and end with C|G, or start G|C and end T|A. With these two ‘orientation’ nt removed, the remaining 115 nt of each segment can be split into the first 100 ‘message’ nt and the remaining 15 ‘indexing’ nt. The indexing nt can be decoded to determine the file identifier  $ID$  and the position index  $i_3$  and hence  $i$ , and errors may be detected by testing the parity trit  $P$ . The value of  $i$  permits the identification of the correct random keystream, the randomizing effect of which can be removed by applying the keystream encoding (step 1.7) a second time, but *subtracting* the keystream trits. The file identifier and position index then permit reconstruction of the DNA-encoded files, which can then be converted to base-3 using the reverse of the encoding table in step 1.4 above and then to the original bytes using the given Huffman code. Precise details of the decoding procedure are left as an exercise for the reader.

Errors introduced during DNA synthesis, storage or sequencing could lead to various artefacts, particularly nt insertion, deletion or substitution. Recovery of information from fragments with such errors may be possible (details left as exercise)—we have not found this to be necessary due to the large numbers of perfectly-sequenced fragments available via high-throughput sequencing.

## 4 Version history

**Version 2.0:** Updated to reflect changes to representation of file length (step 1.3 above) and introduction of keystream encoding (step 1.7 above), with corresponding changes elsewhere.

**Version 1.0:** First complete, finalized description, as published with the *Nature* paper.

**Versions 0.1–0.4:** Preliminary versions of original description.





```

'false' for subtracting (decrypting)

$L=length($keystream);
if (length($DNA) != $L) {die "\&encrypt error";}

$DNA =~ tr/acgtACGT/01230123/;
@bases=split(//,$DNA);
@kstrits=split(//,$keystream);

foreach (1..$L-1) {
    $i=$L-$_;
    $bases[$i]=($bases[$i]-$bases[$i-1]) % 4; #
    # compute differences
    # (note: this is performed right to left
    # in order to do it in-place)
}

foreach (1..$L-1) {
    $bases[$_]=(($bases[$_-1]) % 3;
    # subtract 1 to recover base-3 trits
}

foreach (1..$L-1) {
    $bases[$_]=(($add) ? ($bases[$_]+$kstrits[$_] % 3 : ($bases[$_-1]-$kstrits[$_] % 3);
    # add/subtract keystream trits as appropriate
}

foreach (1..$L-1) {
    $bases[$_]=(($bases[$_]+1) % 4;
    # add 1 to generate base-4 differences
}

foreach (1..$L-1) {
    $bases[$_]=(($bases[$_]+$bases[$_-1]) % 4;
    # add successively to generate 0,1,2,3 coding
}

$DNA=join("", @bases);
# return to string format

```

```
$DNA =~ tr/0123/ACGT/;
# convert 0,1,2,3 to A,C,G,T
$DNA;
# return result
}
```

1; # need to end with a true value if subroutine to be used via 'require'

---