

ARTICLE TYPE

Indel Error Correction Codes for DNA Digital Data Storage and Retrieval

NallappaBhavithran G | Selvakumar R*

¹ Vellore Institute of Technology,
Vellore, Tamil Nadu, India

Correspondence

*Selvakumar R.

Email: rselvakumar@vit.ac.in

Abstract

A procedure for storage and retrieval of Digital information in DNA strings is discussed by constructing an error correction code which is capable of correcting indel errors and maintaining the stability of DNA strings. For correcting indel errors, Varshamov–Tenengolts algorithm is used. In the VT algorithm, syndrome is calculated using the position values of the retrieved word. Then, the wrongly added or deleted positions can be found by the syndrome and weight of the retrieved word. The stability of DNA strings depends on the proportionate occurrence of guanine (G) and cytosine (C). Reverse complementary property is essential for preventing de-hybridization. To ensure these properties, VT encoded information is mapped to DNA strings by appropriate choices of group homomorphisms. Code for any desired length(n) can be created using the construction methods presented in this paper. The code's reverse complement distance is calculated by using the value of n .

KEYWORDS:

Mathematical biology in general, coding theorem, Kernel code, DNA code, VT codes

1 | INTRODUCTION

In recent decades, the volume of data has expanded exponentially, and this trend will continue. So, to store this much data, a large number of non-biodegradable and well-maintained silica-based storage systems are used. Since DNA-based storage systems are long-lasting and bio-degradable, they can be an alternative for these problems. Humans have successfully recovered data from fossilised sources, while silica-based storage systems have been only effective for almost a century. Moreover, enormous amounts of storage are found in DNA, and according to Castillo, every piece of information on the Internet could be stored in a gram of DNA. Although DNA synthesis and sequencing are now expensive, according to the Molecular Information Storage Program, they will become cheaper in the future¹.

DNA is in helical structure made up of two long strands of nucleotides. Adenine(A), Guanine(G), Cytosine(C), and Thymine(T) are the four bases found in each nucleotide, along with the phosphate and sugar particles. Those two strands of nucleotides are held together by hydrogen bonds between the complementary bases $A = T$ (double bond) and $G \equiv C$ (triple Bonds) and denoted as $A^c = T$, $G^c = C$ and vice versa. Oligos are single-stranded DNA fragments that are synthesised using DNA synthesisers. With phosphate as the backbone, oligos are synthesised in a growing chain of nucleotides by passing one basis at a time. These oligos are commonly synthesised with hundreds of base pairs to optimise the cost and errors. The error rate of the synthesis depends on the number of base pairs present in an oligo. Clelland² found a way to hide information on the human genome. He suggested that it is hard to find the location of the gene to decrypt the information. This paved the way for the creation of DNA based storage system. To create a DNA-based storage system, the binary information has to be encoded into

four base pairs³. Following that, oligomers are made based on those base pair blocks and kept in e. coli form. For the retrieval process, the data held inside is extracted by breaking up the e. coli into smaller pieces. A sequencer can then read the base pairs in each piece and merge them^{4,5}. Here, the primary goal is to develop a reliable coding scheme to minimise errors that might arise while storing and retrieving the information from DNA nucleotides.

Starting from kilobytes on 2012^{6,7} researchers created coding schemes that managed to store and retrieve megabytes^{8,9} of information from DNA. In 2016, Yazdi¹⁰ devised a system to store six universities' budgets that could be randomly accessed and modified. To get random access, they placed addresses on both ends of 27 DNA strands, which were used to store the university information. These addresses contain 20 nucleotide bases and are self-uncorrelated to each other. Additionally, each piece of information is encoded into three nucleotide bases (A, T and C) dependent on their address, and the modification is done using OE-PCR and g-block technologies.

Kiah¹¹ created profile vectors from a field of q elements to store non-binary information on the DNA. For this, the set of all possible values of length l is taken and ordered in lexicographical order. The information is then cut into pieces of size l and mapped to the corresponding profile vector. But, the problem here is to choose the profile vector set efficiently so that no two elements decode to the same number. For decoding, Varshmov's algorithm is applied to the corresponding profile vector of the codeword.

A wide range of algebraic coding properties has been applied to reduce errors in DNA storage. Liu¹² created a cyclic code on a quotient field to inherit reversibility property to DNA code. In^{13,14,15}, they define bounds for different algebraic characteristics of DNA code. In this paper, we proposed a two-layered coding technique. The first layer addresses the Indel errors discussed in section 2, and the second layer is used to satisfy the constraints stated in section 3. In section 4, kernel codes are introduced for creating our DNA code. In section 5, an algorithm is proposed to construct indel error-free code.

2 | TYPES OF ERRORS

2.1 | Synthesizing Errors

Since the primer is created artificially, the occurrence of errors is unavoidable. The basic errors are:

- *Insertion*: When a nucleotide is wrongly inserted into the primer is called insertion.
- *Deletion*: When a nucleotide is deleted from the primer is called a deletion.
- *Substitution*: When a nucleotide is replaced by another in the primer is called substitution.

The place of these error occurrences is unidentifiable and in some instances, insertions and deletions affect the length of the primer. But in the long run, these errors co-occur so that advantage is also removed. The Levenshtein distance is a new distance that has been proposed to connect two strings. This distance specifies the number of steps required to get from one string to another, where each step is an insertion, deletion, or substitution. Eg: $d_l(GATCTG, GTCAG) = 2$. The second letter 'A' has to be inserted/deleted based on the requirement, and the 5th word has to be substituted from 'T' to 'A'. Since calculating this is an NP-hard problem, approaches based on the hamming distance(3.1) and uncorrelation have been developed.

2.2 | Sequencing Errors

Tantum errors are a type of sequencing error that occurs when reconstructing a whole sequence from fragments. The error is caused when a single base pair or substring of the primer is replicated and attached adjacent to the original substring or base pair. For example, AGACAGTG if this is changed to AGAGACAGTG then the tantum error occurred is of length 2 and it happened on the first two base pairs. Here, note that AG also repeated in 5th, 6th places of the original sequences, and that is not considered a tantum errors since it's not following the original immediately. This error rectification is mostly used for treatments and other medical purposes. But when we speak about the durability of the DNA and other things to store for a long time we need to consider this also¹⁶. Next, the hybridization property of DNA helps to stabilize the DNA but for that, the preferred GC-content(3.4) has to be maintained. Sometimes due to external factors wrong or unwanted hybridization occurs. One of these is the formation of secondary structures (folds) in the long strand of bps. This happens because the complement of prefix or suffix is a sub-string of that string itself. Since small strand primers can be sequenced to generate a long stable strand, complement

substrings should not be included in the concatenation of two strings.

i.e, Let $x, y, z \in C_{DNA}(n, M, d)$, C be a DNA-code and θ be an automorphism on $\sum_{DNA}^n = \{A, T, G, C\}^n$, such that

$$\theta(x) = \text{Reverse compliment of } x, \forall x \in \sum_{DNA}^n$$

Then, $\theta(x)$ should not be a substring of ZY . To rectify all of these errors, several constraints stated in section 3 are implemented on DNA code words.

3 | DEFINITION AND CONSTRAINTS DNA CODES

A DNA code $C_{DNA}(n, M, d) \subset \sum_{DNA}^n = \{A, T, G, C\}^n$ (nucleotide bases) with each DNA code-word of length n and size M and minimum distance d . The *Hamming distance* $H(x, y)$ between two codewords is the number of distinct elements in those two codewords. The *reverse* of a codeword $x = (x_1, x_2, \dots, x_n)$ is $(x_n, x_{n-1}, \dots, x_2, x_1)$ and denoted by x^R . Similarly the *reverse compliment* of a codeword $x = (x_1, x_2, \dots, x_n)$ is $(x_n^c, x_{n-1}^c, \dots, x_2^c, x_1^c)$ and denoted by x^{RC} . There are basically four constraints one should see while constructing a code. They are as follows:

3.1 | Hamming Distance constraint

Let C be a DNA code, $H(x, y) \geq d, \forall x, y \in C$ with $x \neq y$. We call this distance d as the minimum distance for the Code C . This constraint is inherited from coding theory and this helps in determining the number of errors that can be corrected by this code.

3.2 | Reverse constraint

Let C be a DNA code, $H(x^R, y) \geq d, \forall x, y \in C$. Also considering the case $x = y$. This step acts as a bridge for constructing the reverse compliment constraint. One can use the bound of this constraint and implement on the other.

3.3 | Reverse Complement constraint

Let C be a DNA code, $H(x^{RC}, y) \geq d, \forall x, y \in C$. Also considering the case $x = y$. This constraint helps in the reduction of unwanted hybridization errors.

3.4 | GC-content constraint

Let C be a DNA code, $x \in C$. The number of bits in x that has G or C is called the GC-weight of the codeword. The GC weight constraint is that for a fixed weight w , $w_{GC}(x) = w, \forall x \in C$. The use of this constraint comes from the biological strand point. The GC-content of the codeword is the percentage of $\frac{\text{GC-Weight}}{\text{length of the codeword}}$.

Example: Let $AAGCT \in C$ $w_{GC}(AAGCT) = 2$ and GC-content = $\frac{2}{5} \times 100\% = 40\%$.

Since the G, C has three hydrogen bonds it is more stable than the A, T bond. If the GC content is high, PCR amplification is difficult, whereas a low GC content results in a non-stable gene. So, this forces to fix a precise amount of GC content based on the usage of the code. Average amount of GC-content ranges from 40% – 60%

3.5 | Correlation constraint

Let C be a DNA code of length n , $x, y \in C$. The correlation of x, y is denoted as $x \circ y$ and contains n bits and

$$(x \circ y)[i] = \begin{cases} 1 & \text{if } x[i : n] = y[1 : n - i] \\ 0 & \text{otherwise} \end{cases}$$

(i.e.), For example: Let $X = CATCGT$, $Y = TCGTAC$, $x, y \in C$, $X \circ Y = 001001$

$$\begin{array}{rcccccc}
 X = & C & A & T & C & G & T \\
 Y = & T & C & G & T & A & C & 0 \\
 & & T & C & G & T & A & C & 0 \\
 & & & T & C & G & T & A & C & 1 \\
 & & & & T & C & G & T & A & C & 0 \\
 & & & & & T & C & G & T & A & C & 0 \\
 & & & & & & T & C & G & T & A & C & 1
 \end{array}$$

These constrains are being implemented to our DNA code by using the group homomorphism and the error correcting properties of the kernel codes¹⁷.

4 | KERNEL CODES

Definition 1: Let $\mathcal{G} = \mathcal{G}_1 \times \mathcal{G}_2 \times \dots \times \mathcal{G}_n$, where \mathcal{G}_i 's are groups and $(S, *)$ be an abelian group with identity element e . Then $\mu : \mathcal{G} \rightarrow S$ such that $\mu(g_1, g_2, \dots, g_k) = \mu_1(g_1) * \mu_2(g_2) * \dots * \mu_n(g_n)$ such that $g_i \in \mathcal{G}_i$ and $\mu_i : \mathcal{G}_i \rightarrow S$ is a homomorphism $\forall i = 1$ to n . The construction states that μ is a homomorphism. The kernel of the homomorphism $K = \{g \in \mathcal{G} / \mu(g) = e\}$ is called the Kernel Code¹⁷. **Examples:** Let us consider (\mathcal{Z}_2, \cdot_2) , (\mathcal{Z}_4, \cdot_4) , $(S_3, *)$, $(\mathcal{Z}_2, \cdot_{\mathcal{Z}_2})$ as $\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3, S$ respectively, and $\mu : \mathcal{Z}_2 \times \mathcal{Z}_4 \times S_3 \rightarrow \mathcal{Z}_2$ such that $\mu_1(0) = 0, \mu_1(1) = 1, \mu_2(0) = 0, \mu_2(1) = 1, \mu_2(2) = 0, \mu_2(3) = 1, \mu_3((12)(12) = e') = 0, \mu_3((12) = a_1) = 1, \mu_3((13) = a_2) = 1, \mu_3((23) = a_3) = 1, \mu_3((123) = b_1) = 0, \mu_3((132) = b_2) = 0$. The kernel Code of this homo-morphism is $\{00e', 00b_1, 00b_2, 01a_1, 01a_2, 01a_3, 02e', 02b_1, 02b_2, 03a_1, 03a_2, 03a_3, 10a_1, 10a_2, 10a_3, 11e', 11b_1, 11b_2, 12a_1, 12a_2, 12a_3, 13e', 13b_1, 13b_2\}$.

4.1 | Concatenated Kernel Code

The codes are constructed in double-layer outer and inner kernel code. Inner code is the kernel code as (4). But here \mathcal{G}_i 's are the same abelian group G and the mapping is from

$$\mu : G^n \rightarrow S$$

where $(S, *)$ is an abelian group and $\mu(g_1, g_2, \dots, g_n) = \mu_1(g_1) * \mu_2(g_2) * \dots * \mu_n(g_n)$ such that $g_i \in G$ and $\mu_i : G \rightarrow S$ is a homomorphism.¹⁸

Outer Code: Let G be an abelian group and $\mu' : G^k \rightarrow G^n$ defined as $\mu'(g_1, g_2, \dots, g_k) = (g_1, g_2, \dots, g_k, h_1(g_1, g_2, \dots, g_k), \dots, h_{n-k}(g_1, g_2, \dots, g_k))$ h_i 's are homomorphism can be defined depending upon the error correction. The outer code is the image of μ' .

The kernel of $\mu \circ \mu'$ from the images of μ' is called the concatenated kernel code.

Example: Let's take the G as \mathcal{Z}_3 and the information is mapped to S as \mathcal{Z}_3 . and $k = 3$ and $n = 5$. Let μ_i be identity homomorphisms (i.e), $\mu_i(0) = 0, \mu_i(1) = 1, \mu_i(2) = 2$ for $i = 1, 2, 3$ and $h_1(g_1, g_2, g_3) = g_1 + g_3$ and $h_2(g_1, g_2, g_3) = g_2 + g_3$.

The kernel Code is $\{000, 012, 021, 102, 111, 120, 201, 210, 222\}$

For $\mu'(102) = (1, 0, 2, 1 + 2, 0 + 2) = (1, 0, 2, 0, 2) = 10202$

The concatenated Kernel code is $\{00000, 01220, 02110, 10202, 11122, 12012, 20101, 21021, 22211\}$. Here, The kernel set of the μ is only considered for the μ' mapping which helps for the error correction.

5 | CONSTRUCTION OF DNA CODES

Here, a code of length n is constructed with message block of length:

$$l = n - \log_2(2 * n - 1) - 1 \quad (1)$$

Each element of the block is from the same finite abelian group \mathcal{G} . The corresponding DNA code for each block is generated with GC-constraint(3.4), Reverse Complement constraint(3.3), and less correlation(3.5). For this, VT code and kernel code have been used to encode the information. The algorithm goes as follows:

1. The set of all information of length l is generated.
2. Each of that information is encoded using VT code to restrict indel errors.
3. Each VT-encoded information is then mapped to a unique element in the Kernel code of length $n+1$.
4. Now, the mapped elements of the sets are again outer encoded using the proposed homomorphism to address the constraints(3).
5. Then the obtained word is mapped onto the corresponding base pairs.

The algorithm begins by encoding the information using VT codes. VT codes are used for correcting single deletion, insertion or substitution(indel) errors.¹⁹ has suggested a way for correcting multiple indel errors by segmenting codewords into smaller lengths.

For a fixed codeword length(say n), the algorithm(1) calculates the maximum length for the message bits(l) (1). All possible information on length l is encoded using the VT code. Then, the obtained set of all VT encoded information(say Y) is sent as output. Example: For $n = 8$, the value l is 3. Since, $n - 1 = 7$ is not a power of 2, the par_pos value is $[2^0, 2^1, 2^2, 7]$.

Let $i = 2$, then b_2 is 010 and Y_2 is initialised as 0000100. The syndrome of Y_2 (i.e 10) is partitioned as $7+2+1$. So, Y_2 is finalised to 1100101.

Algorithm 1 Encoding using VT codes

Require: Finite Group and n

▷ n is the length of the codeword

Ensure: VT encoded information bit

▷ l is the length of the message

```

1:  $l \leftarrow n - \log_2(n - 1) - 1$ 
2: if  $n - 1 = 2^k$  then
3:    $\text{par\_pos} \leftarrow 2^0, 2^1, \dots, 2^k - 1$ 
4: else
5:    $\text{par\_pos} \leftarrow 2^0, 2^1, \dots, 2^k - 1, 2^k$ 
6: end if
7: for  $i = 0$  to  $2^l - 1$  do
8:    $Y_i \leftarrow$  all zeros of length  $n-1$ 
9:    $b_i \leftarrow$  Make the possible binary.
10:   $Y_i \leftarrow b_i$  in positions other than  $\text{par\_pos}$ 
11:   $\text{syndrome} \leftarrow \text{mod}(\text{Negative sum of all values in } Y_i, 2^{n+1})$ 
12:  if  $\text{syndrome} \neq 0$  then
13:    for  $\text{value in par\_pos}$  reverse do
14:      if  $\text{syndrome} \geq \text{value}$  then
15:         $Y_i[\text{value}] \leftarrow 1$ 
16:         $\text{syndrome} -= \text{value}$ 
17:      end if
18:    end for
19:  end if
20: end for
21: return  $Y$ 
```

For the length(n) codeword, algorithm 2 generates the Kernel code(K) of length $n+1$. All the elements from K that starts with 1 are collected into a subset(S). For example: For $n = 3$, $\mathcal{G} = \mathbb{Z}_2$, $K = \{ 0000, 0001, 0010, 0011, 0101, 0110, 0111, 1000, 1001,$

Algorithm 2 Getting the required subset of kernel code**Require:** Finite Group and n $\triangleright //n$ is the length of the codeword**Ensure:** Subset of Kernel Code

```

1: for  $i = 1$  to  $n + 1$  do
2:    $\mu_i =$  Make the possible Homomorphism.
3: end for
4: Compute  $C$  the set of all elements in the Cartesian product of  $\underbrace{G \times G \times \dots \times G}_{n \text{ times}}$ 
5: for  $j = 1$  to  $\|C\|$  do
6:   if  $\mu_1(g_1)\mu_2(g_2) \dots \mu_n(g_n) = 0$  and  $\mu(g_1) = 1$  then
7:     Add to set  $S$ 
8:   end if
9: end for
10: return  $S$ 

```

 $\triangleright //$ The required subset of kernel codes

1010, 1011, 1101, 1110, 1111} and $S = \{1001, 1010, 1111\}$.

Then algorithm 3 maps each element of S to an element of Y . Let a be an element in S , $a[i : j]$ represents the sub-sequence starting from i^{th} position to the j^{th} position of the sequence.

(i.e) if $a = 01101$, $a[2 : 4] = 110$

. The element $a \in S$ is mapped to $y \in Y$ if $a[2:n-1] = y$. Considering the previous example. Let $n = 8$ and $y = 1100101$ then kernel map of y is 111001011. Clearly, 111001011 is in K .

For the final part of the encoding, $n - 1$ redundancy bits are added to each element a of the kernel subset S . Here, i^{th} redundancy

Algorithm 3 Mapping information to the subset of kernel code**Require:** K, y, n $\triangleright //y$ is from Y **Ensure:** The map of the word(y) to K

```

1: Map( $y$ )
2:  $m \leftarrow \text{length}(y)$ 
3: if  $m + 1 = n$  then
4:   for element in  $K$  do
5:     if element[2: $m$ ] ==  $y$  then
6:       return element
7:     break
8:   end if
9: end for
10: else
11: return Map(0 $y$ )
12: end if

```

bit is obtained using the homomorphism h_i . The homomorphism goes as follows

- The h'_i 's is depending upon the i^{th} information bit for $i = 1, 2, \dots, n$. except
- If n is even the bit $\frac{n}{2}$ depends upon the the $(n + 1)^{th}$ bit also.

Here, for $i = 1, 2, \dots, \left\lfloor \frac{n-1}{2} \right\rfloor$ the homomorphism of h_i is $h_i(g_{i+1}) = g_{i+1}$ and for $i = \left\lfloor \frac{n+1}{2} \right\rfloor, \dots, n - 1$ $h_i(g_1, g_{i+1}) = g_1 + g_{i+1}$ and if n is even $h_{\frac{n}{2}}(g_1, g_{\frac{n}{2}}, g_{n+1}) = g_1 + g_{\frac{n}{2}} + g_{n+1}$ (4.1).

Example: For the infomation $b_2 = 010$ $y = 1100101$, the kernel mapping is $a = 111001011$ and the encoding is $E_a = 1110010111101010$.

The binary codeword is being separated into small strings of length two with the idea portrayed in the algorithm 4. This binary conversion is dependent on an element from each half. In this way, 1 length information is mapped onto the n-length bps.

Example: For $E_a = 1110010111101010$, the DNA encoded codeword be *GGGCATAT*.

This way of encoding preserves the GC-content of 50% for even length n and for odd length, it ranges from 40% to 60%. This determines the stability of the DNA code.

This also give a greater distance for reverse compliment constraint. For a n-length DNA code the RC-distance is

Algorithm 4 Binary code to DNA strings

Require: E_a

Ensure: encoded information with n bps

```

1:  $n \leftarrow \frac{\text{length}(E_a)}{2}$ 
2: for  $i = 1$  to  $n$  do
3:    $x[i] \leftarrow \text{join}(E_a[i], E_a[n + i])$ 
4: end for
5: MAP 00  $\Rightarrow$  C, 01  $\Rightarrow$  A 10  $\Rightarrow$  T, 11  $\Rightarrow$  G
6: for  $i = 1$  to  $n$  do
7:    $y[i] \leftarrow \text{MAP}(x[i])$ 
8: end for
9: return y

```

$$d_{RC} = 2 \times \left\lfloor \frac{n-3}{2} \right\rfloor$$

This construction also has a less correlation which mostly happens at the end. This acts as an improvised coding scheme for¹⁴ where the 3^l and other S_l methods have been used. This meets the constraints described in 3 and serves as a solution to the DNA storage problem.

The decoding is done using an algorithm 5 that considers the second to last element of the word and decodes it using the Map given. The map is defined from base pair to \mathcal{Z}_2 . Suppose, the DNA codeword from the example(*GGGCATAT*) is recived as

Algorithm 5 Decoding Kernel codes to binary

Require: d

▷ DNA codeword

Ensure: Binary conversion

```

1: MAP C, A  $\Rightarrow$  0 and T, G  $\Rightarrow$  1
2:  $n \leftarrow \text{length}(y)$ 
3: for  $i = 2$  to  $n$  do
4:    $a[i-1] \leftarrow \text{MAP}(y[i])$ 
5: end for
6: return  $\hat{a}$ 

```

GGGATAT. Then the decoding is as follows

the 'T, G' as 1 and 'A', 'C' as 0. So, $\hat{a} = 110101$ calculated from algorithm 5 is given as input to the 6. Then, the value of y is initialized to 110101, syndrome = 2, weight $w = 4 > 2 = \text{syndrome}$. Since weight is less than syndrome 0 has to be added, after the two 1's 0. Therefore y becomes 1100101. Then, the information is taken without parity positions from y(i.e) 010 which is b_2 .

Figure 1 depicts how the encoding has been done for an information set of length(l) 3 and codeword(n) of length 8 using python3.

Algorithm 6 VT decoding to get information**Require:** \hat{a}

▷ DNA codeword

Ensure: information

```

 $y \leftarrow \hat{a}$ 
syndrome  $\leftarrow \text{mod}(\text{Negative sum of all values in } Y, 2*n+1)$ 
 $w \leftarrow \text{weight}(y)$ 
 $lv \leftarrow n - 1$ 
if  $\text{len}(\hat{a}) == lv-1$  then
    if syndrome = 0 then
         $y = y0$ 
    else if syndrome  $\leq w$  then
        0 has to be added to y
        count the ones until it equals syndrome value
        add 0 to next place
    else
        1 has to be added to y
        count zero until it reaches one less than s-w
        Add 1 to next place
    end if
else if  $\text{len}(\hat{a}) == lv+1$  then
     $m \leftarrow 2*lv + 1$ 
    if syndrome ==  $m - lv - 1$  or  $s == 0$  then
        remove last entry from y
    else if syndrome ==  $m - w$  then
        remove first entry
    else if syndrome  $> m - w$  then
        1 has to be deleted from y
        when number of ones reaches m-s
        remove next one
    else
        0 has to be removed from y
        when number of ones reaches m-s-w
        remove next zero
    end if
else
return Not a codeword
return value of y that are not in par_pos
end if

```

6 | CONCLUSION

As DNA has high chances for the occurrence of mutation errors and instability, algorithms have been proposed for the construction of stable DNA code. This procedure ensures uncorrelatedness and balanced GC- content. Further, Vt algorithm used here corrects single indel error that occurs in DNA strings during retrieval.


```

File - man (1)
C:\Users\Admin\PycharmProjects\pythonProject4\venv\Scripts\python.exe C:/Users/Admin/Pyck
Enter the length of the codeword: 8
Binary_number      VT_encoded      Kernel_code \
0      [0, 0, 0]    [0, 0, 0, 0, 0, 0, 0, 0]    [1, 0, 0, 0, 0, 0, 0, 0, 1]
1      [0, 0, 1]    [0, 1, 0, 0, 0, 1, 1]    [1, 0, 1, 0, 0, 0, 1, 1, 0]
2      [0, 1, 0]    [1, 1, 0, 0, 1, 0, 1]    [1, 1, 1, 0, 0, 1, 0, 1, 1]
3      [0, 1, 1]    [0, 0, 0, 1, 1, 1, 0]    [1, 0, 0, 0, 1, 1, 1, 0, 0]
4      [1, 0, 0]    [1, 0, 1, 1, 0, 0, 1]    [1, 1, 0, 1, 1, 0, 0, 1, 1]
5      [1, 0, 1]    [0, 1, 1, 1, 0, 1, 0]    [1, 0, 1, 1, 1, 0, 1, 0, 1]
6      [1, 1, 0]    [0, 0, 1, 0, 1, 0, 1]    [1, 0, 0, 1, 0, 1, 0, 1, 0]
7      [1, 1, 1]    [1, 0, 1, 0, 1, 1, 0]    [1, 1, 0, 1, 0, 1, 1, 0, 1]

Concatenated Kernel Code      DNA_codeword
0      [1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1]    [G, C, C, C, A, A, A, A]
1      [1, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0]    [T, C, G, C, C, A, T, T]
2      [1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0]    [G, G, G, C, A, T, A, T]
3      [1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1]    [T, C, C, C, G, T, T, A]
4      [1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0]    [G, G, C, G, T, A, A, T]
5      [1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1]    [G, C, G, G, T, A, T, A]
6      [1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1]    [T, C, C, G, C, T, A, T]
7      [1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1]    [G, G, C, G, A, T, T, A]

Process finished with exit code 0

```

Page 1 of 1

FIGURE 1 Encoding for three-length information set

Source: It is the python code output which was executed on the Pycharm version: 2021.2.2 (open source software) on the computer with the processor Intel(R)Xeon(R) E3-1225 v5@ 3.30GHz.

References

1. Dong Y, Sun F, Ping Z, Ouyang Q, Qian L. DNA storage: research landscape and future prospects. *National Science Review* 2020; 7(6): 1092-1107. doi: 10.1093/nsr/nwaa007
2. Clelland C, Risca V, Bancroft C. Hiding messages in DNA microdots. *Nature* 1999; 399: 533-4. doi: 10.1038/21092
3. Berlekamp ER. *Algebraic coding theory (revised edition)*. Newyork: World Scientific . 2015.
4. Levenshtein VI. Efficient reconstruction of sequences. *IEEE Transactions on Information Theory* 2001; 47(1): 2–22.
5. Kannan S, McGregor A. More on reconstructing strings from random traces: insertions and deletions. 2005: 297-301. doi: 10.1109/ISIT.2005.1523342
6. Church GM, Gao Y, Kosuri S. Next-Generation Digital Information Storage in DNA. *Science* 2012; 337(6102): 1628-1628. doi: 10.1126/science.1226355
7. Blawat M, Gaedke K, Hütter I, et al. Forward Error Correction for DNA Data Storage. *Procedia Computer Science* 2016; 80: 1011-1022. International Conference on Computational Science 2016, ICCS 2016, 6-8 June 2016, San Diego, California, USA doi: <https://doi.org/10.1016/j.procs.2016.05.398>
8. Organick L, Ang SD, Chen YJ, et al. Scaling up DNA data storage and random access retrieval. *BioRxiv* 2017: 114553.
9. Organick L, Ang SD, Chen YJ, et al. Random access in large-scale DNA data storage. *Nature biotechnology* 2018; 36(3): 242–248.

10. Tabatabaei Yazdi S, Yuan Y, Ma J, Zhao H, Milenkovic O. A rewritable, random-access DNA-based storage system. *Scientific reports* 2015; 5(1): 1–10.
11. Kiah HM, Puleo GJ, Milenkovic O. Codes for DNA Sequence Profiles. *IEEE Transactions on Information Theory* 2016; 62(6): 3125–3146. doi: 10.1109/TIT.2016.2555321
12. Liu J, Liu H. DNA Codes Over the Ring $\mathcal{F}_4[U]/\langle U^3 \rangle$. *IEEE Access* 2020; 8: 77528–77534. doi: 10.1109/ACCESS.2020.2989203
13. Aboluiou N, Smith DH, Perkins S. Linear and nonlinear constructions of DNA codes with Hamming distance d , constant GC-content and a reverse-complement constraint. *Discrete Mathematics* 2012; 312(5): 1062–1075. doi: <https://doi.org/10.1016/j.disc.2011.11.021>
14. Goldman N, Bertone P, Chen S, et al. Towards practical, high-capacity, low-maintenance information storage in synthesized DNA. *nature* 2013; 494(7435): 77–80.
15. Lenz A, Siegel PH, Wachter-Zeh A, Yaakobi E. Coding Over Sets for DNA Storage. *IEEE Transactions on Information Theory* 2020; 66(4): 2331–2351. doi: 10.1109/TIT.2019.2961265
16. Jain S, Farnoud Hassanzadeh F, Schwartz M, Bruck J. Duplication-Correcting Codes for Data Storage in the DNA of Living Organisms. *IEEE Transactions on Information Theory* 2017; 63(8): 4996–5010. doi: 10.1109/TIT.2017.2688361
17. Ramachandran S, P Balasubramanie P. Construction of Kernel codes and its trellis. *ACCST research Journal* 2003; I: 94–96.
18. Selvakumar R, Pavan Kumar C. Concatenated kernel codes. *Discrete Mathematics, Algorithms and Applications* 2020; 12(03): 2050044.
19. Liu Z, Mitzenmacher M. Codes for deletion and insertion channels with segmented errors. *IEEE Transactions on Information Theory* 2009; 56(1): 224–232.

