# DNA Codec with Polar Codes, Huffman Compression and PGP Encryption

## IT495 - DNA Storage

# Team Members

| Nilay Shah | 201901026 |
|---|---|
| Kathan Sanghavi | 201901053 |
| Aayush Prajapati | 201901099 |
| Yash Sakaria | 201901165 |
| Ronak Jethava | 201901174 |
| Parth Prajapati | 201901429 |

# Literature Review

# Literature Review

- **Channel polarization: A method for constructing capacity-achieving codes for symmetric binary-input memoryless channels Erdal Arıkan, Senior Member, IEEE**

  This paper by Erdal Arıkan gives an idea about the construction of Error-correcting aka Capacity-achieving Polar Codes. We understood this paper's basic definition and implementation of Polar Encoder and Decoder. Based on the base formulas given in this, we are able to design efficient algorithms, that can be run in Python.

- **Introduction to Polar Codes, NPTEL - NOC IITM**

  This YouTube playlist is by Prof. Andrew Thangaraj. Through this playlist, he explains What are Polar Codes, Polar Transforms, Polar-Encoder and Polar-Decoder the help of an algorithm. We took reference of this video to get familiar with the concepts of Polar Transform, Polar Encoder and Polar Decoder for carrying out our research in DNA Codec.

- **A Brief Introduction to Polar Codes, Henry D. PfisterY.**

  This paper by Henry D. Pfister has basic details on kronecker delta products, designing efficient encoding and decoding algorithm and scheming polar code based on effective channel error rates. Also provides details on analysis on error rate of codes using monte carlo simulation of BSC, BEC channels.

# Literature Review

- **Kurniawan, A. Albone and H. Rahyuwibowo, "The design of mini PGP security," Proceedings of the 2011 International Conference on Electrical Engineering and Informatics.**

  This paper gives an idea about creating a PGP application which has the conventional functions of encryption, decryption and digital signature. With help, we can plan to design the algorithm with some changes made in reference to considering the original system implementation.

- **"A Characterization of the DNA Data Storage Channel", Reinhard Heckel1, Gediminas Mikutis2 & Robert N.Grass**

  This study describes the DNA data storage channel and how DNA can be utilised as an archival storage devices. It also underlines the limitations and imperfections in DNA synthesis, sequencing, and manipulation. We learned from others how vital it is to have a qualitative understanding of errors and molecule loss while designing a DNA storage system. In this paper, the error probabilities have also been described by analysing experimental data.
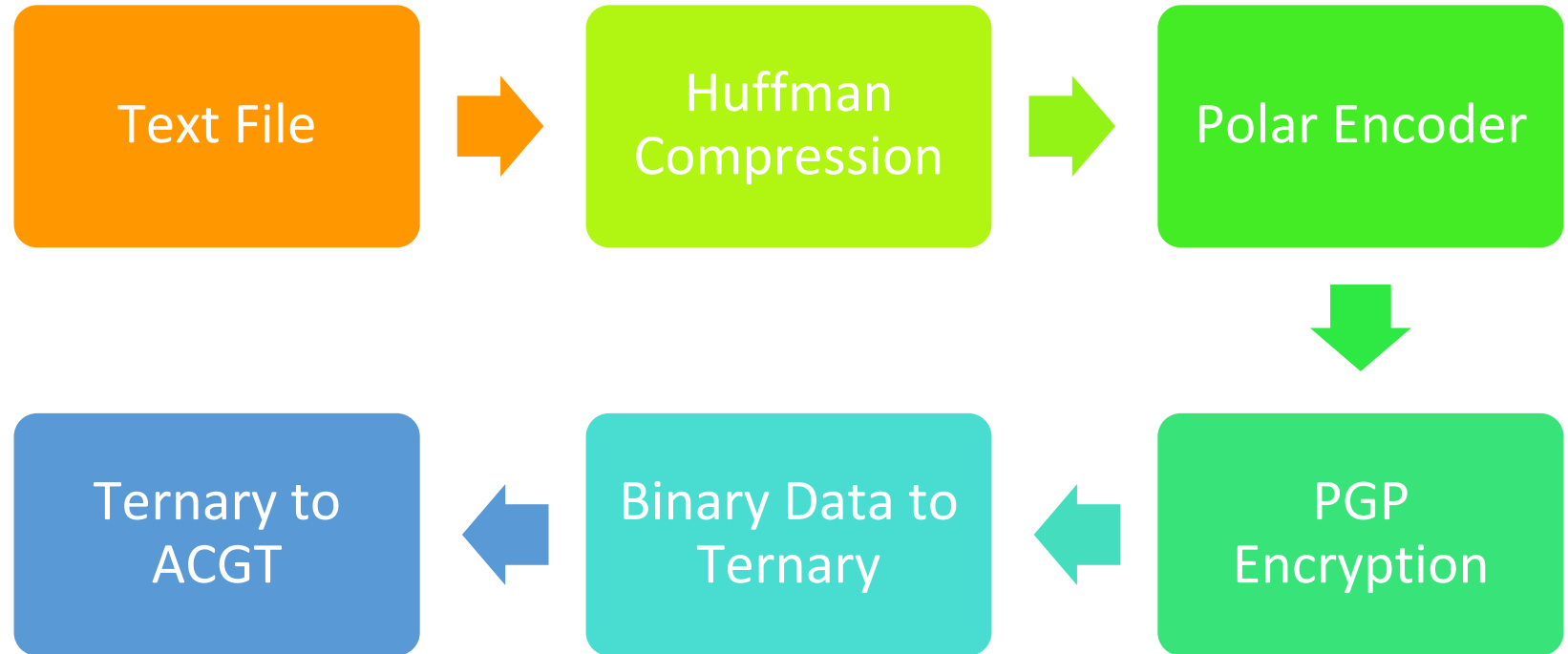
# Encrypted Error Correcting - DNA Codec System

# DNA Codec

**Huffman Compression** - Compressing initial size of data is the best choice for Archival Storage like DNA!

**Polar Codes** - Linear Block Error Correcting Code, going to be used in 5G! Will provide the capability of error correctness for our DNA storage.

**PGP Encryption** – Laptops/PCs are now coming with encrypted Hard Disks, so why not have similar security for DNA storage? PGP Encryption on DNA leverages the capability of DNA of storing huge amounts of data only using a few grams. We are using the GnuPG library of python to do encryption-decryption. Kleopatra software is used for key – certificate generation.

# Encoder

# Text File:

Hello World!

# Text File to Binary Data Conversion

ASCII to Binary: 01001000 01100101 01101100 01101100 01101111 00100000 01010111 01101111 01110010 01101100 01100100 00100001

# Huffman Compression

Huffman Compression: 1100 1101 01 01 101 1110 1111 101 000 01 001 100

# Polar Codes - Encoding

For Polar Codes Generator Matrix GN,

$$G_N = (I_{N/2} \otimes F) R_N (I_2 \otimes G_{N/2}), for\ N \geq 2$$

Where, $G_1 = I_1$ . Using (1) we can get $G_2$ , $G_3$ , ... recursively.

We design an efficient encoding algorithm based on Generator Matrix.

# Polar Encoding Calculation

$G_1 = [1]$

$$G_2 = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$$

$$G_3 = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}$$

We get generator matrix using formula described. For coding part, we are using recursion and N can vary as we are not hardcoding the generator matrix.

# Polar Encoding Pseudocode

```python
def polar_encoder(u):
    if (len(u)==1):
        x = u
    else:
        u1u2 = (u[1::2]+u[2::2]) % 2
        u2 = u[2::2]
        x = [polar_encoder(u1u2), polar_encoder(u2)]
    return x
```

# Polar Encoding

Polar Encoder, N = 7

1 1 0 1 1 0 0 0 1 1 1 0 0 1 0 0 1 0 0 0 0 1 0 0 1 1 1 0 0 1 0 0 1 0 0 1 1 1 0 1
0 1 1 1 1 1 0 1 1 0 1 1 0 0 0 1 1 1 0 0 1 0 0 1 0 0 0 0 1 0 0 0 0 0 1 0 0 1 0 1 0
1 1 1 0 0 0 1 1 1 0 0 0 1 0 1 0 0 0 0 1 0 0 1 1 0 1 1 1 1 0 0 1 1 1 0 1 0 0 1 1 1
0 0 0 1 0

# PGP Encryption Block

Polar Encoded message → Custom Mapping to ASCII Characters → PGP Encryption → ASCII to Binary Mapping

# PGP Encryption

011010000100011000110100010001000110001000110000010001010110011001000110101011001011010001100010101010111011011100101000101010100110100000101010001011001000100000101110100001100110011010000110010011000100011000001011000001100000111101001010111011011101110101010101011100001110110011100010011000100110010101100101101111010100000100110100010001100110101010011000100101111011001000111001101101011110110011001001111011110010111001101101010000010100010110001101001110010100010101100101100100001011111011000110100.....

# Binary to Ternary Conversion

$$y_1 = x_1$$

$$y_i = x_i + x_{i-1} \; Where \; i = 2,3,...,N$$

| x(i-1) \ x(i) | 0 | 1 |
|---|---|---|
| 0 | 0 | 1 |
| 1 | 1 | 2 |

# Ternary to ACGT Conversion

$Mapping\ of\ y_1 : 0 \to A,\ 1 \to C,\ \ 2 \to G$

| x(i) / y(i-1) | 0 | 1 | 2 |
|---|---|---|---|
| A | C | G | T |
| C | G | T | A |
| G | T | A | C |
| T | A | C | G |

# Ternary to ACGT

ACTGAGACGTCTACTGACGATCTCGTCTACTCGTCAGTAGACGATCGTACTCG
TCTCTCAGTCAGTCTACGATCTCTCTGACTCTGAGACGATCGTCTCTCTCATCT
GAGCAGTCTCTACTCTCTCGATCTCGTACTCTCTCGTCTCAGTCTACTCGTACT
CTGCTCTACGATCGATCGATCTCGTAGCTAGACTGACGAGTAGCTACGTCTCA
GTACGATCGTACTGCAGAGTCTCTCATCTGCTCATCTCTCTCTCTGCAGTACTG
CTCAGTCATCGTCTAGCTACTCTCTGACTCTGAGCATCTCGTACTCGATCTCGT
CTACTGACTGAGCTCTAGCTACTCGAGATGCTCAGTCTACTGCTAGCTCAGAT
GCTCAGTCAGTCTAGCATCTGCAGTCTCATCGATCTCTCGTACTCTCGTCTCAG
TAGCTCTAGCAGTCTCTACTCTCTGACTCTGACTCGTAGAGCATCTGACGATCT
CGTCAGTCAGAGTCTCAGAGATCTGCTCATCGTAGAGACTGAGCTACGATCTC
GTCTACTCGTAGCTCAGTAGCTCATCTGAGAGCTCTACGAGTCAGATCT..

# Decoder

# ACGT Data

ACTGAGACGTCTACTGACGATCTCGTCTACTCGTCAGTAGACGATCGTACTC
GTCTCTCAGTCAGTCTACGATCTCTCTGACTCTGAGACGATCGTCTCTCTCAT
CTGAGCAGTCTCTACTCTCTCGATCTCGTACTCTCTCGTCTCAGTCTACTCGT
ACTCTGCTCTACGATCGATCGATCTCGTAGCTAGACTGACGAGTAGCTACGT
CTCAGTACGATCGTACTGCAGAGTCTCTCATCTGCTCATCTCTCTCTCTGCAG
TACTGCTCAGTCATCGTCTAGCTACTCTCTGACTCTGAGCATCTCGTACTCGA
TCTCGTCTACTGACTGAGCTCTAGCTACTCGAGATGCTCAGTCTACTGCTAG
CTCAGATGCTCAGTCAGTCTAGCATCTGCAGTCTCATCGATCTCTCGTACTCT
CGTCTCAGTAGCTCTAGCAGTCTCTACTCTCTGACTCTGACTCGTAGAGCATC
TGACGATCTCGTCAGTCAGAGTCTCAGAGATCTGCTCATCGTAGAGACTGAG
CTACGATCTCGTCTACTCGTAGCTCAGTAGCTCATCTGAGAGCTCTACGAGT
CAGATCT..

# ACGT to Ternary Conversion

Mapping of $y_1$ : $0 \rightarrow A$, $1 \rightarrow C$, $2 \rightarrow G$

| y(i-1) \ y(i) | A | C | G | T |
|---|---|---|---|---|
| A | - | 0 | 1 | 2 |
| C | 0 | - | 1 | 2 |
| G | 0 | 1 | - | 2 |
| T | 0 | 1 | 2 | - |

# ACGT to Ternary Conversion

012111000110012100121110011001100121001100121000011001111210
121011000121111121011121100121001111111221121122101111100111
111012111000011111100111210110011000011122111000121012101 2111
000121011012100110012100001112100001210000122211101111221122
112211111111122210001221121012210011012100111112101112112221
110000110121110011001210121121110121001101112221121011001 2210
121121122211210121011012221122210111221012111100001111001112
100121110122101111001111121011210110001112221121001211100121
012111011121111211221122100011110121121000121110011001100121
121001211221121111211100011012112111012210001211121110110001
112100011001210012101211210012111121011012100111101110110122
101210001101111121121000111221011011000121122112211211210 0111
211101101211101211100121001101210121122210001101221011000011
122100111101221010000110121122112221110121122100121012101210
110111100111100012111012221 12..

# Ternary to Binary Conversion

$$x_1 = y_1$$

$$x_i = (y_i + x_{i-1}) \bmod 2, \qquad Where\ i = 2,3,\dots,N$$

| x(i-1) \ y(i) | 0 | 1 | 2 |
|---|---|---|---|
| 0 | 0 | 1 | – |
| 1 | – | 0 | 1 |

# Ternary to Binary Conversion

0110100001000110001101000100010001100010001100000100010101
1001100100001101010110010110100011000101010101110110111100101
0001010100110100000101010001011001000100000101110100001100
1100110100001100100110001000110000010110000011000001111010
0101011101110111010101010101110000111011001110001001100010 1
0110010110111101000001001101000100011001101101001100010010
1111011001000111001101101011110110011001001111011110010111 00
1101010000010100010110001101001110010100010101100101100100
00101111011000110100.....

# PGP Decryption Block

# PGP Decryption

1101100011100100100000100111001001001110
1011111011011000111001001000010000010010
1011100011100010100001001101111001110100
11100010

# Polar Codes - Decoding

Worst Case time complexity for decoding block-length N code is O(NlogN).

Decoder can be seen as N decision elements, one for each source element ui. They are activated

in the order 1 to N according to the likelihood ratio,

$$L_N^{(i)}\left(y_1^N, \hat{u}_1^{i-1}\right) \triangleq \frac{W_N^{(i)}\left(y_1^N, \hat{u}_1^{i-1}|0\right)}{W_N^{(i)}\left(y_1^N, \hat{u}_1^{i-1}|1\right)}$$

# Polar Codes - Decoding

Decision at $\hat{u}_i$ is generated as,

$$\hat{u}_i = \begin{cases} 0, & if\ L_N^{(i)}(y_1^N, \hat{u}_1^{i-1}) \geq 1 \\ 1, & otherwise \end{cases}$$

Efficient Decoding Algorithm is made based on equation (1) and (2).

# Polar Decoding

1100 1101 01 01 101 1110 1111 101 000 01 001 100

# Huffman Decompression

01001000 01100101 01101100 01101100 01101111 00100000
01010111 01101111 01110010 01101100 01100100 00100001

## Binary to ASCII

**Hello World!**

# User Application - UI Walkthrough

# Encoder

# Decoder

# Applications

Encoder Web App: http://hn8qenyu4brw.app.pywebio.online/

Decoder Web App: http://hc2zz2pry0cj.app.pywebio.online/

GitHub Repository: https://github.com/KathanS/DNA_Codec

# User Application - System Design Overview

# Encoding App Web interface

# Blockwise Output [Encoder]

# Blockwise Output [Encoder]

# Decoding App Web Interface

# Blockwise Output [Decoder]

# Blockwise Output [Decoder]

# Microservices like Architecture

- Encoder Web App outputs DNA string and Secret Code, which will be required for Decoding.

- We are using Python Pickles to store outputs of intermediate blocks, this way system gets decoupled and one block can pass the message having secret code to next block after finishing its task.

- We are able to scale each block independently, this increases robustness and overall performance of the system.

- **Information Density**

  Without PGP: ~2.5

  With PGP: ~6

# References

1. Channel polarization: A method for constructing capacity-achieving codes for symmetric binary-input memoryless channels Erdal Arıkan, Senior Member, IEEE, https://arxiv.org/pdf/0807.3917.pdf
2. Introduction to Polar Codes, NPTEL - NOC IITM, (1225) Introduction to Polar Codes: Polar Transform - YouTube
3. A Brief Introduction to Polar Codes, Henry D. Pfister, polar.pdf (duke.edu)
4. A few concepts of DNA Codec - https://youtube.com/playlist?list=PLZ3QKRH1yB54D8HioxHVe7mrEx85yArOR
5. PGP Encryption/Decryption - What is PGP Encryption and How Does It Work? (varonis.com)
6. Huffman Compression - Huffman Encoding & Python Implementation | by Yağmur Çiğdem Aktaş | Towards Data Science
7. The design of mini PGP security, https://ieeexplore.ieee.org/document/6021726?arnumber=6021726