THE MINISTRY OF EDUCATION AND SCIENCE OF THE REPUBLIC OF KAZAKHSTAN

INTERNATIONAL INFORMATION TECHNOLOGIES UNIVERSITY

Faculty of Information technologies

Department of Computer Systems, Software Engineering and Telecommunications

**Project report**

**Course: Information theory**

Students: Arsen Yerbol, Yersain Makazhanov

Group: CSSE - 1608K

Accepted : senior-lecturer L.A. Kozina

Almaty 2019

**CONTENT**

**INTRODUCTION**

This project is done specially for "Information Theory" course by 3$^{rd}$ year students at International Information Technologies University.

The aim was to write a code of Communication System divided into 6 parts. Project required the knowledge of one programming language(data structures, basic functions) and encoding algorithms(Shannon-Fano, Huffman, Hamming(7, 4)). The last was given as a requisite of a course itself.

During the half of semester we developed a whole structured project as shown in **Figure 1**.
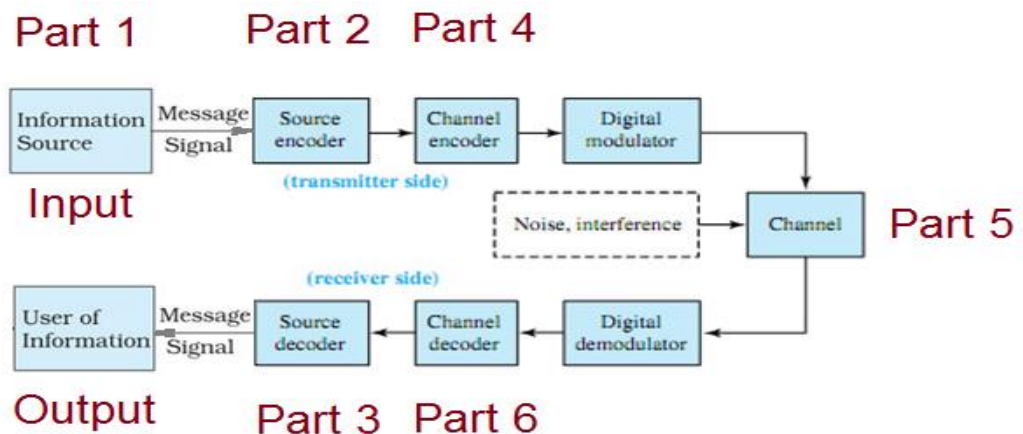


**Figure – 1 – Communication system.**

We chose **Python** programming language, because all team members were familiar and confident with it. As an IDE we chose **jupyter notebook** due to the same reason.

# PART 1

The first part of a project required "Text.txt" file as an input and probabilities of each unique symbol as an output.

Example: A – 0,014, B – 0,012.

Figure 2 shows the beginning of a project. Firstly, we imported dependencies that we need.

Numpy for arrrays and random for error generating.

We made a closer look on text. Then created a dictionary with all unique symbols in it, and finally counted probabilities of each.

**PART 1**

```
In [1]: import numpy as np
        import random

In [2]: f = open("text.txt", "r")

In [3]: text = f.read()[:-1]

In [4]: text

Out[4]: 'In 1951, David A. Huffman and his MIT information theory classmates were given the choice of a term paper or a final exam. The professor, Robert M.Fano, assigned a term paper on the problem of finding the most efficient binary code. Huffman, unable to prove any codes were the most efficient, was about to give up and start studying for the final when he hit upon the idea of using a frequency-sorted binary tree and quickly proved this method the most efficient. In doing so, the student outdid his professor, who had worked with information theory inventor Claude Shannon to develop a similar code. Huffman avoided the major flaw of the suboptimal Shannon-Fano coding by building the tree from the bottom up instead of from the top down.'

In [5]: keys = list(set(text))

In [6]: values = np.zeros((len(keys),), dtype = int)

In [7]: d = dict(zip(keys, values))

In [8]: for char in text:
            d[char]+=1

In [9]: for key in d:
            d[key] = d[key] / len(text)

In [10]: d = sorted(d.items(), key=lambda x: x[1], reverse = True)
```

**Figure – 2 – Code for part1.**

```
In [11]: d

Out[11]: [(' ', 0.17185385656292287),
         ('e', 0.08795669824086604),
         ('o', 0.07713125845737483),
         ('t', 0.06765899864682003),
         ('n', 0.06359945872801083),
         ('a', 0.056833558863328824),
         ('i', 0.05548037889039242),
         ('r', 0.044654939106901215),
         ('h', 0.03924221921515562),
         ('f', 0.03924221921515562),
         ('d', 0.037889039242219216),
         ('s', 0.035182679296346414),
         ('m', 0.02706359945872801),
         ('u', 0.023004059539918808),
         ('p', 0.02029769959404601),
         ('c', 0.016238159675236806),
         ('l', 0.016238159675236806),
         ('b', 0.013531799729364006),
         ('y', 0.012178619756427604),
         ('w', 0.012178619756427604),
         ('g', 0.012178619756427604),
         ('v', 0.010825439783491205),
         ('.', 0.009472259810554804),
         (',', 0.009472259810554804),
         ('H', 0.0040595399188092015),
         ('I', 0.0040595399188092015),
         ('S', 0.0027063599458728013),
         ('T', 0.0027063599458728013),
         ('k', 0.0027063599458728013),
         ('-', 0.0027063599458728013),
         ('M', 0.0027063599458728013),
         ('1', 0.0027063599458728013),
         ('q', 0.0027063599458728013),
         ('F', 0.0027063599458728013),
         ('j', 0.0013531799729364006),
         ('5', 0.0013531799729364006),
         ('D', 0.0013531799729364006),
         ('C', 0.0013531799729364006),
         ('A', 0.0013531799729364006),
         ('9', 0.0013531799729364006),
         ('R', 0.0013531799729364006),
         ('x', 0.0013531799729364006)]
```

**Figure – 3 – Output of dictionary.**

## PART 2

The goal of 2<sup>nd</sup> part was to encode the input text to the sequence of a binary digits according to our probabilities.

Here we chose Huffman encoding, because it seemed more interesting than Shannon-Fano and of course, it is more reliable.

We developed our own algorithm of Huffman encoding without using of Trees and Nodes. We started giving 0s and 1s just from the beginning and did not wait until we get all symbols combined.

For example, lets say that we have A – 0.5 and B – 0.3. Our algorithm in such case, firstly gives 1 to A and 0 to B. The main advantage of such algorithm is that we are winning in terms of memory, because it is not using hashes and Trees.

As a result we got encoded text in binary digits.

**PART 2**

```
In [12]: d2 = dict(d.copy())

In [13]: i = len(d) - 2

In [14]: for x in d2.keys():
             d2[x] = ''

In [15]: while i!=-1:
             first = d[i][0]+d[i+1][0]
             second = d[i][1]+d[i+1][1]
             a = (first, second)
             d.append(a)
             if d[i][1]>d[i+1][1]:
                 for key, value in d2.items():
                     if key in d[i][0]:
                         d2[key] += '1'
                     elif key in d[i+1][0]:
                         d2[key] += '0'
             else:
                 for key, value in d2.items():
                     if key in d[i][0]:
                         d2[key] += '0'
                     elif key in d[i+1][0]:
                         d2[key] += '1'
             d = dict(d)
             d = sorted(d.items(), key=lambda x: x[1], reverse = True)
             del d[-1]
             del d[-1]
             i-=1
         print(d)

         [(' pDCj5RxA9.hfdos,HIctnalbmiywgvruk-STqFM1e', 1.0)]
```

**Figure – 4 – Huffman.**

```
In [16]: d2

Out[16]: {' ': '111',
          'e': '000',
          'o': '1101',
          't': '1001',
          'n': '0001',
          'a': '1110',
          'i': '1010',
          'r': '0100',
          'h': '01011',
          'f': '10011',
          'd': '00011',
          's': '10101',
          'm': '00110',
          'u': '11100',
          'p': '011011',
          'c': '000101',
          'l': '110110',
          'b': '010110',
          'y': '010010',
          'w': '110010',
          'g': '100010',
          'v': '000010',
          '.': '0111011',
          ',': '1100101',
          'H': '00100101',
          'I': '10100101',
          'S': '01001100',
          'T': '11001100',
          'k': '00001100',
          '-': '10001100',
          'M': '01101100',
          '1': '11101100',
          'q': '00101100',
          'F': '10101100',
          'j': '0101111011',
          '5': '1101111011',
          'D': '0001111011',
          'C': '1001111011',
          'A': '0111111011',
          '9': '1111111011',
          'R': '0011111011',
          'x': '1011111011'}
```

**Figure – 5 – Binary representation.**

```
In [17]: for key in d2.keys():
             d2[key] = d2[key][::-1]
```

```
In [18]: output = ''
         for word in text:
             for char in word:
                 output+=d2[char]
```

```
In [19]: output
```

Out[19]: '1010010110001110011011111011111111101111011001101111010011111110111100011101000001011100011110111111011011101111010010
00011111001110010110001110011100111101111000110011111101000101101011101011001101101010010010011111101011000110011011001001100011
110010101011011000111100111010000101100100100101111010000110110111101010101011000111001000101011101001100000100001110010
00101010100000001000111100111010000111101000110101010010101000000011110111100111101111110010000010011001111101100111101
10000001011110110010111011111110010101010000111011100101110110110011001101101011100101011110100011101001101101111001110011111101
110010001010101010110110010101001111111011111100101011110100000101001110010110110110001010101111000101101010011111101111110
101101010101010010001100000001100011100111111110010000010011001111101100111110110000001011110111000111100111010000111110110100010
101101100100110110010011000011110010101001111101001010101100000010001110011110110011100010101010000101111100111001101011000010
01010101100001001000100100111010100101100001100100100100101111010001001110000101101101100000101011011101010000111110011100101100011
100010100101111001110000011011010010110110001110011011111110110000101010100000001101011110001001011110100011011110000001010
1110100100011000100000001110011011100111011001110100000100111011001011110110010001011011100001000011101110011001101111010010
11110110101010010011100111100110111101000010101000000001100111101011011011110001000111101010010011011001010011111010
11001001111100001001001010000010001111100110110100101110011101000011110010101000011011011111010011101000000001001111010
0001111101010011100111101111010101110001111011001110011100101010100000001111110111111001111011100111100011100101110010111011011011
111110010010100000011010000110011000100010000011001101100010010010101110010110100010110000111001001001011110011
0000000111011110001100011100110100000110101101000001100000110110100010111101100010101011000000001100011100111101001010110101
11101110000010010101011100011100111001101000011011001011101011001110001100111001010101000000100010001001001110110111110100
101100011111000101010111010100111100111001101111100111011001100111001110100001110101010010001110000111001100111000111100
010101100011110101001010101011111011000010111001000100010101011101011011001001001110101001110101011111110101100111110001101110100111
011001000110000000110001110100110101101100110101110101110001100110110001001100011100010101011100011110011101000001011001001001
10111010010100001000000100011010011010110110101111011011001101010011100110101010101110010100111011001110110011001011101110011011
1111100000001000000011011101101101101101101111101010101010100010101011011001011101010001011100000010110110111101010010000011
1110011100101100011100011011101101000010101011000000011000111100111010000110110001111101111010101100101111001011011011
0100111111011111011110011111101010001110100100110111011100010111100101010100011110010100101011000011110
001100010010110101011011110001011111010000101110000101100010001111011010010010011011101000110101101011100000101100000100011111
10011101000011110010010000000111100100101011011001110011011000011011010101011001100110110100110011111101101111010110001
0101100100000111100011110111100111111001001010110110011110011011000011110011011111101101111000101101010011100011011110'
```

**Figure – 6 – Binary representation of text.**

## PART 3

The aim of this part was to decode back sequence of binary digits into initial text.  It was pretty easy, we were adding each digit until it matches the first value in our dictionary. You can see the output on Figure N.

```
In [22]: char = ''
         sentence = ''
         for i in output:
             char+=i
             if char in d2.values():
                 sentence+=list(d2.keys())[list(d2.values()).index(char)]
                 char = ''

In [23]: sentence

Out[23]: 'In 1951, David A. Huffman and his MIT information theory classmates were given the choice of a term paper or a final exa
         m. The professor, Robert M.Fano, assigned a term paper on the problem of finding the most efficient binary code. Huffman,
         unable to prove any codes were the most efficient, was about to give up and start studying for the final when he hit upon
         the idea of using a frequency-sorted binary tree and quickly proved this method the most efficient. In doing so, the stud
         ent outdid his professor, who had worked with information theory inventor Claude Shannon to develop a similar code. Huffm
         an avoided the major flaw of the suboptimal Shannon-Fano coding by building the tree from the bottom up instead of from t
         he top down.'
```

**Figure – 7 – Decoding text back.**

## PART 4

In this part we encoded our sequence of binary digits from part 2 with the use of Hamming(7, 4) code. We divided the sequence into blocks by 4 and the left symbols saved in a variable, because it is out of encoding (if the length of sequence cannot be divided by 4 with the mod 0).

Then we calculated the sets of parity bits with Formula 1 and added them into our blocks.

Finally we concatenated all blocks and left symbols to get the output.

$$r1 = i1 \ XOR \ i2 \ XOR \ i3 \quad (1)$$

$$r2 = i2 \ XOR \ i3 \ XOR \ i4 \quad (2)$$

$$r3 = i1 \ XOR \ i2 \ XOR \ i4 \quad (3)$$

where r1,r2,r3 – parity bits , i1,i2,i3,i4 – data bits.

```
In [28]: trios = []
         for i in range(len(blocks_4)):
             t_1 = int(blocks_4[i][0]) ^ int(blocks_4[i][1]) ^ int(blocks_4[i][2])
             t_2 = int(blocks_4[i][1]) ^ int(blocks_4[i][2]) ^ int(blocks_4[i][3])
             t_3 = int(blocks_4[i][0]) ^ int(blocks_4[i][1]) ^ int(blocks_4[i][3])
             t = (str(t_1)+str(t_2)+str(t_3))
             blocks_4[i]+=t

In [29]: blocks_4

Out[29]: ['1010011',
          '0101100',
          '1000101',
          '1110100',
          '0110001',
          '1111111',
          '1011000',
          '1111111',
          '1110100',
          '1111111',
          '0110001',
          '0110001',
          '1111111',
          '0100111',
          '1111111',
          '1110100',
          '1111111',
          '0000000',
          '1110100',
```

**Figure – 8 – Appending parity bits.**

```python
In [25]: blocks_4 = []
         b_4 = str()
         for i in range(len(output)):
             b_4 += output[i]
             if (i+1)%4==0:
                 blocks_4.append(b_4)
                 b_4 = str()
```

```python
In [26]: left_4 = output[len(blocks_4)*4:]
```

```python
In [27]: blocks_4
```

```
'1101',
'0111',
'1001',
'1011',
'0101',
'0010',
'1001',
'1001',
'1111',
'0101',
'1000',
'1100',
'1101',
'1001',
'0011',
'0001',
'1110',
'0101',
'0110',
'1110',
```

**Figure – 9 – Dividing into blocks.**

```python
In [30]: output_4 = str()
         for i in blocks_4:
             output_4+=i
         output_4+=left_4
```

```python
In [31]: output_4
```



**Figure – 10 – Binary representation of text with parity bits.**

## PART 5

Here we added 1 error to each block. It was required to make a random index and put the reverse number on the place of that index in block. The output looks like this.

```python
In [33]: for i in range(len(blocks_5)):
             index = random.randint(0, 6)
             part = list(blocks_5[i])
             if part[index]=='0':
                 part[index] = '1'
             else:
                 part[index] = '0'
             blocks_5[i] = ''.join(part)

In [34]: blocks_5

Out[34]: ['1010111',
          '0100100',
          '1010101',
          '1111100',
          '0110000',
          '1111101',
          '1011001',
          '1111101',
          '1100100',
          '1011111',
          '0110000',
          '0110011',
          '1111011',
          '0100110',
          '0111111',
          '1110101',
          '1011111',
          '0010000',
          '0110100',
```

**Figure – 11 – Adding errors.**

```python
In [35]: output_5 = str()
         for i in blocks_5:
             output_5+=i
         output_5+=left_5

In [36]: output_5
```

Out[36]: '1010111010010010101011111100011000011111011011001111110110010010111101100000110011111011011010011001111111101011011110...'

**Figure – 12 – Binary representation of text with errors.**

## PART 6

Final part of a project had a purpose of decoding back everything. We firstly calculated error syndromes with Formula 2. Then with the use of that syndromes we came to the first blocks of 4 digits without parity bits and errors. Furthermore, we combined all blocks and left digits to generate a sequence of binary digits as it was in part 2. Finally we decoded it with the algorithm wrote in part 3. As a result we got our initial text.

$$S1 = r1 \text{ XOR } i1 \text{ XOR } i2 \text{ XOR } i3 \quad (4)$$

$$S2 = r2 \text{ XOR } i2 \text{ XOR } i3 \text{ XOR } i4 \quad (5)$$

$$S3 = r3 \text{ XOR } i1 \text{ XOR } i2 \text{ XOR } i4 \quad (6)$$

where S1-3 – error syndromes, r1,r2,r3 – parity bits, i1,i2,i3,i4 – data bits.

```
In [39]: syndromes = []
         for i in range(len(blocks_6)):
             s_1 = int(blocks_6[i][0]) ^ int(blocks_6[i][1]) ^ int(blocks_6[i][2]) ^ int(blocks_6[i][4])
             s_2 = int(blocks_6[i][1]) ^ int(blocks_6[i][2]) ^ int(blocks_6[i][3]) ^ int(blocks_6[i][5])
             s_3 = int(blocks_6[i][0]) ^ int(blocks_6[i][1]) ^ int(blocks_6[i][3]) ^ int(blocks_6[i][6])
             s = (str(s_1)+str(s_2)+str(s_3))
             syndromes.append(s)

In [40]: syndromes

Out[40]: ['100',
          '011',
          '110',
          '011',
          '001',
          '010',
          '001',
          '010',
          '110',
          '111',
          '001',
          '010',
          '100',
          '001',
          '101',
          '001',
          '111',
          '110',
          '101',
```

**Figure – 13 – Calculating error syndromes.**

```
In [41]: vals = [6, 5, 3, 4, 0, 2, 1]
         keys = ['001', '010', '011', '100', '101', '110', '111']
         check = dict(zip(keys, vals))

In [42]: check

Out[42]: {'001': 6, '010': 5, '011': 3, '100': 4, '101': 0, '110': 2, '111': 1}

In [43]: for i in range(len(blocks_6)):
             part_6 = list(blocks_6[i])
             if part_6[check[syndromes[i]]] == '0':
                 part_6[check[syndromes[i]]] = '1'
             else:
                 part_6[check[syndromes[i]]] = '0'
             blocks_6[i] = ''.join(part_6)

In [44]: for i in range(len(blocks_6)):
             blocks_6[i] = blocks_6[i][:4]

In [45]: blocks_6

Out[45]: ['1010',
          '0101',
          '1000',
          '1110',
          '0110',
          '1111',
          '1011',
          '1111',
          '1110',
          '1111',
          '0110',
          '0110',
          '1111',
          '0100',
          '1111',
          '1110',
          '1111',
          '0000',
          '1110',
```

**Figure – 14 – Finding errors.**

```
In [46]: output_6 = str()
         for i in blocks_6:
             output_6 += i
         output_6+=left_6
```

```
In [47]: output_6
```

Out[47]: '101001011000111001101111110111111111110111101100110111110100111111110111100001110100000101110001111011111110110111101111010010
0001111100111100101110001111100011101111000110001111101001011101011110011101101010010010011001111110101100011001101100100110011
11001010110111000111100111010000101100100100101111010000110110111101011010101100011100100010101110100110001000010000111010
001010101000000010001111001110100001111010001101010101010110101010100000011110111100111110111100100000100010011111101100111101
100000010111101101010111011111110010101000101101011110011011110011010011001101100111100110011110100011111011000101011
11001000101010110010110010110101011111110111110010101011011011011100000010010011110110011011000101010111101011110011111011110
1011101010101010001000001000111011111111001010000100110011110110110010011110100100011111011100011111101100010
101101101001011000011001111011110011111100101011001000100011001100011111001110100001110110011011101011100110011100110
0101011010001010100000100111110110101001010110000111001001001000011101011110010101011100100010010011100011110100001011111010100
100001010011111100111100001110110010011011000011110011011011111110100000010110111110101110100010000111110011011100011
11000100100001000111100111011010010011010000110100111000010111110100000010100100101110100011110100011101101111010
11110111101101010101000111100011110011111001011011111100100001010101000000001100111111111011011110010010011000011111010
11001001111100001001001011101101110010000100011111100110110110010011110100101010110100110111101101110001110011110010
00011111010010011001110011110101010110001111001101001000011010110000011111101101110011110110101010101000001000011100110011110
1111110010001000000111011110110010000010001111010111011101101011100111110011110100001111010110010101010100000101000100010010100
1110110100000100111101010110010111011010000111101100011100111011001111110110100100010101010010001000010001010010000100000101110
0101110001110001011100101010101101101111010010101010111100111001101001000001100111101001100101010010110011100111011000010
0110010001100000001000111010010110110110101011010111111110101010101001100011100011000110110010010110110110111001010111001001
11111001001000010000001101110111011011011110111111110101010101010010110100010100100001100111010110101011000011100100000101011001001
00000011101111100011100111001101010000110101010000001100000111011011010101000000011000111011101101111001110011110100110101
11101110000010011101010111100011110011101101000011101001101010010011110011001101010101010100000100001000100010010111011011101100
101100011111100011010101011000010010111101011101010010111111100111010100011110101011111111010100111101100001011000011
0101110001111101001010101011110001010101011010111010010101010110110010010011110100111011010011111111010101011001110100111
01100100010000001110011100101001001010111011101101110110011110101010110010110110101011010001101110001011001010001101011010010
111110000001010000001101100011110110111010110110111011111110101000101001001111011010010001001011011111101001001011010011010
1100111001001100011011011110110011110011101011101010100011110111111101010101001100011110000011000100010001000010001111
010011111101111100111100110110001111010010011010101011110100100101010001001101110110011011000101011110100010111100010010001111000
0011000100101010101110010011111101010010011100001010000101100000010100001001110101001000100110101010010101011100010011000100001111
10011101010000111100100010000001111001100101011011001111001110001100111010101011001001100110111001101110010100110011101000
010110110010001111000111010010111111100011011001010101100110011110011101000011110010110101100001101110001101110'

**Figure – 15 – Binary representation of text without errors.**

 'M': '00110110',
 '1': '00110111',
 'q': '00110100',
 'F': '00110101',
 'j': '1101111010',
 '5': '1101111011',
 'D': '1101111000',
 'C': '1101111001',
 'A': '1101111110',
 '9': '1101111111',
 'R': '1101111100',
 'x': '1101111101'}

```
In [49]: char = ''
         sentence = ''
         for i in output_6:
             char+=i
             if char in d2.values():
                 sentence+=list(d2.keys())[list(d2.values()).index(char)]
                 char = ''
```

```
In [50]: sentence
```

Out[50]: 'In 1951, David A. Huffman and his MIT information theory classmates were given the choice of a term paper or a final exam. The professor, Robert M.Fano, assigned a term paper on the problem of finding the most efficient binary code. Huffman, unable to prove any codes were the most efficient, was about to give up and start studying for the final when he hit upon the idea of using a frequency-sorted binary tree and quickly proved this method the most efficient. In doing so, the student outdid his professor, who had worked with information theory inventor Claude Shannon to develop a similar code. Huffman avoided the major flaw of the suboptimal Shannon-Fano coding by building the tree from the bottom up instead of from the top down.'

**Figure – 16 – Decoded text.**

**CONCLUSION**

During this project we gained the knowledge of encoding algorithms practically. We dealt with Huffman and Hamming encodings and realized them in Python language.

We faced several troubles on 2$^{nd}$ part where we were asked to write an algorithm of Huffman encoding. However, after some research we developed our own and maybe state-of-art algorithm without using of Trees and Nodes. The solution for other parts came up really fast in one breath.

Project also revealed some forgotten functions and features of programming language. Moreover, it required real participation and interaction with the course material, which made it more interesting. Researching, critical thinking and team work were the main three objectives of a project for us. These skills worked and fit our personalities the best and we came up with solution without any difficulty.

# REFERENCES

- «Information Theory» course, L. Kozina, N. Duzbayev, International Information Technologies University

- Information Measures: Information and its Description in Science and Engineering, **Arndt C.**

- Elements Of Information Theory, **Thomas Cover.**

# APPENDIX

**This Project is done by**

**Yersain Makazhanov**

**Arsen Yerbol**

PART 1

```python
import numpy as np

import random


f = open("text.txt", "r")


text = f.read()


keys = list(set(text))


values = np.zeros((len(keys),), dtype = int)


d = dict(zip(keys, values))


for char in text:

    d[char]+=1


for key in d:

    d[key] = d[key] / len(text)
```

```
d = sorted(d.items(), key=lambda x: x[1], reverse = True)


PART 2

d2 = dict(d.copy())


i = len(d) – 2


for x in d2.keys():

    d2[x] = ''


while i!=-1:

    first = d[i][0]+d[i+1][0]

    second = d[i][1]+d[i+1][1]

    a = (first, second)

    d.append(a)

    if d[i][1]>d[i+1][1]:

        for key, value in d2.items():

            if key in d[i][0]:

                d2[key] += '1'

            elif key in d[i+1][0]:

                d2[key] += '0'

    else:

        for key, value in d2.items():

            if key in d[i][0]:
```

```python
            d2[key] += '0'

        elif key in d[i+1][0]:

            d2[key] += '1'

    d = dict(d)

    d = sorted(d.items(), key=lambda x: x[1], reverse = True)

    del d[-1]

    del d[-1]

    i-=1
print(d)


for key in d2.keys():

    d2[key] = d2[key][::-1]


output = ''
for word in text:

    for char in word:

        output+=d2[char]
PART 3
char = ''

sentence = ''

for i in output:

    char+=i

    if char in d2.values():

        sentence+=list(d2.keys())[list(d2.values()).index(char)]
```

```python
        char = ''
```

PART 4

```python
blocks_4 = []

b_4 = str()

for i in range(len(output)):

    b_4 += output[i]

    if (i+1)%4==0:

        blocks_4.append(b_4)

        b_4 = str()


left_4 = output[len(blocks_4)*4:]


trios = []

for i in range(len(blocks_4)):

    t_1 = int(blocks_4[i][0]) ^ int(blocks_4[i][1]) ^ int(blocks_4[i][2])

    t_2 = int(blocks_4[i][1]) ^ int(blocks_4[i][2]) ^ int(blocks_4[i][3])

    t_3 = int(blocks_4[i][0]) ^ int(blocks_4[i][1]) ^ int(blocks_4[i][3])

    t = (str(t_1)+str(t_2)+str(t_3))

    blocks_4[i]+=t


output_4 = str()

for i in blocks_4:

    output_4+=i

output_4+=left_4
```

```
PART 5
blocks_5 = []

b_5 = str()

for i in range(len(output_4)):

    b_5 += output_4[i]

    if (i+1)%7==0:

        blocks_5.append(b_5)

        b_5 = str()

left_5 = output_4[len(blocks_5)*7:]


for i in range(len(blocks_5)):

    index = random.randint(0, 6)

    part = list(blocks_5[i])

    if part[index]=='0':

        part[index] = '1'

    else:

        part[index] = '0'

    blocks_5[i] = ''.join(part)


output_5 = str()

for i in blocks_5:

    output_5+=i

output_5+=left_5
```

PART 6

```python
blocks_6 = []

b_6 = str()

for i in range(len(output_5)):

    b_6 += output_5[i]

    if (i+1)%7==0:

        blocks_6.append(b_6)

        b_6 = str()

left_6 = output_5[len(blocks_6)*7:]


syndromes = []

for i in range(len(blocks_6)):

    s_1 = int(blocks_6[i][0]) ^ int(blocks_6[i][1]) ^ int(blocks_6[i][2]) ^ int(blocks_6[i][4])

    s_2 = int(blocks_6[i][1]) ^ int(blocks_6[i][2]) ^ int(blocks_6[i][3]) ^ int(blocks_6[i][5])

    s_3 = int(blocks_6[i][0]) ^ int(blocks_6[i][1]) ^ int(blocks_6[i][3]) ^ int(blocks_6[i][6])

    s = (str(s_1)+str(s_2)+str(s_3))

    syndromes.append(s)


vals = [6, 5, 3, 4, 0, 2, 1]

keys = ['001', '010', '011', '100', '101', '110', '111']

check = dict(zip(keys, vals))
```

```python
for i in range(len(blocks_6)):
    part_6 = list(blocks_6[i])
    if part_6[check[syndromes[i]]] == '0':
        part_6[check[syndromes[i]]] = '1'
    else:
        part_6[check[syndromes[i]]] = '0'
    blocks_6[i] = ''.join(part_6)


for i in range(len(blocks_6)):
    blocks_6[i] = blocks_6[i][:4]


output_6 = str()
for i in blocks_6:
    output_6 += i
output_6+=left_6


char = ''
sentence = ''
for i in output_6:
    char+=i
    if char in d2.values():
        sentence+=list(d2.keys())[list(d2.values()).index(char)]
        char = ''
```