

ГУАП

КАФЕДРА № 42

ОТЧЕТ
ЗАЩИЩЕН С ОЦЕНКОЙ _____
ПРЕПОДАВАТЕЛЬ

Доцент, канд. техн. наук
должность, уч. степень, звание

подпись, дата

В.А. Миклуш
инициалы, фамилия

ОТЧЕТ О ЛАБОРАТОРНОЙ РАБОТЕ №3

Словарно-ориентированные методы кодирования

по курсу: Теория информации, данные, знания

РАБОТУ ВЫПОЛНИЛ

СТУДЕНТ ГР. № _____ 4329

подпись, дата

Д.С. Шаповалова
инициалы, фамилия

Санкт-Петербург 2025

1. Цель работы:

Изучение словарно-ориентированных методов кодирования (LZ77, LZ78, LZW).

2. Задание:

В соответствии с вариантом:

1. Построить алгоритм (блок-схему)
2. Написать программу, реализующую заданный метод кодирования
3. Провести ручную трассировку
4. Сравнить полученные результаты между собой
5. Рассчитать среднее число элементарных сигналов

3. Исходные данные

Вариант 23:

Текст: косарь касьян косой косит косо. не скосит косарь касьян покоса

Метод кодирования: LZ78

4. Теоретические сведения:

Алгоритм LZ78:

При старте алгоритма этот словарь содержит одну пустую строку длины ноль. Алгоритм считывает символы сообщения до тех пор, пока накапливаемая строка входит целиком в одну из фраз словаря. Как только эта строка перестанет соответствовать хотя бы одной фразе словаря, алгоритм генерирует код, состоящий из индекса строки в словаре, которая до последнего введенного символа содержала входную строку, и символа, нарушившего совпадение. Затем в словарь добавляется введенная строка. Если словарь уже заполнен, то из него предварительно удаляют менее всех используемую в сравнениях фразу.

Длина полученного двоичного кода будет округленным в большую сторону:

$$L_k = \log(\text{размер словаря}) + 8, \quad (1)$$

(8 битами кодируются символы, например, ASCII+)

Расчет среднего числа элементарных сигналов:

Среднее число элементарных сигналов на один входной символ:

$$N_{\text{ср}} = \frac{\text{Длина исходного текста (в символах)}}{\text{Число фраз в выходе LZ78}}, \quad (2)$$

Средняя длина одной фразы в символах (обратная величина):

$$L_{\text{ср}} = \frac{\text{Число фраз}}{\text{Длина исходного текста}}, \quad (3)$$

4. Ход работы:

Для выполнения работы был выбран язык программирования высокого уровня Python. Была разработана программа (Приложение А), реализующая алгоритм LZ78, кодирующая поданную на вход последовательность.

Для наглядности работы программы и алгоритма были сделаны блок схемы (рисунки 1-2):



Рисунок 1 – Блок схема основной программы

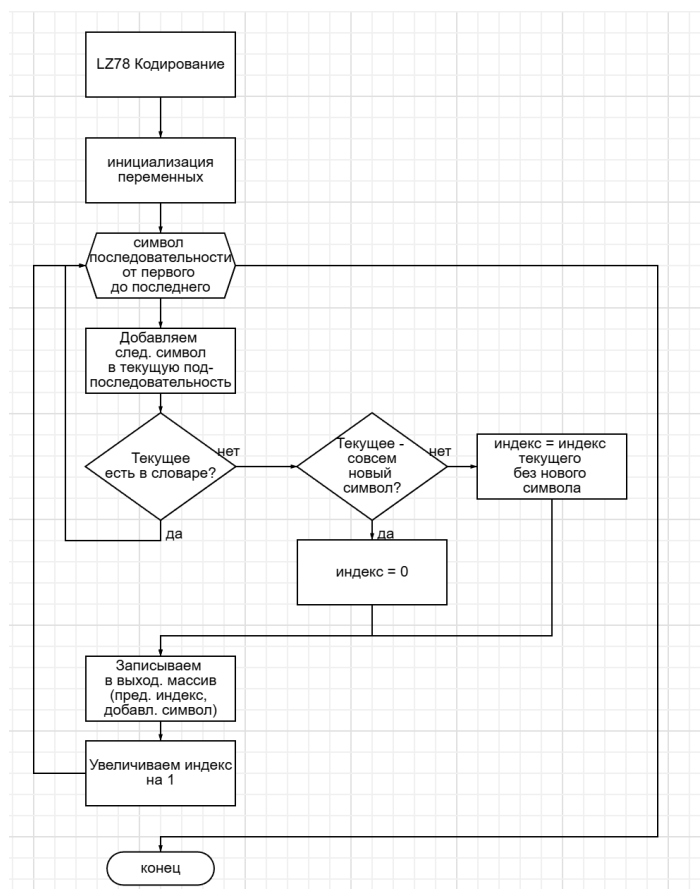


Рисунок 2 – Блок схема подпрограммы – алгоритма LZ78

Результат кодирования текста программой представлен на рисунке 3:

```
(0, 'к')
(0, 'о')
(0, 'с')
(0, 'а')
(0, 'р')
(0, 'ь')
(0, ' ')
(1, 'а')
(3, 'ь')
(0, 'я')
(0, 'н')
(7, 'к')
(2, 'с')
(2, 'й')
(12, 'о')
(3, 'и')
(0, 'т')
(15, 'с')
(2, '.')
(7, 'н')
(0, 'е')
(7, 'с')
(1, 'о')
(16, 'т')
(18, 'а')
(5, 'ь')
(12, 'а')
(9, 'я')
(11, ' ')
(0, 'п')
(2, 'к')
(13, 'а')
```

Рисунок 3 – Закодированная последовательность

Для проверки работы написанной программы и правильности результата была выполнена ручная трассировка, представленная на рисунке 4:

Косарь, Касяк, косой, косит, Косо, Кос, скосит, косарь, Касяк, покоса

| Вход в словарь «>>» | Код | Позиция словаре |
|------------------------|---------|--------------------|
| К | (0, К) | 1 |
| О | (0, О) | 2 |
| С | (0, С) | 3 |
| А | (0, А) | 4 |
| Р | (0, Р) | 5 |
| Б | (0, Б) | 6 |
| Л | (0, Л) | 7 |
| КА | (1, А) | 8 |
| СБ | (3, Б) | 9 |
| А | (0, А) | 10 |
| К | (0, К) | 11 |
| ЛК | (1, К) | 12 |
| ОС | (2, С) | 13 |
| ОИ | (2, И) | 14 |
| ЛКО | (12, О) | 15 |
| СИ | (3, И) | 16 |
| Т | (0, Т) | 17 |
| ЛКОС | (15, С) | 18 |
| О. | (2, .) | 19 |
| ЛИ | (12, И) | 20 |
| Е | (0, Е) | 21 |
| ЛС | (12, С) | 22 |
| КО | (1, О) | 23 |
| СИТ | (16, Т) | 24 |
| ЛКОА | (18, А) | 25 |
| РБ | (5, Б) | 26 |
| ЛКА | (12, А) | 27 |
| СБА | (9, А) | 28 |
| ИЛ | (11, Л) | 29 |
| П | (0, П) | 30 |
| ОК | (2, К) | 31 |
| ОСА | (13, А) | 32 |

Рисунок 4 – Ручная трассировка LZ78

Как видно, расхождения в результатах отсутствуют, что подтверждает правильность работы написанной функции.

Расчёт среднего числа элементарных сигналов:

Входных символов: 62

Элементарных сигналов (фраз): 32

Среднее число сигналов на символ: $N_{\text{ср}} = 0.5161$ по формуле (1)

Средняя длина фразы: $L_{\text{ср}} = 1.94$ символов по формуле (2)

Длина двоичного кода: $L_{\text{к}} = 32 * (\log(33) + 8) = 32 * (6 + 8) = 448$ бит по формуле (3)

Вывод

В ходе выполнения лабораторной работы был изучен метод кодирования текстовой информации алгоритмом LZ78. Была написана программа, реализующая кодирование по методу LZ78, а также для проверки её работы была проведена ручная трассировка и сравнение результатов: почти идентичны. Также была рассчитана длина двоичного кода закодированной последовательности = 448 бит.

ПРИЛОЖЕНИЕ А

Листинг программы

```
def lz78_encode(data):
    dictionary = {}          # словарь: {фраза: индекс}
    phrases = []             # результат: список кортежей (индекс, символ)
    current = ""
    index = 1                 # индексы начинаются с 1

    for char in data:
        current_with_char = current + char
        if current_with_char in dictionary:
            current = current_with_char
        else:
            # current есть в словаре (или пуст), char — новый символ
            prev_index = dictionary[current] if current != "" else 0
            phrases.append((prev_index, char))
            dictionary[current_with_char] = index
            index += 1
            current = ""

    return phrases

def read_input(filename):
    with open(filename, 'r', encoding='utf-8') as f:
        return f.read()

def write_output(phrases, filename='output.txt'):
    with open(filename, 'w', encoding='utf-8') as f:
        for idx, char in phrases:
            # Записываем в формате: (индекс, символ)
            f.write(f"({idx}, '{char}')\n")

if __name__ == '__main__':
    data = read_input('text.txt')
    encoded = lz78_encode(data)
    write_output(encoded, 'output.txt')
    print("Сжатие LZ78 завершено. Результат в output.txt")
```