

**МИНОБРНАУКИ РОССИИ**  
**САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ**  
**ЭЛЕКТРОТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ**  
**«ЛЭТИ» ИМ. В.И. УЛЬЯНОВА (ЛЕНИНА)**  
**Кафедра МО ЭВМ**

**ОТЧЕТ**  
**по лабораторной работе №5**  
**по дисциплине «Построение и анализ алгоритмов»**  
**Тема: Поиск набора подстрок в строке. Ахо-Корасик.**

Студент гр. 3343

\_\_\_\_\_

Никишин С.А.

Преподаватель

\_\_\_\_\_

Жангиров Т.Р.

Санкт-Петербург

2025

### **Цель работы.**

Изучить алгоритм поиска набора подстрок в строке. Разработать программу, находящую все вхождения подстрок в строку с помощью алгоритма Ахо-Корасика.

### **Задание.**

Вариант 2. Подсчитать количество вершин в автомате; вывести список найденных образцов, имеющих пересечения с другими найденными образцами в строке поиска.

### **Задание 6.1.**

Разработайте программу, решающую задачу точного поиска набора образцов.

#### **Вход:**

Первая строка содержит текст ( $T$ ,  $1 \leq |T| \leq 100000$ ).

Вторая - число  $n$  ( $1 \leq n \leq 3000$ ), каждая следующая из  $n$  строк содержит шаблон из набора  $P = \{p_1, \dots, p_n\}$   $1 \leq |p_i| \leq 75$

Все строки содержат символы из алфавита  $\{A, C, G, T, N\}$

#### **Выход:**

Все вхождения образцов из  $P$  в  $T$ .

Каждое вхождение образца в текст представить в виде двух чисел -  $i$   $p$

Где  $i$  - позиция в тексте (нумерация начинается с 1), с которой начинается вхождение образца с номером  $p$  (нумерация образцов начинается с 1).

Строки выхода должны быть отсортированы по возрастанию, сначала номера позиции, затем номера шаблона.

---

#### **Sample Input:**

```
NTAG
3
TAGT
TAG
T
```

---

#### **Sample Output:**

```
2 2
2 3
```

## Задание 6.2.

Используя реализацию точного множественного поиска, решите задачу точного поиска для одного образца с *джокером*.

В шаблоне встречается специальный символ, именуемый джокером (wild card), который "совпадает" с любым символом. По заданному содержащему шаблон образцу  $P$  необходимо найти все вхождения  $P$  в текст  $T$ .

Например, образец  $ab??c?$  с джокером  $?$  встречается дважды в тексте  $xabvccbababcah$ .

Символ джокер не входит в алфавит, символы которого используются в  $T$ . Каждый джокер соответствует одному символу, а не подстроке неопределённой длины. В шаблон входит хотя бы один символ не джокер, т.е. шаблоны вида  $???$  недопустимы.

Все строки содержат символы из алфавита  $\{A, C, G, T, N\}$

### Вход:

Текст ( $T$ ,  $1 \leq |T| \leq 100000$ )

Шаблон ( $P$ ,  $1 \leq |P| \leq 40$ )

Символ джокера

### Выход:

Строки с номерами позиций вхождений шаблона (каждая строка содержит только один номер).

Номера должны выводиться в порядке возрастания.

---

### Sample Input:

ACTANCA

A\$A\$A\$

\$

---

### Sample Output:

1

## Описание алгоритма.

Алгоритм Ахо-Корасика — это эффективный алгоритм поиска множества подстрок в тексте. Он сочетает в себе идеи использования бора (префиксного дерева) и автомата.

Этапы алгоритма (используется в заданиях 6.1 и 6.2):

1. Построение бора (Trie). Создаём бор для всех искомых строк:

1.1. Корень соответствует пустой строке.

1.2. Для каждого символа строки переходим по соответствующему ребру (если его нет — создаем).

1.3. В конце помечаем терминальную вершину (храним длину образца или его идентификатор).

## 2. Построение суффиксных и конечных ссылок:

2.1. Строим суффиксные ссылки (перебираем узлы-детей и их символы):

2.1.1. Суффиксная ссылка корня и суффиксная ссылка его потомков указывает на сам корень.

2.1.2. Иначе идем по суффиксным ссылкам родителя, пока не найдем узел-ребенка с текущим символом.

2.2. Построение конечных ссылок:

2.2.1. Если суффиксная ссылка ведет в терминальный узел, она является терминальной

2.2.2. Иначе продолжаем идти по суффиксным ссылкам, пока не найдем терминальный узел.

## 3. Поиск шаблонов в тексте:

3.1. Начинаем с корня.

3.2. Для каждого символа в тексте:

3.2.1. Пока символ и символ в рассматриваемом узле не совпали, идем по суффиксным ссылкам.

3.2.2. Если символы совпали, переходим в дочерний узел.

3.3. Проверяем является ли данный узел терминальным и идем по конечным ссылкам.

3.4. Сохраняем найденные вхождения.

## 4. В итоге получаем найденные вхождения шаблонов.

Поиск при использовании шаблона с символом-джокером (используется в задании 6.2):

1. Делим шаблон на подшаблоны игнорируя символы-джокеры.

2. Ищем вхождения подшаблонов алгоритмом Ахо-Корасика.

3. Ищем вхождения шаблона:

3.1. Создается массив, который хранит количество найденных подстрок для каждой позиции текста.

3.2. Перебираем позиции найденных вхождений подшаблонов. Увеличиваем значение ячейки массива с индексов позиции вхождения.

3.3. Если число в ячейки массива соответствует числу подшаблонов, индекс ячейки — начало вхождения полного шаблона.

Описание алгоритма индивидуализации:

1. Ищем вхождения шаблонов алгоритмом Ахо-Корасика.

2. Рассматриваем найденный массив вхождений. Рассматриваем каждое вхождение, с остальными, с которыми еще не рассматривалось это вхождение ранее.

3. Вычисляем концы рассматриваемых вхождений и проверяем не оканчивается ли хотя бы одно из них раньше, чем начинается другое.

3.1. Если оканчивается рассматриваем следующую пару

3.2. Иначе находим точку пересечения и заносим информацию о вхождениях и точке пересечения в массив найденных пересечений

4. В конце получаем массив найденных пересечений.

	Сложность по времени
Построение бора	$O(L)$
Построение ссылок	$O(A * L)$
Поиск	$O(T + L + n)$
Деление шаблона на подшаблоны	$O(P)$
Поиск пересечений	$O(n^2)$

	Сложность по памяти
--	---------------------

Хранение автомата	$O(A * L)$
-------------------	------------

, где  $L = \Sigma$  длин всех шаблонов,  $T$  – длина текста,  $A$  – число символов в алфавите (число разных символов среди всех шаблонов),  $P$  – длина шаблона с джокерами,  $n$  – число вхождений

Исходные коды реализующие алгоритмы из заданий 6.1 и 6.2 представлены в приложении (приложении 1 и 2 соответственно).

### Описание функций и структур данных.

Структуры данных:

1. *TrieNode* – структура, представляющая из себя узел бора.

Содержит поля:

*id* - уникальный идентификатор узла (для вывода)

*isTerminal* - флаг, является ли вершина терминалом

*patternIndices* - номера шаблонов, заканчивающихся в этом узле

*patternLength* - длина шаблона (для определения позиции)

*childrens* - указатели на дочерние узлы

*suffixLink* - суффиксная ссылка

*finalLink* - конечная ссылка

2. *AhoCorasicAlgorithm* – класс, который хранит автомат (бор с ссылками) реализует методы для работы с ним, а так же сам алгоритм Ахо-Корасика.

3.  $C = [0] * (\text{len}(\text{text}) + 1)$  - массив, где  $C[i]$  - количество встретившихся в тексте безмасочных подстрок шаблона, который начинается в тексте на позиции  $i$ .

Функции и методы:

1. *def \_\_init\_\_(self, patterns: list):* метод класса *AhoCorasicAlgorithm*, инициализирует алгоритм, создает автомат.

Содержит:

*root* – корень бора (тип *TrieNode*)

*nodeCount* – счетчик числа узлов в боре (тип *int*)

Параметры:

*patterns* – массив шаблонов, добавляемых в бор (тип *list*).

Сложность по времени:  $O(L)$ , где  $L = \Sigma$  длин всех шаблонов

2. *def \_\_add(self, pattern: str, index: int):* метод класса *AhoCorasicAlgorithm*, добавляет новый шаблон в бор.

Параметры:

*pattern* – массив шаблон, добавляемый в бор (тип *str*).

*index* – индекс добавляемого шаблона.

Сложность по времени:  $O(L)$ , где  $L = \Sigma$  длин всех шаблонов

3. *def \_\_makeLinks(self):* метод класса *AhoCorasicAlgorithm*, создает суффиксные и конечные ссылки.

Сложность по времени:  $O(A * L)$ , где  $L = \Sigma$  длин всех шаблонов,  $A$  – число символов в алфавите (число разных символов среди всех шаблонов)

4. *def search(self, text):* метод класса *AhoCorasicAlgorithm*, ищет все вхождения шаблонов в тексте.

Параметры:

*text* – текст, в котором осуществляется поиск (тип *str*).

Возвращает:

*results* – отсортированный по позициям в текст массив найденных вхождений (элемент массива имеет вид (позиция начала вхождения, индекс шаблона) – *(int, int)*)

Сложность по времени:  $O(T + L + n)$ , где  $L = \Sigma$  длин всех шаблонов,  $T$  – длина текста,  $n$  – число вхождений

5. *def getNodeCount(self)*: метод класса *AhoCorasicAlgorithm*, позволяет получить число узлов в автомате.

Возвращает:

*nodeCount* – счетчик числа узлов в боре (тип *int*)

Сложность по времени:  $O(1)$

6. *def \_\_printAutomat(self)*: метод класса *AhoCorasicAlgorithm*, выводит информацию о построенном автомате.

Сложность по времени:  $O(L)$ , где  $L = \Sigma$  длин всех шаблонов

7. *def findIntersectingPatterns(results: list, patterns: list)*: ищет пересечения вхождений.

Параметры:

*results* – массив найденных вхождений (тип *list*).

*patterns* – массив шаблонов

Возвращает:

*intersectionPairs* – массив с найденными пересечениями (тип *list*).

Сложность по времени:  $O(n^2)$ , где  $n$  – число элементов массива *results*.

8. *def findSubPatterns(pattern: str, jokerSymbol: str)*: делит шаблон с символом-джокером на отдельные части.

Параметры:

*pattern* – шаблон с джокер-символами, который требуется разделить (тип *str*).

*jokerSymbol* – джокер-символ (тип *str*).

Возвращает:



*subPatterns, positions* – список подшаблонов и список их позиций  
(типы *list*)

Сложность по времени:  $O(P)$ , где  $P$  – длина шаблона с джокерами.

### Тестирование.

Тестирование программ представлено на таблице.

Таблица 1. Задание 6.1

№ Теста	Входные данные	Выходные данные	Комментарий
1	NTAG 3 TAGT TAG T	2 2 2 3	Верно
2	shershe 4 he she his hers	1 2 2 1 2 4 5 2 6 1	Верно
3	TTTTTTTTTT 2 NAN GNG		Верно
4	TATTATAT 3 TATT TAT TATTAT	1 1 1 2 1 3 4 2 6 2	Верно
5	ABBBAAABAB 3 ABB ABB ABA	1 1 1 2 6 3	Верно

Таблица 2. Задание 6.2

№ Теста	Входные данные	Выходные данные	Комментарий
1	ACTANCA A\$\$A\$ \$	1	Верно
2	ABBBACBBABVBC A\$BB\$ \$	1 5 9	Верно
3	ABBBACBBABVBC \$C\$ \$	5	Верно
4	ABBA C\$ \$		Верно

Таблица 3. Индивидуализация

№ Теста	Входные данные	Выходные данные	Комментарий
1	NTAG 3 TAGT TAG T	TAG и T пересекаются в точке 2	Верно
2	shershe 4 he she his hers	she и he пересекаются в точке 2 she и her пересекаются в точке 2 he и her пересекаются в точке 2 her и she пересекаются в точке 5 she и he пересекаются в точке 6	Верно
3	ABBBACBBABBB A\$BB\$ \$	A\$BB\$ и A\$BB\$ пересекаются в точке 5 A\$BB\$ и A\$BB\$ пересекаются в точке 9	Верно
4	ABCMMMCAB 2 ABC CAB	Пересечений нет	Верно

Примеры работ программ с выводом промежуточных результатов при введённом значении указаны в приложении (приложение 3 и приложение 4 для программ 1 и 2 соответственно)

### **Выводы.**

Был изучен алгоритм Ахо-Корасика поиска всех вхождений набора подстрок в строку. Разработана программа, реализующая алгоритм Ахо-Корасика, а также алгоритм поиска шаблона с джокерами.

# ПРИЛОЖЕНИЕ

## ПРИЛОЖЕНИЕ 1

```
FROM COLLECTIONS IMPORT DEQUE
IMPORT GRAPHVIZ

CLASS TRIENODE:
    """КЛАСС УЗЛА ВОРА"""
    DEF __INIT__(SELF, NODEID: INT):
        SELF.ID = NODEID          # УНИКАЛЬНЫЙ ИДЕНТИФИКАТОР УЗЛА

        SELF.ISTERMINAL = FALSE    # ФЛАГ, ЯВЛЯЕТСЯ ЛИ ВЕРШИНА ТЕРМИНАЛОМ
        SELF.PATTERNINDICES = []   # НОМЕРА ШАБЛОНОВ, ЗАКАНЧИВАЮЩИХСЯ В ЭТОМ УЗЛЕ
        SELF.PATTERNLENGTH = 0     # ДЛИНА ШАБЛОНА

        SELF.CHILDRENS = {}        # УКАЗАТЕЛИ НА ДОЧЕРНИЕ УЗЛЫ
        SELF.SUFFIXLINK = NONE     # СУФФИКСНАЯ ССЫЛКА
        SELF.FINALLINK = NONE     # КОНЕЧНАЯ ССЫЛКА

CLASS AHOCORASICALGORITHM:
    """КЛАСС, РЕАЛИЗУЮЩИЙ АЛГОРИТМ АХО-КОРАСИКА"""
    DEF __INIT__(SELF, PATTERNS: LIST):
        """ИНИЦИАЛИЗИРУЕТ АЛГОРИТМ, СОЗДАЕТ АВТОМАТ"""

        SELF.ROOT = TRIENODE(0)    # КОРНЕВОЙ УЗЕЛ
        SELF.ROOT.SUFFIXLINK = SELF.ROOT # СУФФИКСНАЯ ССЫЛКА КОРНЯ
        SELF.NODECOUNT = 1        # СЧЕТЧИК УЗЛОВ
        SELF.PATTERNS = PATTERNS   # СОХРАНЯЕМ ШАБЛОНЫ
        SELF.MAX_ARCS = 0          # МАКСИМАЛЬНОЕ КОЛИЧЕСТВО ДУГ ИЗ ОДНОЙ ВЕРШИНЫ

        # СТРОИМ ВОР
        FOR INDEX IN RANGE(LEN(PATTERNS)):
            SELF.__ADD(PATTERNS[INDEX], INDEX)

        # СОЗДАЕМ ССЫЛКИ, СТРОИМ АВТОМАТ
        SELF.__MAKELINKS()

    DEF __ADD(SELF, PATTERN: STR, INDEX: INT):
        """ДОБАВЛЯЕТ НОВЫЙ ШАБЛОН В ВОР"""
        CURRENTNODE = SELF.ROOT
        # ПЕРЕБИРАЕМ СИМВОЛЫ ШАБЛОНА
        FOR CHAR IN PATTERN:
            IF CHAR NOT IN CURRENTNODE.CHILDRENS:
                NEWNODE = TRIENODE(SELF.NODECOUNT)
                CURRENTNODE.CHILDRENS[CHAR] = NEWNODE
                SELF.NODECOUNT += 1
                # ОБНОВЛЯЕМ МАКСИМАЛЬНОЕ КОЛИЧЕСТВО ДУГ
                SELF.MAX_ARCS = MAX(SELF.MAX_ARCS, LEN(CURRENTNODE.CHILDRENS))
            CURRENTNODE = CURRENTNODE.CHILDRENS[CHAR]

        CURRENTNODE.ISTERMINAL = TRUE
        CURRENTNODE.PATTERNINDICES.APPEND(INDEX)
        CURRENTNODE.PATTERNLENGTH = LEN(PATTERN)

    DEF __MAKELINKS(SELF):
        """СОЗДАЕТ СУФФИКСНЫЕ И КОНЕЧНЫЕ ССЫЛКИ"""
        QUEUE = DEQUE()
        QUEUE.APPEND(SELF.ROOT)

        WHILE QUEUE:
            CURRENTNODE = QUEUE.POPLEFT()
```

```

FOR CHAR, CHILDNODE IN CURRENTNODE.CHILDRENS.ITEMS():
    QUEUE.APPEND(CHILDNODE)

    # ДЛЯ ДЕТЕЙ КОРНЯ СУФФИКСНАЯ ССЫЛКА ВЕДЕТ В КОРЕНЬ
    IF CURRENTNODE == SELF.ROOT:
        CHILDNODE.SUFFIXLINK = SELF.ROOT
    ELSE:
        # ИЩЕМ ПЕРВУЮ ВОЗМОЖНУЮ СУФФИКСНУЮ ССЫЛКУ
        TEMP = CURRENTNODE.SUFFIXLINK
        WHILE (TEMP != SELF.ROOT) AND (CHAR NOT IN TEMP.CHILDRENS):
            TEMP = TEMP.SUFFIXLINK

        # УСТАНАВЛИВАЕМ НАЙДЕННУЮ ССЫЛКУ ИЛИ ССЫЛКУ НА КОРЕНЬ
        IF CHAR IN TEMP.CHILDRENS:
            CHILDNODE.SUFFIXLINK = TEMP.CHILDRENS[CHAR]
        ELSE:
            CHILDNODE.SUFFIXLINK = SELF.ROOT

    # ПОСТРОЕНИЕ КОНЕЧНОЙ ССЫЛКИ
    IF CHILDNODE.SUFFIXLINK.ISTERMINAL:
        CHILDNODE.FINALLINK = CHILDNODE.SUFFIXLINK
    ELSE:
        CHILDNODE.FINALLINK = CHILDNODE.SUFFIXLINK.FINALLINK

DEF SEARCH(SELF, TEXT):
    """Ищет все вхождения шаблонов в тексте"""
    RESULTS = []
    CURRENTNODE = SELF.ROOT

    FOR POSITION IN RANGE(LEN(TEXT)):
        CHAR = TEXT[POSITION]

        # ИСПОЛЬЗУЕМ СУФФИКСНЫЕ ССЫЛКИ ПРИ ОТСУТСТВИИ ПЕРЕХОДА
        WHILE (CURRENTNODE != SELF.ROOT) AND (CHAR NOT IN CURRENTNODE.CHILDRENS):
            CURRENTNODE = CURRENTNODE.SUFFIXLINK

        # ПЕРЕХОДИМ ПО СИМВОЛУ, ЕСЛИ ПЕРЕХОД СУЩЕСТВУЕТ
        IF CHAR IN CURRENTNODE.CHILDRENS:
            CURRENTNODE = CURRENTNODE.CHILDRENS[CHAR]

        # ПРОВЕРЯЕМ ТЕРМИНАЛЬНЫЕ УЗЛЫ
        IF CURRENTNODE.ISTERMINAL:
            FOR PATTERNINDEX IN CURRENTNODE.PATTERNINDICES:
                STARTPOSITION = POSITION - CURRENTNODE.PATTERNLENGTH + 1
                RESULTS.APPEND((STARTPOSITION, PATTERNINDEX, CURRENTNODE.PATTERNLENGTH))

        # ПРОВЕРЯЕМ КОНЕЧНЫЕ ССЫЛКИ ДЛЯ НАХОЖДЕНИЯ ВСЕХ ВЛОЖЕННЫХ ШАБЛОНОВ
        TEMP = CURRENTNODE.FINALLINK
        WHILE TEMP:
            FOR PATTERNINDEX IN TEMP.PATTERNINDICES:
                STARTPOSITION = POSITION - TEMP.PATTERNLENGTH + 1
                RESULTS.APPEND((STARTPOSITION, PATTERNINDEX, TEMP.PATTERNLENGTH))
            TEMP = TEMP.FINALLINK

    RETURN SORTED(RESULTS)

DEF GETNODECOUNT(SELF):
    """Позволяет получить число узлов в автомате"""
    RETURN SELF.NODECOUNT

DEF GETMAXARCS(SELF):
    """Возвращает максимальное количество дуг из одной вершины"""
    RETURN SELF.MAX_ARCS

```

```

DEF VISUALIZEBOR (SELF, FILENAME="BOR_TREE") :
    """"Визуализирует бор в PNG файл""""
    DOT = GRAPHVIZ.DIGRAPH (COMMENT='BOR TREE')
    DOT.ATTR ('NODE', SHAPE='CIRCLE')

    QUEUE = DEQUE ([SELF.ROOT])
    VISITED = SET ([SELF.ROOT])

    # ДОБАВЛЯЕМ КОРНЕВОЙ УЗЕЛ
    ROOT_LABEL = "0"
    IF SELF.ROOT.ISTERMINAL:
        ROOT_LABEL += "\\n" + ",".JOIN (STR (I+1) FOR I IN SELF.ROOT.PATTERNINDICES)
    DOT.NODE ('0', ROOT_LABEL, STYLE='FILLED', FILLCOLOR='LIGHTBLUE')

    WHILE QUEUE:
        NODE = QUEUE.POPLEFT ()

        FOR CHAR, CHILD IN NODE.CHILDRENS.ITEMS () :
            IF CHILD NOT IN VISITED:
                VISITED.ADD (CHILD)
                QUEUE.APPEND (CHILD)

            # СОЗДАЕМ МЕТКУ ДЛЯ УЗЛА
            NODE_LABEL = F"{CHILD.ID}"
            IF CHILD.ISTERMINAL:
                PATTERNS = ",".JOIN (STR (I+1) FOR I IN CHILD.PATTERNINDICES)
                NODE_LABEL += F"\\n{PATTERNS}"

            DOT.NODE (STR (CHILD.ID), NODE_LABEL)

            # ДОБАВЛЯЕМ РЕБРО
            DOT.EDGE (STR (NODE.ID), STR (CHILD.ID), LABEL=CHAR)

    # СОХРАНЯЕМ В ФАЙЛ
    DOT.RENDER (FILENAME, FORMAT='PNG', CLEANUP=TRUE)
    PRINT (F"ДЕРЕВО БОРА СОХРАНЕНО В ФАЙЛ: {FILENAME}.PNG")

DEF PRINTAUTOMATINFO (SELF) :
    """"Выводит краткую информацию об автомате""""
    PRINT (F"Информация об автомате:")
    PRINT (F"Количество узлов: {SELF.GETNODECOUNT ()}")
    PRINT (F"Количество шаблонов: {LEN (SELF.PATTERNS) }")
    PRINT (F"Максимальное количество дуг из одной вершины: {SELF.GETMAXARCS ()}")

    # ПОДСЧИТЫВАЕМ ДОПОЛНИТЕЛЬНУЮ СТАТИСТИКУ
    TERMINAL_COUNT = 0
    QUEUE = DEQUE ([SELF.ROOT])

    WHILE QUEUE:
        NODE = QUEUE.POPLEFT ()

        IF NODE.ISTERMINAL:
            TERMINAL_COUNT += 1

        FOR CHILD IN NODE.CHILDRENS.VALUES () :
            QUEUE.APPEND (CHILD)

    PRINT (F"Терминальных узлов: {TERMINAL_COUNT}")

DEF REMOVEFOUND PATTERNS (ТЕХТ, RESULTS, PATTERNS) :
    """"Вырезает найденные образцы из строки и возвращает остаток""""
    IF NOT RESULTS:

```

```

    RETURN TEXT

# Создаем массив для отметки позиций, которые нужно удалить
REMOVE_MASK = [False] * len(TEXT)

# Помечаем позиции, которые попадают в найденные образцы
for start_pos, pattern_idx, pattern_length in results:
    pattern = patterns[pattern_idx]
    end_pos = start_pos + pattern_length

    # Проверяем границы и помечаем позиции для удаления
    for i in range(max(0, start_pos), min(len(TEXT), end_pos)):
        REMOVE_MASK[i] = True

# Собираем остаток строки
REMAINDER = []
for i, char in enumerate(TEXT):
    if not REMOVE_MASK[i]:
        REMAINDER.append(char)

RETURN ' '.join(REMAINDER)

def findPatternRanges(results, patterns):
    """Находит диапазоны найденных образцов для наглядного вывода"""
    RANGES = []
    for start_pos, pattern_idx, pattern_length in results:
        pattern = patterns[pattern_idx]
        end_pos = start_pos + pattern_length
        RANGES.append((start_pos, end_pos, pattern))
    RETURN SORTED(RANGES)

def main():
    # Ввод данных
    TEXT = input("Введите текст: ")
    NUM = int(input("Введите количество шаблонов: "))
    PATTERNS = [input(f"Введите шаблон {i + 1}: ") for i in range(NUM)]

    # Создаем автомат Ахо-Корасика
    ahoCorasicAlgorithm = AhoCorasicAlgorithm(PATTERNS)

    # Визуализируем бор (опционально)
    try:
        ahoCorasicAlgorithm.visualizeBOR("BOR_visualization")
    except:
        print("Для визуализации бора установите graphviz: pip install graphviz")

    # Выводим информацию об автомате
    ahoCorasicAlgorithm.printAutomatInfo()

    # Пункт 1 варианта 5: Максимальное количество дуг из одной вершины
    max_arcs = ahoCorasicAlgorithm.getMaxArcs()
    print(f"\n=== РЕЗУЛЬТАТЫ ДЛЯ ВАРИАНТА 5 ===")
    print(f"1. Максимальное количество дуг, исходящих из одной вершины: {max_arcs}")

    # Поиск вхождений
    results = ahoCorasicAlgorithm.search(TEXT)

    # Вывод результатов поиска
    print(f"\nНайдено вхождений: {len(results)}")
    if results:
        print("Найденные вхождения:")
        RANGES = findPatternRanges(results, PATTERNS)

```

```

FOR START_POS, END_POS, PATTERN IN RANGES:
    PRINT(F"    Позиции {START_POS}-{END_POS-1}: '{PATTERN}'")

# ПУНКТ 2 ВАРИАНТА 5: ВЫРЕЗАЕМ НАЙДЕННЫЕ ОБРАЗЦЫ
REMAINDER = REMOVEFOUND PATTERNS(ТЕКСТ, RESULTS, PATTERNS)
PRINT(F"\n2. РЕЗУЛЬТАТ ВЫРЕЗАНИЯ НАЙДЕННЫХ ОБРАЗЦОВ:")
PRINT(F"    ИСХОДНЫЙ ТЕКСТ: '{ТЕКСТ}'")
PRINT(F"    ОСТАТОК СТРОКИ: '{REMAINDER}'")

# ДОПОЛНИТЕЛЬНАЯ ИНФОРМАЦИЯ О ВЫРЕЗАНИИ
ORIGINAL_LENGTH = LEN(ТЕКСТ)
REMAINDER_LENGTH = LEN(REMAINDER)
REMOVED_LENGTH = ORIGINAL_LENGTH - REMAINDER_LENGTH
PRINT(F"    СТАТИСТИКА: УДАЛЕНО {REMOVED_LENGTH} СИМВОЛОВ ИЗ {ORIGINAL_LENGTH}")

IF REMAINDER:
    PRINT(F"    СОХРАНЕНО: {REMAINDER_LENGTH} СИМВОЛОВ")
ELSE:
    PRINT(F"    ВСЕ ТЕКСТ БЫЛ УДАЛЕН (СОСТОИТ ТОЛЬКО ИЗ ШАБЛОНОВ)")
ELSE:
    PRINT("ШАБЛОНЫ НЕ НАЙДЕНЫ В ТЕКСТЕ")
    PRINT(F"ОСТАТОК СТРОКИ: '{ТЕКСТ}'")

IF __NAME__ == "__MAIN__":
    MAIN()

```

## Приложение 2

```

FROM COLLECTIONS IMPORT DEQUE
IMPORT GRAPHVIZ

CLASS TRIENode:
    """Класс узла бора"""
    DEF __INIT__(SELF, NODEID: INT):
        SELF.ID = NODEID # УНИКАЛЬНЫЙ ИДЕНТИФИКАТОР УЗЛА

        SELF.ISTERMINAL = FALSE # ФЛАГ, ЯВЛЯЕТСЯ ЛИ ВЕРШИНА ТЕРМИНАЛОМ
        SELF.PATTERNINDICES = [] # НОМЕРА ШАБЛОНОВ, ЗАКАНЧИВАЮЩИХСЯ В ЭТОМ УЗЛЕ
        SELF.PATTERNLENGTH = 0 # ДЛИНА ШАБЛОНА

        SELF.CHILDRENS = {} # УКАЗАТЕЛИ НА ДОЧЕРНИЕ УЗЛЫ
        SELF.SUFFIXLINK = NONE # СУФФИКСНАЯ ССЫЛКА
        SELF.FINALLINK = NONE # КОНЕЧНАЯ ССЫЛКА

CLASS AhoCorasickAlgorithm:
    """Класс, реализующий алгоритм Ахо-Корасика"""
    DEF __INIT__(SELF, PATTERNS: LIST):
        """ИНИЦИАЛИЗИРУЕТ АЛГОРИТМ, СОЗДАЕТ АВТОМАТ"""

        SELF.ROOT = TRIENode(0) # КОРНЕВОЙ УЗЕЛ
        SELF.ROOT.SUFFIXLINK = SELF.ROOT # СУФФИКСНАЯ ССЫЛКА КОРНЯ
        SELF.NODECOUNT = 1 # СЧЕТЧИК УЗЛОВ
        SELF.PATTERNS = PATTERNS # СОХРАНЯЕМ ШАБЛОНЫ
        SELF.MAX_ARCS = 0 # МАКСИМАЛЬНОЕ КОЛИЧЕСТВО ДУГ ИЗ ОДНОЙ ВЕРШИНЫ

        # СТРОИМ БОР
        FOR INDEX IN RANGE(LEN(PATTERNS)):
            SELF.__ADD(PATTERNS[INDEX], INDEX)

        # СОЗДАЕМ ССЫЛКИ, СТРОИМ АВТОМАТ
        SELF.__MAKELINKS()

```

```

DEF _ADD (SELF, PATTERN: STR, INDEX: INT):
    """ДОБАВЛЯЕТ НОВЫЙ ШАБЛОН В БОР"""
    CURRENTNODE = SELF.ROOT
    # ПЕРЕБИВАЕМ СИМВОЛЫ ШАБЛОНА
    FOR CHAR IN PATTERN:
        IF CHAR NOT IN CURRENTNODE.CHILDRENS:
            NEWNODE = TRIENODE (SELF.NODECOUNT)
            CURRENTNODE.CHILDRENS[CHAR] = NEWNODE
            SELF.NODECOUNT += 1
            # ОБНОВЛЯЕМ МАКСИМАЛЬНОЕ КОЛИЧЕСТВО ДУГ
            SELF.MAX_ARCS = MAX (SELF.MAX_ARCS, LEN (CURRENTNODE.CHILDRENS))
            CURRENTNODE = CURRENTNODE.CHILDRENS[CHAR]

    CURRENTNODE.ISTERMINAL = TRUE
    CURRENTNODE.PATTERNINDICES.APPEND (INDEX)
    CURRENTNODE.PATTERNLENGTH = LEN (PATTERN)

DEF _MAKELINKS (SELF):
    """СОЗДАЕТ СУФФИКСНЫЕ И КОНЕЧНЫЕ ССЫЛКИ"""
    QUEUE = DEQUE ()
    QUEUE.APPEND (SELF.ROOT)

    WHILE QUEUE:
        CURRENTNODE = QUEUE.POPLEFT ()

        FOR CHAR, CHILDNODE IN CURRENTNODE.CHILDRENS.ITEMS ():
            QUEUE.APPEND (CHILDNODE)

            # ДЛЯ ДЕТЕЙ КОРНЯ СУФФИКСНАЯ ССЫЛКА ВЕДЕТ В КОРЕНЬ
            IF CURRENTNODE == SELF.ROOT:
                CHILDNODE.SUFFIXLINK = SELF.ROOT
            ELSE:
                # ИЩЕМ ПЕРВУЮ ВОЗМОЖНУЮ СУФФИКСНУЮ ССЫЛКУ
                TEMP = CURRENTNODE.SUFFIXLINK
                WHILE (TEMP != SELF.ROOT) AND (CHAR NOT IN TEMP.CHILDRENS):
                    TEMP = TEMP.SUFFIXLINK

                # УСТАНАВЛИВАЕМ НАЙДЕННУЮ ССЫЛКУ ИЛИ ССЫЛКУ НА КОРЕНЬ
                IF CHAR IN TEMP.CHILDRENS:
                    CHILDNODE.SUFFIXLINK = TEMP.CHILDRENS[CHAR]
                ELSE:
                    CHILDNODE.SUFFIXLINK = SELF.ROOT

            # ПОСТРОЕНИЕ КОНЕЧНОЙ ССЫЛКИ
            IF CHILDNODE.SUFFIXLINK.ISTERMINAL:
                CHILDNODE.FINALLINK = CHILDNODE.SUFFIXLINK
            ELSE:
                CHILDNODE.FINALLINK = CHILDNODE.SUFFIXLINK.FINALLINK

DEF SEARCH (SELF, TEXT):
    """ИЩЕТ ВСЕ ВХОЖДЕНИЯ ШАБЛОНОВ В ТЕКСТЕ"""
    RESULTS = []
    CURRENTNODE = SELF.ROOT

    FOR POSITION IN RANGE (LEN (TEXT)):
        CHAR = TEXT[POSITION]

        # ИСПОЛЬЗУЕМ СУФФИКСНЫЕ ССЫЛКИ ПРИ ОТСУТСТВИИ ПЕРЕХОДА
        WHILE (CURRENTNODE != SELF.ROOT) AND (CHAR NOT IN CURRENTNODE.CHILDRENS):
            CURRENTNODE = CURRENTNODE.SUFFIXLINK

        # ПЕРЕХОДИМ ПО СИМВОЛУ, ЕСЛИ ПЕРЕХОД СУЩЕСТВУЕТ
        IF CHAR IN CURRENTNODE.CHILDRENS:

```



```

        CURRENTNODE = CURRENTNODE.CHILDRENS[CHAR]

# ПРОВЕРЯЕМ ТЕРМИНАЛЬНЫЕ УЗЛЫ
IF CURRENTNODE.ISTERMINAL:
    FOR PATTERNINDEX IN CURRENTNODE.PATTERNINDICES:
        STARTPOSITION = POSITION - CURRENTNODE.PATTERNLENGTH + 1
        RESULTS.APPEND((STARTPOSITION, PATTERNINDEX, CURRENTNODE.PATTERNLENGTH))

# ПРОВЕРЯЕМ КОНЕЧНЫЕ ССЫЛКИ ДЛЯ НАХОЖДЕНИЯ ВСЕХ ВЛОЖЕННЫХ ШАБЛОНОВ
TEMP = CURRENTNODE.FINALLINK
WHILE TEMP:
    FOR PATTERNINDEX IN TEMP.PATTERNINDICES:
        STARTPOSITION = POSITION - TEMP.PATTERNLENGTH + 1
        RESULTS.APPEND((STARTPOSITION, PATTERNINDEX, TEMP.PATTERNLENGTH))
    TEMP = TEMP.FINALLINK

RETURN SORTED(RESULTS)

DEF GETNODECOUNT(SELF):
    """Позволяет получить число узлов в автомате"""
    RETURN SELF.NODECOUNT

DEF GETMAXARCS(SELF):
    """Возвращает максимальное количество дуг из одной вершины"""
    RETURN SELF.MAX_ARCS

DEF VISUALIZEBOR(SELF, FILENAME="BOR_TREE"):
    """Визуализирует бор в PNG файл"""
    DOT = GRAPHVIZ.DIGRAPH(COMMENT='BOR TREE')
    DOT.ATTR('NODE', SHAPE='CIRCLE')

    QUEUE = DEQUE([SELF.ROOT])
    VISITED = SET([SELF.ROOT])

    # ДОБАВЛЯЕМ КОРНЕВОЙ УЗЕЛ
    ROOT_LABEL = "0"
    IF SELF.ROOT.ISTERMINAL:
        ROOT_LABEL += "\\N" + ",".JOIN(STR(I+1) FOR I IN SELF.ROOT.PATTERNINDICES)
    DOT.NODE('0', ROOT_LABEL, STYLE='FILLED', FILLCOLOR='LIGHTBLUE')

    WHILE QUEUE:
        NODE = QUEUE.POPLEFT()

        FOR CHAR, CHILD IN NODE.CHILDRENS.ITEMS():
            IF CHILD NOT IN VISITED:
                VISITED.ADD(CHILD)
                QUEUE.APPEND(CHILD)

            # СОЗДАЕМ МЕТКУ ДЛЯ УЗЛА
            NODE_LABEL = F"{CHILD.ID}"
            IF CHILD.ISTERMINAL:
                PATTERNS = ",".JOIN(STR(I+1) FOR I IN CHILD.PATTERNINDICES)
                NODE_LABEL += F"\\N{PATTERNS}"

            DOT.NODE(STR(CHILD.ID), NODE_LABEL)

            # ДОБАВЛЯЕМ РЕБРО
            DOT.EDGE(STR(NODE.ID), STR(CHILD.ID), LABEL=CHAR)

    # СОХРАНЯЕМ В ФАЙЛ
    DOT.RENDER(FILENAME, FORMAT='PNG', CLEANUP=TRUE)
    PRINT(F"ДЕРЕВО БОРА СОХРАНЕНО В ФАЙЛ: {FILENAME}.PNG")

DEF PRINTAUTOMATINFO(SELF):

```

```

    """Выводит краткую информацию об автомате"""
    PRINT (F"Информация об автомате:")
    PRINT (F"Количество узлов: {SELF.GETNODECOUNT()}")
    PRINT (F"Количество шаблонов: {LEN(SELF.PATTERNS)}")
    PRINT (F"Максимальное количество дуг из одной вершины: {SELF.GETMAXARCS()}")

    # Подсчитываем дополнительную статистику
    TERMINAL_COUNT = 0
    QUEUE = DEQUE([SELF.ROOT])

    WHILE QUEUE:
        NODE = QUEUE.POPLEFT()

        IF NODE.ISTERMINAL:
            TERMINAL_COUNT += 1

        FOR CHILD IN NODE.CHILDRENS.VALUES():
            QUEUE.APPEND(CHILD)

    PRINT (F"Терминальных узлов: {TERMINAL_COUNT}")

DEF REMOVEFOUND PATTERNS (TEXT, RESULTS, PATTERNS):
    """Вырезает найденные образцы из строки и возвращает остаток"""
    IF NOT RESULTS:
        RETURN TEXT

    # Создаем массив для отметки позиций, которые нужно удалить
    REMOVE_MASK = [FALSE] * LEN(TEXT)

    # Помечаем позиции, которые попадают в найденные образцы
    FOR START_POS, PATTERN_IDX, PATTERN_LENGTH IN RESULTS:
        PATTERN = PATTERNS[PATTERN_IDX]
        END_POS = START_POS + PATTERN_LENGTH

        # Проверяем границы и помечаем позиции для удаления
        FOR I IN RANGE(MAX(0, START_POS), MIN(LEN(TEXT), END_POS)):
            REMOVE_MASK[I] = TRUE

    # Собираем остаток строки
    REMAINDER = []
    FOR I, CHAR IN ENUMERATE(TEXT):
        IF NOT REMOVE_MASK[I]:
            REMAINDER.APPEND(CHAR)

    RETURN ' '.JOIN(REMAINDER)

DEF FINDPATTERNRANGES (RESULTS, PATTERNS):
    """Находит диапазоны найденных образцов для наглядного вывода"""
    RANGES = []
    FOR START_POS, PATTERN_IDX, PATTERN_LENGTH IN RESULTS:
        PATTERN = PATTERNS[PATTERN_IDX]
        END_POS = START_POS + PATTERN_LENGTH
        RANGES.APPEND((START_POS, END_POS, PATTERN))
    RETURN SORTED(RANGES)

DEF MAIN():
    # Ввод данных
    TEXT = INPUT("Введите текст: ")
    NUM = INT(INPUT("Введите количество шаблонов: "))
    PATTERNS = [INPUT(F"Введите шаблон {I + 1}: ") FOR I IN RANGE(NUM)]

```

```

# СОЗДАЕМ АВТОМАТ АХО-КОРАСИКА
AHO_CORASIC_ALGORITHM = AHO_CORASIC_ALGORITHM(PATTERNS)

# ВИЗУАЛИЗИРУЕМ БОР (ОПЦИОНАЛЬНО)
TRY:
    AHO_CORASIC_ALGORITHM.VISUALIZEBOR("BOR_VISUALIZATION")
EXCEPT:
    PRINT("Для визуализации бора установите GRAPHVIZ: PIP INSTALL GRAPHVIZ")

# ВЫВОДИМ ИНФОРМАЦИЮ ОБ АВТОМАТЕ
AHO_CORASIC_ALGORITHM.PRINTAUTOMATINFO()

# ПУНКТ 1 ВАРИАНТА 5: МАКСИМАЛЬНОЕ КОЛИЧЕСТВО ДУГ ИЗ ОДНОЙ ВЕРШИНЫ
MAX_ARCS = AHO_CORASIC_ALGORITHM.GETMAXARCS()
PRINT(F"\n=== РЕЗУЛЬТАТЫ ДЛЯ ВАРИАНТА 5 ===")
PRINT(F"1. МАКСИМАЛЬНОЕ КОЛИЧЕСТВО ДУГ, ИСХОДЯЩИХ ИЗ ОДНОЙ ВЕРШИНЫ: {MAX_ARCS}")

# ПОИСК ВХОЖДЕНИЙ
RESULTS = AHO_CORASIC_ALGORITHM.SEARCH(TEXT)

# ВЫВОД РЕЗУЛЬТАТОВ ПОИСКА
PRINT(F"\nНайдено вхождений: {LEN(RESULTS)}")
IF RESULTS:
    PRINT("Найденные вхождения:")
    RANGES = FINDPATTERNRANGES(RESULTS, PATTERNS)
    FOR START_POS, END_POS, PATTERN IN RANGES:
        PRINT(F"    Позиции {START_POS}-{END_POS-1}: '{PATTERN}'")

# ПУНКТ 2 ВАРИАНТА 5: ВЫРЕЗАЕМ НАЙДЕННЫЕ ОБРАЗЦЫ
REMAINDER = REMOVEFOUNDPATTERNS(TEXT, RESULTS, PATTERNS)
PRINT(F"\n2. РЕЗУЛЬТАТ ВЫРЕЗАНИЯ НАЙДЕННЫХ ОБРАЗЦОВ:")
PRINT(F"    Исходный текст: '{TEXT}'")
PRINT(F"    Остаток строки: '{REMAINDER}'")

# ДОПОЛНИТЕЛЬНАЯ ИНФОРМАЦИЯ О ВЫРЕЗАНИИ
ORIGINAL_LENGTH = LEN(TEXT)
REMAINDER_LENGTH = LEN(REMAINDER)
REMOVED_LENGTH = ORIGINAL_LENGTH - REMAINDER_LENGTH
PRINT(F"    СТАТИСТИКА: УДАЛЕНО {REMOVED_LENGTH} СИМВОЛОВ ИЗ {ORIGINAL_LENGTH}")

IF REMAINDER:
    PRINT(F"    СОХРАНЕНО: {REMAINDER_LENGTH} СИМВОЛОВ")
ELSE:
    PRINT(F"    ВСЬ ТЕКСТ БЫЛ УДАЛЕН (СОСТОИТ ТОЛЬКО ИЗ ШАБЛОНОВ)")
ELSE:
    PRINT("ШАБЛОНЫ НЕ НАЙДЕНЫ В ТЕКСТЕ")
    PRINT(F"Остаток строки: '{TEXT}'")

IF __NAME__ == "__MAIN__":
    MAIN()

```