



Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное
учреждение высшего образования
«Московский государственный технический университет имени Н.
Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н. Э. Баумана)

ФАКУЛЬТЕТ «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

Отчет по лабораторной работе № 1 по курсу "Анализ алгоритмов"

Тема Расстояние Левенштейна и Дамерау-Левенштейна

Студент Калашков П. А.

Группа ИУ7-56Б

Оценка (баллы) _____

Преподаватель Волкова Л. Л.

Содержание

Введение	3
1 Аналитическая часть	4
1.1 Расстояние Левенштейна	4
1.1.1 Матричный алгоритм нахождения расстояния	5
1.2 Расстояние Дамерау — Левенштейна	6
1.2.1 Рекурсивный алгоритм нахождения расстояния	7
1.2.2 Матричный алгоритм нахождения расстояния	8
1.2.3 Рекурсивный алгоритм нахождения расстояния с использо- ванием кеша	9
1.3 Вывод	9
2 Конструкторская часть	10
2.1 Описание используемых типов данных	10
2.2 Сведения о модулях программы	10
2.3 Схемы алгоритмов	11
2.4 Классы эквивалентности тестирования	15
2.5 Использование памяти	15
2.6 Вывод	16
3 Технологическая часть	17
3.1 Средства реализации	17
3.2 Листинги кода	17
3.3 Функциональные тесты	21
3.4 Вывод	21
4 Исследовательская часть	22
4.1 Технические характеристики	22
4.2 Демонстрация работы программы	22
4.3 Время выполнения алгоритмов	24
4.4 Вывод	26

Заключение	27
Список литературы	28

Введение

Операции работы со строками являются очень важной частью всего программирования. Часто возникает потребность в использовании строк для различных задач - обычные статьи, записи в базу данных и так далее. Отсюда возникает несколько важных задач, для решения которых нужны алгоритмы сравнения строк. Об этих алгоритмах и пойдет речь в данной работе. Подобные алгоритмы используются при:

- исправлении ошибок в тексте, предлагая заменить введенное слово с ошибкой на наиболее подходящее;
- поиске слова в тексте по подстроке (например, в поисковых системах или в биоинформатике для сравнения цепочек молекул);
- сравнении целых текстовых файлов.

Целью данной работы является изучение, реализация и исследование алгоритмов нахождения расстояний Левенштейна и Дamerau-Левенштейна. Для достижения поставленной цели необходимо выполнить следующие задачи:

- изучить и реализовать алгоритмы Левенштейна и Дamerau-Левенштейна;
- провести тестирование по времени и по памяти для алгоритмов Левенштейна и Дamerau-Левенштейна;
- провести сравнительный анализ по времени алгоритмов нахождения расстояния Левенштейна и Дamerau-Левенштейна;
- провести сравнительный анализ по времени матричной, рекурсивной, а также рекурсивной с использованием кэша реализаций алгоритма нахождения расстояния Дamerau-Левенштейна;
- описать и обосновать полученные результаты в отчете о выполненной лабораторной работе, выполненного как расчётно-пояснительная записка к работе.

1 Аналитическая часть

В данном разделе будут разобраны алгоритмы нахождения расстояния - алгоритмы Левенштейна и Дамерау-Левенштейна.

1.1 Расстояние Левенштейна

Расстояние Левенштейна [1] между двумя строками - метрика, позволяющая определить «схожесть» двух строк — минимальное количество операций вставки одного символа, удаления одного символа и замены одного символа на другой, необходимых для превращения одной строки в другую (каждая операция имеет свою цену - штраф).

Редакционное предписание - последовательность действий, необходимых для получения из первой строки вторую, и минимизирующую суммарную цену (и является расстоянием Левенштейна).

Пусть S_1 и S_2 - две строки, длиной N и M соответственно. Введем следующие обозначения:

- I (англ. Insert) - вставка символа в произвольной позиции ($w(\lambda, b) = 1$);
- D (англ. Delete) - удаление символа в произвольной позиции ($w(\lambda, b) = 1$);
- R (англ. Replace) - замена символа на другой ($w(a, b) = 1, a \neq b$);
- M (англ. Match) - совпадение двух символов ($w(a, a) = 0$).

С учетом введенных обозначений, расстояние Левенштейна может быть подсчитано по следующей рекуррентной формуле:

$$D(i, j) = \begin{cases} 0 & i = 0, j = 0 \\ i & j = 0, i > 0 \\ j & i = 0, j > 0 \\ \min\{ \\ \quad D(i, j - 1) + 1 \\ \quad D(i - 1, j) + 1 & i > 0, j > 0 \\ \quad D(i - 1, j - 1) + m(a[i], b[j]) & (1.2) \\ \} \end{cases} \quad (1.1)$$

Функция 1.2 определена как:

$$m(a, b) = \begin{cases} 0 & \text{если } a = b, \\ 1 & \text{иначе} \end{cases}. \quad (1.2)$$

1.1.1 Матричный алгоритм нахождения расстояния

Рекурсивный алгоритм вычисления расстояния Левенштейна может быть не эффективен при больших i и j , так как множество промежуточных значений $D(i, j)$ вычисляются не один раз, что сильно замедляет время выполнения программы.

В качестве оптимизации можно использовать *матрицу* для хранения промежуточных значений. Матрица имеет размеры:

$$(length(S1) + 1) * ((length(S2) + 1)), \quad (1.3)$$

где $length(S)$ – длина строки S

Значение в ячейке $[i, j]$ равно значению $D(S1[1...i], S2[1...j])$. Первая строка и первый столбец тривиальны.

Всю таблицу (за исключением первого столбца и первой строки) заполня-

ем в соответствии с формулой 1.8.

$$A[i][j] = \min \begin{cases} A[i-1][j] + 1 \\ A[i][j-1] + 1 \\ A[i-1][j-1] + m(S1[i], S2[j]) \end{cases} \quad (1.4)$$

Функция 1.9 определена как:

$$m(S1[i], S2[j]) = \begin{cases} 0, & \text{если } S1[i-1] = S2[j-1], \\ 1, & \text{иначе} \end{cases} \quad (1.5)$$

Результат вычисления расстояния Левенштейна будет ячейка матрицы с индексами $i = \text{length}(S1)$ и $j = \text{length}(S2)$.

1.2 Расстояние Дамерау — Левенштейна

Расстояние Дамерау-Левенштейна [2] между двумя строками, состоящими из конечного числа символов — это минимальное число операций вставки, удаления, замены одного символа и транспозиции двух соседних символов, необходимых для перевода одной строки в другую.

Является модификацией расстояния Левенштейна - добавлена операции *транспозиции*, то есть перестановки, двух символов.

Расстояние Дамерау — Левенштейна может быть найдено по формуле 1.6, которая задана как

$$d_{a,b}(i, j) = \begin{cases} \max(i, j), & \text{если } \min(i, j) = 0, \\ \min\{ \\ \quad d_{a,b}(i, j - 1) + 1, \\ \quad d_{a,b}(i - 1, j) + 1, \\ \quad d_{a,b}(i - 1, j - 1) + m(a[i], b[j]), & \text{иначе} \\ \quad \left[\begin{array}{ll} d_{a,b}(i - 2, j - 2) + 1, & \text{если } i, j > 1; \\ & a[i] = b[j - 1]; \\ & b[j] = a[i - 1] \\ & \infty, & \text{иначе} \end{array} \right. \\ \} \end{cases}, \quad (1.6)$$

Формула выводится по тем же соображениям, что и формула (1.1).

1.2.1 Рекурсивный алгоритм нахождения расстояния

Рекурсивный алгоритм вычисления расстояния Дамерау-Левенштейна реализует формулу 1.6

Минимальная цена преобразования - минимальное значение приведенных вариантов.

Если полагать, что a' , b' - строки a и b без последнего символа соответственно, а a'' , b'' - строки a и b без двух последних символов, то цена преобразования из строки a в b может быть выражена так:

1. сумма цены преобразования строки a' в b и цены проведения операции удаления, которая необходима для преобразования a' в a ;
2. сумма цены преобразования строки a в b' и цены проведения операции вставки, которая необходима для преобразования b' в b ;
3. сумма цены преобразования из a' в b' и операции замены, предполагая, что a и b оканчиваются на разные символы;

4. сумма цены преобразования из a'' в b'' и операции перестановки, предполагая, что длины a'' и b'' больше 1 и последние два символа a'' , поменянные местами, совпадут с двумя последними символами b'' ;
5. цена преобразования из a' в b' , предполагая, что a и b оканчиваются на один и тот же символ.

1.2.2 Матричный алгоритм нахождения расстояния

Рекурсивный алгоритм вычисления расстояния Дамерау-Левенштейна может быть не эффективен при больших i и j , так как множество промежуточных значений $D(i, j)$ вычисляются не один раз, что сильно замедляет время выполнения программы.

В качестве оптимизации можно использовать *матрицу* для хранения промежуточных значений. Матрица имеет размеры:

$$(length(S1) + 1) * ((length(S2) + 1)), \quad (1.7)$$

где $length(S)$ – длина строки S

Значение в ячейке $[i, j]$ равно значению $D(S1[1...i], S2[1...j])$. Первая строка и первый столбец тривиальны.

Всю таблицу (за исключением первого столбца и первой строки) заполняем в соответствии с формулой 1.8.

$$A[i][j] = \min \begin{cases} A[i-1][j] + 1 \\ A[i][j-1] + 1 \\ A[i-1][j-1] + m(S1[i], S2[j]) \\ A[i-2][j-2] + 1, \text{ если } i > 1, j > 1 \text{ и} \\ \quad S1[i-2] = S2[j-1], S2[i-1] = S2[j-2] \end{cases} \quad (1.9) \cdot (1.8)$$

Функция 1.9 определена как:

$$m(S1[i], S2[j]) = \begin{cases} 0, & \text{если } S1[i-1] = S2[j-1], \\ 1, & \text{иначе} \end{cases}. \quad (1.9)$$

Результат вычисления расстояния Дамерау-Левенштейна будет ячейка матрицы с индексами $i = \text{length}(S1)$ и $j = \text{length}(S2)$.

1.2.3 Рекурсивный алгоритм нахождения расстояния с использованием кеша

В качестве оптимизации рекурсивного алгоритма заполнения можно использовать *кеш*, который будет представлять собой матрицу.

Суть оптимизации - при выполнении рекурсии происходит параллельное заполнение матрицы.

Если рекурсивный алгоритм выполняет прогон для данных, которые еще не были обработаны, то результат нахождения заносится в матрицу. Иначе, если обработанные данные встречаются снова, то для них расстояние не находится и алгоритм переходит к следующему шагу.

1.3 Вывод

В данном разделе были теоретически разобраны формулы Левенштейна и Дамерау-Левенштейна, которые являются рекуррентными, что позволяет реализовать их как рекурсивно, так и итерационно.

В качестве входных данных в программу будет подаваться две строки, также реализовано меню для вызова алгоритмов и замеров времени. Ограничением для работы программного продукта является то, что программе на вход может подаваться строка на английском или русском языке, а также программа должна корректно обрабатывать случай ввода пустых строк.

Реализуемое ПО будет работать в двух режимах - пользовательский, в котором можно выбрать алгоритм и вывести для него посчитанное значение, а также экспериментальный режим, в котором можно произвести сравнение алгоритмов по времени работы на различных входных данных.

2 Конструкторская часть

В этом разделе будут представлено описание используемых типов данных, а также схемы алгоритмов вычисления расстояния Левенштейна и Дамерау-Левенштейна.

2.1 Описание используемых типов данных

При реализации алгоритмов будут использованы следующие структуры данных:

- две строки типа *str*;
- длина строки - целое число типа *int*;
- в матричной реализации алгоритмов поиска расстояний Левенштейна и Дамерау-Левенштейна, а также рекурсивной реализации с кешем - матрица, которая является двумерным списком типа *int*.

2.2 Сведения о модулях программы

Программа состоит из шести модулей:

- *main.py* - файл, содержащий точку входа;
- *menu.py* - файл, содержащий код меню программы;
- *test.py* - файл, содержащий код тестирования алгоритмов;
- *utils.py* - файл, содержащий служебные алгоритмы;
- *constants.py* - файл, содержащий константы программы;
- *algorithms.py* - файл, содержащий код всех алгоритмов.

2.3 Схемы алгоритмов

На рисунках 2.1-2.2 представлены схемы алгоритмов вычисления расстояния Левенштейна и Дамерау-Левенштейна.

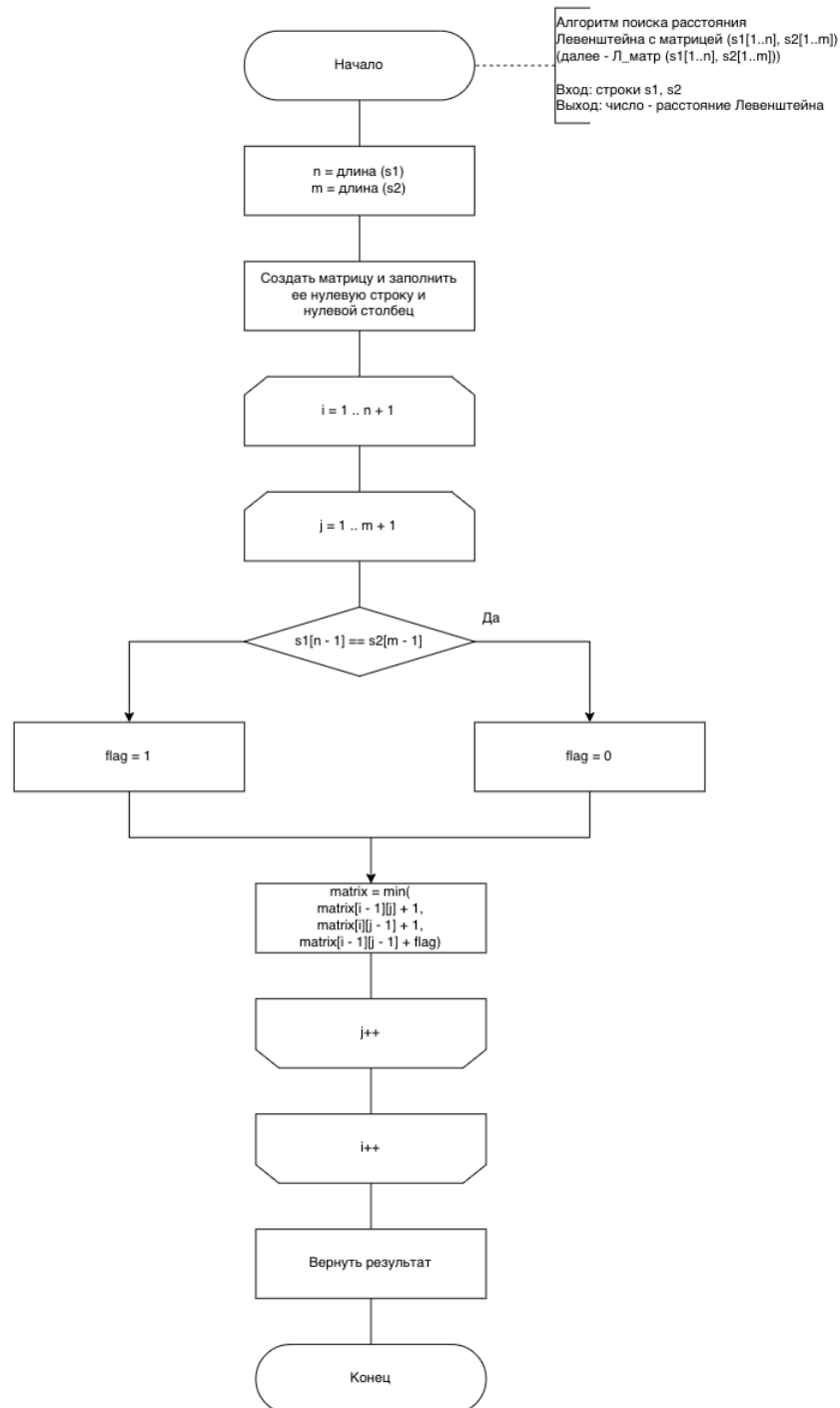


Рисунок 2.1 – Схема матричного алгоритма нахождения расстояния Левенштейна

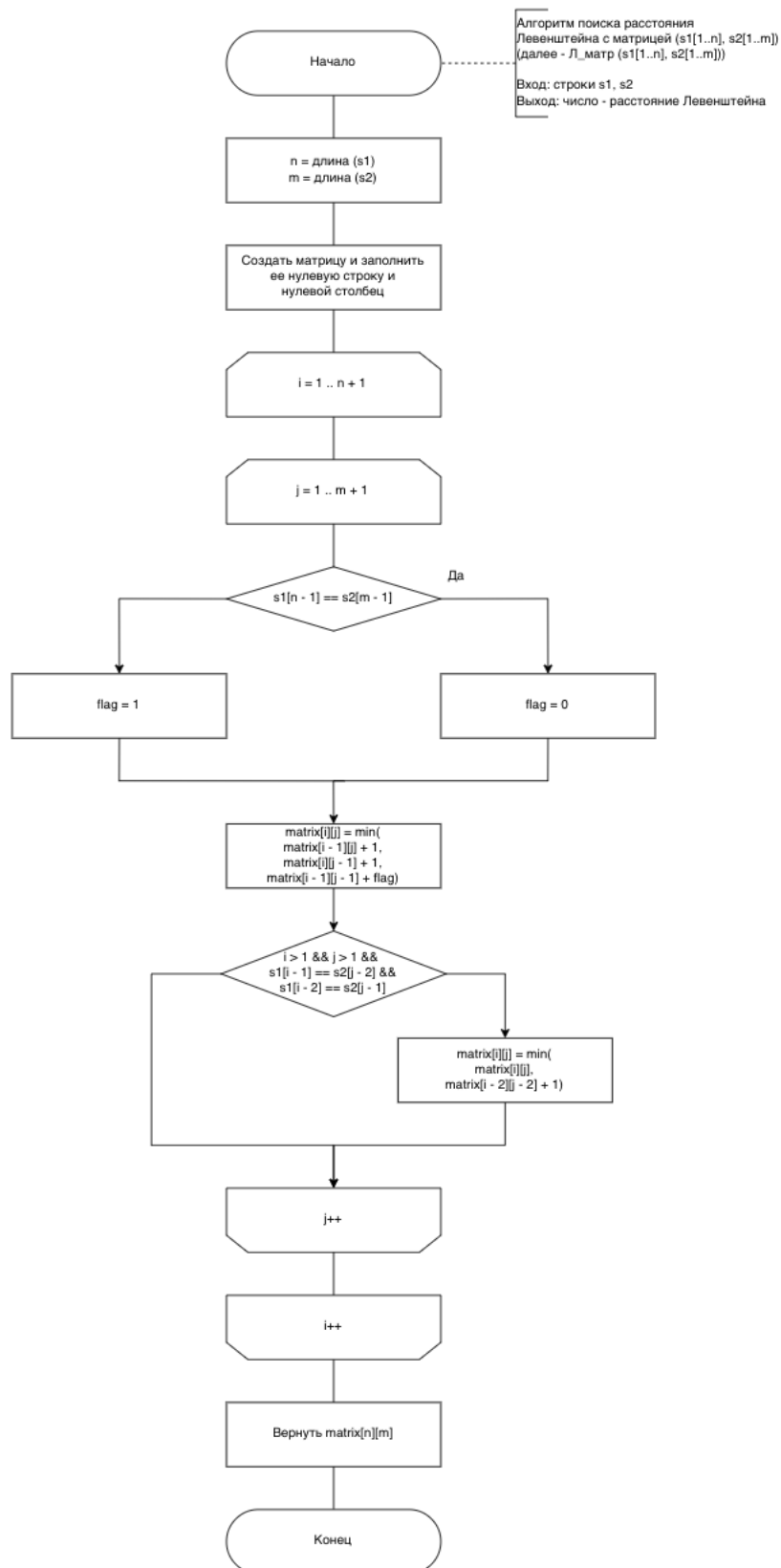


Рисунок 2.2 – Схема матричного алгоритма нахождения расстояния Дамерау-Левенштейна

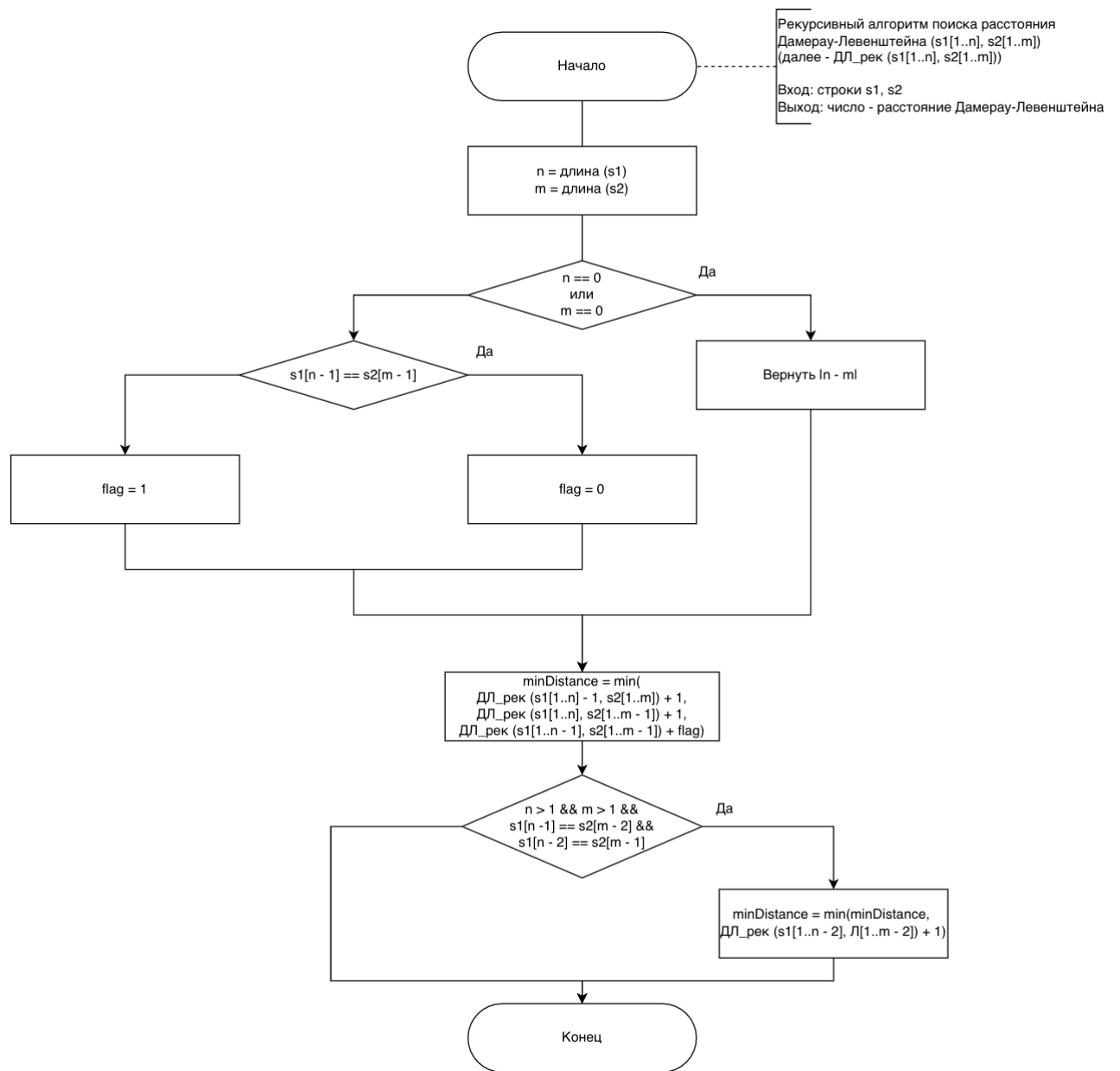


Рисунок 2.3 – Схема рекурсивного алгоритма нахождения расстояния Дамерау-Левенштейна

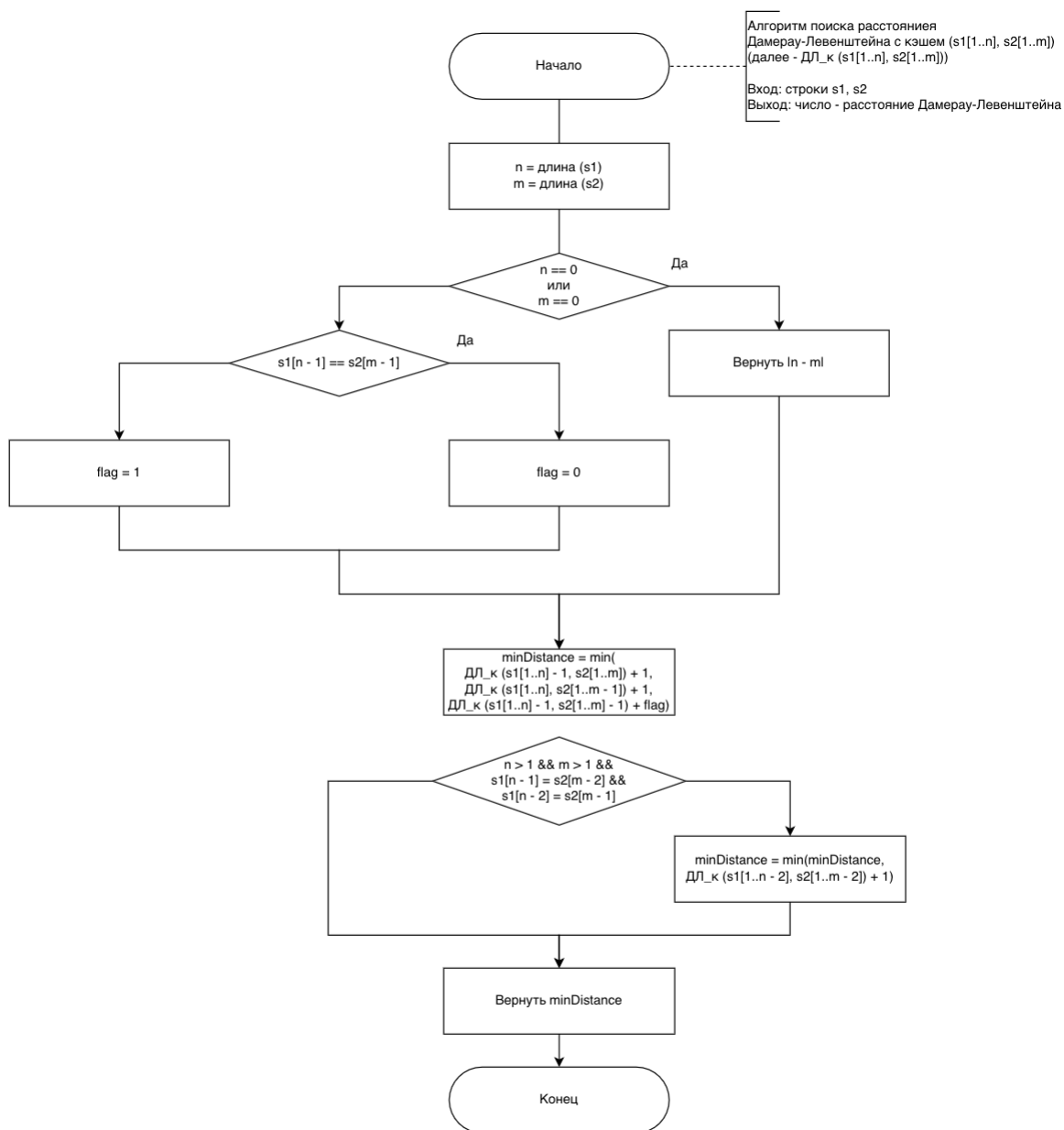


Рисунок 2.4 – Схема рекурсивного алгоритма нахождения расстояния Дамерау-Левенштейна с использованием кеша (матрицы)

2.4 Классы эквивалентности тестирования

Для тестирования выделены классы эквивалентности, представленные ниже.

1. Ввод двух пустых строк.
2. Одна из строк - пустая.
3. Расстояния, которые вычислены алгоритмами Левенштейна и Дамерау-Левенштейна, равны.
4. Расстояния, которые вычислены алгоритмами Левенштейна и Дамерау-Левенштейна, дают разные результаты.

2.5 Использование памяти

С точки зрения замеров и сравнения используемой памяти, алгоритмы Левенштейна и Дамерау-Левенштейна не отличаются друг от друга. Тогда рассмотрим только рекурсивную и матричную реализации данных алгоритмов.

Пусть:

- n - длина строки S_1
- m - длина строки S_2

Тогда затраты по памяти будут такими:

- алгоритм нахождения расстояния Левенштейна (матричный):
 - для матрицы - $((n + 1) * (m + 1)) * \text{sizeof}(\text{int})$
 - для S_1, S_2 - $(n + m) * \text{sizeof}(\text{char})$
 - для n, m - $2 * \text{sizeof}(\text{int})$
 - доп. переменные - $3 * \text{sizeof}(\text{int})$
 - адрес возврата

- алгоритм нахождения расстояния Дамерау-Левенштейна (матричный):
 - для матрицы - $((n + 1) * (m + 1)) * \text{sizeof}(\text{int})$
 - для S1, S2 - $(n + m) * \text{sizeof}(\text{char})$
 - для n, m - $2 * \text{sizeof}(\text{int})$
 - доп. переменные - $4 * \text{sizeof}(\text{int})$
 - адрес возврата

- алгоритм нахождения расстояния Дамерау-Левенштейна (рекурсивный), где для каждого вызова:
 - для S1, S2 - $(n + m) * \text{sizeof}(\text{char})$
 - для n, m - $2 * \text{sizeof}(\text{int})$
 - доп. переменные - $2 * \text{sizeof}(\text{int})$
 - адрес возврата

- алгоритм нахождения расстояния Дамерау-Левенштейна с использованием кеша в виде матрицы (память на саму матрицу: $((n + 1) * (m + 1)) * \text{sizeof}(\text{int})$) (рекурсивный), где для каждого вызова:
 - для S1, S2 - $(n + m) * \text{sizeof}(\text{char})$
 - для n, m - $2 * \text{sizeof}(\text{int})$
 - доп. переменные - $2 * \text{sizeof}(\text{int})$
 - ссылка на матрицу - 8 байт
 - адрес возврата

2.6 Вывод

В данном разделе были представлено описание используемых типов данных, а также схемы алгоритмов, рассматриваемых в лабораторной работе.

3 Технологическая часть

В данном разделе будут рассмотрены средства реализации, а также представлены листинги алгоритмов определения расстояния Левенштейна и Дамерау-Левенштейна.

3.1 Средства реализации

В данной работе для реализации был выбран язык программирования *Python*[3]. В текущей лабораторной работе требуется замерить процессорное время для выполняемой программы, а также построить графики. Все эти инструменты присутствуют в выбранном языке программирования.

Время работы было замерено с помощью функции *process_time(...)* из библиотеки *time*[4].

3.2 Листинги кода

В листингах 3.1-3.4 представлены реализации алгоритмов нахождения расстояния Левенштейна и Дамерау-Левенштейна.

Листинг 3.1 – Алгоритм нахождения расстояния Левенштейна (матричный)

```
1 def levenshteinDistance(str1: str, str2: str, output: bool = True)
  -> int:
2     n = len(str1)
3     m = len(str2)
4     matrix = createLevenshteinMatrix(n + 1, m + 1)
5
6     for i in range(1, n + 1):
7         for j in range(1, m + 1):
8             add = matrix[i - 1][j] + 1
9             delete = matrix[i][j - 1] + 1
10            change = matrix[i - 1][j - 1]
11
12            if (str1[i - 1] != str2[j - 1]):
13                change += 1
14
15            matrix[i][j] = min(add, delete, change)
16
17    if (output):
18        printMatrix(matrix, str1, str2)
19
20    return matrix[n][m]
```

Листинг 3.2 – Алгоритм нахождения расстояния Дамерау-Левенштейна (матричный)

```
1 def damerauLevenshteinDistance(str1: str, str2: str, output: bool =
  True) -> int:
2     n = len(str1)
3     m = len(str2)
4     matrix = createLevenshteinMatrix(n + 1, m + 1)
5
6     for i in range(1, n + 1):
7         for j in range(1, m + 1):
8             add = matrix[i - 1][j] + 1
9             delete = matrix[i][j - 1] + 1
10            change = matrix[i - 1][j - 1]
11            if (str1[i - 1] != str2[j - 1]):
12                change += 1
13            swap = n
14            if (i > 1 and j > 1 and
```

```

15         str1[i - 1] == str2[i - 2] and
16         str1[i - 2] == str2[i - 1]):
17             swap = matrix[i - 2][j - 2] + 1
18
19         matrix[i][j] = min(add, delete, change, swap)
20
21     if (output):
22         printMatrix(matrix, str1, str2)
23
24     return matrix[n][m]

```

Листинг 3.3 – Алгоритм нахождения расстояния Дамерау-Левенштейна (рекурсивный)

```

1 def damerauLevenshteinDistanceRecursive(str1: str, str2: str,
2     output: bool = True) -> int:
3     n = len(str1)
4     m = len(str2)
5     flag = 0
6     if ((n == 0) or (m == 0)):
7         return abs(n - m)
8
9     if (str1[-1] != str2[-1]):
10         flag = 1
11
12     minDistance = min(
13         damerauLevenshteinDistanceRecursive(str1[:-1], str2) + 1,
14         damerauLevenshteinDistanceRecursive(str1, str2[:-1]) + 1,
15         damerauLevenshteinDistanceRecursive(str1[:-1], str2[:-1]) +
16         flag
17     )
18     if (n > 1 and m > 1 and str1[-1] == str2[-2] and str1[-2] ==
19         str2[-1]):
20         minDistance = min(
21             minDistance,
22             damerauLevenshteinDistanceRecursive(str1[:-2],
23                 str2[:-2]) + 1
24         )
25
26     return minDistance

```

Листинг 3.4 – Алгоритм нахождения расстояния Дамерау-Левенштейна
(рекурсивный с использованием кэша)

```
1 def recursiveWithCache(  
2     str1: str, str2: str, n: int, m: int, matrix: Matrix) -> None:  
3     if (matrix[n][m] != -1):  
4         return matrix[n][m]  
5     if (n == 0):  
6         matrix[n][m] = m  
7         return matrix[n][m]  
8     if ((n > 0) and (m == 0)):  
9         matrix[n][m] = n  
10        return matrix[n][m]  
11  
12    add = recursiveWithCache(str1, str2, n - 1, m, matrix) + 1  
13    delete = recursiveWithCache(str1, str2, n, m - 1, matrix) + 1  
14    change = recursiveWithCache(str1, str2, n - 1, m - 1, matrix)  
15    if (str1[n - 1] != str2[m - 1]):  
16        change += 1 # flag  
17  
18    matrix[n][m] = min(add, delete, change)  
19    if (n > 1 and m > 1 and  
20        str1[n - 1] == str2[m - 2] and  
21        str1[n - 2] == str2[m - 1]):  
22  
23        matrix[n][m] = min(  
24            matrix[n][m],  
25            recursiveWithCache(str1, str2, n - 2, m - 2, matrix) + 1  
26        )  
27  
28    return matrix[n][m]  
29  
30 def damerauLevenshteinDistanceRecursiveCache(  
31     str1: str, str2: str, output: bool = True) -> int:  
32     n = len(str1)  
33     m = len(str2)  
34     matrix = createLevenshteinMatrix(n + 1, m + 1)  
35  
36     for i in range(n + 1):  
37         for j in range(m + 1):  
38             matrix[i][j] = -1  
39
```

```

40 recursiveWithCache(str1 , str2 , n , m, matrix)
41
42 if (output):
43     printMatrix(matrix , str1 , str2)
44
45 return matrix[n][m]

```

3.3 Функциональные тесты

В таблице 3.1 приведены тесты для функций, реализующих алгоритмы нахождения расстояния Левенштейна и Дамерау-Левенштейна. Тесты *для всех алгоритмов* пройдены успешно.

Таблица 3.1 – Функциональные тесты

№	Строка 1	Строка 2	Ожидаемый результат	
			Левенштейн	Дамерау-Л.
1	"пустая строка"	"пустая строка"	0	0
2	"пустая строка"	слово	5	5
3	проверка	"пустая строка"	8	8
4	ремонт	емонт	1	1
5	гигиена	иена	3	3
6	нисан	автоваз	6	6
7	спасибо	пожалуйста	9	9
8	что	кто	1	1
9	ты	тыква	3	3
10	есть	кушать	4	4
11	abba	baab	3	2
12	abcba	bacab	4	2

3.4 Вывод

Были представлены всех алгоритмов нахождения расстояния Левенштейна и Дамерау-Левенштейна, которые были описаны в предыдущем разделе. Также в данном разделе была приведена информации о выбранных средствах для разработки алгоритмов.

4 Исследовательская часть

В данном разделе будут приведены примеры работы программа, а также проведен сравнительный анализ алгоритмов при различных ситуациях на основе полученных данных.

4.1 Технические характеристики

Технические характеристики устройства, на котором выполнялось тестирование представлены далее:

- операционная система: Mac OS Monterey Версия 12.5.1 (21G83) [5] x86_64;
- память: 16 GB;
- процессор: 2,7 GHz 4-ядерный процессор Intel Core i7 [6].

При тестировании ноутбук был включен в сеть электропитания. Во время тестирования ноутбук был нагружен только встроенными приложениями окружения, а также системой тестирования.

4.2 Демонстрация работы программы

На рисунке 4.1 представлен результат работы программы.

Меню

1. Расстояние Левенштейна
2. Расстояние Дameraу-Левенштейна
3. Расстояние Дameraу-Левенштейна (рекурсивно)
4. Расстояние Дameraу-Левенштейна (рекурсивно с кешем)
5. Измерить время
0. Выход

Выбор: 1

Введите 1-ую строку: нирыба

Введите 2-ую строку: нимясо

Матрица, с помощью которой происходило вычисление расстояния Левенштейна:

	∅	н	и	м	я	с	о
∅	0	1	2	3	4	5	6
н	1	0	1	2	3	4	5
и	2	1	0	1	2	3	4
р	3	2	1	1	2	3	4
ы	4	3	2	2	2	3	4
б	5	4	3	3	3	3	4
а	6	5	4	4	4	4	4

Результат: 4

Меню

1. Расстояние Левенштейна
2. Расстояние Дameraу-Левенштейна
3. Расстояние Дameraу-Левенштейна (рекурсивно)
4. Расстояние Дameraу-Левенштейна (рекурсивно с кешем)
5. Измерить время
0. Выход

Выбор: █

Рисунок 4.1 – Пример работы программы

4.3 Время выполнения алгоритмов

Как было сказано выше, используется функция замера процессорного времени `process_time(...)` из библиотеки `time` на Python. Функция возвращает пользовательское процессорное время типа `float`.

Использовать функцию приходится дважды, затем из конечного времени нужно вычесть начальное, чтобы получить результат.

Замеры проводились для длины слова от 0 до 9 по 100 раз на различных входных данных.

Результаты замеров приведены в таблице 4.1 (время в мс).

Таблица 4.1 – Результаты замеров времени

Длина	Л.(матр)	Д-Л.(матр.)	Д-Л.(рек.)	Д.-Л.(рек. с кэшем)
0	0.0033	0.0074	0.0089	0.0032
1	0.0138	0.0130	0.0091	0.0216
2	0.0154	0.0169	0.0326	0.0363
3	0.0225	0.0227	0.1430	0.0521
4	0.0284	0.0331	0.6516	0.0751
5	0.0410	0.0472	3.1557	0.1328
6	0.0509	0.0633	16.7735	0.1755
7	0.0653	0.0854	89.8081	0.2375
8	0.0866	0.1064	496.2408	0.3080
9	0.1049	0.1354	2724.1102	0.3792

Также на рисунках 4.2, 4.3 приведены графические результаты замеров.

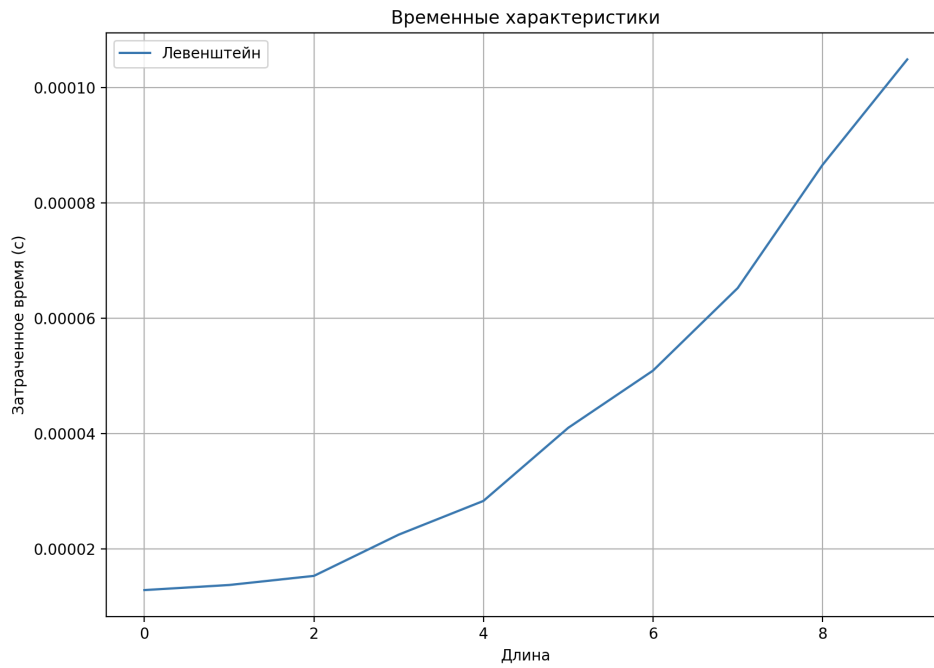


Рисунок 4.2 – Результат работы алгоритма нахождения расстояния Левенштейна (матричного)

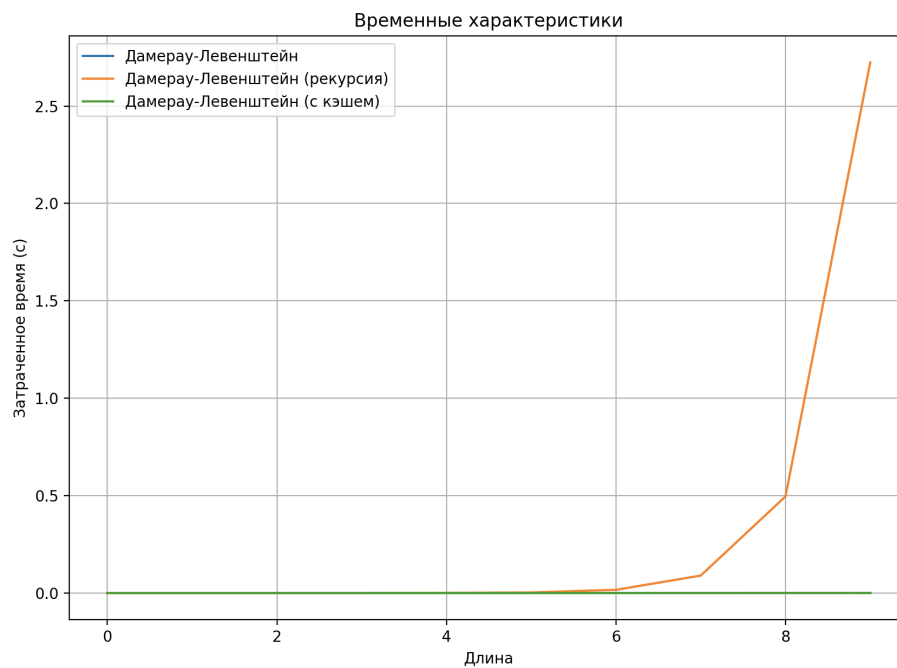


Рисунок 4.3 – Сравнение алгоритмов нахождения расстояния Дamerau-Левенштейна (матричного, рекурсивного и рекурсивного и использованием кэша)

4.4 Вывод

Исходя из замеров по памяти, итеративные алгоритмы проигрывают рекурсивным, потому что максимальный размер памяти в них растет, как произведение длин строк, а в рекурсивных - как сумма длин строк.

В результате эксперимента было получено, что обычно матричный алгоритм нахождения расстояния Дамерау-Левенштейна быстрее рекурсивного алгоритма с использованием кэша, однако занимает он намного больше памяти. Так, для длины слова в 9 символов матричная реализация быстрее рекурсивной в 2 раза, однако занимает в 5 раз больше памяти.

Также при проведении эксперимента было выявлено, что на длине строк в 4 символа рекурсивная реализация алгоритма Дамерау-Левенштейна в уже в 21 раз медленнее матричной реализации алгоритма. При увеличении длины строк в геометрической прогрессии растет и время работы рекурсивной реализации. Следовательно, для строк длиной более 4 символов стоит использовать матричную реализацию.

Заключение

В результате исследования было определено, что время алгоритмов нахождения расстояний Левенштейна и Дамерау-Левенштейна растет в геометрической прогрессии при увеличении длин строк. Лучшие показатели по времени дает матричная реализация алгоритма нахождения расстояния Дамерау-Левенштейна и его рекурсивная реализация с кешем, использование которых приводит к 21-кратному превосходству по времени работы уже на длине строки в 4 символа за счет сохранения необходимых промежуточных вычислений. При этом матричные реализации занимают довольно много памяти при большой длине строк.

Цель, которая была поставлена в начале лабораторной работы была достигнута, а также в ходе выполнения лабораторной работы были решены следующие задачи:

- были изучены и реализованы алгоритмы нахождения расстояния Левенштейна и Дамерау-Левенштейна;
- были также изучены матричная реализация, а также реализация с использованием кэша в виде матрицы для алгоритма нахождения расстояния Дамерау-Левенштейна;
- проведен сравнительный анализ алгоритмов нахождения расстояний Дамерау-Левенштейна в матричной, рекурсивной и рекурсивной и использованием кэша реализациях;
- подготовлен отчет о лабораторной работе.

Список литературы

- [1] Вычисление редакционного расстояния [Электронный ресурс]. Режим доступа: <https://habr.com/ru/post/117063/> (дата обращения: 17.09.2022).
- [2] Нечёткий поиск в тексте и словаре [Электронный ресурс]. Режим доступа: <https://habr.com/ru/post/114997/> (дата обращения: 17.09.2022).
- [3] Welcome to Python [Электронный ресурс]. Режим доступа: <https://www.python.org> (дата обращения: 17.09.2022).
- [4] time — Time access and conversions [Электронный ресурс]. Режим доступа: <https://docs.python.org/3/library/time.html#functions> (дата обращения: 17.09.2022).
- [5] macOS Monterey [Электронный ресурс]. Режим доступа: <https://www.apple.com/macos/monterey/> (дата обращения: 17.09.2022).
- [6] Процессор Intel® Core™ i7 [Электронный ресурс]. Режим доступа: <https://www.intel.com/processors/core/i7/docs> (дата обращения: 17.09.2022).