



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 дисциплине "Анализ Алгоритмов"

Тема

Студент Козырных А.Д.

Группа ИУ7-52Б

Преподаватель Волкова Л. Л., Строганов Д.В.

2024 г.

СОДЕРЖАНИЕ

ВВЕДЕНИЕ	3
1 Аналитическая часть	4
1.1 Описание алгоритмов	4
1.1.1 Расстояние Левенштейна	4
1.1.2 Расстояние Дамерау-Левенштейна	5
2 Конструкторская часть	6
2.1 Схемы алгоритмов	6
3 Технологическая часть	10
3.1 Требования к программному обеспечению	10
3.2 Средства реализации	10
3.3 Реализация алгоритмов	10
4 Исследовательская часть	13
4.1 Технические характеристики	13
4.2 Описание используемых типов данных	13
4.3 Оценка памяти	13
4.4 Время выполнения алгоритмов	15
4.5 Вывод	18
ЗАКЛЮЧЕНИЕ	19
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	20

ВВЕДЕНИЕ

Расстояние Левенштейна и Дамерау-Левенштейна - это минимальное количество операций (вставка, удаление и замена символов), требуемых для преобразования одной строки в другую. Расстояние Левенштейна используется для исправления ошибок в словах, поиска дубликатов текстов, сравнения геномов и прочих полезных операций с символьными последовательностями.

Цель лабораторной работы — сравнение алгоритмов нахождения расстояния Левенштейна и Дамерау-Левенштейна. Для достижения поставленной цели необходимо выполнить следующие задачи:

- реализовать указанные алгоритмы поиска расстояний (один алгоритм с использованием рекурсии, два алгоритма с использованием динамического программирования);
- проанализировать рекурсивную и матричную реализации алгоритмов по затраченному процессорному времени и памяти на основе экспериментальных данных;
- описать и обосновать полученные результаты в отчете.

1 Аналитическая часть

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

1.1 Описание алгоритмов

Расстояние Левенштейна — это минимальное число односимвольных преобразований (удаления, вставки или замены), необходимых для превращения одной строки в другую [1].

1.1.1 Расстояние Левенштейна

Для двух строк S_1 и S_2 , представленных в виде списков символов, длиной M и N соответственно, расстояние Левенштейна рассчитывается по рекуррентной формуле 1.1:

$$D(S_1, S_2) = \begin{cases} \begin{cases} len(S_1), & \text{если } len(S_2) = 0 \\ len(S_2), & \text{если } len(S_1) = 0 \\ D(tail(S_1), tail(S_2)), & \text{если } head(S_1) = head(S_2) \\ 1 + \min \begin{cases} D(tail(S_1), S_2), \\ D(S_1, tail(S_2)), \\ D(tail(S_1), tail(S_2)), \end{cases} & \text{иначе} \end{cases} \end{cases} \quad (1.1)$$

где:

- $len(S)$ — длина списка S ;
- $head(S)$ — первый элемент списка S ;
- $tail(S)$ — список S без первого элемента;

1.1.2 Расстояние Дамерау-Левенштейна

В алгоритме поиска расстояния Дамерау-Левенштейна, помимо вставки, удаления, и замены присутствует операция перестановки символов. Расстояние Дамерау-Левенштейна может быть вычисленно по рекуррентной формуле 1.2:

$$D(S_1, S_2) = \begin{cases} \begin{cases} len(S_1), & \text{если } len(S_2) = 0 \\ len(S_2), & \text{если } len(S_1) = 0 \\ D(tail(S_1), tail(S_2)), & \text{если } head(S_1) = head(S_2) \end{cases} \\ 1 + \min \begin{cases} D(tail(S_1), S_2), \\ D(S_1, tail(S_2)), \\ D(tail(S_1), tail(S_2)), \end{cases} \\ 1 + \min \begin{cases} D(tail(tail(S_1)), tail(tail(S_2))), & (*) \\ D(tail(S_1), S_2), \\ D(S_1, tail(S_2)), \\ D(tail(S_1), tail(S_2)), \end{cases} \end{cases} \quad (1.2)$$

(*): если $head(S_1) = head(tail(S_2))$ и $head(S_2) = head(tail(S_1))$;

ВЫВОД

В данном разделе были рассмотрены алгоритмы нахождения расстояния Левенштейна и Дамерау-Левенштейна.

2 Конструкторская часть

В данном разделе будут приведены схемы алгоритмов нахождения расстояния Левенштейна и Дамерау-Левенштейна.

2.1 Схемы алгоритмов

Алгоритмы на вход получают строки S_1 и S_2 , а на выходе возвращают целое число, которое показывает расстояние.

На рисунках 2.1 — 2.3 представлены схемы алгоритмов нахождения расстояния Левенштейна и Дамерау-Левенштейна.

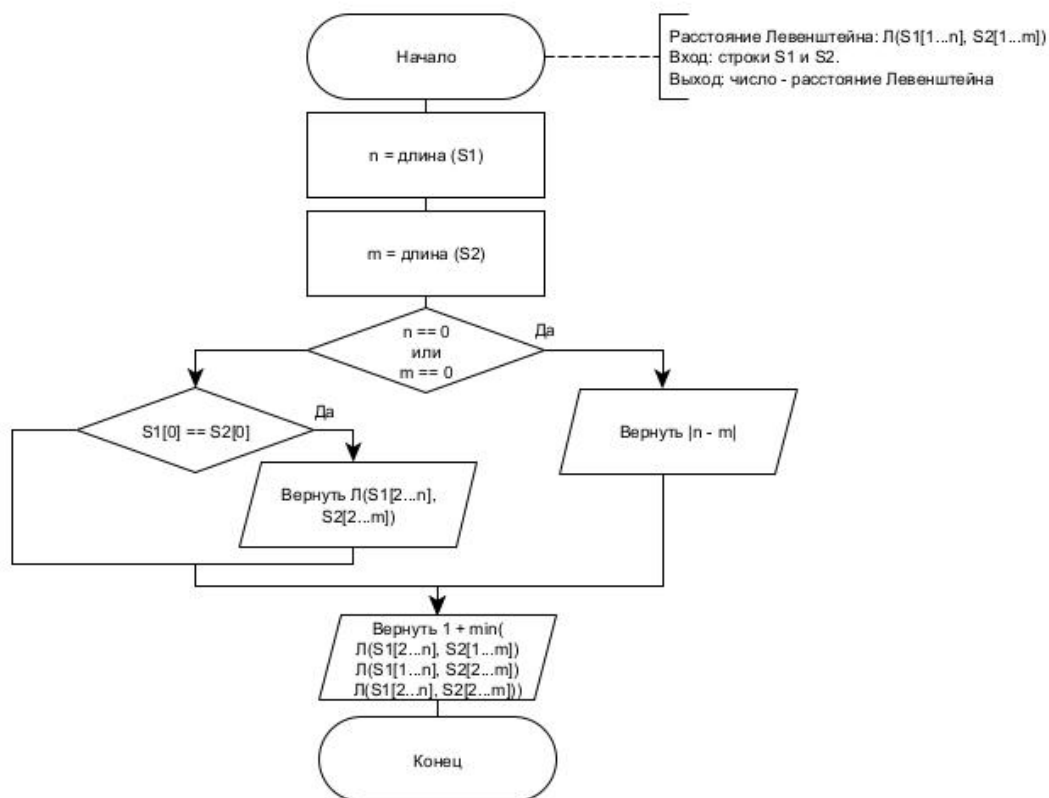


Рисунок 2.1 – Схема рекурсивного алгоритма расстояния Левенштейна

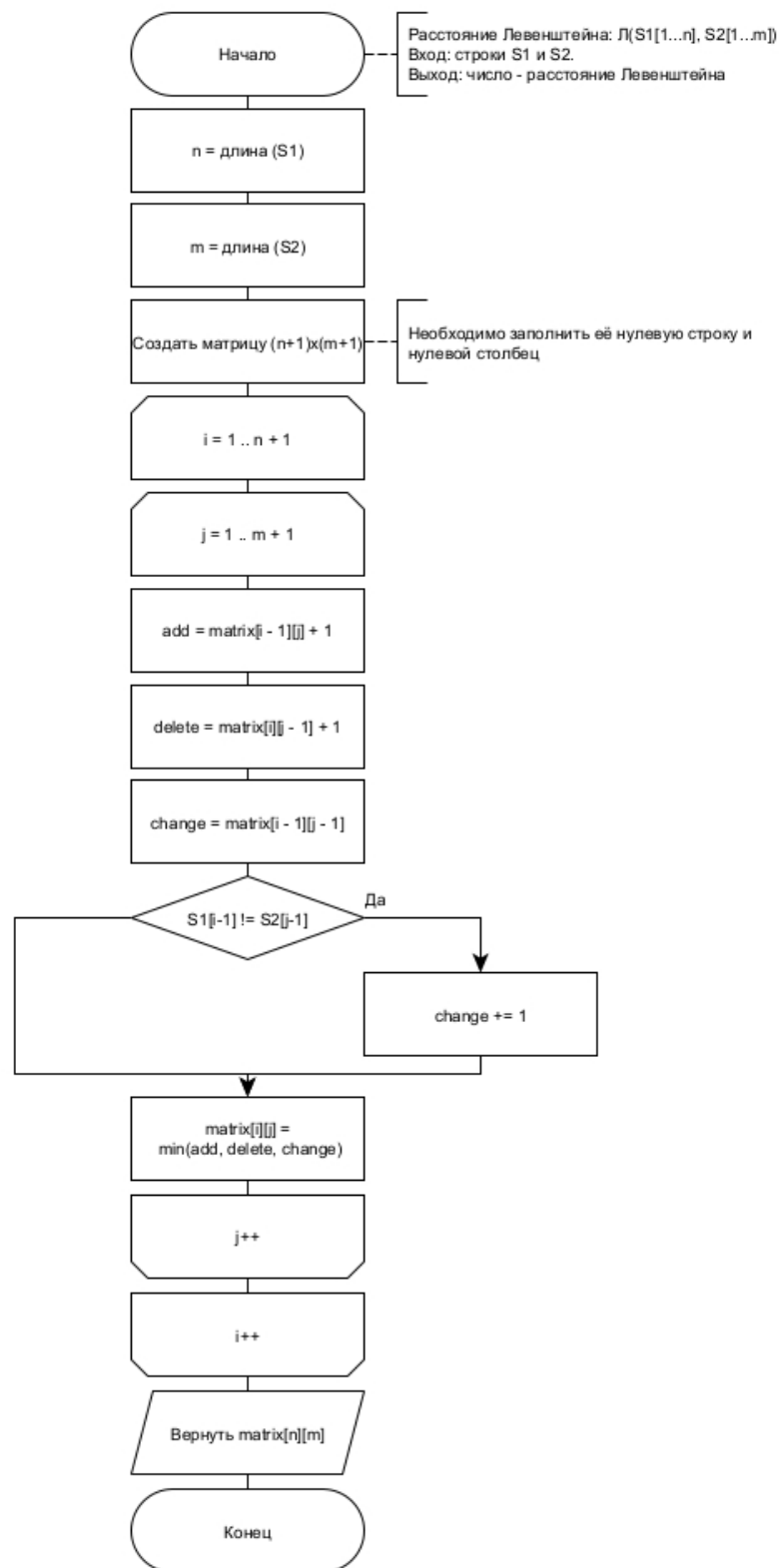


Рисунок 2.2 – Схема динамического алгоритма расстояния Левенштейна

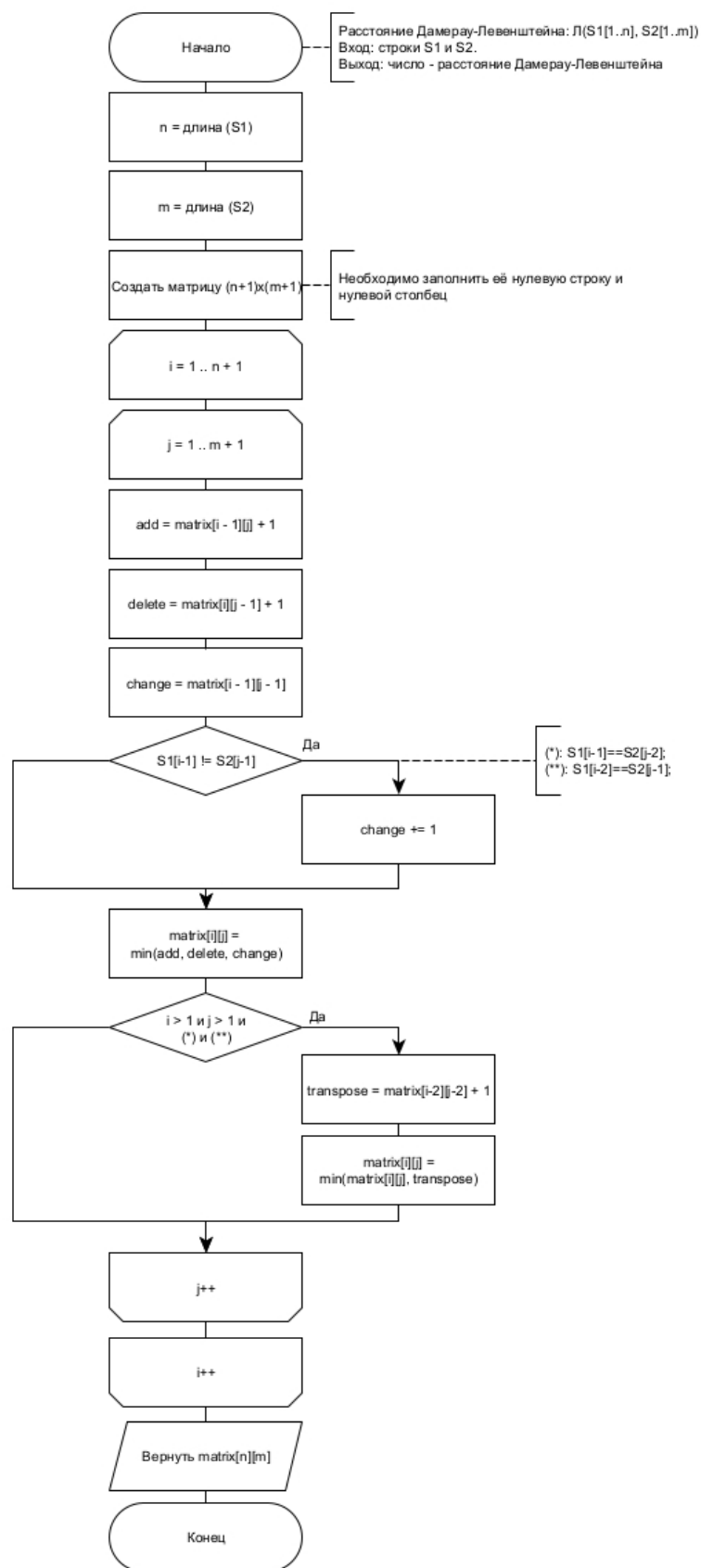


Рисунок 2.3 – Схема динамического алгоритма расстояния Дамерау-Левенштейна

ВЫВОД

В данном разделе были представлены схемы алгоритмов нахождения расстояния Левенштейна и Дамерау-Левенштейна.

3 Технологическая часть

В данном разделе будут приведены требования к ПО, реализация алгоритмов и средства реализации.

3.1 Требования к программному обеспечению

Входные данные: две строки из латинских или кириллических символов;

Выходные данные: целое число, которое является искомым расстоянием между двумя строками.

3.2 Средства реализации

Для реализации был выбран ЯП Python [2]. Выбор обусловлен наличием функции вычисления процессорного времени в библиотеке `time` [3]. Время было замерено с помощью функции `process_time()`.

3.3 Реализация алгоритмов

В листингах 3.1 — 3.3 представлены реализации алгоритмов нахождения расстояний Левенштейна и Дamerau-Левенштейна.

Листинг 3.1 – Реализация рекурсивного алгоритма нахождения расстояния Левенштейна

```
1 def tail(string : str):  
2     return string[1:]  
3  
4 def head(string : str):  
5     return string[0]  
6  
7  
8 def lvsh_recursive(string1 : str, string2 : str):  
9     if len(string1) == 0:  
10         return len(string2)
```

```

11     if len(string2) == 0:
12         return len(string1)
13
14     if head(string1) == head(string2):
15         return lvsh_recursive(tail(string1), tail(string2))
16     else:
17         a = lvsh_recursive(string1, tail(string2)) + 1
18         b = lvsh_recursive(tail(string1), string2) + 1
19         c = lvsh_recursive(tail(string1), tail(string2)) + 1
20
21     return min(a, b, c)

```

Листинг 3.2 – Реализация динамического алгоритма нахождения расстояния Левенштейна

```

1 def lvsh_matrix_min(matrix, i, j):
2     return min(matrix[i - 1][j], matrix[i][j - 1], matrix[i - 1][j - 1])
3
4 def lvsh_dynamic(string1 : str, string2 : str):
5     lenstr1 = len(string1)
6     lenstr2 = len(string2)
7     if lenstr1 == 0:
8         return lenstr2
9     if lenstr2 == 0:
10        return lenstr1
11
12    matrix = [[0 for _ in range(lenstr2 + 1)] for _ in range(lenstr1 +
13        1)]
14    for i in range(1, lenstr1 + 1):
15        matrix[i][0] = i
16    for j in range(1, lenstr2 + 1):
17        matrix[0][j] = j
18
19    for i in range(lenstr1):
20        for j in range(lenstr2):
21            if string1[i] == string2[j]:
22                matrix[i + 1][j + 1] = lvsh_matrix_min(matrix, i + 1, j
23                    + 1)
24            else:
25                matrix[i + 1][j + 1] = lvsh_matrix_min(matrix, i + 1, j
26                    + 1) + 1

```

```

25     test = min(lenstr1 , lenstr2)
26     return matrix[test][test] + abs(lenstr1 - lenstr2)

```

Листинг 3.3 – Реализация динамического алгоритма нахождения расстояния Дамерау-Левенштейна

```

1 def dlsh_dynamic(string1: str , string2: str):
2     n = len(string1)
3     m = len(string2)
4
5
6     matrix = [[0 for _ in range(m + 1)] for _ in range(n + 1)]
7     for i in range(1, n + 1):
8         matrix[i][0] = i
9     for j in range(1, m + 1):
10        matrix[0][j] = j
11
12    for i in range(1, n + 1):
13        for j in range(1, m + 1):
14            A = matrix[i - 1][j] + 1
15            B = matrix[i][j - 1] + 1
16            C = matrix[i - 1][j - 1]
17
18            if string1[i - 1] != string2[j - 1]:
19                C += 1
20
21            matrix[i][j] = min(A, B, C)
22
23            if i > 1 and j > 1 and string1[i - 1] == string2[j - 2] and
24                string1[i - 2] == string2[j - 1]:
25                transpose = matrix[i - 2][j - 2] + 1
26                matrix[i][j] = min(matrix[i][j], transpose)
27
28    return matrix[n][m]

```

ВЫВОД

В данном разделе рассмотрены средства реализации, требования к ПО и реализации алгоритмов.

4 Исследовательская часть

4.1 Технические характеристики

Характеристики используемого оборудования:

- Операционная система — Linux [4]
- Оперативная память — 16 Гб.
- Процессор — AMD Ryzen 7 5800H with Radeon Graphics (16) @ 4.463GHz [5]

4.2 Описание используемых типов данных

Используемые типы данных:

- Строка - последовательность символов типа *str*
- Длина строки - целое число типа *int*
- Матрица - двумерный массив типа *int*

4.3 Оценка памяти

Рекурсивная версия алгоритма Левенштейна не использует явных структур данных для хранения промежуточных вычислений. Вместо этого каждый вызов функции обрабатывает небольшой фрагмент строк, и функция вызывает саму себя несколько раз. Глубина рекурсии в худшем случае составляет:

$$(\text{len}(S_1) + \text{len}(S_2)). \quad (4.1)$$

При этом каждый рекурсивный вызов требует хранения локальных переменных. 3 переменные типа *int* (a, b, c), 2 переменные типа *str*. В результате, максимальный расход памяти составляет:

$$(\text{len}(S_1) + \text{len}(S_2)) \cdot (3 \cdot \text{size}(\text{int}) + 2 \cdot \text{size}(\text{str})), \quad (4.2)$$

где $size$ - функция, вычисляющая размер параметра.

Алгоритм, основанный на динамическом программировании, использует двумерную матрицу размером $len(S_1+1) \times len(S_2+1)$. Эта матрица хранит результаты всех промежуточных вычислений (расстояние Левенштейна для всех подстрок). Дополнительно хранятся локальные переменные: 3 переменных типа int , 2 переменных типа str . По итогу расход памяти в данном случае составляет:

$$(len(S_1 + 1) \cdot len(S_2 + 1) \cdot size(int)) + 3 \cdot size(int) + 2 \cdot size(str)). \quad (4.3)$$

По расходу памяти алгоритм, использующий принцип динамического программирования, проигрывает рекурсивному: максимальный размер используемой памяти в них растёт как произведение длин строк, в то время как у рекурсивного алгоритма — как сумма длин строк.

Алгоритм Дамерау-Левенштейна, также реализованный через динамическое программирование, аналогичен по своей структуре алгоритму Левенштейна. Основное отличие заключается в дополнительной проверке для учёта операций перестановки соседних символов. Для этого используется та же матрица размера $len(S_1 + 1) \times (len(S_2 + 1))$, что и в динамическом алгоритме Левенштейна.

Несмотря на добавление новой операции (перестановка), использование памяти остаётся также на уровне произведение длин строк, поскольку не требуется дополнительное пространство для хранения результатов перестановок. Как и в случае с Левенштейном, каждая клетка матрицы заполняется лишь однажды.

4.4 Время выполнения алгоритмов

Замер времени выполнения каждого из алгоритмов находится в таблице 4.1. Замеры для каждого из алгоритмов для одного и того же размера проводились 500 раз , и результаты замеров усреднялись.

Таблица 4.1 – Время работы алгоритмов (в секундах)

Длина строк	Лев рек.	Лев дин.	Дам-Лев дин.
1	1.23e-07	1.24e-07	5.23e-07
2	7.67e-07	1.14e-06	1.01e-06
3	4.13e-06	2.20e-06	2.57e-06
4	1.67e-05	3.97e-06	5.09e-06
5	9.79e-05	6.27e-06	6.37e-06
6	4.29e-04	9.34e-06	9.26e-06
7	2.12e-03	1.17e-05	1.27e-05
8	1.04e-02	1.54e-05	1.71e-05
9	5.62e-02	2.00e-05	2.21e-05

На рисунке 4.1 показаны графики зависимости времени от количества символов для динамических алгоритмов нахождения расстояния Левенштейна и Дameraу-Левенштейна.

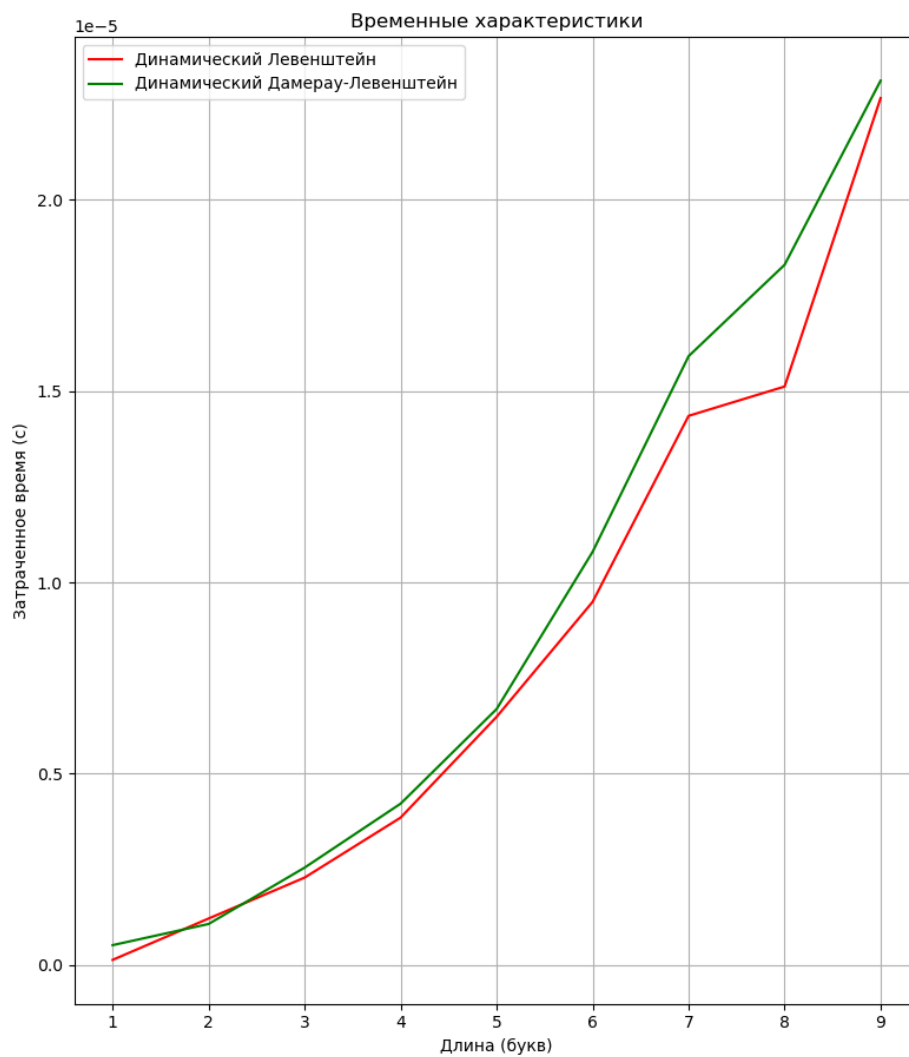


Рисунок 4.1 – Временные показатели динамических алгоритмов Левенштейна и Дameraу-левенштейна

На рисунке 4.2 показаны графики зависимости времени от количества символов для всех реализаций алгоритмов нахождения расстояния Левенштейна и Дameraу-Левенштейна.

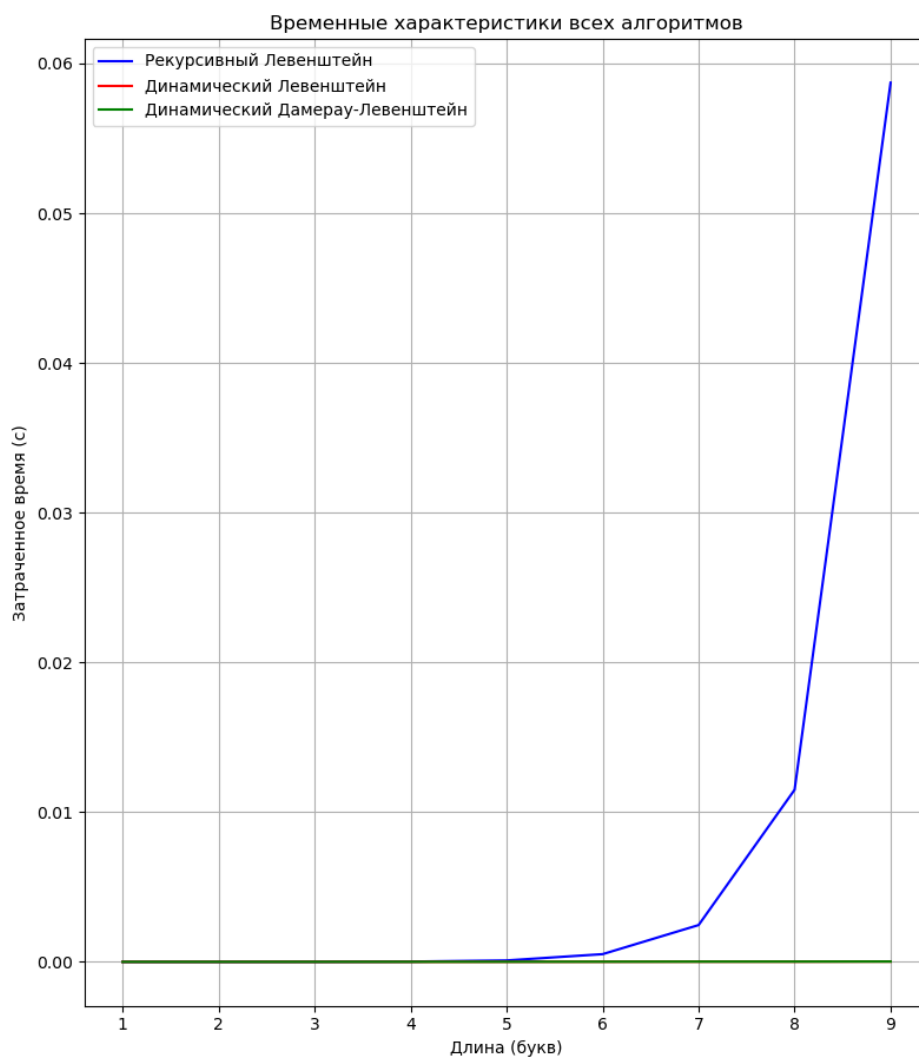


Рисунок 4.2 – Общие временные показатели

Наиболее эффективные реализации алгоритмы нахождения расстояния Левенштейна и Дameraу-левенштейна — это те алгоритмы, где используется динамический подход, так как в рекурсивном подходе идет повторный расчет.

4.5 Вывод

Рекурсивный алгоритм, вычисляющий расстояние Левенштейна, работает по времени на несколько порядков дольше, чем динамический вариант. Также стоит заметить, что динамические алгоритмы вычисления расстояний Левенштейна и Дамерау-Левенштейна сопоставимы между собой по времени выполнения и примерно равны.

Анализ расхода памяти показывает, что рекурсивный алгоритм имеет меньшие требования к памяти по сравнению с алгоритмом, использующим динамическое программирование. В случае динамического варианта алгоритма, считающего расстояние Дамерау-Левенштейна, несмотря на добавление операции перестановки, потребление памяти остается на уровне алгоритма, считающего расстояние Левенштейна, так как не требуется дополнительное пространство для хранения результатов перестановок.

ЗАКЛЮЧЕНИЕ

Было экспериментально подтверждено различие во временной эффективности рекурсивной и нерекурсивной реализаций выбранных алгоритмов нахождения расстояния между строками при помощи разработанного программного обеспечения на материале замеров процессорного времени выполнения реализаций на различных длинах строк.

В результате исследований можно сделать вывод о том, что матричная реализация данных алгоритмов, по сравнению с рекурсивной, заметно выигрывает по времени при росте длин строк, но проигрывает по количеству затрачиваемой памяти.

В ходе выполнения данной лабораторной работы были решены следующие задачи:

- реализованы указанные алгоритмы для нахождения расстояния Левенштейна (в рекурсивной и динамической вариации), Дамерау-Левенштейна (в динамической);
- проанализированы рекурсивная и динамическая реализации алгоритма Левенштейна, динамические реализации алгоритмов Левенштейна и Дамерау-Левенштейна по затрачиваемым ресурсам (времени и памяти);
- описаны и обоснованы полученные результаты в отчете.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- [1] Расстояние Левенштейна для чайников [Электронный ресурс]. URL: <https://habr.com/ru/articles/676858/> (дата обращения: 03.09.2024)
- [2] Welcome to Python [Электронный ресурс]. URL: <https://www.python.org> (дата обращения: 07.09.2024).
- [3] time — Time access and conversions [Электронный ресурс]. URL: <https://docs.python.org/3/library/time.html#functions> (дата обращения: 07.09.2024).
- [4] Arch Linux [Электронный ресурс]. URL: <https://archlinux.org/> (дата обращения: 07.09.2024).
- [5] AMD Ryzen 7 5800H Processor. Бенчмарки [Электронный ресурс] URL: <https://www.notebookcheck-ru.com/AMD-Ryzen-7-5800H.519526.0.html> (дата обращения: 07.09.2024).