



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»
на тему:
«Редакционные расстояния»

Студент Рунов К. А.

Группа ИУ7-54Б

Преподаватели Волкова Л. Л., Строганов Д. В.

2023 г.

СОДЕРЖАНИЕ

ВВЕДЕНИЕ	3
1 Аналитическая часть	5
1.1 Расстояние Левенштейна	5
1.2 Расстояние Дамерау — Левенштейна	6
2 Конструкторская часть	11
2.1 Требования к программному обеспечению	11
2.2 Разработка алгоритмов	13
2.3 Описание типов и структур данных	17
3 Технологическая часть	18
3.1 Средства реализации	18
3.2 Сведения о модулях программы	18
3.3 Реализация алгоритмов	19
3.4 Функциональные тесты	19
4 Исследовательская часть	21
4.1 Технические характеристики	21
4.2 Демонстрация работы программы	21
4.3 Временные характеристики	25
4.4 Теоретические оценки затрачиваемой памяти	30
4.5 Характеристики по памяти	32
ЗАКЛЮЧЕНИЕ	38
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	40
ПРИЛОЖЕНИЕ А	41
ПРИЛОЖЕНИЕ Б	49

ВВЕДЕНИЕ

Редакционные расстояния — расстояние Левенштейна и его модификация — расстояние Дамерау — Левенштейна — метрики сходства между двумя символьными последовательностями. Расстоянием в этих метриках считается минимальное число односимвольных преобразований (удаления, вставки, замены или транспозиции), необходимых для преобразования одной последовательности символов в другую.

Редакционные расстояния применяются

- для исправления ошибок в слове поискового запроса;
- в формах заполнения информации на сайтах;
- для распознавания рукописных символов;
- в базах данных. [1]

Целью данной лабораторной работы является изучение, реализация и исследование алгоритмов поиска расстояний Левенштейна и Дамерау — Левенштейна.

Для достижения поставленной цели нужно решить следующие задачи:

- 1) описать алгоритмы поиска расстояний Левенштейна и Дамерау — Левенштейна;
- 2) обосновать выбор средств реализации алгоритмов;
- 3) реализовать алгоритмы:
 - итеративный алгоритм нахождения расстояния Левенштейна,
 - итеративный алгоритм нахождения расстояния Дамерау — Левенштейна,
 - рекурсивный алгоритм нахождения расстояния Дамерау — Левенштейна,
 - рекурсивный с кэшированием алгоритм нахождения расстояния Дамерау — Левенштейна;
- 4) провести сравнительный анализ алгоритмов по критериям:

- используемое процессорное время,
 - максимальная затрачиваемая память;
- 5) описать и проанализировать полученные результаты в отчёте.

1 Аналитическая часть

1.1 Расстояние Левенштейна

Расстояние Левенштейна — метрика, определяющая понятие расстояния между двумя последовательностями символов, как минимального количества редакторских операций вставки (I , от англ. insert), замены (R , от англ. replace) и удаления (D , от англ. delete), необходимых для преобразования одной строки в другую [2]. Для каждой операции должна быть определена её стоимость. Введём обозначения для стоимостей. Пусть:

- 1) $w(a, b)$ — цена замены символа a на b ;
- 2) $w(\lambda, b)$ — цена вставки символа b ;
- 3) $w(a, \lambda)$ — цена удаления символа a .

Определим стоимости операций:

$$w(a, b) = \begin{cases} 1, & \text{если } a \neq b; \\ 0, & \text{иначе.} \end{cases} \quad (1.1)$$

Отсутствие операций в случае совпадения символов будем обозначать за M (от англ. match).

Введём в рассмотрение функцию $D(S_1[1..i], S_2[1..j])$, значением которой является редакционное расстояние между подстроками $S_1[1..i]$ и $S_2[1..j]$, где $S_1[1..i]$ — подстрока S_1 длины i . Так, если $S_1 = \text{"скат"}$, то $S_1[1..0] = \lambda$, $S_1[1..1] = \text{"с"}$, $S_1[1..2] = \text{"ск"}$. Расстояние Левенштейна между строками S_1 и S_2 длин L_1 и L_2 соответственно вычисляется по рекуррентной формуле:

$$\begin{aligned} D(S_1[1..i], S_2[1..j]) &= \\ &= \begin{cases} \max(i, j), & i \cdot j = 0; \\ \min \begin{cases} D(S_1[1..i], S_2[1..j-1]) + 1, \\ D(S_1[1..i-1], S_2[1..j]) + 1, \\ D(S_1[1..i-1], S_2[1..j-1]) + w(S_1[i], S_2[j]), \end{cases} & i \cdot j \neq 0, \end{cases} \end{cases} \quad (1.2) \end{aligned}$$

где $i = L_1$, $j = L_2$.

1.1.1 Итерационный алгоритм нахождения расстояния Левенштейна

Рекурсивная реализация алгоритма поиска расстояния Левенштейна малоэффективна по времени при больших L_1 и L_2 , так как производится много повторных, лишних вычислений. Реализацию можно оптимизировать с помощью динамического программирования. Например, ввести матрицу размерности $(L_1 + 1) \times (L_2 + 1)$ и заполнять её промежуточными значениями $D(S_1[1..i], S_2[1..j])$, используя их затем по ходу вычислений. Значения в ячейках $[i][j]$ (i -я строка, j -й столбец) матрицы равны значениям $D(S_1[1..i], S_2[1..j])$ соответственно. Можно заметить, что всю матрицу для вычислений хранить не обязательно — двух строк будет достаточно.

1.2 Расстояние Дамерау — Левенштейна

Расстояние Дамерау — Левенштейна — метрика, которая определяет расстояние между двумя последовательностями символов, как и расстояние Левенштейна, но к исходному набору редакторских операций добавляется ещё одна — транспозиция (Т, от англ. transposition). Операция транспозиции меняет местами соседние буквы в строке. Обозначим её стоимость: $w(ab, ba) = 1$.

Расстояние Дамерау — Левенштейна $\mathcal{D}(S_1, S_2)$ между строками S_1 и S_2 длин L_1 и L_2 соответственно может быть вычислено по рекуррентной фор-

муле:

$$\begin{aligned}
& \mathcal{D}(S_1[1..i], S_2[1..j]) = \\
& = \begin{cases} \max(i, j), & i \cdot j = 0; \\ \min \begin{cases} \mathcal{D}(S_1[1..i], S_2[1..j-1]) + 1, \\ \mathcal{D}(S_1[1..i-1], S_2[1..j]) + 1, \\ \mathcal{D}(S_1[1..i-1], S_2[1..j-1]) + w(S_1[i], S_2[j]), \\ \begin{cases} \mathcal{D}(S_1[1..i-2], S_2[1..j-2]) + 1, & \text{если } i > 1, j > 1, \\ S_1[i] = S_2[j-1], \\ S_1[i-1] = S_2[j]; \\ \infty, & \text{иначе,} \end{cases} \end{cases} & i \cdot j \neq 0, \end{cases}
\end{aligned} \tag{1.3}$$

где $i = L_1, j = L_2$.

1.2.1 Рекурсивный алгоритм нахождения расстояния Дамерау — Левенштейна

Рекурсивный алгоритм нахождения расстояния Дамерау — Левенштейна реализует рекуррентную формулу (1.3). Таким образом, верно следующее:

- 1) $\mathcal{D}(\lambda, \lambda) = 0$, — для преобразования пустой строки в пустую строку требуется 0 операций вставки, замены, удаления и транспозиции;
- 2) $\mathcal{D}(S_1, \lambda) = |S_1|$ (длина S_1), — для преобразования строки S_1 в пустую строку требуется $|S_1|$ операций (удаления);
- 3) $\mathcal{D}(\lambda, S_2) = |S_2|$, — для преобразования пустой строки в строку S_2 требуется $|S_2|$ операций (вставки);
- 4) $\mathcal{D}(c_1, c_2) = \begin{cases} 1, c_1 \neq c_2; \\ 0, c_1 = c_2, \end{cases}$ — для преобразования одного символа в другой требуется 1 операция (замены), если символы отличаются, и 0 операций, если символы совпадают;

5) Для преобразования одной пары символов в другую:

$$\mathcal{D}(c_1c_2, c_3c_4) = \begin{cases} 0, c_1 = c_3, c_2 = c_4; // MM \\ 1, c_1 = c_3, c_2 \neq c_4; // MR \\ 1, c_1 \neq c_3, c_2 = c_4; // RM \\ 1, c_1 = c_4, c_2 = c_3; // T \\ 2, \text{ иначе}; // RR \end{cases}$$

6) Пусть $S_1 = S'_1c_2 = S''_1c_1c_2$, $S_2 = S'_2c_4 = S''_2c_3c_4$, где S'_1 и S'_2 — строки S_1 и S_2 без последних символов, S''_1 и S''_2 — строки S_1 и S_2 без двух последних символов, а c_1c_2 и c_3c_4 — пары их последних символов соответственно.

$$\mathcal{D}(S_1, S_2) = \min \begin{cases} \mathcal{D}(S'_1c_2, S'_2) + 1, \\ \mathcal{D}(S'_1, S'_2c_4) + 1, \\ \mathcal{D}(S'_1, S'_2) + w(c_2, c_4), \\ \mathcal{D}(S''_1, S''_2) + 1, \end{cases} \quad \text{если } c_2 = c_3, c_1 = c_4.$$

Можно заметить, что $\mathcal{D}(S_1, S_2)$ вычисляется как минимальная длина последовательности редакторских операций, которыми можно преобразовать строку S_1 в S'_2 плюс цена вставки последнего символа из S_2 , строку S'_1 в S_2 плюс цена удаления последнего символа из S_1 , строку S'_1 в S'_2 плюс цена замены последних символов строк S_1 и S_2 , строку S''_1 в S''_2 плюс цена транспозиции двух последних символов строк, если она возможна. Для любых подстрок \mathcal{S}_1 и \mathcal{S}_2 строк S_1 и S_2 , $\mathcal{D}(\mathcal{S}_1, \mathcal{S}_2)$ вычисляет минимальное количество редакторских операций. Следовательно, $\mathcal{D}(S_1, S_2)$ действительно считает расстояние Дамерау — Левенштейна для произвольных строк S_1 и S_2 .

1.2.2 Рекурсивный с кэшированием алгоритм нахождения расстояния Дамерау — Левенштейна

Рекурсивный алгоритм нахождения расстояния Дамерау — Левенштейна прост, но его реализация на ЭВМ без дополнительных оптимизаций неэффективна, так как по несколько раз считаются значения, которые уже были вычислены, а вычисление их может оказаться достаточно трудоёмким.

Идея возможной оптимизации состоит в следующем: хранить получаемые

по ходу выполнения алгоритма значения в матрице, а перед вычислением очередного — проверять, было ли оно посчитано ранее (заполнена ли соответствующая ячейка матрицы), и, если да, — брать его оттуда, не прибегая к повторным вычислениям. Для того, чтобы узнать, заполнена ячейка матрицы или нет, нужно изначально проинициализировать все ячейки матрицы какими-то значениями, которые на практике не могут быть получены в результате расчёта расстояния Дамерау — Левенштейна, например, -1 .

1.2.3 Итерационный алгоритм нахождения расстояния Дамерау — Левенштейна

Как алгоритм поиска расстояния Левенштейна, так и алгоритм поиска расстояния Дамерау — Левенштейна можно реализовать нерекурсивно с помощью динамического программирования, используя матрицу расстояний.

Процесс вычисления значения ячейки матрицы показан на рисунке 1.1.

	s_{21}	s_{22}	s_{23}
s_{11}	N_T	\dots	\dots
s_{12}	\vdots	N_R	N_D
s_{13}	\vdots	N_I	\mathbf{N}

$$\mathbf{N} = \min \begin{cases} N_I + 1, \\ N_D + 1, \\ N_R + w(s_{13}, s_{23}), \\ \begin{cases} N_T + 1, & \text{если } s_{13} = s_{22}, s_{12} = s_{23} \text{ и } s_{11}, s_{21} \text{ существуют,} \\ \infty, & \text{иначе.} \end{cases} \end{cases}$$

Рисунок 1.1 – Вычисление расстояния Дамерау — Левенштейна с использованием матрицы

Таким образом, для нахождения расстояния Дамерау — Левенштейна хранить всю матрицу не обязательно — трёх строк будет достаточно. Заполнив последнее значение очередной строки, можно выполнить «циклическую прокрутку» строк матрицы вверх, после чего изменить значение первой ячейки

последней теперь строки матрицы на значение первой ячейки предпоследней строки матрицы плюс 1, а затем продолжить заполнение матрицы по алгоритму, представленному на рисунке выше, перезаписывая значения ячеек.

Вывод

В данном разделе были рассмотрены алгоритмы нахождения расстояний Левенштейна и Дамерау — Левенштейна. Поскольку данные расстояния могут быть вычислены с помощью рекуррентных формул, то алгоритмы могут быть реализованы как рекурсивно, так и итерационно.

2 Конструкторская часть

В этой части будут приведены требования к программе и схемы алгоритмов нахождения расстояний Левенштейна и Дameraу — Левенштейна.

2.1 Требования к программному обеспечению

Программа должна поддерживать два режима работы:

- 1) ручной ввод;
- 2) массивный замер времени.

В режиме ручного ввода пользователь должен иметь возможность

- ввести два слова через пробел с клавиатуры;
- использовать символы в кодировке Unicode;
- ввести слова размером до 25 000 символов в кодировке Unicode.

В режиме ручного ввода программа должна

- запустить каждую из реализаций четырёх алгоритмов (итерационного алгоритма нахождения расстояния Левенштейна, итерационного алгоритма нахождения расстояния Дameraу — Левенштейна, рекурсивного алгоритма нахождения расстояния Дameraу — Левенштейна, рекурсивного с кешированием алгоритма нахождения расстояния Дameraу — Левенштейна);
- вывести заполненные матрицы расстояний для итерационных алгоритмов нахождения расстояний Левенштейна и Дameraу — Левенштейна;
- вывести последовательность редакторских операций, которая привела к нахождению расстояний Левенштейна и Дameraу — Левенштейна, для итерационных алгоритмов нахождения расстояний Левенштейна и Дameraу — Левенштейна;
- вывести результат работы каждого алгоритма с названием этого алгоритма.

В режиме массированного замера времени пользователь должен иметь возможность

- ввести начальную длину строк, конечную длину строк, шаг изменения длин строк, которые будут использоваться в замерах времени;
- ввести количество замеров времени, которое требуется произвести для каждой длины строк (в результате будет выдано среднее значение по каждой длине);
- посмотреть графики, построенные на основе данных из последнего замера времени.

В режиме массированного замера времени программа должна

- на основе данных о длинах строк, полученных от пользователя, сгенерировать строки указанных длин из фиксированного набора символов в кодировке Unicode; конкретный символ строки должен выбираться в результате генерации псевдослучайного целого числа по модулю количества символов в фиксированном наборе — полученное число есть индекс символа в наборе;
- каждую пару сгенерированных строк подать на вход всем алгоритмам, время работы реализаций которых требуется измерить;
- произвести замеры времени работы для каждой реализации алгоритмов без учёта времени на создание и инициализацию матрицы в тех алгоритмах, где матрица используется;
- записать данные о каждом замере времени в файл;
- на основе данных из файла построить четыре графика (по графику на каждый алгоритм), где значения по оси Ox — длины строк, значения по оси Oy — время работы реализации алгоритма в наносекундах, и вывести результат пользователю.

2.2 Разработка алгоритмов

На вход алгоритмам подаются строки $S1$, $S2$ и их длины $L1$, $L2$ соответственно.

На рисунке 2.1 представлена схема алгоритма нахождения расстояния Левенштейна. На рисунках 2.2 – 2.5 представлены схемы алгоритмов нахождения расстояния Дamerau — Левенштейна.

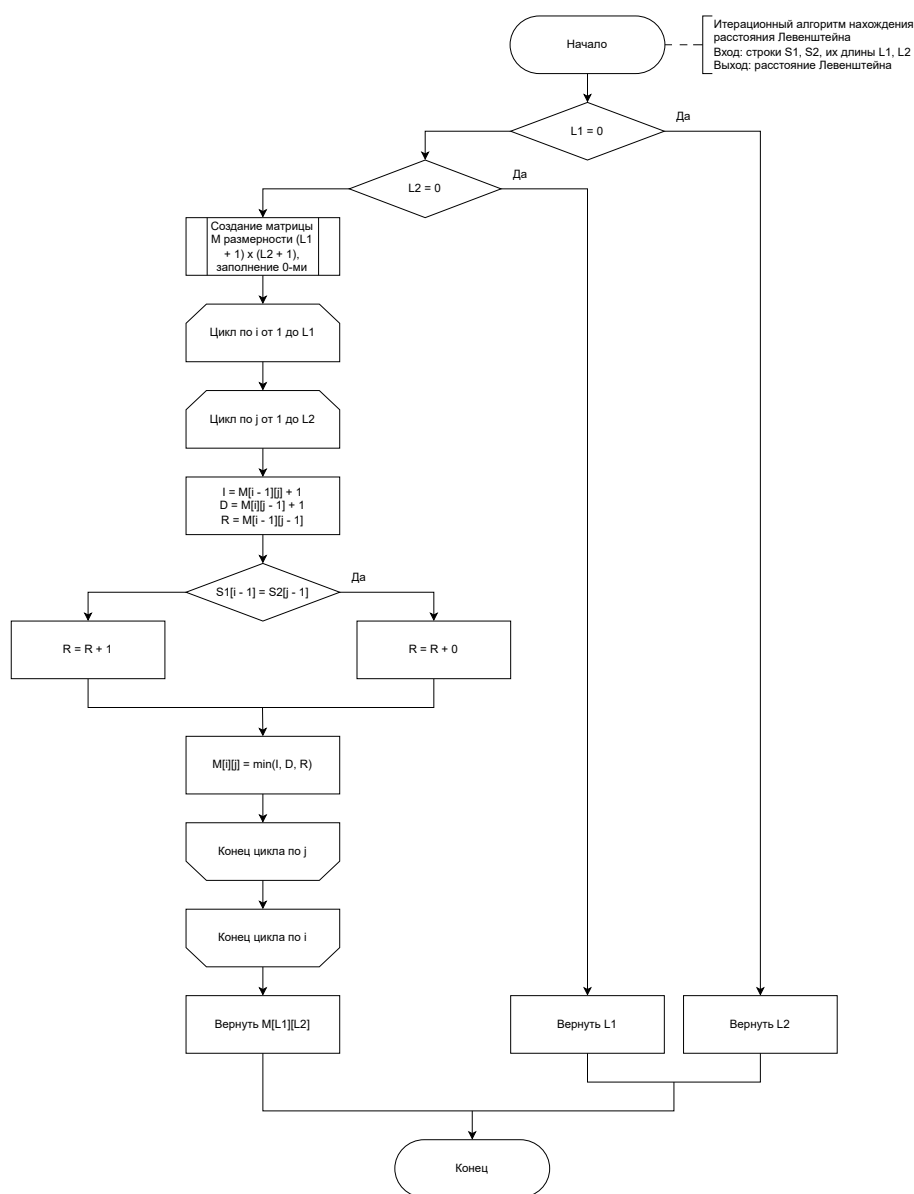


Рисунок 2.1 – Итерационный алгоритм нахождения расстояния Левенштейна

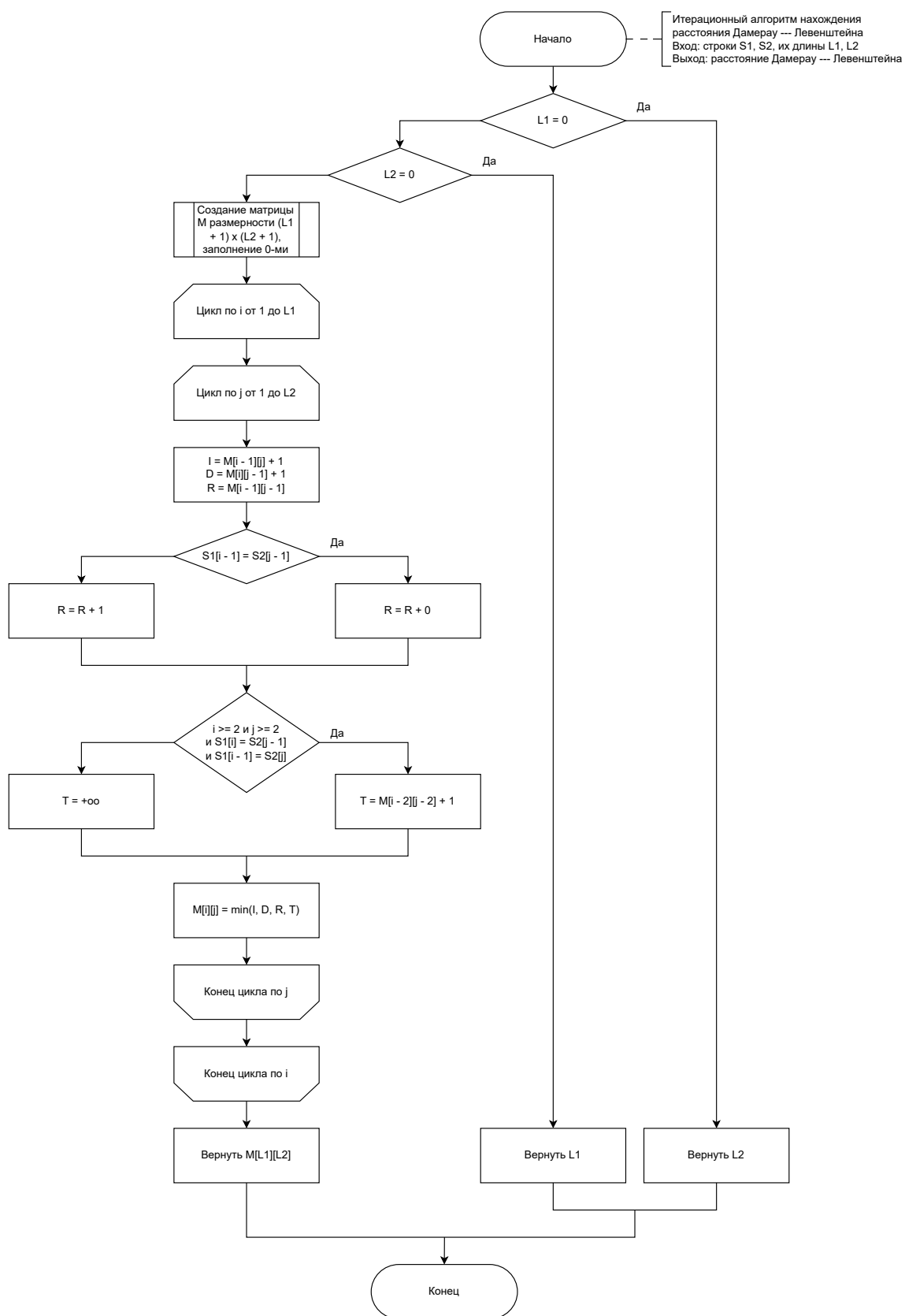


Рисунок 2.2 – Итерационный алгоритм нахождения расстояния Дамерау — Левенштейна

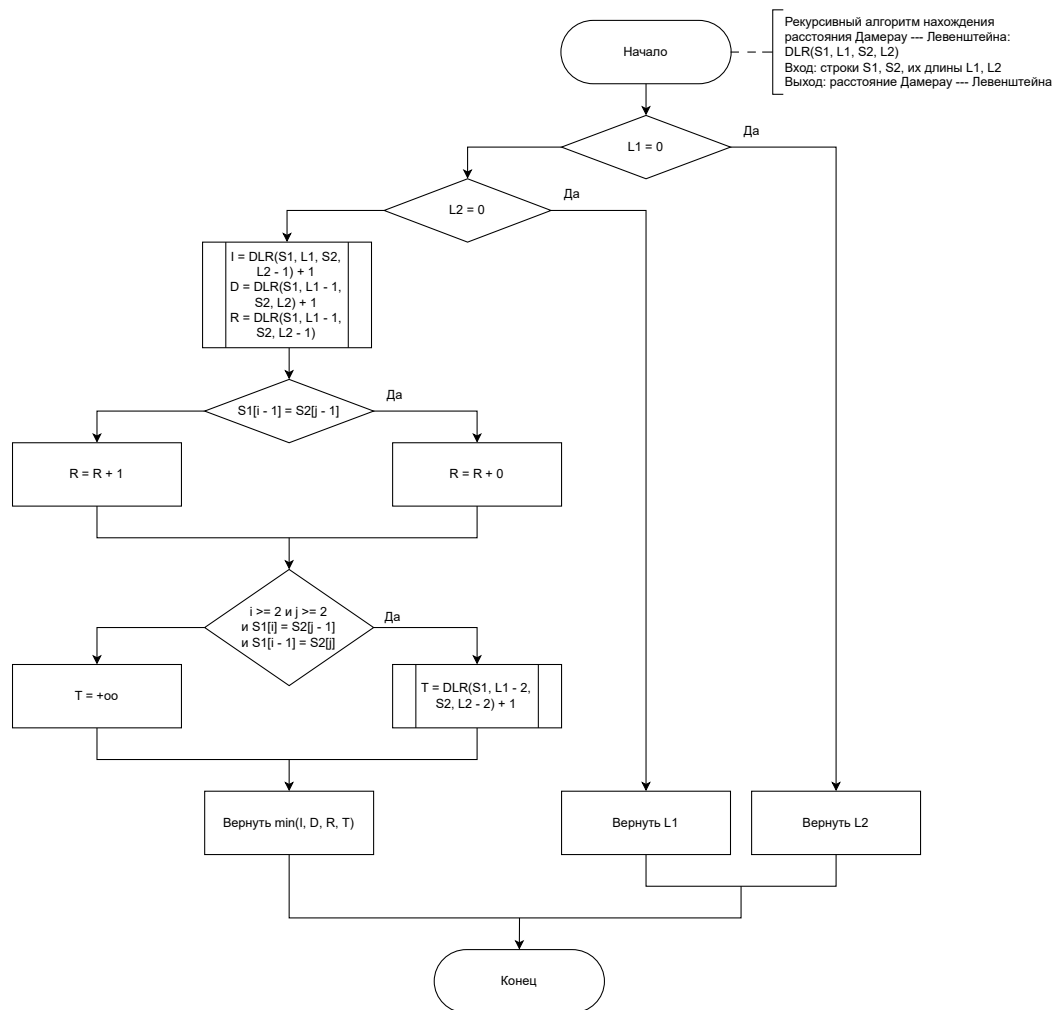


Рисунок 2.3 – Рекурсивный алгоритм нахождения расстояния Дамерау — Левенштейна

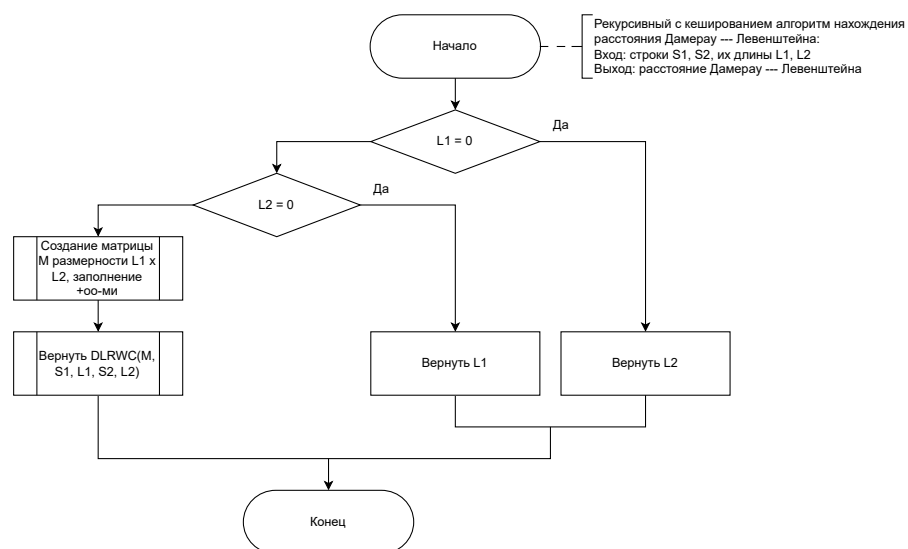


Рисунок 2.4 – Рекурсивный с кешированием алгоритм нахождения расстояния Дамерау — Левенштейна

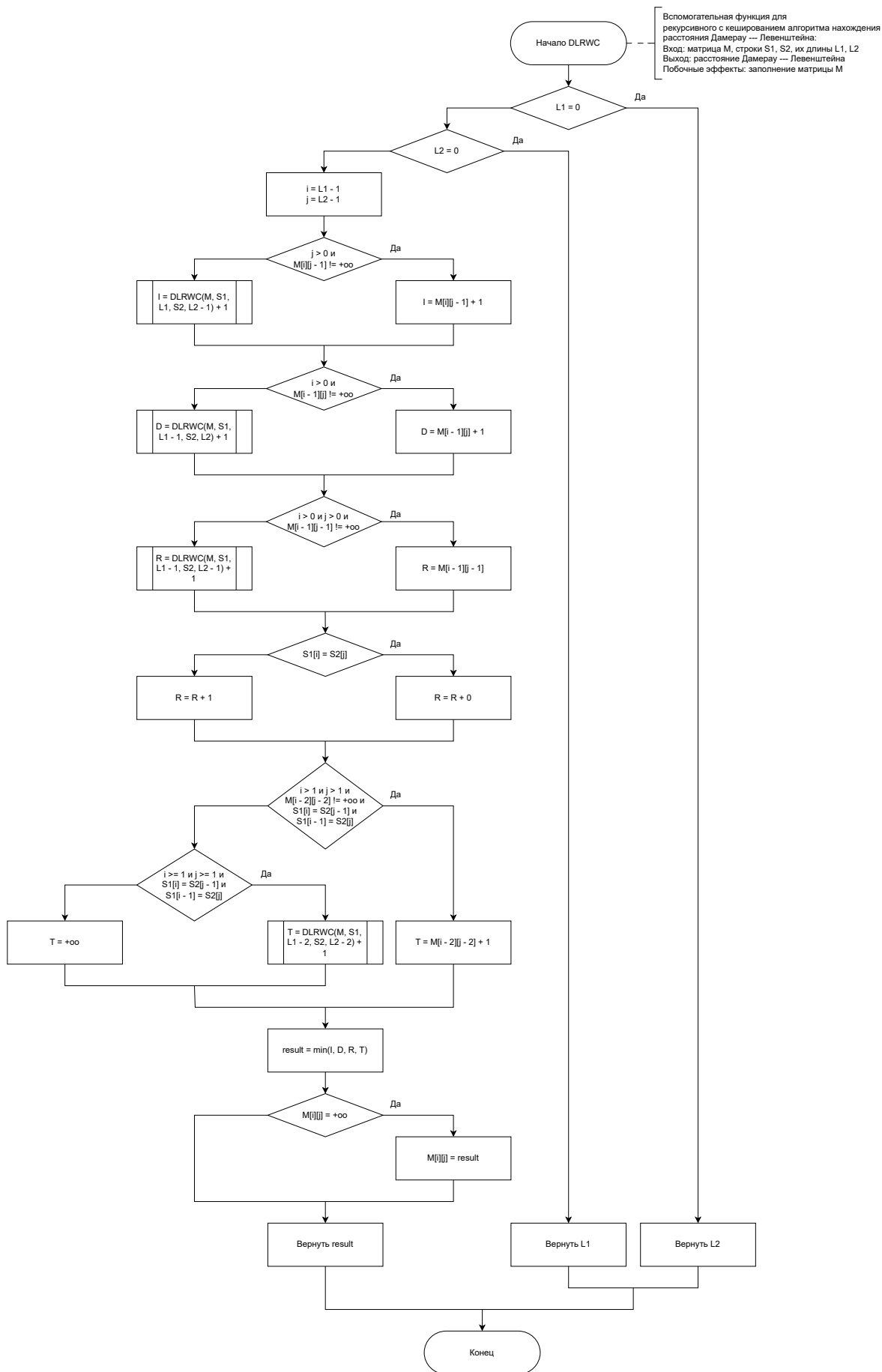


Рисунок 2.5 – Рекурсивная функция нахождения расстояния Дамерау — Левенштейна с использованием матрицы в качестве кеша

2.3 Описание типов и структур данных

Для реализации алгоритмов будут использованы следующие типы и структуры данных:

- `wchar_t` — для символов в кодировке Unicode, из которых состоят строки;
- `wchar_t*` — для строк, подаваемых на вход;
- `size_t` (`unsigned long int`) — для длин строк; для вычисленных расстояний Левенштейна и Дамерау — Левенштейна (как возвращаемых из функций, так и хранящихся в матрице); для размеров матрицы при её создании;
- `size_t*` — для строк матрицы;
- `size_t **` — для матрицы (массива указателей на строки);
- `struct BenchmarkData` (см. 2.1) — для временного хранения данных о замерах времени реализаций алгоритмов.

Листинг 2.1 – Структура для временного хранения данных о замерах времени реализаций алгоритмов

```
1 struct BenchmarkData
2 {
3     size_t string_length;
4     long t_lev_iter_ns;
5     long t_damlev_iter_ns;
6     long t_damlev_rec_ns;
7     long t_damlev_rec_cache_ns;
8 };
```

Вывод

В данном разделе на основе теоретических данных были построены схемы алгоритмов, выбраны типы и структуры данных, которые будут использоваться при реализации алгоритмов.

3 Технологическая часть

В данном разделе будут описаны требования к программному обеспечению, средства реализации и функциональные тесты. Также будут приложены листинги кода.

3.1 Средства реализации

Для реализации данной работы был выбран язык C++ [3] по следующим причинам:

- имеется опыт разработки на данном языке;
- в стандартной библиотеке языка присутствует функция `clock_gettime` [4], которая позволяет рассчитать процессорное время конкретного потока;
- наличие библиотеки для проведения автоматизированного тестирования `gtest` [5];
- язык поддерживается отладчиком `gdb` [6].

Для построения графиков был выбран язык Python [7], так как в нём есть все необходимые для этого библиотеки [8] [9] [10].

В качестве среды разработки был выбран Neovim [11].

3.2 Сведения о модулях программы

Программа поделена на следующие модули:

- `main.cpp` — файл, содержащий точку входа в программу, из которой происходит вызов алгоритмов по разработанному интерфейсу;
- `func.cpp` — файл, содержащий реализации алгоритмов нахождения расстояний Левенштейна и Дамерау — Левенштейна;
- `test.cpp` — файл, содержащий модульные тесты и точку входа в программу для их запуска;

- `benchmark.cpp` — файл, содержащий функции для проведения замеров времени работы реализаций алгоритмов;
- `globals.cpp` — файл, содержащий глобальные переменные для замеров количества рекурсивных вызовов и максимального размера стека;
- `plot.py` — файл, содержащий программу для построения графиков по данным о времени выполнения реализаций алгоритмов.

3.3 Реализация алгоритмов

Листинги исходных кодов программ 4.1 – 4.11 приведены в приложении.

3.4 Функциональные тесты

В таблице 3.1 приведены функциональные тесты для алгоритмов нахождения расстояний Левенштейна и Дameraу — Левенштейна. Все тесты пройдены успешно.

Таблица 3.1 – Функциональные тесты

Входные данные		Расстояние и алгоритм нахождения расстояния			
Строка 1	Строка 2	Левенштейна	Дameraу — Левенштейна		
		Итерационный	Итерационный	Рекурсивный	
				Без кеша	С кешем
λ	λ	0	0	0	0
<i>abc</i>	λ	3	3	3	3
λ	<i>a</i>	1	1	1	1
<i>a</i>	<i>a</i>	0	0	0	0
<i>a</i>	<i>b</i>	1	1	1	1
$\xi\eta$	$\eta\xi$	2	1	1	1
<i>abcab</i>	<i>bacba</i>	4	2	2	2
<i>abcabcab</i>	<i>bacbacba</i>	6	3	3	3
<i>python</i>	<i>ptyhon</i>	2	1	1	1
скат	кот	3	2	2	2
<i>ls</i>	<i>sl</i>	2	1	1	1
<i>make</i>	<i>mkae</i>	2	1	1	1

Вывод

Были реализованы и протестированы алгоритмы:

- Итерационный алгоритм нахождения расстояния Левенштейна;
- Итерационный алгоритм нахождения расстояния Дамерау — Левенштейна;
- Рекурсивный алгоритм нахождения расстояния Дамерау — Левенштейна;
- Рекурсивный с кешированием алгоритм нахождения расстояния Дамерау — Левенштейна.

4 Исследовательская часть

В данном разделе будут приведены примеры работы программ и сравнительный анализ алгоритмов на основе полученных данных.

4.1 Технические характеристики

Технические характеристики устройства, на котором выполнялись замеры времени:

- Процессор: AMD Ryzen™ 7 4700U 2.0 ГГц [12], 8 физических ядер, 8 потоков;
- Оперативная память: 8 ГБ, DDR4, 3200 МГц;
- Операционная система: NixOS 23.05.4448.5550a [13];
- Версия ядра: 6.1.59.

4.2 Демонстрация работы программы

На рисунках 4.1 – 4.3 представлена демонстрация работы программы: ручной ввод слов, проведение замеров времени и построение графиков.

```

human in ~/University/aa on main ● ● λ make run
cd build && ./main
Выберите действие:
  1 – Ввести два слова
  2 – Провести замеры времени
  3 – Провести замеры времени (ручной ввод)
  4 – Построить графики
  5 – Выйти
> 1
abcsab басба
Левенштейн итерационный:
  λ  b  a  c  b  a
λ  0  1  2  3  4  5
a  1  1  1  2  3  4
b  2  1  2  2  2  3
c  3  2  2  2  3  3
a  4  3  2  3  3  3
b  5  4  3  3  3  4
Редакторские операции: RRMRR
Результат: 4

Дамерау–Левенштейн итерационный:
  λ  b  a  c  b  a
λ  0  1  2  3  4  5
a  1  1  1  2  3  4
b  2  1  1  2  2  3
c  3  2  2  1  2  3
a  4  3  2  2  2  2
b  5  4  3  3  2  2
Редакторские операции: TMT
Результат: 2

Левенштейн Итерационный : 4
Дамерау–Левенштейн Итерационный : 2
Дамерау–Левенштейн Рекурсивный : 2
Дамерау–Левенштейн Рекурсивный с кешированием : 2

```

Рисунок 4.1 – Демонстрация работы программы, ввод двух слов

```
human in ~/University/aa on main ● ● λ make run
cd build && ./main
Выберите действие:
  1 – Ввести два слова
  2 – Провести замеры времени
  3 – Провести замеры времени (ручной ввод)
  4 – Построить графики
  5 – Выйти
> 3
Введите начальную длину слов, конечную, шаг изменения
длины, и сколько замеров времени требуется провести
для каждой длины:
1 10 1 10
Производится замер времени...
Длина слов = 1...
Длина слов = 2...
Длина слов = 3...
Длина слов = 4...
Длина слов = 5...
Длина слов = 6...
Длина слов = 7...
Длина слов = 8...
Длина слов = 9...
Длина слов = 10...
Данные записаны в файл.
```

Рисунок 4.2 – Демонстрация работы программы, проведение замеров времени

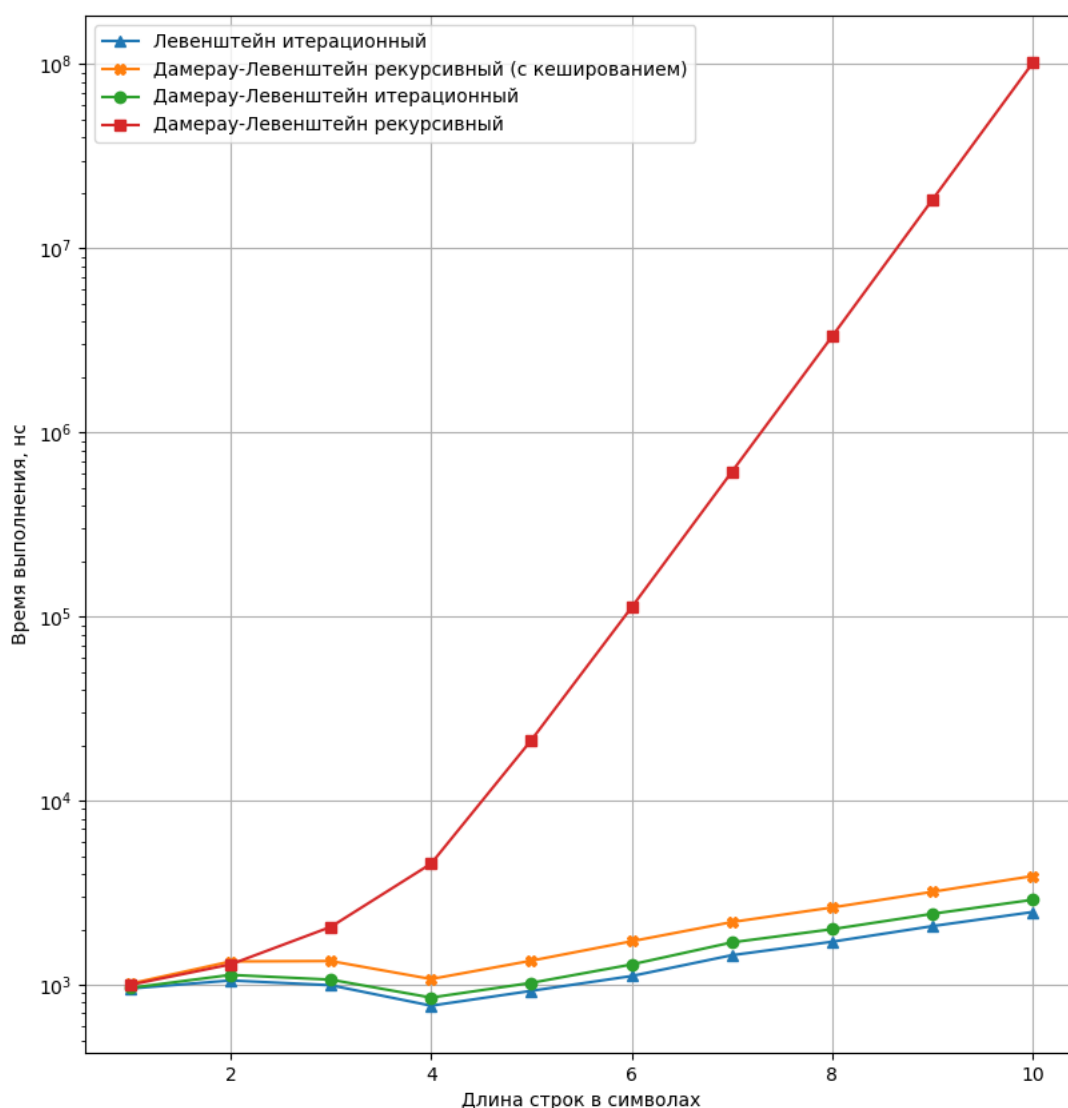


Рисунок 4.3 – Демонстрация работы программы, построение графиков на основе последних замеров времени (создаётся дочерний процесс, который запускает Python-скрипт)

4.3 Временные характеристики

На графиках 4.4 – 4.8 представлены результаты замеров времени выполнения реализаций алгоритмов нахождения расстояний Левенштейна и Дамерау — Левенштейна. Замеры времени проводились для строк одинаковой длины. Отображённое на графиках время является усреднённым для каждой длины строк.

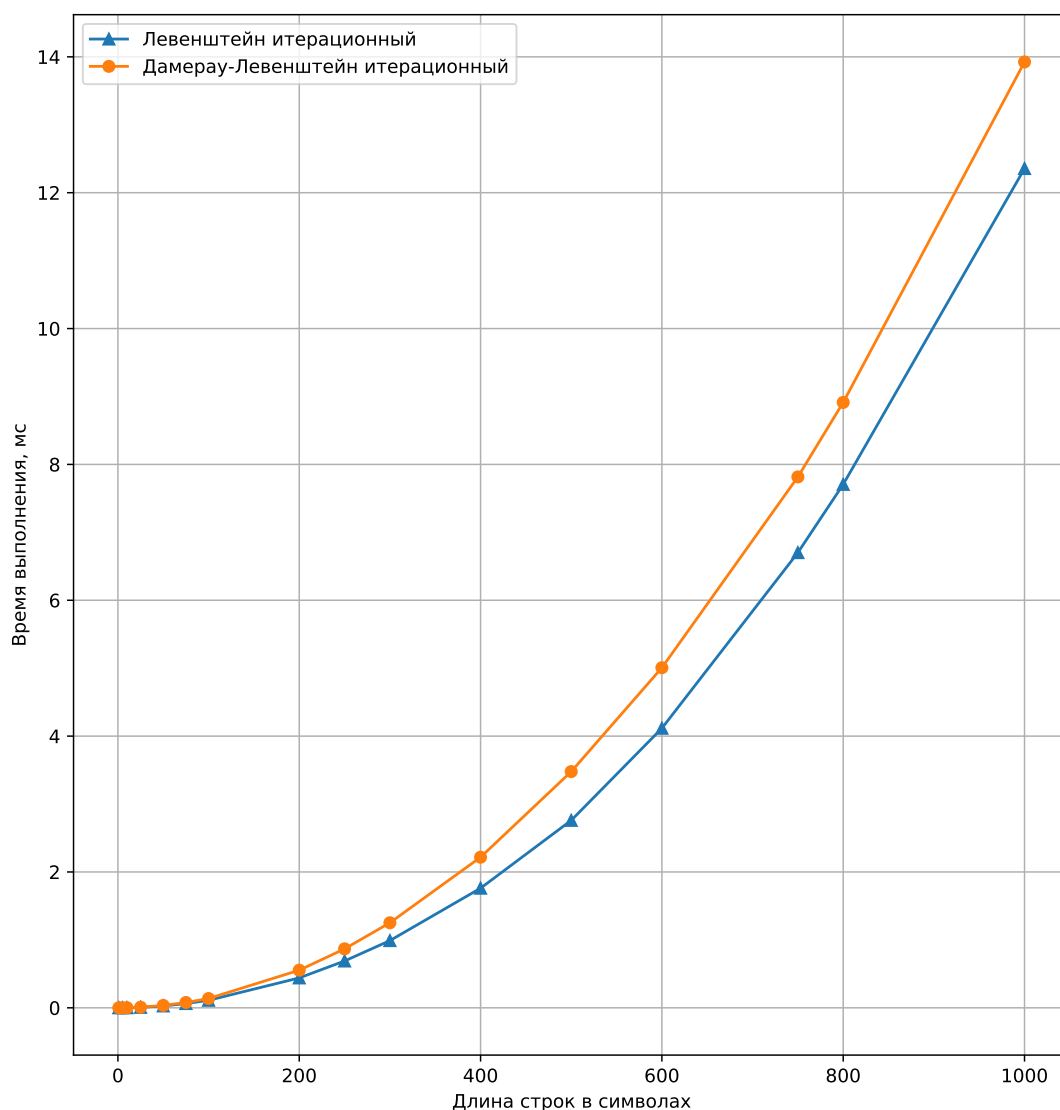


Рисунок 4.4 – Сравнение реализаций итерационных алгоритмов по времени выполнения (среднее время из 100 замеров)

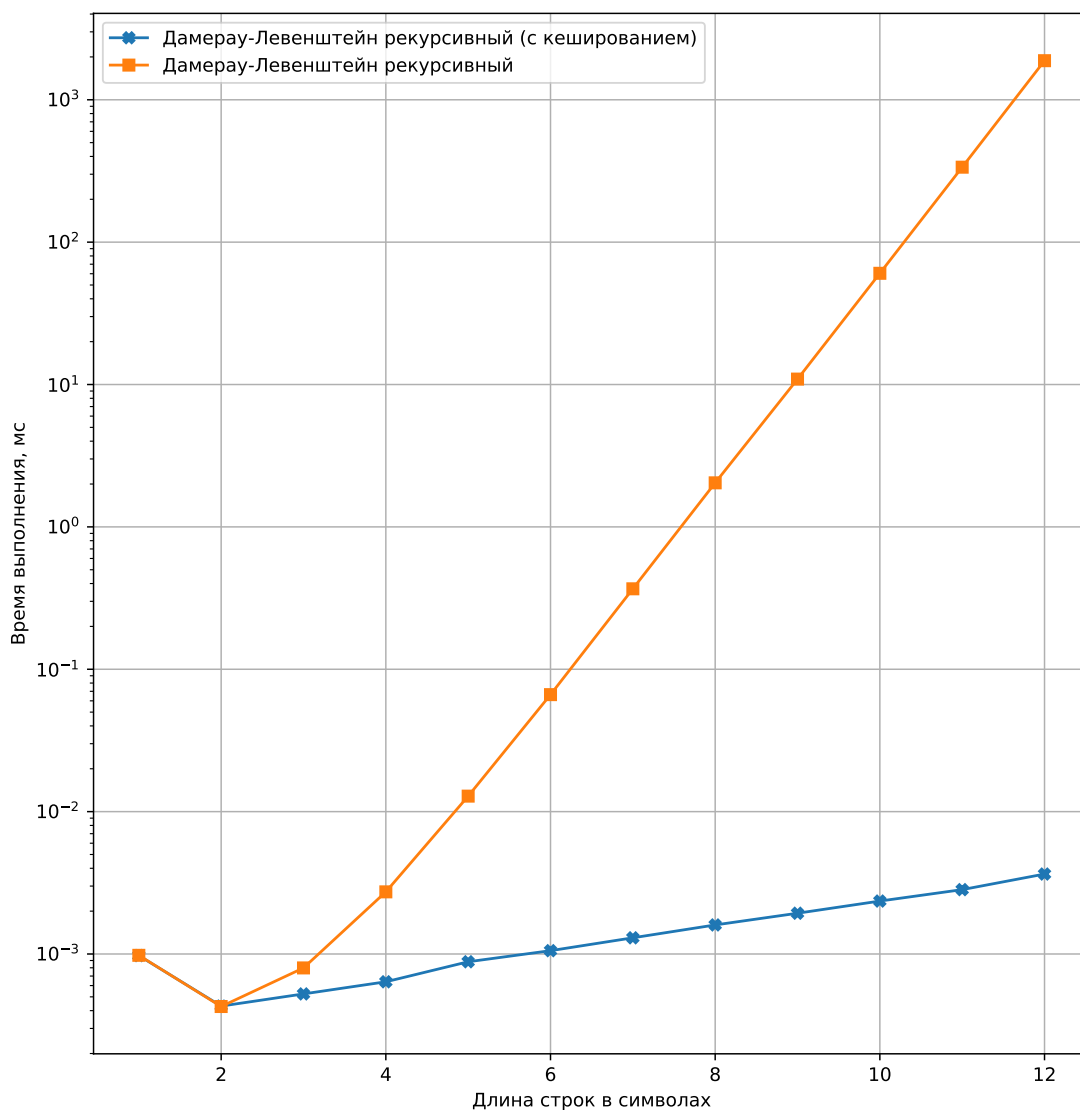


Рисунок 4.5 – Сравнение реализаций рекурсивных алгоритмов по времени выполнения (среднее время из 10 замеров)

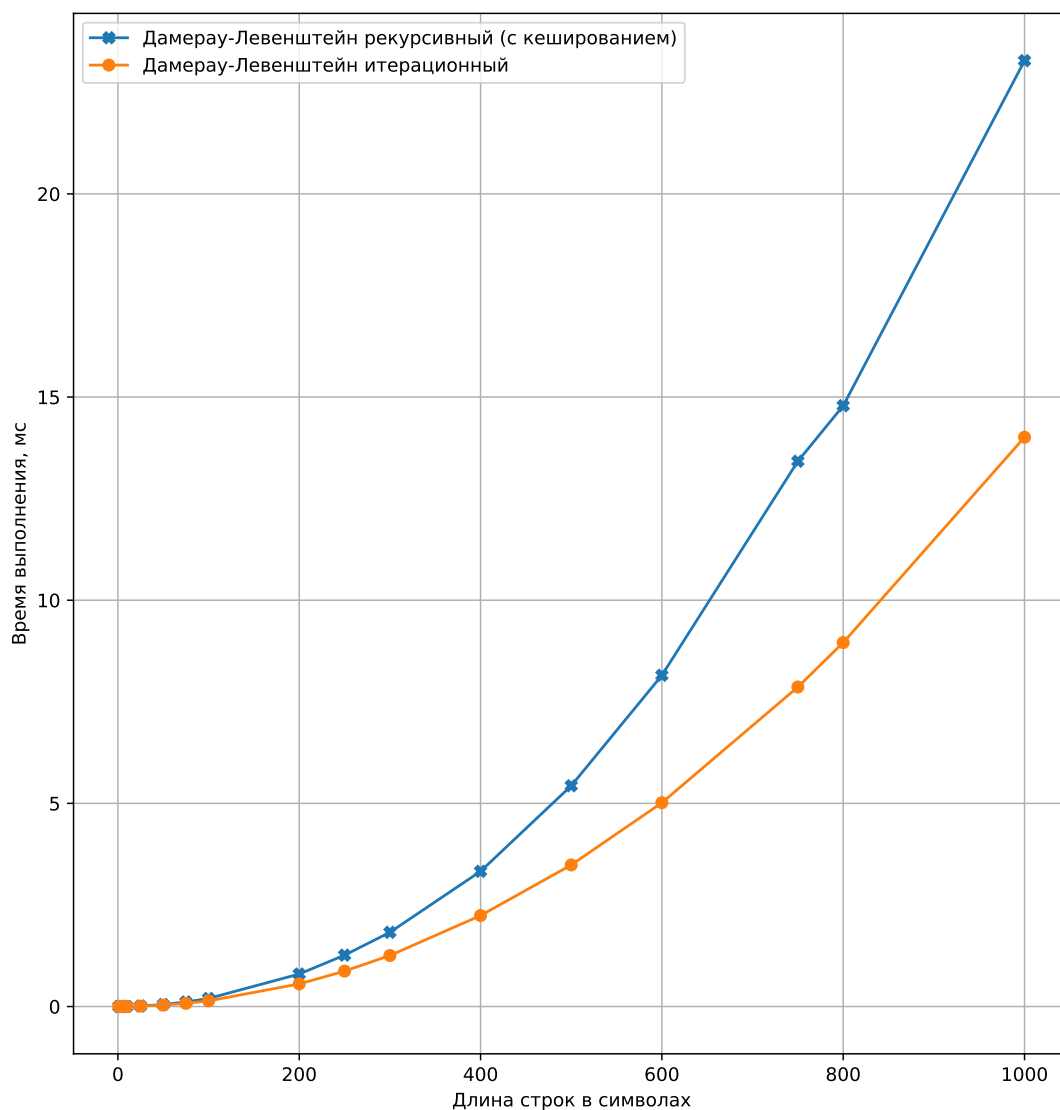


Рисунок 4.6 – Сравнение итерационной и рекурсивной с кешированием реализации алгоритмов нахождения расстояния Дамерау — Левенштейна по времени выполнения (среднее время из 100 замеров)

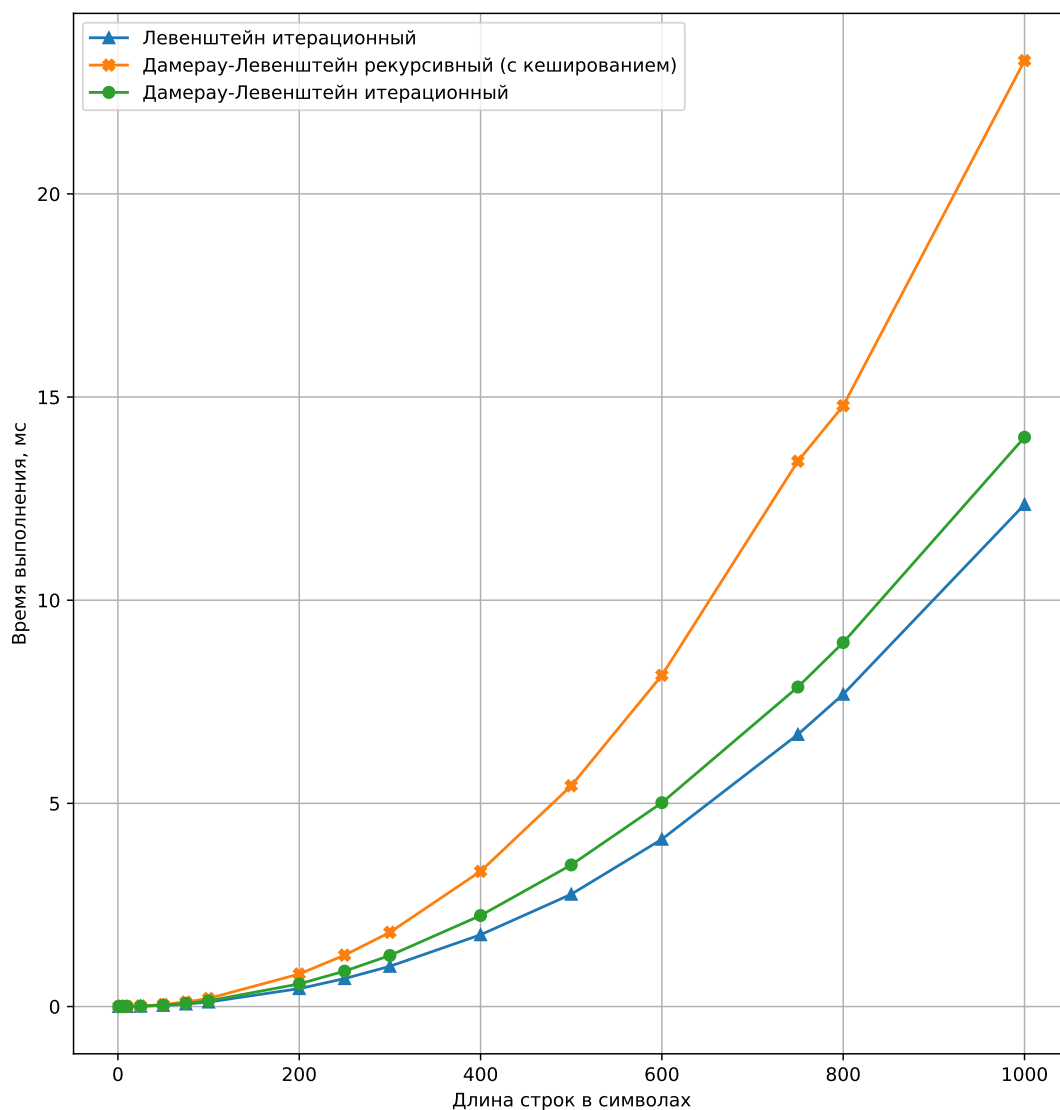


Рисунок 4.7 – Сравнение итерационных и рекурсивной с кешированием реализаций алгоритмов нахождения расстояний Левенштейна и Дамерау — Левенштейна по времени выполнения (среднее время из 100 замеров)

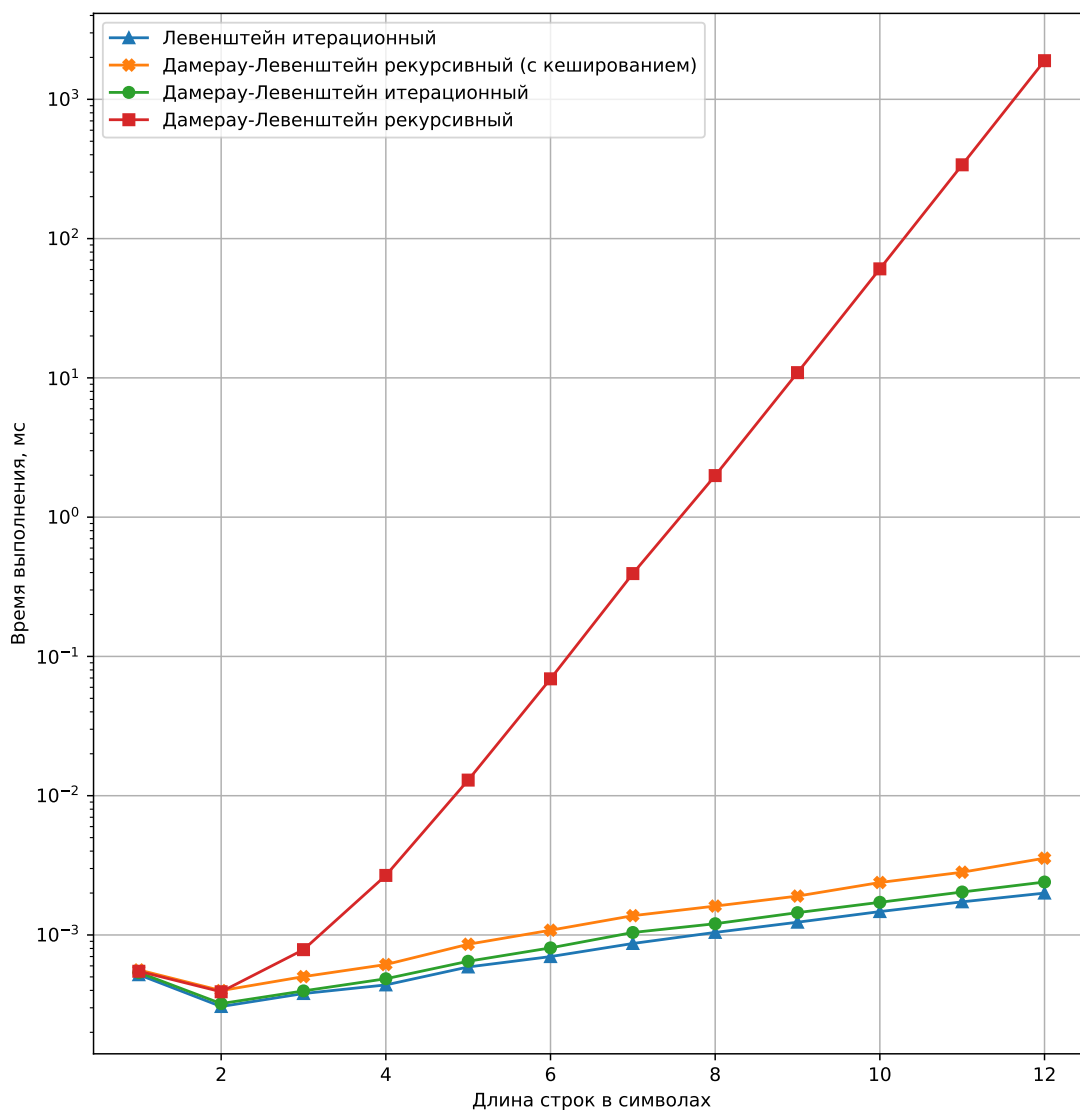


Рисунок 4.8 – Сравнение всех реализаций алгоритмов по времени выполнения (среднее время из 10 замеров)

4.4 Теоретические оценки затрачиваемой памяти

Введём следующие обозначения:

- L_1 — длина первой строки;
- L_2 — длина второй строки;
- *char* — символьный тип данных;
- $*$ — тип данных — указатель;
- *int* — целочисленный тип данных;
- *size()* — функция, вычисляющая размер типа данных в байтах, т.е. на практике $size(char) \sim sizeof(wchar_t)$;

На основе введённых обозначений проведём теоретическую оценку памяти, используемой реализациями алгоритмов.

4.4.1 Оценка памяти реализаций итерационных алгоритмов нахождения расстояний Левенштейна и Дамерау — Левенштейна

- Первая строка: $L_1 \cdot size(char)$;
- Вторая строка: $L_2 \cdot size(char)$;
- Размеры строк: $2 \cdot size(int)$;
- Указатели на строки и адрес возврата: $3 \cdot size(*)$;
- Матрица: $(L_1 + 1) \cdot (L_2 + 1) \cdot size(int) + (L_1 + 1) \cdot size(*) + size(*)$;
- Вспомогательные переменные: $c \cdot size(int)$, где c — количество вспомогательных переменных;

Итого, функция суммарного объёма используемой памяти, имеет вид:

$$\begin{aligned} V(L_1, L_2) = & L_1 \cdot size(char) + L_2 \cdot size(char) \\ & + 2 \cdot size(int) + (L_1 + 1) \cdot (L_2 + 1) \cdot size(int) \\ & + 3 \cdot size(*) + (L_1 + 1) \cdot size(*) + size(*) + c \cdot size(int) \end{aligned} \quad (4.1)$$

Просуммировав все значения и приняв $size(\dots)$ за константы, можно сделать вывод, что $V(L_1, L_2) \in \Theta(L_1 \cdot L_2)$.

4.4.2 Оценка памяти реализации рекурсивного алгоритма нахождения расстояния Дамерау — Левенштейна

- Первая строка: $L_1 \cdot size(char)$;
- Вторая строка: $L_2 \cdot size(char)$;
- Размеры строк: $2 \cdot size(int)$;
- Указатели на строки и адрес возврата: $3 \cdot size(*)$;
- Вспомогательные переменные: $c \cdot size(int)$, где c — количество вспомогательных переменных;

Максимальная глубина стека вызовов при рекурсивной реализации алгоритма нахождения расстояния Дамерау — Левенштейна равна сумме длин входящих строк, соответственно, максимальный расход памяти рассчитывается по:

$$V(L_1, L_2) = (L_1 + L_2) \cdot ((2 + c) \cdot size(int) + 3 \cdot size(*)) \quad (4.2)$$

Просуммировав все значения и приняв $size(\dots)$ за константы, можно сделать вывод, что $V(L_1, L_2) \in \Theta(L_1 + L_2) = \Theta(L)$, то есть, максимальный расход памяти линейно зависит от длины входящих строк.

4.4.3 Оценка памяти реализации рекурсивного с кешированием алгоритма нахождения расстояния Дамерау — Левенштейна

- Первая строка: $L_1 \cdot size(char)$;
- Вторая строка: $L_2 \cdot size(char)$;

- Размеры строк: $2 \cdot \text{size}(\text{int})$;
- Указатели на строки и адрес возврата: $3 \cdot \text{size}(*)$;
- Матрица: $L_1 \cdot L_2 \cdot \text{size}(\text{int}) + L_1 \cdot \text{size}(*) + \text{size}(*)$;
- Вспомогательные переменные: $c \cdot \text{size}(\text{int})$, где c — количество вспомогательных переменных;

Итого, под матрицу используется:

$$V_M(L_1, L_2) = L_1 \cdot L_2 \cdot \text{size}(\text{int}) + L_1 \cdot \text{size}(*) + \text{size}(*) \in \Theta(L_1 \cdot L_2) \quad (4.3)$$

Под стек будет использоваться меньше памяти, чем в исключительно рекурсивной реализации, так как функция не будет углублять стек вызовов, если встретится уже вычисленное ранее значение. Следовательно, $V_S(L_1, L_2) \in \Theta(L)$.

Тогда суммарный объём используемой памяти имеет вид:

$$V(L_1, L_2) = V_M(L_1, L_2) + V_S(L_1, L_2) \in \Theta(L_1 \cdot L_2) + \Theta(L) = \Theta(L_1 \cdot L_2) \quad (4.4)$$

Таким образом, рекурсивная без кеширования реализация алгоритма нахождения Дамерау — Левенштейна расходует меньше памяти при достаточно длинных строках по сравнению с остальными реализациями. Её расход памяти растёт пропорционально $L_1 + L_2$, в то время как расход памяти остальных реализаций алгоритмов растёт пропорционально $L_1 \cdot L_2$.

4.5 Характеристики по памяти

На графиках 4.9 – 4.12 представлены результаты измерения затрачиваемой памяти реализациями алгоритмов нахождения расстояний Левенштейна и Дамерау — Левенштейна. Замеры памяти проводились для строк одинаковой длины. Отображённая на графиках память является усреднённой для каждой длины строк.

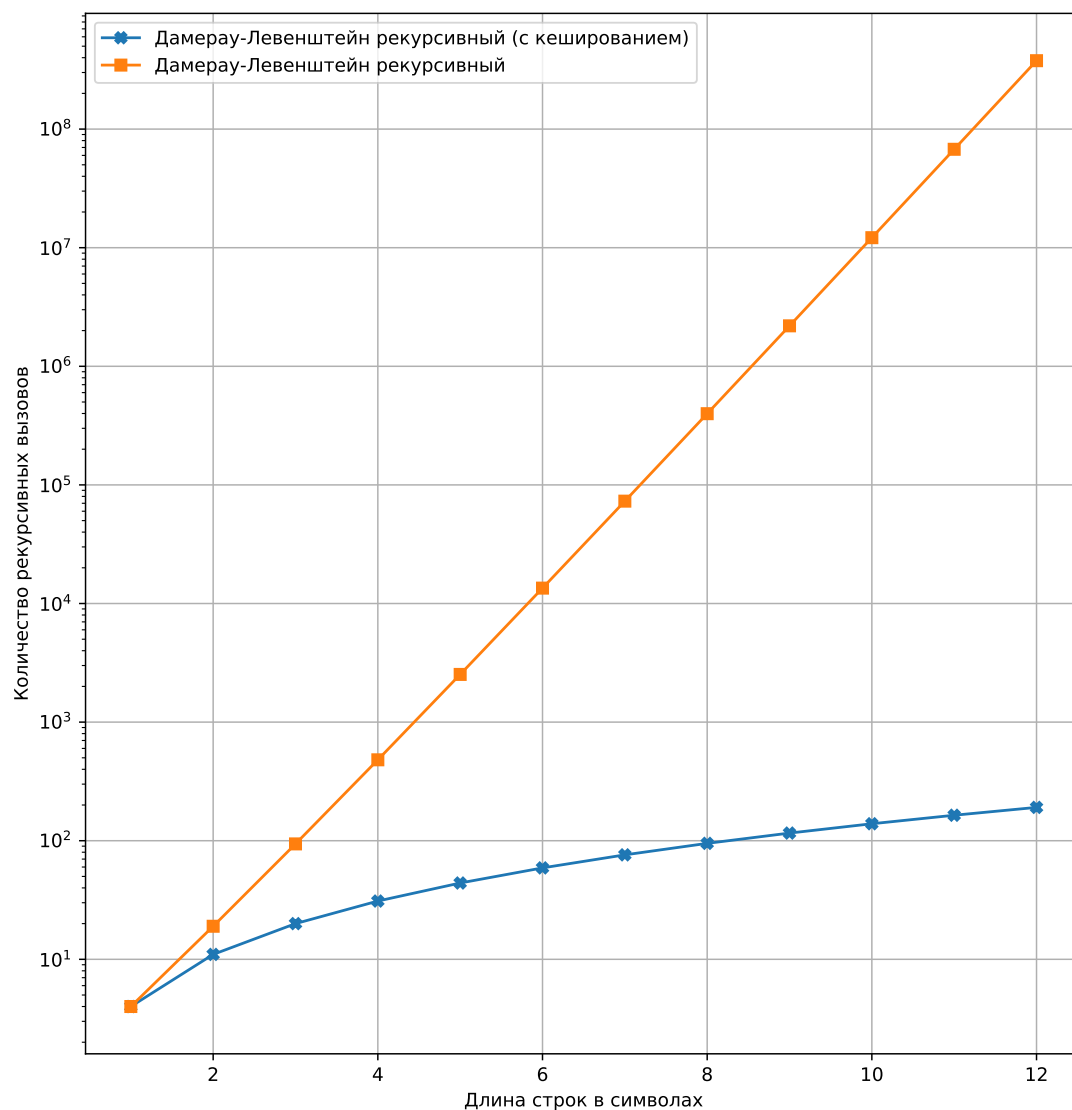


Рисунок 4.9 – Сравнение количества рекурсивных вызовов у рекурсивных реализаций алгоритмов нахождения расстояния Дамерау — Левенштейна

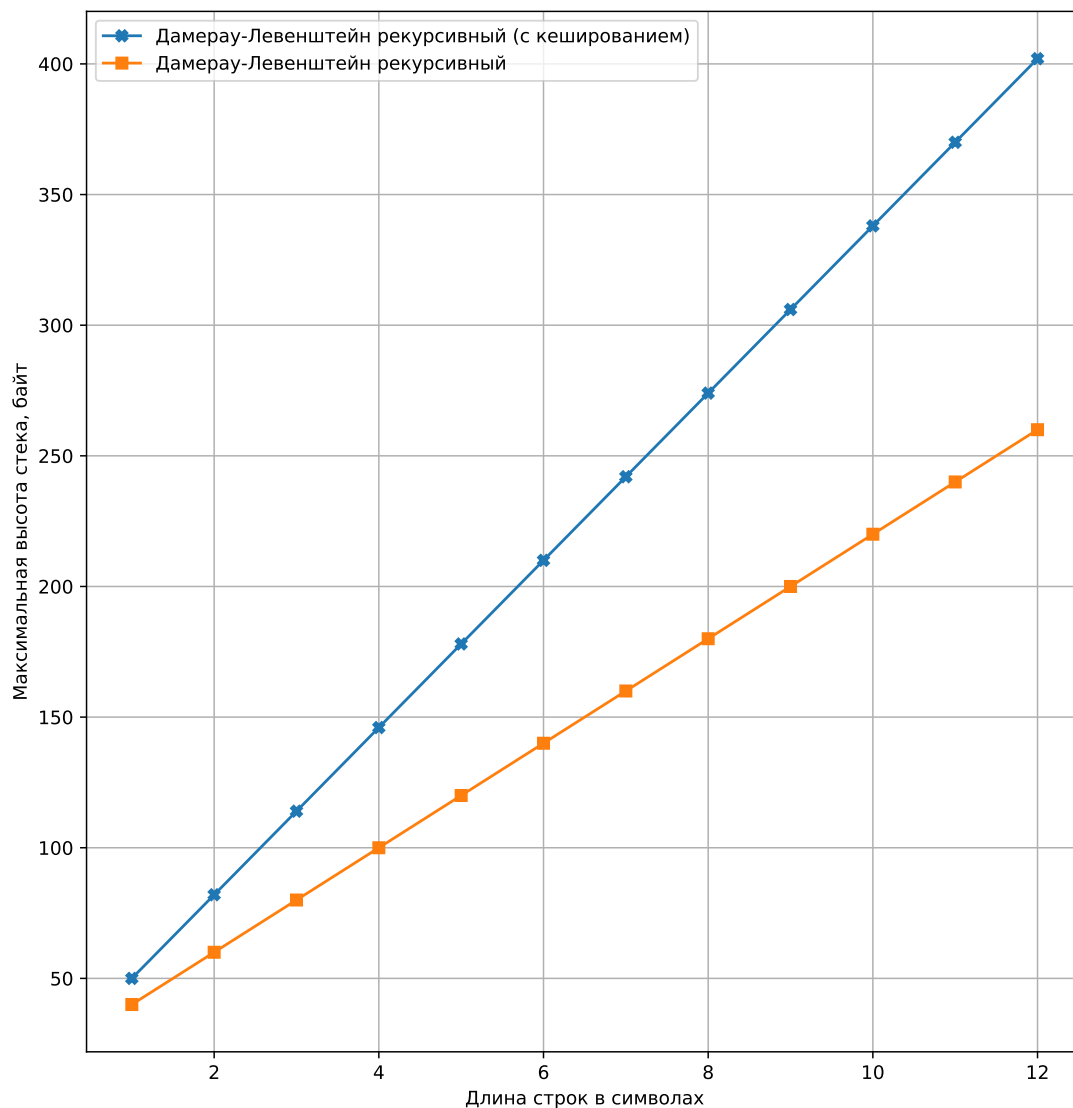


Рисунок 4.10 – Сравнение максимальной высоты стека вызовов у рекурсивных реализаций алгоритмов нахождения расстояния Дамерау — Левенштейна

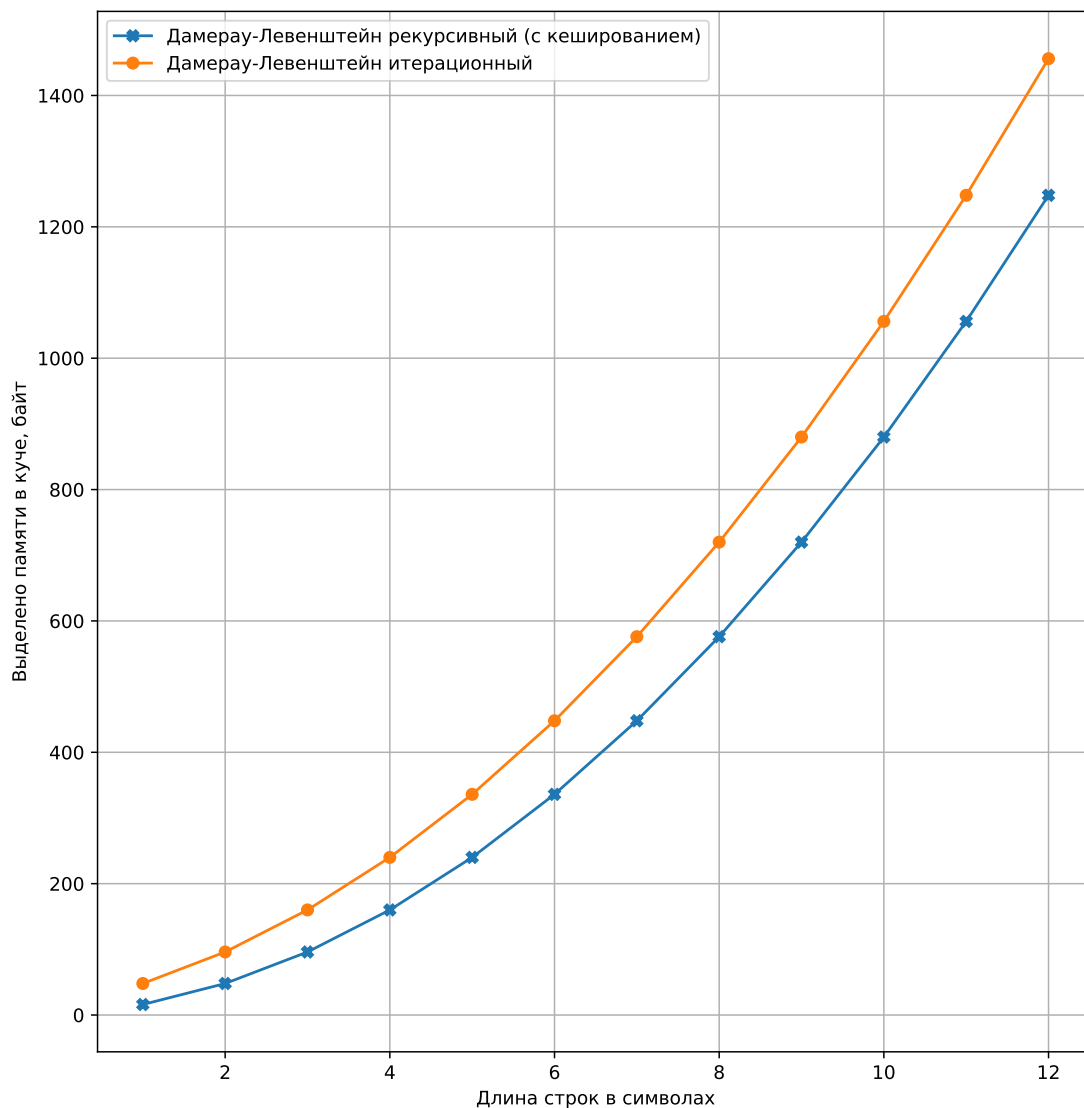


Рисунок 4.11 – Сравнение выделяемой памяти в куче для итерационного и рекурсивного с кешированием реализаций алгоритмов нахождения расстояния Дамерау — Левенштейна

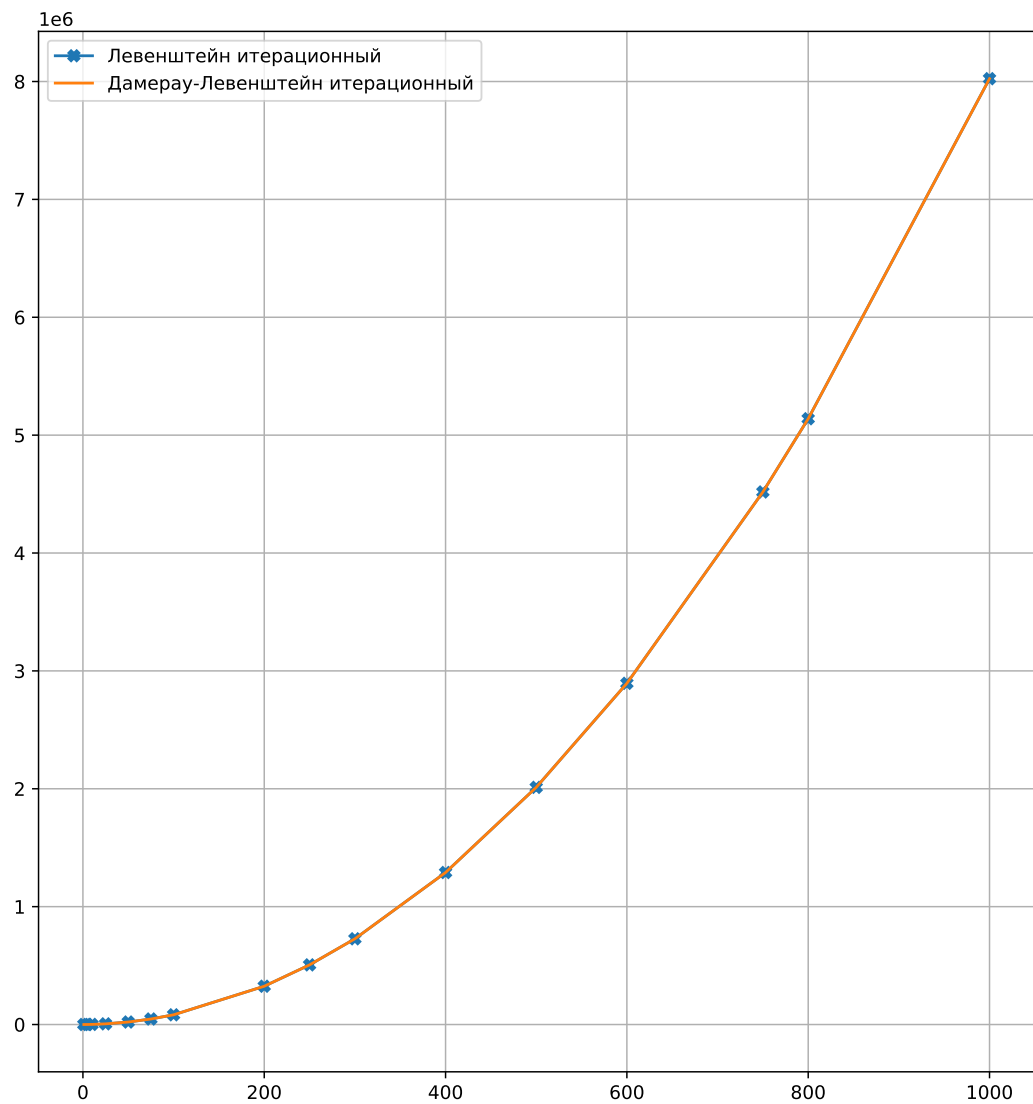


Рисунок 4.12 – Сравнение выделяемой памяти в куче для итерационных реализаций алгоритмов нахождения расстояний Левенштейна и Дameraу — Левенштейна

Вывод

В данном разделе было произведено сравнение количества затраченного времени и памяти реализаций алгоритмов поиска расстояний Левенштейна и Дамерау — Левенштейна. Наименее затратной по времени оказалась итерационная реализация алгоритма нахождения расстояния Левенштейна.

Проанализировав использование памяти в алгоритмах, можно сделать вывод, что итеративные алгоритмы и рекурсивные алгоритмы с кешированием требуют больше памяти по сравнению с рекурсивным алгоритмом без кеширования. В реализациях, использующих матрицы, максимальный используемый объем памяти увеличивается пропорционально произведению длин строк. С другой стороны, для рекурсивного алгоритма без кеширования потребление памяти увеличивается пропорционально сумме длин строк.

ЗАКЛЮЧЕНИЕ

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. А. Погорелов Д., М. Таразанов А. Сравнительный анализ алгоритмов редакционного расстояния Левенштейна и Дамерау-Левенштейна // Синергия Наук. 2019. — Режим доступа: <https://elibrary.ru/item.asp?id=36907767> (дата обращения: 27.10.2023).
2. Левенштейн В. И. Двоичные коды с исправлением выпадений, вставок и замещений символов. — М.: «Наука», Доклады АН СССР, 1965. Т. 163. С. 845–848.
3. ISO C++ [Электронный ресурс]. — Режим доступа: <https://isocpp.org/std/the-standard> (дата обращения: 01.11.2023).
4. Linux man page — clock_gettime(3) [Электронный ресурс]. — Режим доступа: https://linux.die.net/man/3/clock_gettime (дата обращения: 01.11.2023).
5. GoogleTest User's Guide [Электронный ресурс]. — Режим доступа: <https://google.github.io/googletest> (дата обращения: 01.11.2023).
6. GDB: The GNU Project Debugger [Электронный ресурс]. — Режим доступа: <https://www.sourceware.org/gdb> (дата обращения: 01.11.2023).
7. Python [Электронный ресурс]. — Режим доступа: <https://www.python.org> (дата обращения: 01.11.2023).
8. Pandas [Электронный ресурс]. — Режим доступа: <https://pandas.pydata.org> (дата обращения: 01.11.2023).
9. NumPy [Электронный ресурс]. — Режим доступа: <https://numpy.org> (дата обращения: 01.11.2023).
10. Matplotlib [Электронный ресурс]. — Режим доступа: <https://matplotlib.org> (дата обращения: 01.11.2023).
11. Neovim [Электронный ресурс]. — Режим доступа: <https://neovim.io> (дата обращения: 01.11.2023).

12. AMD Ryzen™ 7 4700U [Электронный ресурс]. — Режим доступа: <https://www.amd.com/en/product/9096> (дата обращения: 02.11.2023).
13. NixOS [Электронный ресурс]. — Режим доступа: <https://nixos.org> (дата обращения: 02.11.2023).

ПРИЛОЖЕНИЕ А

Листинг 4.1 – Создание матрицы для последующего использования в реализациях алгоритмов нахождения расстояний Левенштейна и Дamerau – Левенштейна

```
1 size_t **create_matrix(size_t n_rows, size_t n_columns)
2 {
3     if (n_rows < 1 || n_columns < 1)
4         return NULL;
5
6     size_t *mem = (size_t*) malloc(n_rows * n_columns *
7                                     sizeof(size_t));
8
9     if (mem == NULL)
10         return NULL;
11
12     size_t **matrix = (size_t**) malloc(n_rows * sizeof(size_t*));
13
14     if (matrix == NULL)
15     {
16         free(mem);
17         return NULL;
18     }
19
20     for (size_t i = 0; i < n_rows; i++)
21     {
22         matrix[i] = mem + i * n_columns;
23         matrix[i][0] = i;
24     }
25
26     for (size_t j = 1; j < n_columns; j++)
27         matrix[0][j] = j;
28
29     return matrix;
30 }
```

Листинг 4.2 – Освобождение памяти из-под созданной матрицы

```
1 void free_matrix(size_t **matrix, size_t *first_row)
2 {
3     free(first_row);
4     free(matrix);
5 }
```

```
5 }
```

Листинг 4.3 – Нахождение минимального числа из трёх и возврат указателя на него

```
1 size_t *min3(size_t *a, size_t *b, size_t *c)
2 {
3     size_t *result;
4     size_t min = std::min(*a, std::min(*b, *c));
5
6     if (min == *a)
7         result = a;
8     else if (min == *b)
9         result = b;
10    else
11        result = c;
12
13    return result;
14 }
```

Листинг 4.4 – Нахождение минимального числа из четырёх и возврат указателя на него

```
1 size_t *min4(size_t *a, size_t *b, size_t *c, size_t *d)
2 {
3     size_t *result;
4     size_t min = std::min(std::min(*a, *b), std::min(*c, *d));
5
6     if (min == *a)
7         result = a;
8     else if (min == *b)
9         result = b;
10    else if (min == *c)
11        result = c;
12    else
13        result = d;
14
15    return result;
16 }
```

Листинг 4.5 – Часть реализации итерационного алгоритма нахождения расстояния Левенштейна, участвующая в замерах времени выполнения реализаций исследуемых алгоритмов

```

1 size_t lev_ifm_helper(size_t **matrix, const wchar_t *s1, size_t
  len1, const wchar_t *s2, size_t len2)
2 {
3     size_t result = 0;
4     bool replace_skip_cond;
5     size_t insert_cost, delete_cost, replace_cost, *who;
6
7     for (size_t i = 1; i < len1; i++)
8     {
9         for (size_t j = 1; j < len2; j++)
10        {
11            insert_cost = matrix[i - 1][j] + 1;
12            delete_cost = matrix[i][j - 1] + 1;
13            replace_skip_cond = (s1[i] == s2[j]);
14            replace_cost = matrix[i - 1][j - 1] +
              (replace_skip_cond ? 0 : 1);
15            who = min3(&insert_cost, &delete_cost, &replace_cost);
16            matrix[i][j] = *who;
17        }
18    }
19
20    result = matrix[len1 - 1][len2 - 1];
21
22    return result;
23 }

```

Листинг 4.6 – Реализация итерационного алгоритма нахождения расстояния Левенштейна

```

1 size_t levenshtein_iterative_full_matrix(const wchar_t *str1,
  size_t len1, const wchar_t *str2, size_t len2)
2 {
3     if (len1 == 0) return len2;
4     if (len2 == 0) return len1;
5
6     ++len1;
7     ++len2;
8     const wchar_t *s1 = str1 - 1;
9     const wchar_t *s2 = str2 - 1;

```

```

10
11     size_t **matrix = create_matrix(len1 , len2);
12
13     if (matrix == NULL) return -1;
14
15     size_t result = lev_ifm_helper(matrix , s1 , len1 , s2 , len2);
16
17     free_matrix(matrix , matrix[0]);
18
19     return result;
20 }

```

Листинг 4.7 – Часть реализации итерационного алгоритма нахождения расстояния Дамерау — Левенштейна, участвующая в замерах времени выполнения реализаций исследуемых алгоритмов

```

1 size_t damlev_ifm_helper(size_t **matrix , const wchar_t *s1 ,
   size_t len1 , const wchar_t *s2 , size_t len2)
2 {
3     size_t result = 0;
4     bool replace_skip_cond , swap_cond;
5     size_t insert_cost , delete_cost , replace_cost , swap_cost ,
        *who;
6
7     for (size_t i = 1; i < len1; i++)
8     {
9         for (size_t j = 1; j < len2; j++)
10        {
11            insert_cost = matrix[i - 1][j] + 1;
12            delete_cost = matrix[i][j - 1] + 1;
13            replace_skip_cond = (s1[i] == s2[j]);
14            replace_cost = matrix[i - 1][j - 1] +
                (replace_skip_cond ? 0 : 1);
15            if (i >= 2 && j >= 2) [[likely]]
16            {
17                swap_cond = (s1[i] == s2[j - 1] && s1[i - 1] ==
                    s2[j]);
18                swap_cost = swap_cond ? matrix[i - 2][j - 2] + 1
                    : U_INF; // U_INF = -1
19                who = min4(&insert_cost , &delete_cost ,
                    &replace_cost , &swap_cost);
20            }

```

```

21         else
22         {
23             who = min3(&insert_cost , &delete_cost ,
24                       &replace_cost );
25         }
26         matrix[i][j] = *who;
27     }
28 }
29
30 result = matrix[len1 - 1][len2 - 1];
31
32 return result;
33 }

```

Листинг 4.8 – Реализация итерационного алгоритма нахождения расстояния
Дамерау — Левенштейна

```

1 size_t damerau_levenshtein_iterative_full_matrix(const wchar_t
2   *str1 , size_t len1 , const wchar_t *str2 , size_t len2)
3 {
4     if (len1 == 0) return len2;
5     if (len2 == 0) return len1;
6
7     ++len1;
8     ++len2;
9     const wchar_t *s1 = str1 - 1;
10    const wchar_t *s2 = str2 - 1;
11
12    size_t **matrix = create_matrix(len1 , len2);
13    if (matrix == NULL) return -1;
14
15    size_t result = damlev_ifm_helper(matrix , s1 , len1 , s2 , len2);
16
17    free_matrix(matrix , matrix[0]);
18
19    return result;
20 }

```

Листинг 4.9 – Часть реализации рекурсивного с кешированием алгоритма
нахождения расстояния Дамерау — Левенштейна, участвующая в замерах
времени выполнения реализаций исследуемых алгоритмов

```

1 size_t damlev_rwc_helper(size_t **matrix, const wchar_t *str1,
2 size_t len1, const wchar_t *str2, size_t len2)
3 {
4     if (len1 == 0) return len2;
5     if (len2 == 0) return len1;
6
7     size_t i = len1 - 1;
8     size_t j = len2 - 1;
9
10    size_t insert = (((j > 0) && (matrix[i][j - 1] != U_INF))
11        ? matrix[i][j - 1]
12        : damlev_rwc_helper(matrix, str1, len1, str2, len2 - 1))
13        + 1;
14
15    size_t del = (((i > 0) && (matrix[i - 1][j] != U_INF))
16        ? matrix[i - 1][j]
17        : damlev_rwc_helper(matrix, str1, len1 - 1, str2, len2))
18        + 1;
19
20    size_t replace = (((i > 0 && j > 0) && (matrix[i - 1][j - 1]
21        != U_INF))
22        ? matrix[i - 1][j - 1]
23        : damlev_rwc_helper(matrix, str1, len1 - 1, str2, len2 -
24        1))
25        + (str1[i] == str2[j] ? 0 : 1);
26
27    size_t swap = U_INF;
28    if (i > 1 && j > 1 && matrix[i - 2][j - 2] != U_INF &&
29        (str1[i] == str2[j - 1] && str1[i - 1] == str2[j]))
30    {
31        swap = matrix[i - 2][j - 2] + 1;
32    }
33    else if (i >= 1 && j >= 1 && (str1[i] == str2[j - 1] &&
34        str1[i - 1] == str2[j]))
35    {
36        swap = damlev_rwc_helper(matrix, str1, len1 - 2, str2,
37            len2 - 2) + 1;
38    }
39
40    size_t result = *min4(&insert, &del, &replace, &swap);

```

```

36     if (matrix[i][j] == U_INF) matrix[i][j] = result;
37
38     return result;
39 }

```

Листинг 4.10 – Реализация рекурсивного с кешированием алгоритма нахождения расстояния Дамерау — Левенштейна

```

1 size_t damerau_levenshtein_recursive_with_cache(const wchar_t
    *str1, size_t len1, const wchar_t *str2, size_t len2)
2 {
3     if (len1 == 0) return len2;
4     if (len2 == 0) return len1;
5
6     size_t **matrix = create_matrix(len1, len2);
7
8     if (matrix == NULL) return -1;
9
10    for (size_t i = 0; i < len1; i++)
11        for (size_t j = 0; j < len2; j++)
12            matrix[i][j] = U_INF;
13
14    size_t result = damlev_rwc_helper(matrix, str1, len1, str2,
        len2);
15
16    free_matrix(matrix, matrix[0]);
17
18    return result;
19 }

```

Листинг 4.11 – Реализация рекурсивного алгоритма нахождения расстояния Дамерау — Левенштейна

```

1 size_t damerau_levenshtein_recursive_no_cache(const wchar_t
    *str1, size_t len1, const wchar_t *str2, size_t len2)
2 {
3     if (len1 == 0) return len2;
4     if (len2 == 0) return len1;
5
6     size_t insert = damerau_levenshtein_recursive_no_cache(str1,
        len1, str2, len2 - 1) + 1;
7     size_t del = damerau_levenshtein_recursive_no_cache(str1,
        len1 - 1, str2, len2) + 1;

```

```

8      size_t replace = damerau_levenshtein_recursive_no_cache(str1 ,
9          len1 - 1, str2 , len2 - 1)
10          + (str1[len1 - 1] == str2[len2 - 1] ? 0 : 1);
11      size_t swap = (len1 >= 2 && len2 >= 2)
12          ? (
13              (str1[len1 - 1] == str2[len2 - 2] && str1[len1 -
14                  2] == str2[len2 - 1])
15              ? damerau_levenshtein_recursive_no_cache(str1 ,
16                  len1 - 2, str2 , len2 - 2) + 1
17              : U_INF
18          )
19          : U_INF;

      return *min4(&insert , &del , &replace , &swap);
}

```


ПРИЛОЖЕНИЕ Б