



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"

Студент _____ Недолужко Денис Вадимович _____
Группа _____ ИУ7-53Б _____
Название предприятия _____ МГТУ им. Н. Э. Баумана, каф. ИУ7 _____
Тема _____ Расстояние Левенштейна _____

Студент	_____	Недолужко Д. В.
	(Подпись, дата)	(И.О. Фамилия)
Руководитель курсовой работы	_____	Волкова Л. Л.
	(Подпись, дата)	(И.О. Фамилия)

Москва
2021 г.

Содержание

Введение	3
1 Аналитическая часть	5
1.1 Расстояние Левенштейна	5
1.2 Рекурсивный алгоритм нахождения расстояния Левенштейна	5
1.3 Матричный алгоритм нахождения расстояния Левенштейна	6
1.4 Рекурсивный алгоритм нахождения расстояния Левенштейна с заполнением матрицы	7
1.5 Расстояния Дамерау — Левенштейна	7
Вывод	8
2 Конструкторская часть	9
2.1 Схемы алгоритмов	9
2.2 Оценка памяти	13
2.3 Структуры данных	13
Вывод	13
3 Технологическая часть	14
3.1 Выбор средств реализации	14
3.2 Требования к программному обеспечению	14
3.3 Листинги кода	15
Вывод	21
4 Экспериментальная часть	22
4.1 Технические характеристики	22
4.2 Тестирование	22
4.3 Временные характеристики	23
4.4 Сравнительный анализ алгоритмов	25
Вывод	25
Заключение	26
Список литературы	27

Введение

Расстояние Левенштейна (редакционное расстояние, дистанция редактирования) — метрика, измеряющая по модулю разность между двумя последовательностями символов. Она определяется как минимальное количество односимвольных операций (а именно вставки, удаления, замены), необходимых для превращения одной последовательности символов в другую. В общем случае, операциям, используемым в этом преобразовании, можно назначить разные цены. Широко используется в теории информации и компьютерной лингвистике.

Впервые задачу поставил в 1965 году советский математик Владимир Левенштейн при изучении последовательностей $0 - 1$ [1], впоследствии более общую задачу для произвольного алфавита связали с его именем.

Расстояние Левенштейна используется:

- в компьютерной лингвистике для проведения автозамены
- в биоинформатике для сравнения генов, хромосом, белков
- в утилите diff для определения изменений в текстовых файлах

Расстояние Дамерау — Левенштейна (названо в честь учёных Фредерика Дамерау и Владимира Левенштейна) — это мера разницы двух строк символов, определяемая как минимальное количество операций вставки, удаления, замены и транспозиции (перестановки двух соседних символов), необходимых для перевода одной строки в другую. Является модификацией расстояния Левенштейна: к операциям вставки, удаления и замены символов, определённых в расстоянии Левенштейна добавлена операция транспозиции (перестановки) символов.

Целью работы является исследование алгоритмов Левенштейна и Дамерау - Левенштейна.

Для достижения поставленной цели необходимо выполнить следующие задачи:

- а) изучение алгоритма Левенштейна и Дамарау - Левенштейна

б) реализация рекурсивных (без кеша) и итеративных (с кешем) алгоритмов Левенштейна и Дамерау - Левенштейна

в) замер и сравнение процессорного времени затрачиваемых реализованными алгоритмами

г) замер и сравнение используемой памяти реализованными алгоритмами

1 Аналитическая часть

В данном разделе описаны математические модели исследуемой области.

1.1 Расстояние Левенштейна

Расстояние редакционное, расстояние Левенштейна - это минимальное количество операций, необходимое для превращений одной строки в другую.

Для определения стоимости преобразований вводятся следующие операции со величиной штрафа - 1:

- I (от англ. insert) - вставка
- D (от англ. delete) - удаление
- R (от англ. replace) - замена

Также вводится обозначение M (от англ. match). Штраф для совпадения двух символом считается равным 0.

1.2 Рекурсивный алгоритм нахождения расстояния Левенштейна

Расстояние Левенштейна между двумя строками a и b может быть вычислено по формуле 1.1, где $|a|$ означает длину строки a; $a[i]$ — i-ый символ строки a, функция $D(i, j)$ определена как:

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

а функция 1.2 определена как:

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

Рекурсивный алгоритм реализует формулу 1.1. Функция D составлена из следующих соображений:

- а) для перевода из пустой строки в пустую требуется ноль операций;
- б) для перевода из пустой строки в строку a требуется $|a|$ операций;
- в) для перевода из строки a в пустую требуется $|a|$ операций;

Для перевода из строки a в строку b требуется выполнить последовательно некоторое количество операций (удаление, вставка, замена) в некоторой последовательности. Последовательность проведения любых двух операций можно поменять, порядок проведения операций не имеет никакого значения. Полагая, что a', b' — строки a и b без последнего символа соответственно, цена преобразования из строки a в строку b может быть выражена как:

- а) сумма цены преобразования строки a в b и цены проведения операции удаления, которая необходима для преобразования a' в a ;
- б) сумма цены преобразования строки a в b и цены проведения операции вставки, которая необходима для преобразования b' в b ;
- в) сумма цены преобразования из a' в b' и операции замены, предполагая, что a и b оканчиваются разные символы;
- г) цена преобразования из a' в b' , предполагая, что a и b оканчиваются на один и тот же символ.

Минимальной ценой преобразования будет минимальное значение приведенных вариантов.

1.3 Матричный алгоритм нахождения расстояния Левенштейна

Прямая реализация формулы 1.1 может быть малоэффективна по времени исполнения при больших i, j , т. к. множество промежуточ-

ных значений $D(i, j)$ вычисляются заново множество раз подряд. Для оптимизации нахождения расстояния Левенштейна можно использовать матрицу в целях хранения соответствующих промежуточных значений. В таком случае алгоритм представляет собой построчное заполнение матрицы $A_{|a|, |b|}$ значениями $D(i, j)$.

1.4 Рекурсивный алгоритм нахождения расстояния Левенштейна с заполнением матрицы

Рекурсивный алгоритм заполнения можно оптимизировать по времени выполнения с использованием матричного алгоритма. Суть данного метода заключается в параллельном заполнении матрицы при выполнении рекурсии. В случае, если рекурсивный алгоритм выполняет прогон для данных, которые еще не были обработаны, результат нахождения расстояния заносится в матрицу. В случае, если обработанные ранее данные встречаются снова, для них расстояние не находится и алгоритм переходит к следующему шагу.

1.5 Расстояния Дамерау — Левенштейна

Расстояние Дамерау — Левенштейна может быть найдено по формуле 1.3, которая задана как

$$d_{a,b}(i, j) = \begin{cases} \max(i, j), & \text{если } \min(i, j) = 0, \\ \min\{ & \\ \quad d_{a,b}(i, j - 1) + 1, & \\ \quad d_{a,b}(i - 1, j) + 1, & \text{иначе} \\ \quad d_{a,b}(i - 1, j - 1) + m(a[i], b[j]), & \\ \quad \left[\begin{array}{ll} d_{a,b}(i - 2, j - 2) + 1, & \text{если } i, j > 1; \\ & a[i] = b[j - 1]; \\ & b[j] = a[i - 1] \\ & \infty, \quad \text{иначе} \end{array} \right. & \\ \} & \end{cases}, \quad (1.3)$$

Формула выводится по тем же соображениям, что и формула (1.1). Как и в случае с рекурсивным методом, прямое применение этой формулы неэффективно по времени исполнения, то аналогично методу из 1.4 производится добавление матрицы для хранения промежуточных значений рекурсивной формулы.

Вывод

В результате были рассмотрены и описаны математические модели следующих алгоритмов:

- Рекурсивный алгоритм нахождения расстояния Левенштейна
- Матричный алгоритм нахождения расстояния Левенштейна
- Рекурсивный алгоритм нахождения расстояния Дameraу - Левенштейна
- Матричный алгоритм нахождения расстояния Дameraу - Левенштейна

В качестве входных данных алгоритмы принимают 2 строки, по которым рассчитывается искомое расстояние. Выходными данными алгоритма является число - найденное расстояние.

Ограничением для реализуемой программы является обработка строк только в кодах `ascii.c`

2 Конструкторская часть

В данном разделе представлены схемы рассматриваемых алгоритмов и их оценка по памяти.

2.1 Схемы алгоритмов

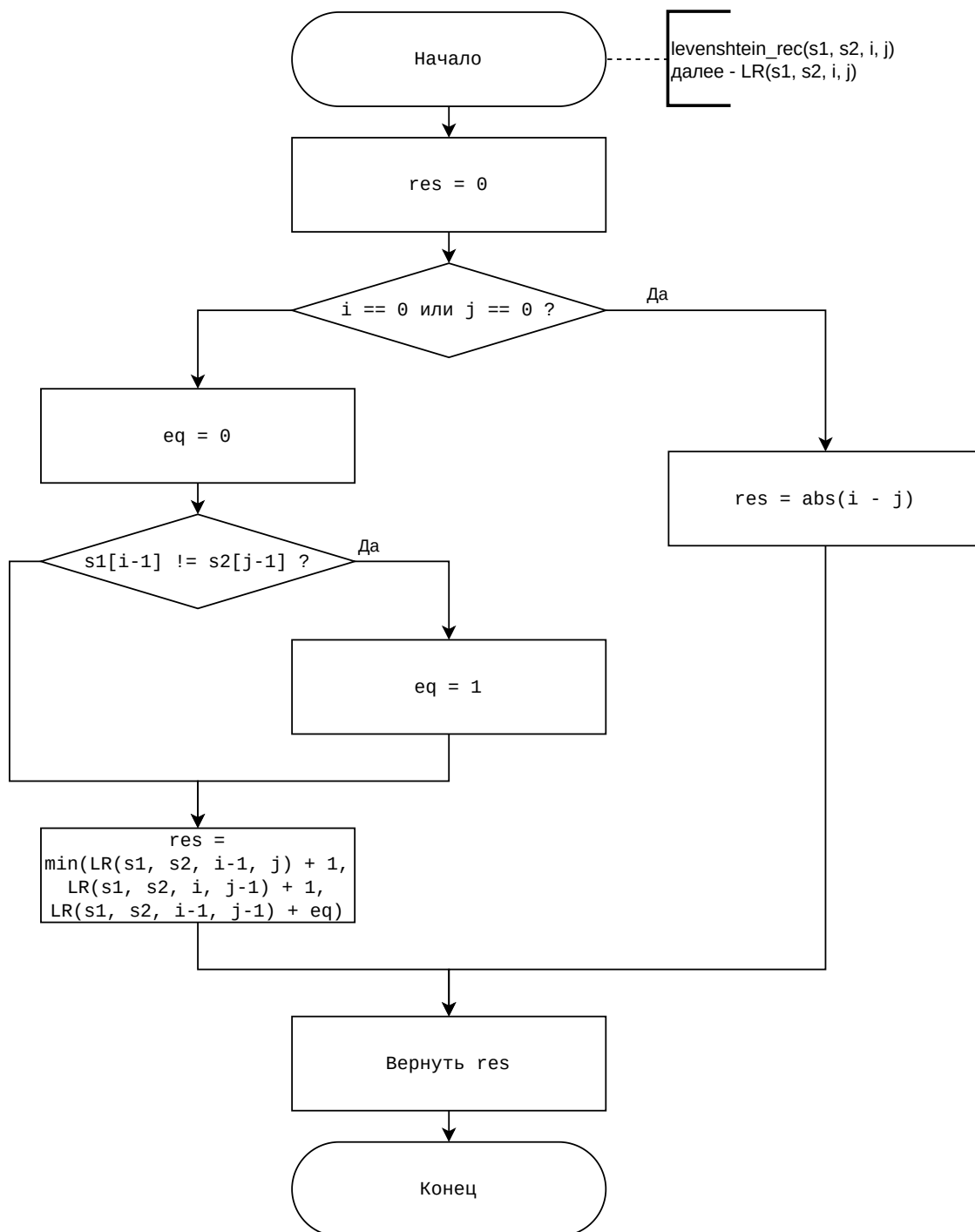


Рисунок 2.1 — Рекурсивный алгоритм Левенштейна

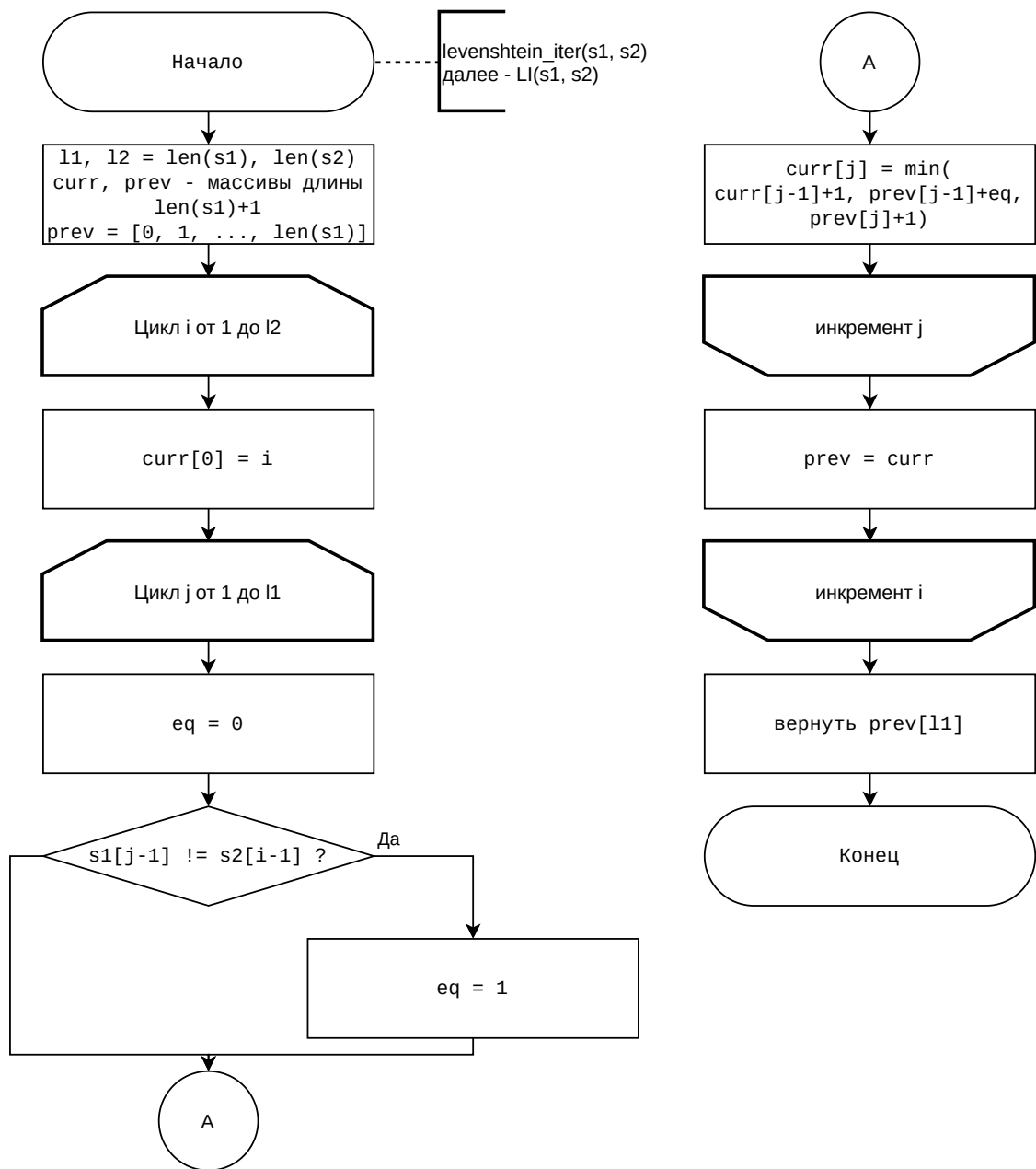


Рисунок 2.2 — Рекурсивный алгоритм Левенштейна

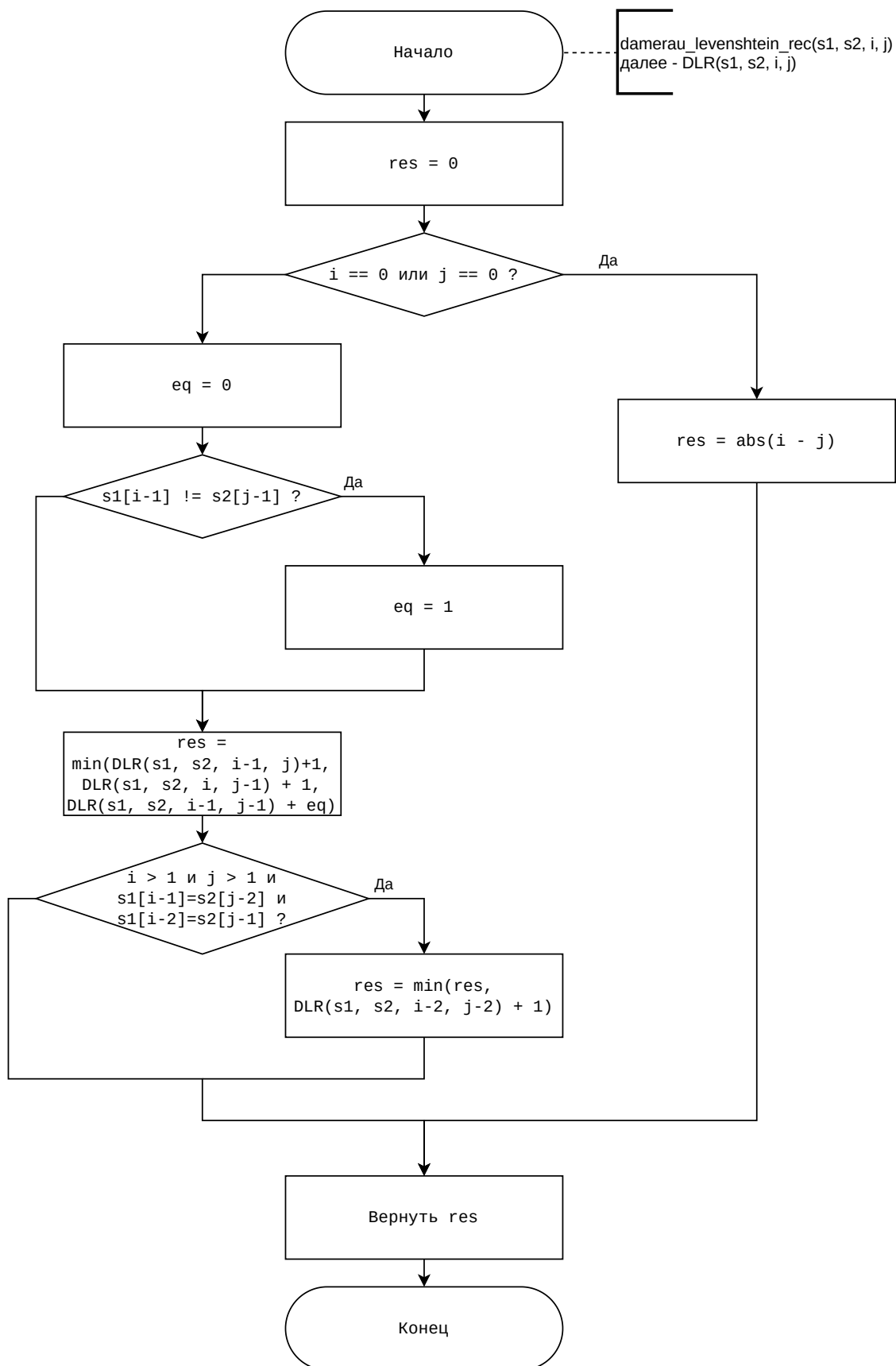


Рисунок 2.3 — Итеративный алгоритм Дамерау - Левенштейна

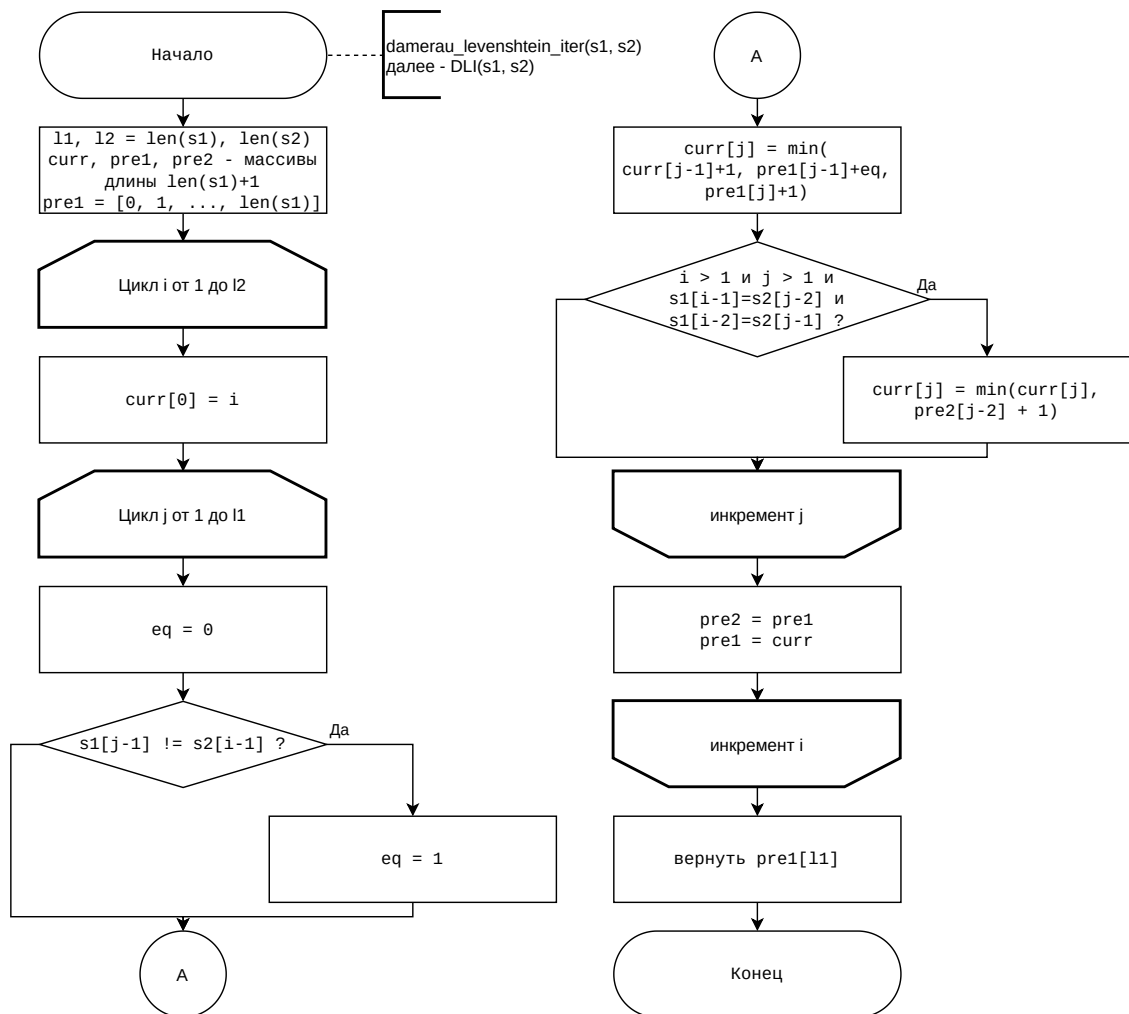


Рисунок 2.4 — Рекурсивный алгоритм Дамерау - Левенштейна

2.2 Оценка памяти

Алгоритмы Левенштейна и Дамерау — Левенштейна не отличаются друг от друга с точки зрения использования памяти, следовательно, достаточно рассмотреть лишь разницу рекурсивной и матричной реализаций этих алгоритмов.

Максимальная глубина стека вызовов при рекурсивной реализации равна сумме длин входящих строк, соответственно, максимальный расход памяти (2.1)

$$(\mathcal{C}(S_1) + \mathcal{C}(S_2)) \cdot (2 \cdot \mathcal{C}(\text{string}) + 3 \cdot \mathcal{C}(\text{int})), \quad (2.1)$$

где \mathcal{C} — оператор вычисления размера, S_1, S_2 — строки, int — целочисленный тип, string — строковый тип.

Использование памяти при итеративной реализации теоретически равно

$$(\mathcal{C}(S_1) + 1) \cdot (\mathcal{C}(S_2) + 1) \cdot \mathcal{C}(\text{int}) + 10 \cdot \mathcal{C}(\text{int}) + 2 \cdot \mathcal{C}(\text{string}). \quad (2.2)$$

2.3 Структуры данных

В качестве входных данных алгоритмы принимают 2 строки, по которым рассчитывается искомое расстояние. Выходными данными алгоритма является число - найденное расстояние.

В качестве строк для реализации алгоритма выберем нуль-терминированные строки. Никаких других структур данных не потребуется.

Вывод

На основе теоретических данных, полученных из аналитического раздела были построены схемы требуемых алгоритмов.

3 Технологическая часть

В данном разделе приведены требования к программному обеспечению, средства реализации и листинги кода.

3.1 Выбор средств реализации

Для реализации библиотеки с реализациями алгоритмов мною был выбран язык C[2] стандарта 1999 года. Он был выбран в силу отсутствия необходимости в использовании высокоуровневых абстракций.

Для реализации клиента был выбран язык C расширения GNU 99 года, так как он предоставляет функции чтения динамических строк.

Для юнит тестирования был выбран язык C++14, для использующая фрейворка GoogleTests[3].

3.2 Требования к программному обеспечению

К программе предъявляется ряд требований:

- На вход подаётся две регистрозависимые строки.
- На выходе — результат выполнения каждого из вышеуказанных алгоритмов.

3.3 Листинги кода

Листинг 3.1 — Основной файл программы main.c

```
1 #include <stdio.h>    // printf, getline
2 #include <stdlib.h>    // free
3 #include <string.h>    // strlen
4
5 #include "levenshtein/leven_dist.h"
6
7 int main(void) {
8     char *s1 = NULL, *s2 = NULL;
9     size_t l1 = 0, l2 = 0;
10
11     int rc = 0;
12
13     printf("Enter line 1: ");
14     if (-1 == (rc = getline(&s1, &l1, stdin))) {
15         printf("Allocation failed");
16         s1 = NULL;
17     } else {
18         printf("Enter line 2: ");
19         if (-1 == (rc = getline(&s2, &l2, stdin))) {
20             printf("Allocation failed");
21         } else {
22             s1[strlen(s1) - 1] = '\0';
23             s2[strlen(s2) - 1] = '\0';
24
25             int res_lr = levenshtein_rec(s1, s2);
26             printf("LR  ('%s', '%s') = %d\n", s1, s2, res_lr);
27
28             int res_li = levenshtein_iter(s1, s2);
29             printf("LI  ('%s', '%s') = %d\n", s1, s2, res_li);
30
31             int res_dlr = damerau_levenshtein_rec(s1, s2);
32             printf("DLR ('%s', '%s') = %d\n", s1, s2, res_dlr);
33
34             int res_dli = damerau_levenshtein_iter(s1, s2);
35             printf("DLI ('%s', '%s') = %d\n", s1, s2, res_dli);
36
37             free(s2);
38             s2 = NULL;
39         }
40         free(s1);
41         s1 = NULL;
42     }
43 }
```

Листинг 3.2 — Рекурсивная функция нахождения расстояния Левенштейна

```
1 static int levenshtein_rec_(char const *const s1, char const *const s2,  
2                             int const i, int const j) {  
3     int res = 0;  
4  
5     if (0 == i || 0 == j) {  
6         res = abs(i - j);  
7     } else {  
8         int eq = 0;  
9  
10        if (s1[i - 1] != s2[j - 1]) eq = 1;  
11  
12        res = min_int(levenshtein_rec_(s1, s2, i - 1, j) + 1,  
13                      min_int(levenshtein_rec_(s1, s2, i, j - 1) + 1,  
14                              levenshtein_rec_(s1, s2, i - 1, j - 1) + eq));  
15    }  
16  
17    return res;  
18 }  
19  
20 int levenshtein_rec(char const *const s1, char const *const s2) {  
21     return levenshtein_rec_(s1, s2, strlen(s1), strlen(s2));  
22 }
```

Листинг 3.3 — Рекурсивная функция нахождения расстояния Дамерау - Левенштейна

```
1 static int damerau_levenshtein_rec_(char const *const s1, char const *const s2,  
2                                     int const i, int const j) {  
3     int res = 0;  
4  
5     if (0 == i || 0 == j) {  
6         res = abs(i - j);  
7     } else {  
8         int eq = 0;  
9  
10        if (s1[i - 1] != s2[j - 1]) eq = 1;  
11  
12        res = min_int(damerau_levenshtein_rec_(s1, s2, i - 1, j) + 1,  
13                      min_int(damerau_levenshtein_rec_(s1, s2, i, j - 1) + 1,  
14                              damerau_levenshtein_rec_(s1, s2, i - 1, j - 1) + eq));  
15  
16        if (i > 1 && j > 1 && s1[i - 1] == s2[j - 2] && s1[i - 2] == s2[j - 1]) {  
17            res = min_int(res, damerau_levenshtein_rec_(s1, s2, i - 2, j - 2) + 1);  
18        }  
19    }  
20  
21    return res;  
22 }  
23  
24 int damerau_levenshtein_rec(char const *const s1, char const *const s2) {  
25     return damerau_levenshtein_rec_(s1, s2, strlen(s1), strlen(s2));  
26 }
```


Листинг 3.4 — Итеративная функция нахождения расстояния Левенштейна

```
1 int levenshtein_iter(char const *const s1, char const *const s2) {
2     size_t l1 = strlen(s1);
3     size_t l2 = strlen(s2);
4
5     int *prev = malloc((l1 + 1) * sizeof(int));
6     int *curr = malloc((l1 + 1) * sizeof(int));
7
8     if (NULL == prev || NULL == curr) {
9         if (NULL != prev) free(prev);
10        if (NULL != curr) free(curr);
11
12        return -1;
13    }
14
15    for (size_t i = 0; i <= l1; ++i) prev[i] = i;
16
17    for (size_t i = 1; i <= l2; ++i) {
18        curr[0] = i;
19
20        for (size_t j = 1; j <= l1; ++j) {
21            int eq = 0;
22
23            if (s1[j - 1] != s2[i - 1]) eq = 1;
24
25            curr[j] =
26                min_int(curr[j - 1] + 1, min_int(prev[j - 1] + eq, prev[j] + 1));
27        }
28
29        swp_ptr(&prev, &curr);
30    }
31
32    int res = prev[l1];
33
34    free(curr);
35    free(prev);
36
37    return res;
38 }
```

Листинг 3.5 — Итеративная функция нахождения расстояния Дамерау - Левенштейна

```
1 int damerau_levenshtein_iter(char const *const s1, char const *const s2) {
2     size_t l1 = strlen(s1);
3     size_t l2 = strlen(s2);
4
5     int *pre2 = malloc((l1 + 1) * sizeof(int));
6     int *pre1 = malloc((l1 + 1) * sizeof(int));
7     int *curr = malloc((l1 + 1) * sizeof(int));
8
9     if (NULL == pre2 || NULL == pre1 || NULL == curr) {
10         if (NULL != pre2) free(pre2);
11         if (NULL != pre1) free(pre1);
12         if (NULL != curr) free(curr);
13
14         return -1;
15     }
16
17     for (size_t i = 0; i <= l1; ++i) pre1[i] = i;
18
19     for (size_t i = 1; i <= l2; ++i) {
20         curr[0] = i;
21
22         for (size_t j = 1; j <= l1; ++j) {
23             int eq = 0;
24
25             if (s1[j - 1] != s2[i - 1]) eq = 1;
26
27             curr[j] =
28                 min_int(curr[j - 1] + 1, min_int(pre1[j - 1] + eq, pre1[j] + 1));
29
30             if (i > 1 && j > 1 && s1[i - 1] == s2[j - 2] && s1[i - 2] == s2[j - 1]) {
31                 curr[j] = min_int(curr[j], pre2[j - 2] + 1);
32             }
33         }
34
35         swp_ptr(&pre2, &pre1);
36         swp_ptr(&pre1, &curr);
37     }
38
39     int res = pre1[l1];
40
41     free(curr);
42     free(pre1);
43     free(pre2);
44
45     return res;
46 }
```

Листинг 3.6 — Вспомогательные функции

```
1 static inline int min_int(const int l, const int r) { return l <= r ? l : r; }
2
3 static inline void swp_pint(int **p1, int **p2) {
4     void *bucket = *p1;
5     *p1 = *p2;
6     *p2 = bucket;
7 }
```

Листинг 3.7 — Юнит тесты для алгоритма поиска расстояния Левенштейна. Часть 1

```
1 #include <cstring>
2 #include <string>
3
4 #include "gtest/gtest.h"
5
6 extern "C" {
7 #include "levenshtein/leven_dist.h"
8 };
9
10 TEST(leven_iter_pos, default) {
11     std::string s1{"CONNECT"};
12     std::string s2{"CONEHEAD"};
13
14     int res = levenshtein_iter(s1.c_str(), s2.c_str());
15
16     ASSERT_EQ(4, res);
17 }
18
19 TEST(leven_iter_pos, empty_one) {
20     std::string s1{""};
21     std::string s2{"CONEHEAD"};
22
23     int res = levenshtein_iter(s1.c_str(), s2.c_str());
24
25     ASSERT_EQ(8, res);
26 }
27
28 TEST(leven_iter_pos, empty_all) {
29     std::string s1{""};
30     std::string s2{""};
31
32     int res = levenshtein_iter(s1.c_str(), s2.c_str());
33
34     ASSERT_EQ(0, res);
35 }
36
37 TEST(leven_iter_pos, same) {
38     std::string s1{"deniska"};
39     std::string s2{"deniska"};
40
41     int res = levenshtein_iter(s1.c_str(), s2.c_str());
42
43     ASSERT_EQ(0, res);
44 }
45
46 TEST(leven_iter_pos, add_first) {
47     std::string s1{"abc"};
48     std::string s2{"bc"};
49
50     int res = levenshtein_iter(s1.c_str(), s2.c_str());
51
52     ASSERT_EQ(1, res);
53 }
```

Листинг 3.8 — Юнит тесты для алгоритма поиска расстояния Левенштейна. Часть 2

```
1 TEST(leven_iter_pos, add_mid) {
2     std::string s1{"abc"};
3     std::string s2{"ac"};
4
5     int res = levenshtein_iter(s1.c_str(), s2.c_str());
6
7     ASSERT_EQ(1, res);
8 }
9
10 TEST(leven_iter_pos, add_last) {
11     std::string s1{"abc"};
12     std::string s2{"ab"};
13
14     int res = levenshtein_iter(s1.c_str(), s2.c_str());
15
16     ASSERT_EQ(1, res);
17 }
18
19 TEST(leven_iter_pos, replace_first) {
20     std::string s1{"abc"};
21     std::string s2{"xbc"};
22
23     int res = levenshtein_iter(s1.c_str(), s2.c_str());
24
25     ASSERT_EQ(1, res);
26 }
27
28 TEST(leven_iter_pos, replace_mid) {
29     std::string s1{"abc"};
30     std::string s2{"axc"};
31
32     int res = levenshtein_iter(s1.c_str(), s2.c_str());
33
34     ASSERT_EQ(1, res);
35 }
36
37 TEST(leven_iter_pos, replace_last) {
38     std::string s1{"abc"};
39     std::string s2{"abx"};
40
41     int res = levenshtein_iter(s1.c_str(), s2.c_str());
42
43     ASSERT_EQ(1, res);
44 }
```

Вывод

Были реализованы алгоритмы: вычисления расстояния Левенштейна рекурсивно, с заполнением кэша, а также вычисления расстояния Дameraу – Левенштейна рекурсивно и вычисления расстояния Дameraу – Левенштейна с заполнением кэша.

4 Экспериментальная часть

В данном разделе будет проведено функциональное тестирование разработанного программного обеспечения. Так же будет произведено измерение временных характеристик и характеристик по памяти каждого из реализованных алгоритмов.

4.1 Технические характеристики

Технические характеристики устройства, на котором выполнялось исследование:

- Процессор: Intel(R) Core(TM) i5-8250U CPU @ 1.60GHz [4]
- Оперативная память: 8 Gb
- Операционная система: elementary OS 6 Odin [5]

4.2 Тестирование

Юнит тестирование проводилось при помощи фреймворка GoogleTests. Были выполнены следующие тесты

Таблица 4.1 — Функциональные тесты

Строка 1	Строка 2	Левенштейн	Дамерау-Левенштейн
CONNECT	CONEHEAD	4	
	CONEHEAD	8	8
		0	0
deniska	deniska	0	0
abc	bc	1	1
abc	ac	1	1
abc	ab	1	1
abc	xbc	1	1
abc	axc	1	1
abc	abx	1	1
CONNECT	CNNETC		2
abcd	badc		2

При проведении функционального тестирования, полученные результаты работы программы совпали с ожидаемыми. Таким образом, функциональное тестирование пройдено успешно.

4.3 Временные характеристики

Результаты замеров по результатам экспериментов приведены в Таблице 4.2. В данной таблице для значений, для которых тестирование не выполнялось, в поле результата находится ” - ”.

Таблица 4.2 — Замер времени для строк разной длины

Длина строк	Время, сек			
	Левенштейн		Дамерау-Левенштейн	
	Рекурсивный	Итеративный	Рекурсивный	Итеративный
1	3e-08	5.3e-08	3.6e-08	6.6e-08
2	6.7e-05	3.1e-07	6.6e-05	3.8e-07
4	8.1e-05	3.2e-07	6.6e-05	3.8e-07
6	6.6e-04	4.2e-07	6.5e-05	3.8e-07
8	0.0019	5.2e-07	0.002	6.7e-07
16	-	2e-06	-	2.7e-06
32	-	8.8e-06	-	1.3e-05
64	-	3.5e-05	-	4.4e-05
128	-	0.00012	-	0.00016
256	-	0.00048	-	0.00064
512	-	0.0019	-	0.0025
1024	-	0.0075	-	0.01

Отдельно сравним итеративные алгоритмы поиска расстояний Левенштейна и Дамерау– Левенштейна. Сравнение будет производиться на основе данных, представленных в Таблице 4.2. Результат можно увидеть на Рисунке 4.1.

При длинах строк менее 64 символов разница по времени между итеративными реализациями незначительна, однако при увеличении длины строки алгоритм поиска расстояния Левенштейна оказывается быстрее вплоть до полутора раз. Это обосновывается тем, что у алгоритма поиска расстояния Дамерау-Левенштейна задействуется дополнительная операция, которая замедляет алгоритм.

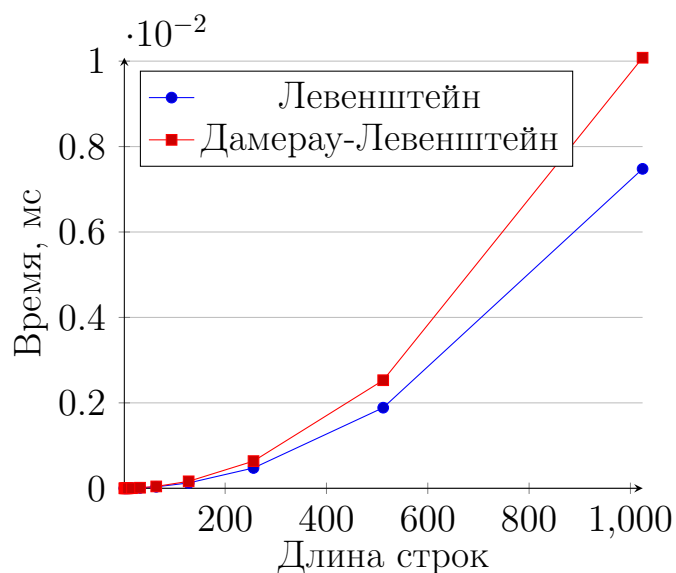


Рисунок 4.1 — Сравнение времени работы итеративных алгоритмов Левенштейна и Дамерау-Левенштейна

Так же сравним рекурсивную и итеративную реализации алгоритма поиска расстояния Левенштейна. Данные представлены в Таблице 4.2 и отображены на Рисунке 4.2.

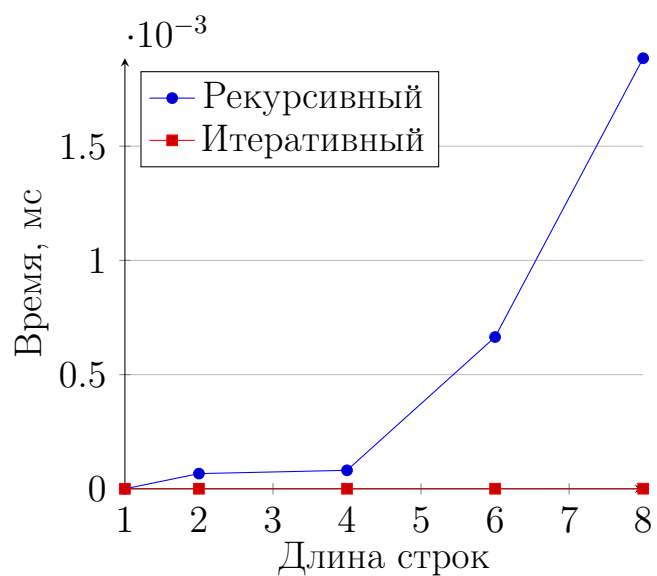


Рисунок 4.2 — Сравнение времени работы итеративного и рекурсивного алгоритмов Левенштейна

4.4 Сравнительный анализ алгоритмов

Приведенные характеристики показывают нам, что рекурсивная реализация алгоритма проигрывает по времени для всех тестовых данных. Например для длины строк длины 1 алгоритм работает в 1.7 раза дольше, для длины строк 4 - в 203 раза дольше, для длины строк 8 - 3500 раз дольше. Поэтому выбор рекурсивного алгоритма перед итеративным не оправдывает себя.

Так как во время печати очень часто возникают ошибки связанные с транспозицией букв, алгоритм поиска расстояния Дамерау – Левенштейна является наиболее предпочтительным, не смотря на то, что он проигрывает по времени и памяти алгоритму Левенштейна.

По аналогии с первым абзацем можно сделать вывод о том, что рекуррентный алгоритм поиска расстояния Дамерау - Левенштейна будет более затратным по времени по сравнению с итеративной реализацией алгоритма поиска расстояния Дамерау – Левенштейна с кешированием.

Вывод

В данном разделе было произведено сравнение количества затраченного времени и памяти вышеизложенных алгоритмов. Наиболее затратным по времени оказался рекурсивный алгоритм нахождения расстояния Дамерау – Левенштейна.

Рекурсивный алгоритм проигрывает по времени перед итеративным тем сильнее, чем больше длина строк. Для строк длины 1 разница незначительна: 1.7. Однако для строк большей длины, разность во времени отличается на порядки: для длины строк 4 - рекурсивный работает в 203 раза больше, а для длины строк 8 - в 3500 раз дольше. Однако, стоит учитывать дополнительные затраты по памяти, возникающие при использовании итеративных алгоритмов.

Заключение

Алгоритмы поиска расстояний Левенштейна и Дамерау-Левенштейна являются самыми популярными алгоритмами, которые помогают найти редакторское расстояние.

В результате выполнения данной лабораторной работы были изучены алгоритмы поиска расстояний Левенштейна 1.1 и Дамерау – Левенштейна 1.3, построены схемы (2.2, 2.4), соответствующие данным алгоритмам, также разобраны рекурсивные алгоритмы (2.1, 2.3). Было проведено сравнение алгоритмов по памяти и по времени. Была написана программа, реализующая данных алгоритмы.

Рекурсивная реализация алгоритма существенного проигрывает перед итеративным вариантом по времени. И эта разница тем сильнее, чем больше длина строк. Для строк длины 1 рекурсивный алгоритм работает в 1.7 раз дольше, для строк длины 4 - в 203 раза дальше, для строк длины 8 - в 3500 раз дольше.

Алгоритм Левенштейна работает быстрее алгоритма в Дамерау - Левенштейна в среднем в 1.3 раза. Однако ошибка перестановки букв в строке так часта, что компенсирует данный проигрыш по времени, поэтому стоит предпочитать итеративный алгоритм поиска расстояния Дамерау - Левенштейна перед другими.

Список литературы

1. *Левенштейн, Владимир Иосифович.* Двоичные коды с исправлением выпадений, вставок и замещений символов / Владимир Иосифович Левенштейн // Доклады Академии наук / Российская академия наук. — Vol. 163. — 1965. — Рр. 845–848.
2. Стандарт С [Электронный ресурс]. — Режим доступа: <http://www.open-std.org/jtc1/sc22/wg14/www/docs/n1256.pdf> (дата обращения: 15.09.2021).
3. Фреймворк Google tests. — Режим доступа: <https://google.github.io/googletest/> (дата обращения: 11.10.2021).
4. Процессор Intel® Core™ i5-8250U. — Режим доступа: <https://ark.intel.com/content/www/ru/ru/ark/products/124967/intel-core-i5-8250u-processor-6m-cache-up-to-3-40-ghz.html> (дата обращения: 11.10.2021).
5. Операционная система elementary os. — Режим доступа: <https://elementary.io/> (дата обращения: 11.10.2021).