



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 «Расстояние Левенштейна и
Дамерау — Левенштейна»

Студент Завойских Евгения Васильевна

Группа ИУ7-53Б

Оценка (баллы) _____

Преподаватели Волкова Л.Л., Строганов Ю.В.

2022 г.

Оглавление

Введение	3
1 Аналитическая часть	4
1.1 Расстояние Левенштейна	4
1.2 Расстояние Дамерау — Левенштейна	4
2 Конструкторская часть	6
2.1 Разработка алгоритмов	6
3 Технологическая часть	13
3.1 Средства реализации	13
3.2 Реализации алгоритмов	13
3.3 Описание тестирования	18
4 Исследовательская часть	19
4.1 Технические характеристики	19
4.2 Постановка эксперимента по замеру времени	19
4.3 Результаты эксперимента	20
4.4 Расчет используемой памяти	22
Заключение	25
Список использованных источников	26

Введение

Расстояние Левенштейна, или редакционное расстояние, – это минимальное количество операций замены, вставки и удаления символа, которое нужно выполнить над одной последовательностью символов, чтобы получить другую. Определение расстояния Дameraу — Левенштейна получается добавлением к указанным действиям операции перестановки соседних символов [1].

Расстояние Левенштейна используется для исправления ошибок в словах, поиска дубликатов текстов, сравнения геномов и прочих операций с символьными последовательностями.

Целью данной работы является изучение и исследование особенностей задач динамического программирования на материале алгоритмов нахождения расстояний Левенштейна и Дameraу — Левенштейна.

Для достижения поставленной цели необходимо выполнить следующие задачи:

- 1) изучить расстояния Левенштейна и Дameraу — Левенштейна;
- 2) разработать алгоритмы нерекурсивного метода поиска расстояния Левенштейна, нерекурсивного метода поиска расстояния Дameraу — Левенштейна, рекурсивного метода поиска расстояния Дameraу — Левенштейна, рекурсивного с кешированием метода поиска расстояния Дameraу — Левенштейна;
- 3) реализовать разработанные алгоритмы;
- 4) выбрать инструменты для замера процессорного времени выполнения реализаций алгоритмов;
- 5) сравнить алгоритмы по процессорному времени работы реализаций;
- 6) сравнить алгоритмы по используемой памяти.

1 Аналитическая часть

1.1 Расстояние Левенштейна

Расстояние Левенштейна – это минимальное количество операций вставки, удаления и замены символа, необходимое для преобразования одной символьной последовательности в другую.

Для решения задачи поиска расстояния Левенштейна нужно найти последовательность операций с минимальной суммарной ценой. Цена каждой операции:

- $w(a, b) = 1$ – для замены одного символа на другой;
- $w(a, \lambda) = 1$ – для удаления одного символа;
- $w(\lambda, b) = 1$ – для вставки одного символа;
- $w(a, a) = 0$ – для двух совпавших символов.

Пусть s_1 и s_2 – две строки над некоторым алфавитом, тогда расстояние Левенштейна можно подсчитать по следующей рекуррентной формуле (1.1)

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & j = 0, i > 0 \\ j, & i = 0, j > 0 \\ \min\{ \\ & D(i, j - 1) + 1, \\ & D(i - 1, j) + 1, \\ & D(i - 1, j - 1) + \\ & \begin{cases} 0, & \text{если } s_1[i] = s_2[j] \\ 1, & \text{иначе} \end{cases} \\ \}, & i > 0, j > 0 \end{cases} \quad (1.1)$$

1.2 Расстояние Дамерау — Левенштейна

Расстояние Дамерау — Левенштейна – это минимальное количество операций вставки, удаления, замены и перестановки соседних символов, необхо-

димое для преобразования одной символьной последовательности в другую.

Алгоритм поиска расстояния Дамерау — Левенштейна аналогичен алгоритму поиска расстояния Левенштейна. Дополнительно необходимо учесть операцию перестановки соседних символов, цена которой равна 1. Тогда расстояние Дамерау — Левенштейна можно подсчитать по формуле (1.2)

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & j = 0, i > 0 \\ j, & i = 0, j > 0 \\ \min\{ \\ \quad D(i, j - 1) + 1, \\ \quad D(i - 1, j) + 1, \\ \quad D(i - 1, j - 1) + \\ \quad \quad \begin{cases} 0, & \text{если } s_1[i] = s_2[j] \\ 1, & \text{иначе} \end{cases} \\ \quad , \\ \quad \begin{cases} D(i - 2, j - 2) + 1, \\ \text{если } i, j > 1, s_1[i] = s_2[j - 1], s_2[j] = s_1[i - 1] \\ \infty, & \text{иначе} \end{cases} \\ \}, & i > 0, j > 0 \end{cases} \quad (1.2)$$

Вывод из аналитической части

В данном разделе были рассмотрены алгоритмы вычисления расстояний Левенштейна и Дамерау — Левенштейна. Формулы вычисления расстояний являются рекуррентными, поэтому каждый из алгоритмов может быть реализован как с использованием рекурсии, так и без.

2 Конструкторская часть

В данном разделе разрабатываются следующие алгоритмы:

- нерекурсивный алгоритм поиска расстояния Левенштейна;
- нерекурсивный алгоритм поиска расстояния Дамерау — Левенштейна;
- рекурсивный алгоритм поиска расстояния Дамерау — Левенштейна без кеширования;
- рекурсивный алгоритм поиска расстояния Дамерау — Левенштейна с кешированием.

В нерекурсивных алгоритмах для сохранения промежуточных результатов вычислений используется матрица, которая построчно заполняется в ходе работы алгоритма. Матрица так же используется в рекурсивном алгоритме поиска с кешированием.

2.1 Разработка алгоритмов

В данном подразделе представлены схемы алгоритмов:

- 1) на рисунке 2.1 – схема нерекурсивного алгоритма поиска расстояния Левенштейна;
- 2) на рисунке 2.2 – схема нерекурсивного алгоритма поиска расстояния Дамерау — Левенштейна;
- 3) на рисунке 2.3 – схема рекурсивного алгоритма поиска расстояния Дамерау — Левенштейна;
- 4) на рисунке 2.4 – схема рекурсивной подпрограммы алгоритма поиска расстояния Дамерау — Левенштейна;
- 5) на рисунке 2.5 – схема рекурсивного алгоритма поиска расстояния Дамерау — Левенштейна с кешированием;
- 6) на рисунке 2.6 – схема рекурсивной подпрограммы алгоритма поиска расстояния Дамерау — Левенштейна с кешированием.

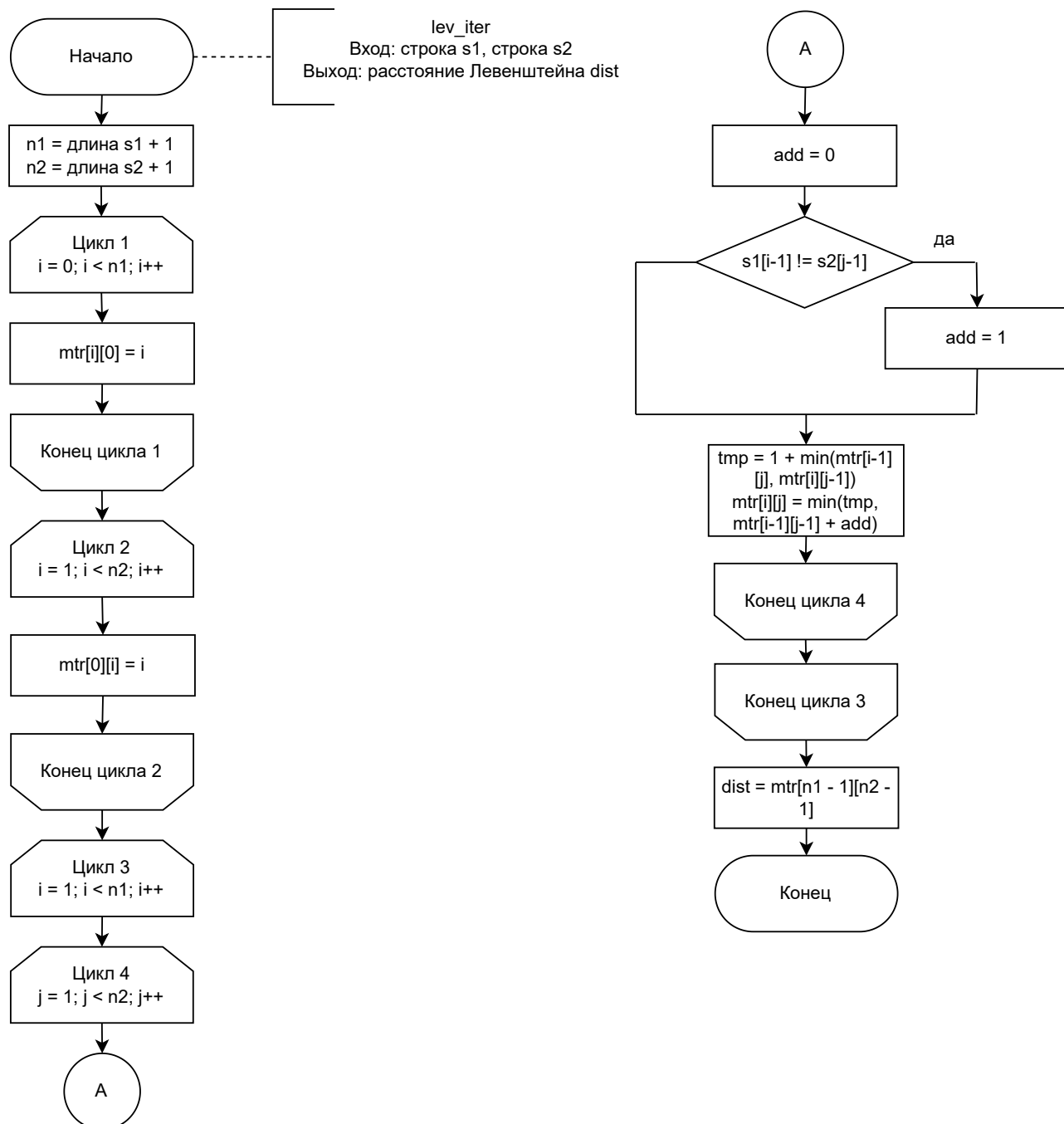


Рисунок 2.1 – Схема нерекурсивного алгоритма поиска расстояния Левенштейна

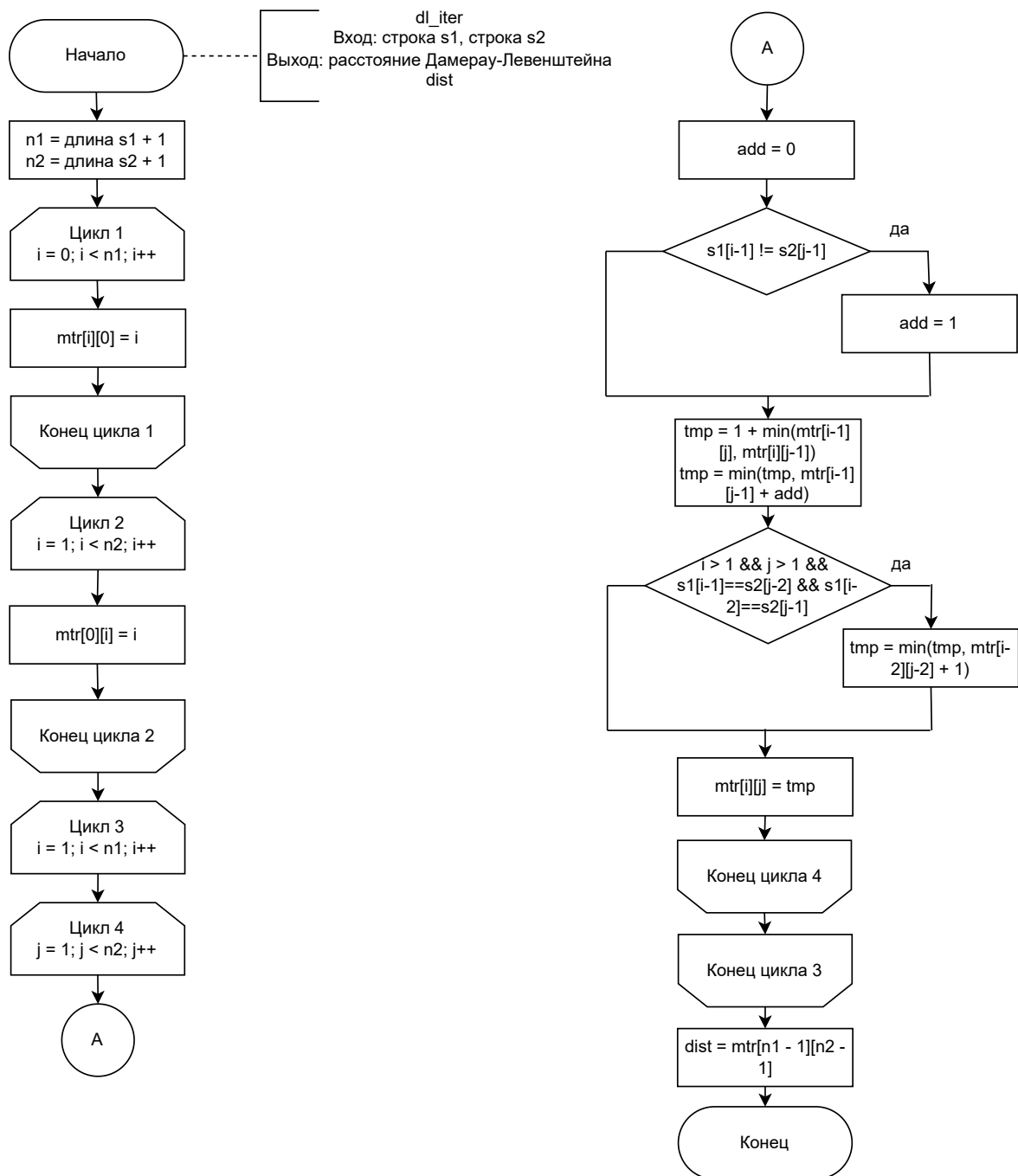


Рисунок 2.2 – Схема нерекурсивного алгоритма поиска расстояния Дамерау — Левенштейна

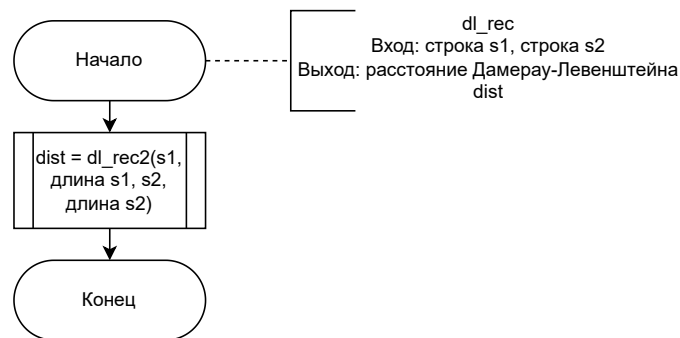


Рисунок 2.3 – Схема рекурсивного алгоритма поиска расстояния Дameraу — Левенштейна

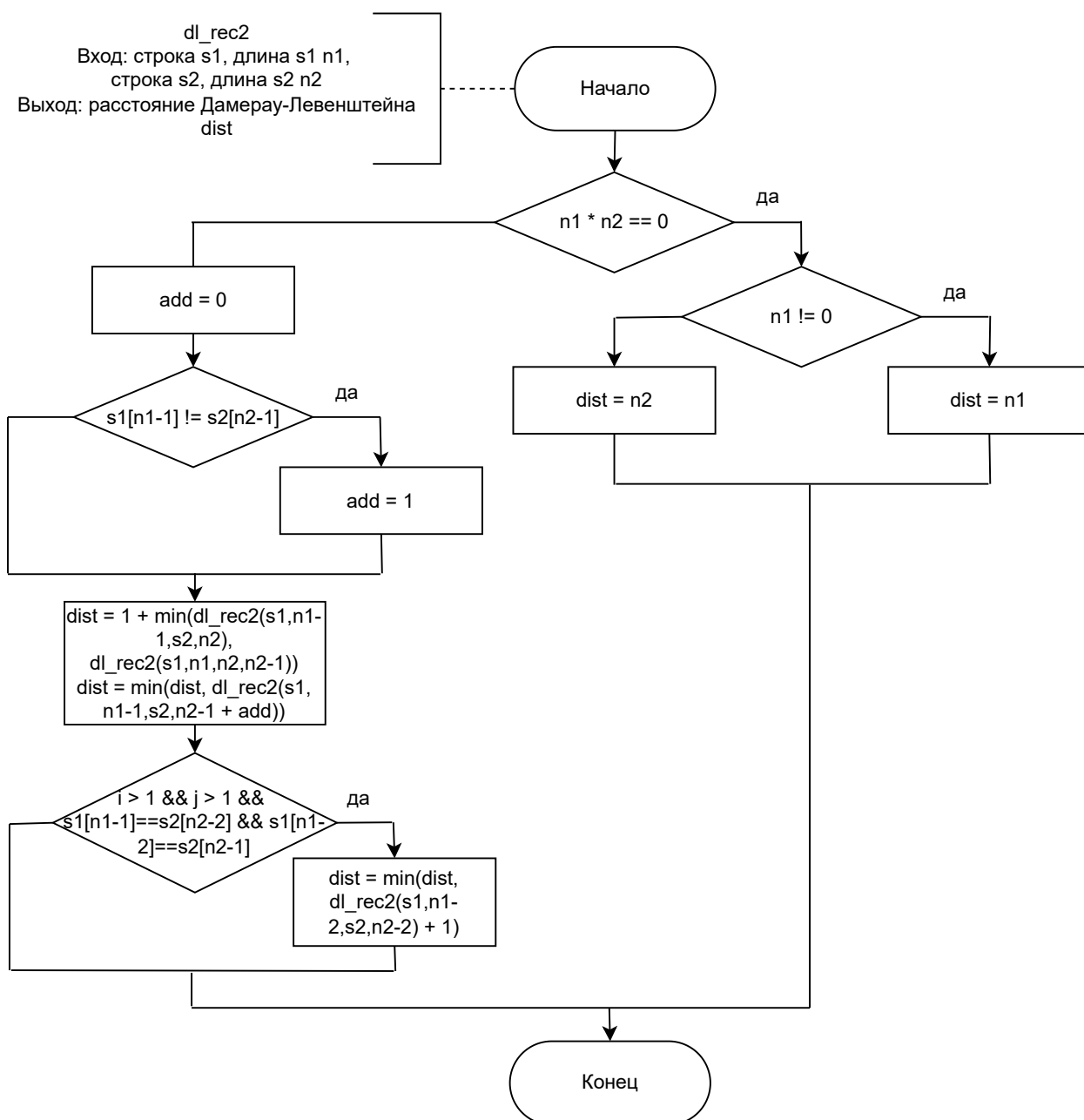


Рисунок 2.4 – Схема рекурсивной подпрограммы алгоритма поиска расстояния Дameraу — Левенштейна

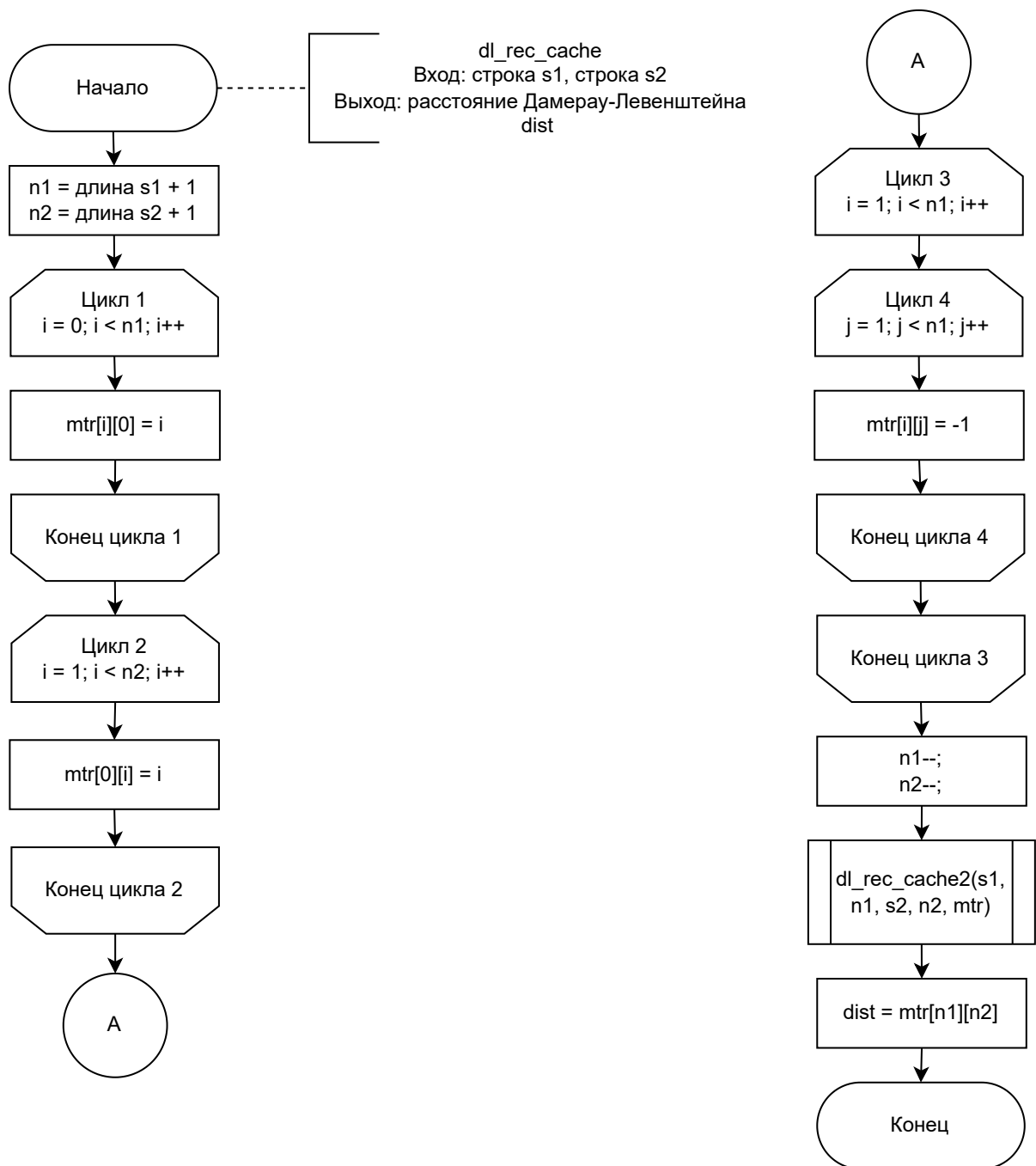


Рисунок 2.5 – Схема рекурсивного алгоритма поиска расстояния Дameraу — Левенштейна с кешированием

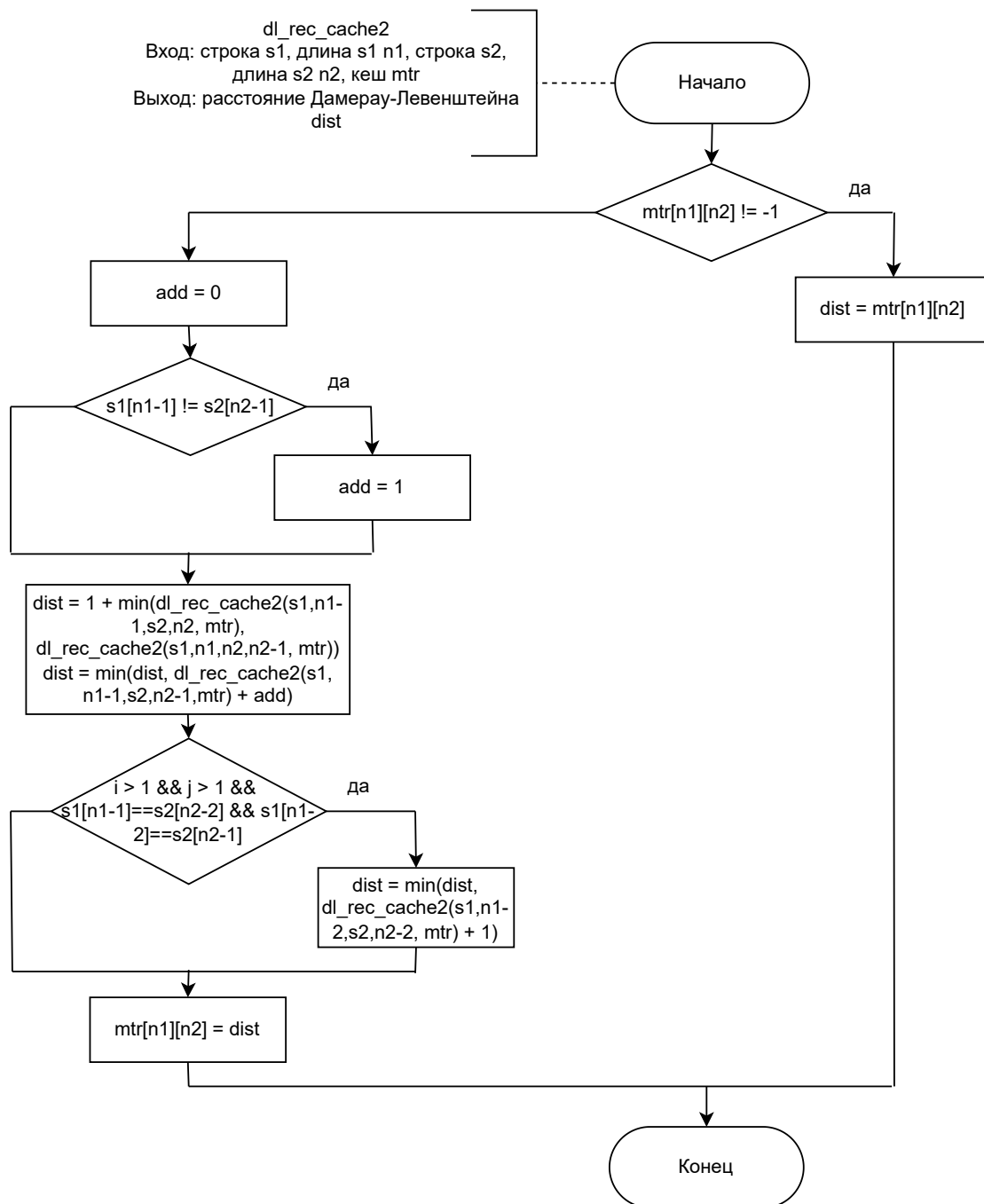


Рисунок 2.6 – Схема рекурсивной подпрограммы алгоритма поиска расстояния Дамерау — Левенштейна с кешированием

Вывод из конструкторской части

В данном разделе были разработаны следующие алгоритмы: нерекурсивный алгоритм поиска расстояния Левентейна, нерекурсивный алгоритм поиска расстояния Дамерау — Левенштейна, рекурсивный алгоритм поиска расстояния Дамерау — Левенштейна с кешированием и без.

3 Технологическая часть

3.1 Средства реализации

Для реализации данной лабораторной работы был выбран язык C++ [2]. В его стандартной библиотеке содержится класс `wstring`, работающий и с латиницей, и с кириллицей [3]. В качестве среды разработки был выбран текстовый редактор VS Code. Замеры процессорного времени проводились при помощи функции `clock` из библиотеки `time.h` [4].

3.2 Реализации алгоритмов

В данном подразделе представлены листинги кода ранее описанных алгоритмов:

- алгоритм нерекурсивного метода поиска расстояния Левенштейна (листинг 3.1);
- алгоритм нерекурсивного метода поиска расстояния Дамерау — Левенштейна (листинг 3.2);
- алгоритм рекурсивного метода поиска расстояния Дамерау — Левенштейна (листинг 3.3);
- алгоритм рекурсивного метода поиска расстояния Дамерау — Левенштейна с кешированием (листинги 3.4).

Листинг 3.1 – Реализация алгоритма нерекурсивного метода поиска расстояния Левенштейна

```
1 int lev_iter(wstring &s1, wstring &s2)
2 {
3     int n1 = s1.length() + 1, n2 = s2.length() + 1;
4     int **mtr = new int * [n1];
5     int add, tmp, dist;
6
7     for (int i = 0; i < n1; i++)
8     {
9         mtr[i] = new int [n2];
10        mtr[i][0] = i;
11    }
12
13    for (int i = 1; i < n2; i++)
14        mtr[0][i] = i;
15
16    for (int i = 1; i < n1; i++)
17        for (int j = 1; j < n2; j++)
18        {
19            add = 0;
20            if (s1[i - 1] != s2[j - 1])
21                add = 1;
22
23            tmp = 1 + min(mtr[i - 1][j], mtr[i][j - 1]);
24            mtr[i][j] = min(tmp, mtr[i - 1][j - 1] + add);
25        }
26
27    dist = mtr[n1 - 1][n2 - 1];
28
29    for (int i = 0; i < n1; i++)
30        delete[] mtr[i];
31    delete[] mtr;
32
33    return dist;
34 }
```

Листинг 3.2 – Реализация алгоритма нерекурсивного метода поиска расстояния Дameraу — Левенштейна

```
1 int dl_iter(wstring &s1, wstring &s2)
2 {
3     int n1 = s1.length() + 1, n2 = s2.length() + 1;
4     int **mtr = new int * [n1];
5     int add, tmp, dist;
6
7     for (int i = 0; i < n1; i++)
8     {
9         mtr[i] = new int [n2];
10        mtr[i][0] = i;
11    }
12
13    for (int i = 1; i < n2; i++)
14        mtr[0][i] = i;
15
16    for (int i = 1; i < n1; i++)
17        for (int j = 1; j < n2; j++)
18        {
19            add = 0;
20            if (s1[i - 1] != s2[j - 1])
21                add = 1;
22
23            tmp = 1 + min(mtr[i - 1][j], mtr[i][j - 1]);
24            tmp = min(tmp, mtr[i - 1][j - 1] + add);
25
26            if (i > 1 && j > 1 && s1[i - 1] == s2[j - 2] && s1[i - 2] == s2[
                j - 1])
27                tmp = min(tmp, mtr[i - 2][j - 2] + 1);
28
29            mtr[i][j] = tmp;
30        }
31
32    dist = mtr[n1 - 1][n2 - 1];
33
34    for (int i = 0; i < n1; i++)
35        delete[] mtr[i];
36    delete[] mtr;
37
38    return dist;
39 }
```

Листинг 3.3 – Реализация алгоритма рекурсивного метода поиска расстояния
Дамерау — Левенштейна

```
1 int dl_rec2(wstring &s1, int n1, wstring &s2, int n2)
2 {
3     int dist;
4     if (n1 * n2 == 0)
5     {
6         if (n1 != 0)
7             dist = n1;
8         else
9             dist = n2;
10    }
11    else
12    {
13        int add = 0;
14        if (s1[n1 - 1] != s2[n2 - 1])
15            add = 1;
16
17        dist = 1 + min(dl_rec2(s1, n1 - 1, s2, n2), dl_rec2(s1, n1, s2, n2 -
18            1));
19        dist = min(dist, dl_rec2(s1, n1 - 1, s2, n2 - 1) + add);
20
21        if (n1 > 1 && n2 > 1 && s1[n1 - 1] == s2[n2 - 2] && s1[n1 - 2] == s2
22            [n2 - 1])
23            dist = min(dist, dl_rec2(s1, n1 - 2, s2, n2 - 2) + 1);
24    }
25    return dist;
26 }
27
28 int dl_rec(wstring &s1, wstring &s2)
29 {
30     return dl_rec2(s1, s1.length(), s2, s2.length());
31 }
```


Листинг 3.4 – Реализация алгоритма рекурсивного с кешированием метода
поиска расстояния Дамерау — Левенштейна

```
1 int dl_rec_cache2(wstring &s1, int n1, wstring &s2, int n2, int **mtr)
2 {
3     int dist;
4     if (mtr[n1][n2] != -1)
5         dist = mtr[n1][n2];
6     else
7     {
8         int add = 0;
9         if (s1[n1 - 1] != s2[n2 - 1])
10             add = 1;
11         dist = 1 + min(dl_rec_cache2(s1, n1 - 1, s2, n2, mtr), dl_rec_cache2
            (s1, n1, s2, n2 - 1, mtr));
12         dist = min(dist, dl_rec_cache2(s1, n1 - 1, s2, n2 - 1, mtr) + add);
13         if (n1 > 1 && n2 > 1 && s1[n1 - 1] == s2[n2 - 2] && s1[n1 - 2] == s2
            [n2 - 1])
14             dist = min(dist, dl_rec_cache2(s1, n1 - 2, s2, n2 - 2, mtr) + 1)
                ;
15         mtr[n1][n2] = dist;
16     }
17     return dist;
18 }
19 int dl_rec_cache(wstring &s1, wstring &s2)
20 {
21     int n1 = s1.length() + 1, n2 = s2.length() + 1;
22     int **mtr = new int * [n1];
23     for (int i = 0; i < n1; i++)
24     {
25         mtr[i] = new int [n2];
26         mtr[i][0] = i;
27     }
28     for (int i = 1; i < n2; i++)
29         mtr[0][i] = i;
30     for (int i = 1; i < n1; i++)
31         for (int j = 1; j < n2; j++)
32             mtr[i][j] = -1;
33     n1--;
34     n2--;
35     dl_rec_cache2(s1, n1, s2, n2, mtr);
36     int dist = mtr[n1][n2];
37
38     for (int i = 0; i <= n1; i++)
39         delete[] mtr[i];
40     delete[] mtr;
41     return dist;
42 }
```

3.3 Описание тестирования

В таблице 3.1 приведены тесты для алгоритмов поиска расстояний Левенштейна и Дамерау — Левенштейна.

Таблица 3.1 – Тесты

Слово №1	Слово №2	Ожидаемый результат	Результат
s	c	1 1 1 1	1 1 1 1
hell	hallo	2 2 2 2	2 2 2 2
corn	cron	2 1 1 1	2 1 1 1
honda	hyundai	3 3 3 3	3 3 3 3
sun	sun	0 0 0 0	0 0 0 0
qwer	wqre	3 2 2 2	3 2 2 2
йцук	цйку	3 2 2 2	3 2 2 2

Вывод из технологической части

В данном разделе был выбран инструмент для замера процессорного времени, были реализованы алгоритмы поиска расстояний Левенштейна и Дамерау — Левенштейна. Также было проведено тестирование реализованных алгоритмов.

4 Исследовательская часть

4.1 Технические характеристики

Технические характеристики устройства, на котором выполнялись замеры времени:

- операционная система Windows 11 64-bit;
- оперативная память 16 ГБ;
- процессор 2.40 ГГц Intel Core i5-1135G7 [5].

4.2 Постановка эксперимента по замеру времени

Для оценки процессорного времени работы реализаций алгоритмов поиска расстояний Левенштейна и Дамерау — Левенштейна был проведен эксперимент, в котором определялось влияние длины символьных последовательностей на время работы каждого из алгоритмов. Тестирование проводилось на словах длиной от 5 до 10 символов с шагом 1 и от 20 до 100 с шагом 10. На последовательностях длиной от 20 символов замеры для рекурсивного алгоритма без кеширования не выполнялись из-за быстрого увеличения времени работы. Поскольку методы замера процессорного времени имеют погрешность и возвращают для достаточно коротких задач константу 0, каждый алгоритм запускался по 500 раз, и для полученных 500 значений определялось среднее арифметическое, которое заносилось в таблицу результатов.

Результаты эксперимента были представлены в виде таблицы и графиков, приведенных в следующем подразделе.

4.3 Результаты эксперимента

Таблица 4.1 – Таблица времени работы реализаций алгоритмов (в мкс)

Длина слова	Лев.	Д.-Л. итер.	Д.-Л. рек.	Д.-Л. рек. кеш
5	0	0	0	0
6	0	0	60.80	0
7	0	0	319.94	0
8	0	0	1807.34	0
9	0	0	9831.58	0
10	0	12.57	55262.60	0
20	0	0	–	0
30	30.00	32.21	–	32.80
40	32.04	35.33	–	64.04
50	34.57	41.03	–	96.28
60	62.36	75.76	–	128.47
70	62.58	94.60	–	188.06
80	94.50	126.73	–	282.39
90	126.47	172.12	–	312.43
100	156.01	190.34	–	360.07

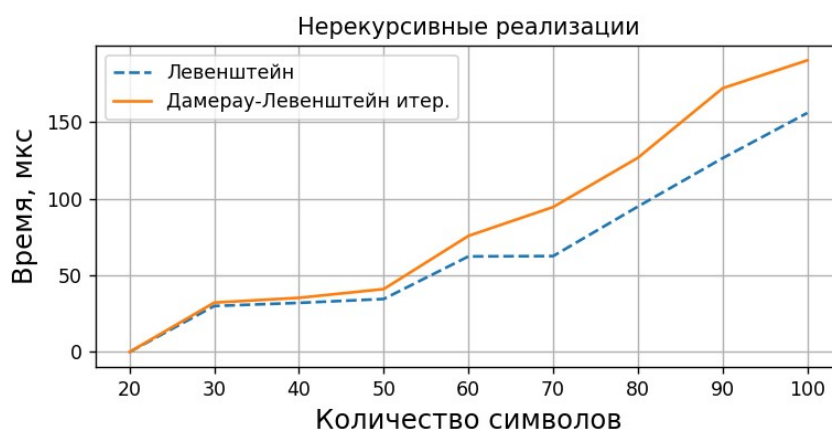


Рисунок 4.1 – Сравнение времени работы реализаций нерекурсивных алгоритмов поиска расстояний Левенштейна и Дамерау — Левенштейна

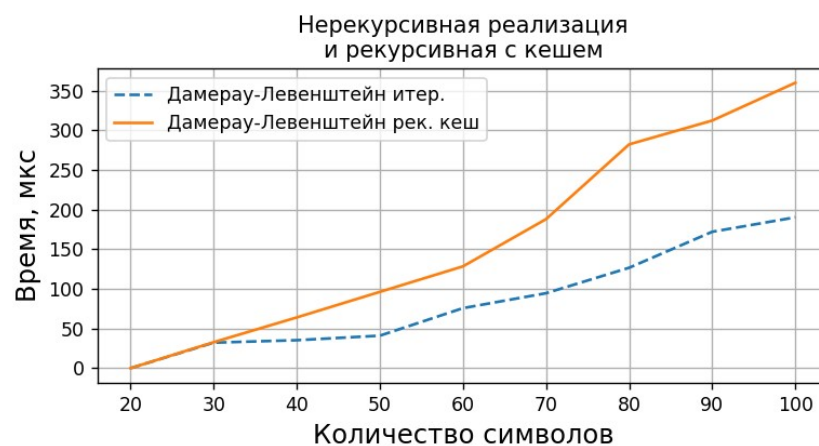


Рисунок 4.2 – Сравнение времени работы реализаций нерекурсивного алгоритма поиска расстояния Дамерау — Левенштейна и рекурсивного алгоритма поиска с кешем

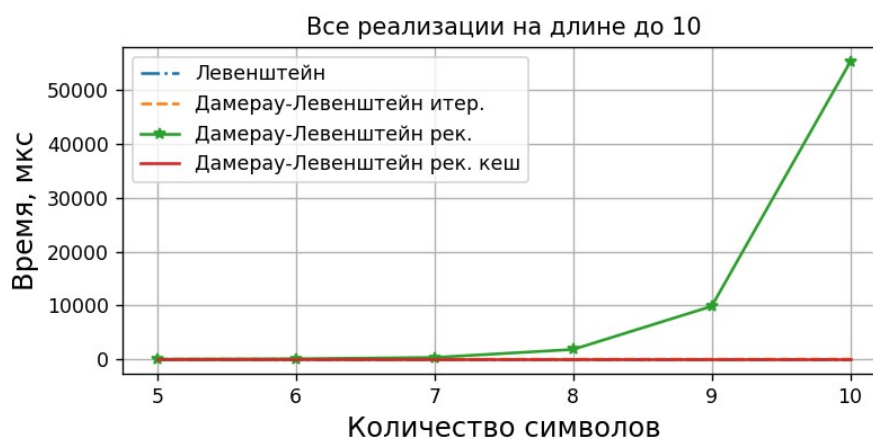


Рисунок 4.3 – Сравнение времени работы реализаций всех алгоритмов на последовательностях малой длины

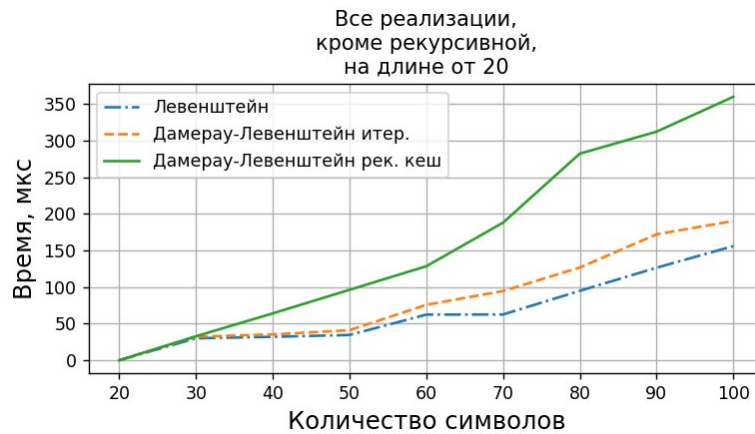


Рисунок 4.4 – Сравнение времени работы реализаций всех алгоритмов, кроме рекурсивного, на последовательностях средней и большой длины

Можно отметить, что на символьных последовательностях длиной до 30 символов реализация нерекурсивного алгоритма поиска расстояния Левенштейна и реализация нерекурсивного алгоритма поиска расстояния Дамерау — Левенштейна отрабатывают приблизительно за одинаковое время. Однако при увеличении длины последовательностей реализация алгоритма поиска расстояния Левенштейна становится более эффективной по времени. Она работает быстрее в 1.2 – 1.3 раза.

Нерекурсивная реализация алгоритма поиска расстояния Дамерау — Левенштейна и рекурсивная реализация с кешем работают в разы быстрее рекурсивной без кеша. Но на той же длине слов нерекурсивная реализация работает примерно в 1.9 раза быстрее рекурсивной реализации с кешем.

4.4 Расчет используемой памяти

Алгоритмы нахождения расстояний Левенштейна и Дамерау — Левенштейна не отличаются друг от друга с точки зрения используемой памяти.

Для каждого вызова рекурсивной реализации алгоритма Дамерау — Левенштейна выделяется память под:

- 2 строки;
- длины строк;
- локальная переменная;
- возвращаемое значение;
- адрес возврата.

Максимальное количество вызовов равно сумме длин двух строк.

Максимальное значение выделяемой памяти выражается формулой (4.1)

$$mem1 = (4 \cdot size(int) + 2 \cdot size(string)) \cdot (|s_1| + |s_2|) \quad (4.1)$$

Для единственного вызова нерекурсивной реализации алгоритма Дамерау — Левенштейна выделяется память под:

- 2 строки;
- длины строк;
- матрицу размерами, равными длинам строк, увеличинным на единицу;
- 4 локальных переменных;
- возвращаемое значение;
- адрес возврата.

Максимальное значение выделяемой памяти выражается формулой (4.2)

$$mem1 = 7 \cdot size(int) + 2 \cdot size(string) + size(int) \cdot (|s_1| + 1) \cdot (|s_2| + 1) \quad (4.2)$$

Память, используемая рекурсивной реализацией, растет пропорционально сумме длин строк, в то время как память, используемая нерекурсивной реализацией, растет пропорционально произведению длин строк, то есть по выделяемой памяти нерекурсивная реализация проигрывает рекурсивной. В рекурсивной реализации с кешем память так же растет пропорционально произведению длин строк.

Вывод из исследовательской части

По времени выполнения нерекурсивная реализация и реализация с кешем гораздо эффективнее рекурсивной без кеша. При этом нерекурсивная реализация работает в 1.9 раза быстрее рекурсивной реализации с кешем. Нерекурсивная реализация алгоритма поиска расстояния Левенштейна работает в 1.2 – 1.3 раза быстрее нерекурсивной реализации алгоритма нахождения расстояния Дамерау — Левенштейна.

По выделяемой памяти реализации, использующие матрицы, проигрывают рекурсивной без кеша: максимальный размер выделяемой памяти в них пропорционален произведению длин строк, в то время как в рекурсивной без

кеша – сумме длин строк.

Заключение

В результате выполнения лабораторной работы при исследовании алгоритмов нахождения расстояний Левенштейна и Дameraу — Левенштейна были изучены и отработаны навыки динамического программирования.

В ходе выполнения лабораторной работы:

1) были изучены расстояния Левенштейна и Дameraу — Левенштейна;
2) были разработаны алгоритмы нерекурсивного метода поиска расстояния Левенштейна, нерекурсивного метода поиска, рекурсивного метода поиска и рекурсивного с кешированием метода поиска расстояния Дameraу — Левенштейна;

3) был реализован каждый из описанных алгоритмов;

4) были выбраны инструменты для замера процессорного времени выполнения реализаций алгоритмов;

5) было проведено сравнение реализованных алгоритмов по времени работы: выявлено, что реализация рекурсивного алгоритма без кеширования уступает по времени всем другим реализациям, нерекурсивная реализация работает в 1.9 раза быстрее рекурсивной реализации с кешем, нерекурсивная реализация алгоритма поиска расстояния Левенштейна работает в 1.2 – 1.3 раза быстрее нерекурсивной реализации алгоритма нахождения расстояния Дameraу — Левенштейна;

6) было проведено сравнение реализованных алгоритмов по используемой памяти: выявлено, что реализации, использующие матрицу, уступают рекурсивной без кеша.

Таким образом, все поставленные задачи были выполнены, а цель достигнута.

Список использованных источников

- [1] Левенштейн В. И. Двоичные коды с исправлением выпадений, вставок и замещений символов. – М.: Издательство «Наука», Доклады АН СССР, 1965. Т. 163. С. 845–848.
- [2] Справочник по языку C++ – Microsoft Learn [Электронный ресурс]. Режим доступа: <https://learn.microsoft.com/ru-ru/cpp/cpp/cpp-language-reference?view=msvc-1701> (дата обращения: 18.11.2022).
- [3] std::wstring – C++ [Электронный ресурс]. Режим доступа: <https://cplusplus.com/reference/string/wstring/> (дата обращения: 15.11.2022).
- [4] clock – ctime – C++ [Электронный ресурс]. Режим доступа: <https://cplusplus.com/reference/ctime/clock/> (дата обращения: 15.11.2022).
- [5] Intel [Электронный ресурс]. Режим доступа: <https://ark.intel.com/content/www/ru/ru/ark/products/208658/intel-core-i51135g7-processor-8m-cache-up-to-4-20-ghz.html> (дата обращения: 18.11.2022).