



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»

Лабораторная работа №1 по дисциплине "Анализ Алгоритмов"

Тема Расстояние Левенштейна и Дамерау-Левенштейна

Студент Светличная А.А.

Группа ИУ7-53Б

Преподаватель Волкова Л. Л., Строганов Ю.В.

Москва — 2022 г.

Оглавление

Введение	3
1 Аналитическая часть	4
1.1 Цель и задачи	4
1.2 Итерационный алгоритм поиска расстояния Левенштейна . .	4
1.3 Итерационный алгоритм поиска расстояния Дамерау-Левенштейна	5
1.4 Рекурсивный алгоритм поиска расстояния Дамерау-Левенштейна	6
1.5 Рекурсивный алгоритм поиска расстояния Дамерау-Левенштейна с использованием кеша	7
2 Конструкторская часть	8
2.1 Описание типов данных	8
2.2 Оценка затрат алгоритмов по памяти	8
2.3 Описания алгоритмов	9
3 Технологическая часть	15
3.1 Требования к программному обеспечению	15
3.2 Выбор языка программирования	15
3.3 Выбор библиотеки и способа для замера времени	15
3.4 Реализации алгоритмов	16
3.5 Тестирование алгоритмов	20
4 Экспериментальная часть	21
4.1 Технические характеристики	21
4.2 Замеры времени	21
Список использованных источников	24

Введение

В программировании большое количество задач связано с обработкой строк. В данной лабораторной работе будут рассмотрены алгоритмы, использующие расстояние Левенштейна – метрика, позволяющая определить «схожесть» двух строк, вычисляя минимальное количество операций вставки, удаления, замены одного символа на другой, необходимых для превращения одной строки в другую.

Расстояние Левенштейна применяется в теории информации и компьютерной лингвистике для:

- 1) исправления ошибок в слове;
- 2) сравнения текстовых файлов утилитой diff;
- 3) для сравнения генов, хромосом и белков в биоинформатике.

1 Аналитическая часть

1.1 Цель и задачи

Целью данной работой является получение навыка динамического программирования на примере разработки алгоритмов поиска редакционных расстояний.

Для достижения поставленной цели требуется решить следующие задачи:

- 1) изучить расстояния Левенштейна и Дamerau-Левенштейна;
- 2) разработать указанные алгоритмы поиска расстояний;
- 3) реализовать разработанные алгоритмы;
- 4) провести сравнительный анализ реализаций алгоритмов по затраченному процессорному времени и памяти на основе экспериментальных данных;
- 5) описать и обосновать полученные результаты в отчете о выполненной лабораторной работе.

1.2 Итерационный алгоритм поиска расстояния Левенштейна

Расстояние Левенштейна [1] – это минимальное количество операций вставки, удаления и замены, необходимых для превращения одной строки в другую.

Цены операций могут зависеть от вида операций (вставка, удаление, замена) и/или от участвующих в ней символов, отражая разную вероятность разных ошибок при вводе текста и т.п. В общем случае

- 1) $w(a, b)$ – цена замены символа a на b , R (от англ. replace);
- 2) $w(\lambda, b)$ – цена вставки символа b , I (от англ. insert);

3) $w(a, \lambda)$ – цена удаления символа a , D (от англ. delete).

Для решения задачи о редакционном расстоянии необходимо найти последовательность замен, минимизирующую суммарную цену. Расстояние Левенштейна является частным случаем этой задачи при

1) $w(a, a) = 0$;

2) $w(a, b) = 1, a \neq b$;

3) $w(\lambda, b) = 1$;

4) $w(a, \lambda) = 1$.

Имеется две строки S_1 и S_2 , длиной m и n соответственно. Расстояние Левенштейна рассчитывается по рекуррентной формуле (1.1).

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

1.3 Итерационный алгоритм поиска расстояния Дамерау-Левенштейна

В алгоритме поиска расстояния Дамерау-Левенштейна, помимо вставки, удаления, и замены присутствует еще одна редакторская операция – транспозиция T (от англ. transposition).

Расстояние Дамерау-Левенштейна может быть вычисленно по рекуррентной формуле (1.2).

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0, \\ i, & j = 0, i > 0, \\ j, & i = 0, j > 0, \\ \min\{ \\ \quad D(i, j - 1) + 1, & i > 0, j > 0, \\ \quad D(i - 1, j) + 1, & S_i = S_{j-1}, \\ \quad D(i - 2, j - 2) + 1, & S_{i-1} = S_j, \\ \quad D(i - 1, j - 1) + \begin{cases} 0, & \text{если } S_i = S_j, \\ 1, & \text{иначе} \end{cases} \\ \}, \\ \min\{ \\ \quad D(i, j - 1) + 1, \\ \quad D(i - 1, j) + 1, & \text{иначе} \\ \quad D(i - 1, j - 1) + \begin{cases} 0, & \text{если } S_i = S_j, \\ 1, & \text{иначе} \end{cases} \\ \} \end{cases} \quad (1.2)$$

1.4 Рекурсивный алгоритм поиска расстояния Дамерау-Левенштейна

Рекурсивный алгоритм нахождения расстояния Дамерау-Левенштейна отличается от его итерационной версии тем, что в нем вместо использования матрицы для хранения предыдущих значений, необходимых для подсчета последующих, эти значения вычисляются каждый раз рекурсивно.

1.5 Рекурсивный алгоритм поиска расстояния Дамерау-Левенштейна с использованием кеша

Данный алгоритм является оптимизацией рекурсивного алгоритма нахождения расстояния Дамерау-Левенштейна. Оптимизация заключается в использовании кеша, который представляется собой матрицу, в которую записываются значения, вычисленные на шагах рекурсии. Таким образом, при необходимости вычисления какого-либо нового значения по рекуррентной формуле, величины, необходимые для этого не вычисляются каждый раз: сначала проверяется, было ли уже вычислено данное значения, и только лишь в случае, если этого не происходило, выполняются рекурсивные вычисления для его получения.

Вывод

В данном разделе были рассмотрены алгоритмы нахождения расстояния Левенштейна и Дамерау-Левенштейна, формулы которых задаются рекуррентно, а следовательно, данные алгоритмы могут быть реализованы рекурсивно и итерационно.

2 Конструкторская часть

2.1 Описание типов данных

Для реализации алгоритмов Левенштейна и Дамерау-Левенштейна были использованы следующие типы данных:

- 1) строка - массив типа *char*;
- 2) длина строки - число типа *int*;
- 3) матрица - структура, содержащая двумерный массив типа *int* и размерности типа *int*.

2.2 Оценка затрат алгоритмов по памяти

Алгоритмы Левенштейна и Дамерау-Левенштейна не отличаются по использованию памяти, поэтому достаточно рассмотреть итерационную и рекурсивную реализации одного алгоритма.

Затраты по памяти для итерационного алгоритма поиска расстояния Левенштейна (Дамерау-Левенштейна):

- 1) длины строк $n, m - 2 \cdot \text{sizeof}(\text{int})$;
- 2) матрица $-(n + 1) \cdot (m + 1) \cdot \text{sizeof}(\text{int}) + 2 \cdot \text{sizeof}(\text{int})$;
- 3) строки $-(n + m + 2) \cdot \text{sizeof}(\text{char})$;
- 4) дополнительные переменные ($i, j, \text{res}, \text{offset}$) $- 4 \cdot \text{sizeof}(\text{int})$;
- 5) адрес возврата.

Суммарные затраты по памяти:

$$(n + 1) \cdot (m + 1) \cdot \text{sizeof}(\text{int}) + 7 \cdot \text{sizeof}(\text{int}) + (n + m + 2) \cdot \text{sizeof}(\text{char})$$

Затраты по памяти для рекурсивного алгоритма поиска расстояния Дамерау-Левенштейна для одного вызова:

- 1) длины строк $n, m - 2 \cdot \text{sizeof}(\text{int})$;
- 2) строки $-(n + m + 2) \cdot \text{sizeof}(\text{char})$;
- 3) дополнительные переменные $(\text{res}, \text{offset}) - 2 \cdot \text{sizeof}(\text{int})$;
- 4) адрес возврата.

Суммарные затраты по памяти (R – количество вызовов рекурсии):
 $(4 \cdot \text{sizeof}(\text{int}) + (n + m + 2) \cdot \text{sizeof}(\text{char})) \cdot R$

Затраты по памяти для рекурсивного алгоритма поиска расстояния Дамерау-Левенштейна с кешем для одного вызова:

- 1) длины строк $n, m - 2 \cdot \text{sizeof}(\text{int})$;
- 2) строки $-(n + m + 2) \cdot \text{sizeof}(\text{char})$;
- 3) дополнительные переменные $(\text{res}, \text{offset}) - 2 \cdot \text{sizeof}(\text{int})$;
- 4) матрица $-(n + 1) \cdot (m + 1) \cdot \text{sizeof}(\text{int}) + 2 \cdot \text{sizeof}(\text{int})$;
- 5) адрес возврата.

Суммарные затраты по памяти (R – количество вызовов рекурсии):
 $(6 \cdot \text{sizeof}(\text{int}) + (n + m + 2) \cdot \text{sizeof}(\text{char})) \cdot R + (n + 1) \cdot (m + 1) \cdot \text{sizeof}(\text{int})$

2.3 Описания алгоритмов

На рисунках ниже показаны схемы алгоритмов Левенштейна и Дамерау-Левенштейна.

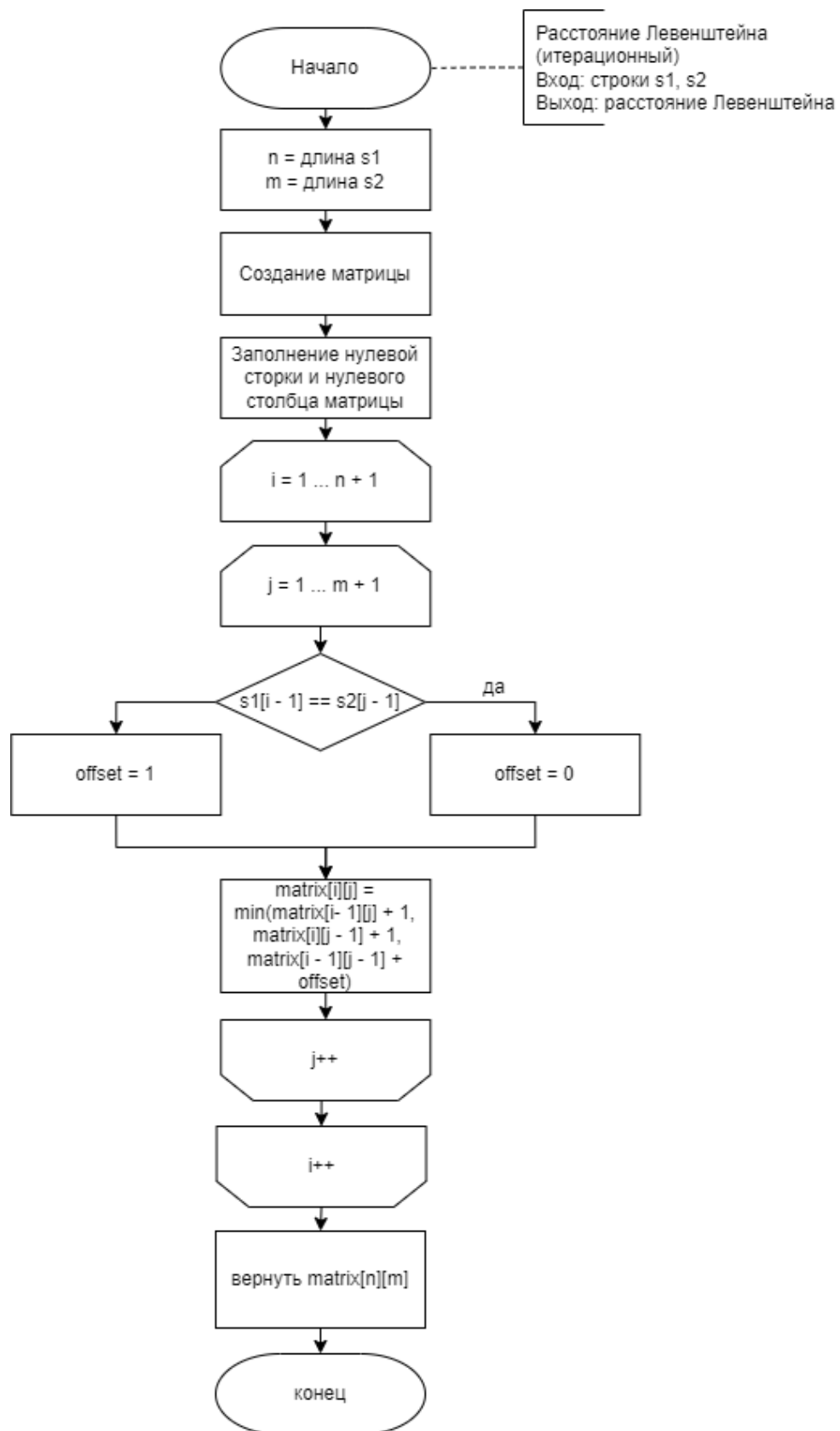


Рисунок 2.1 – Схема итерационного алгоритма поиска расстояния Левенштейна

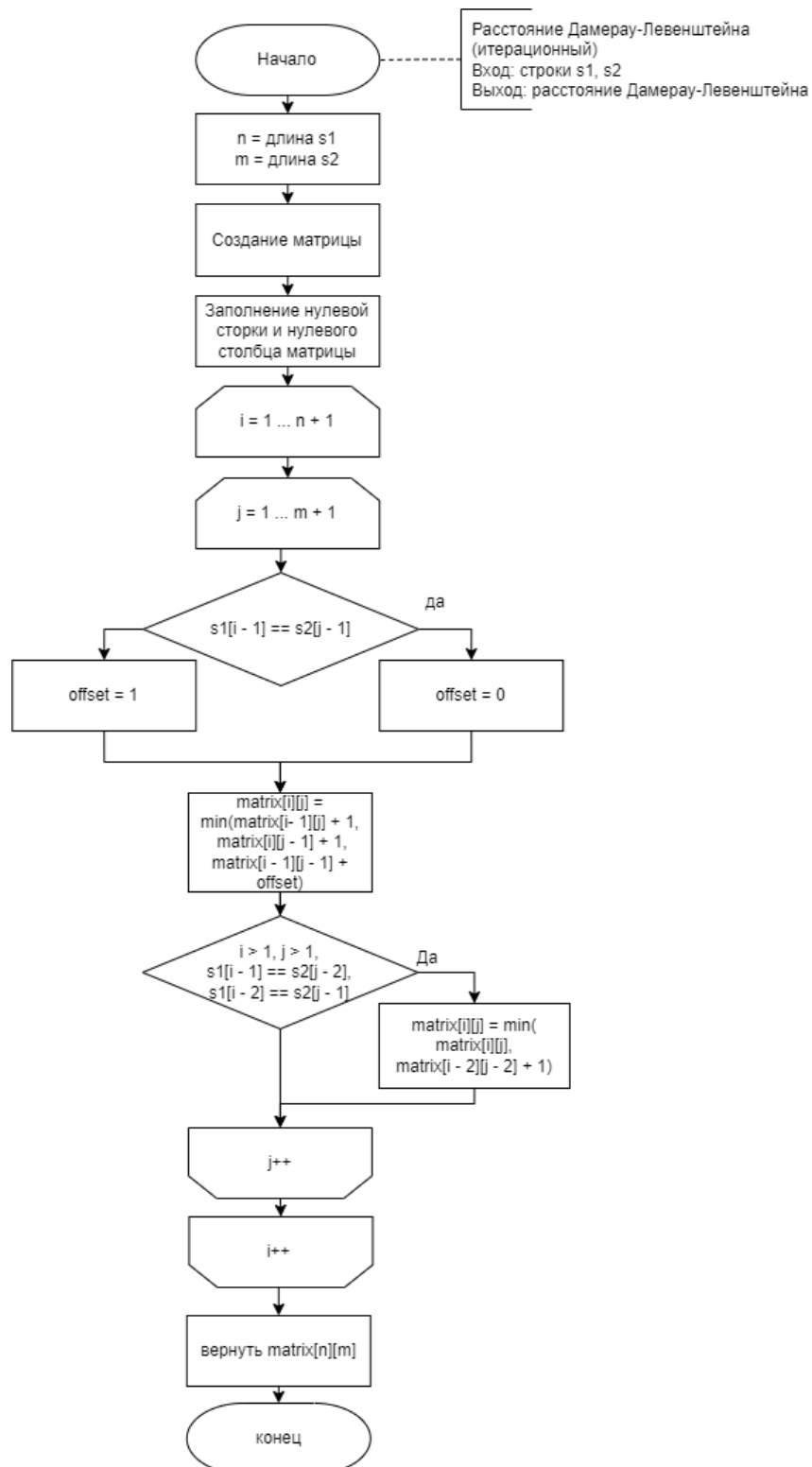


Рисунок 2.2 – Схема итерационного алгоритма поиска расстояния Дамерау-Левенштейна

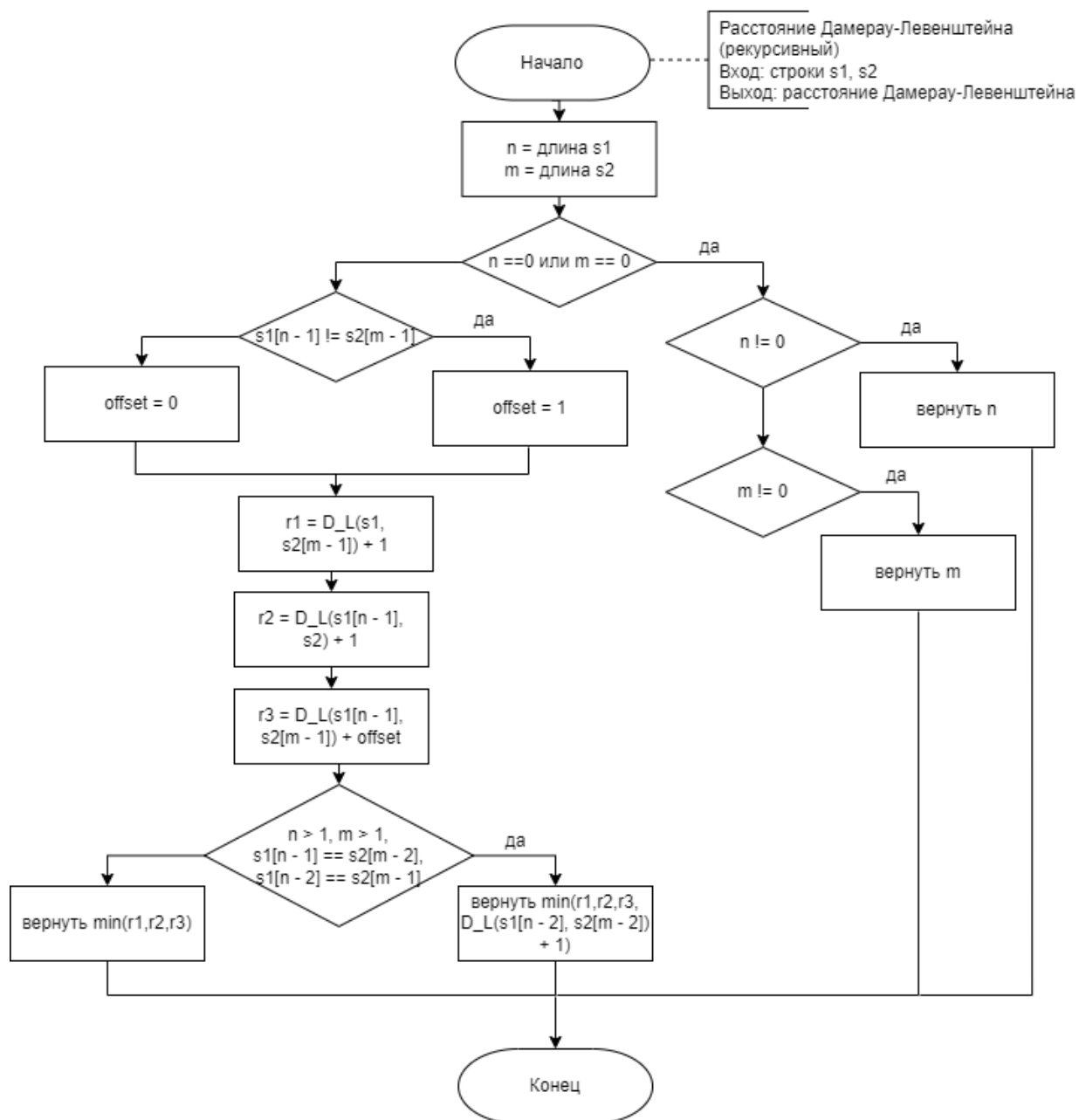


Рисунок 2.3 – Схема рекурсивного алгоритма поиска расстояния Дамерау-Левенштейна

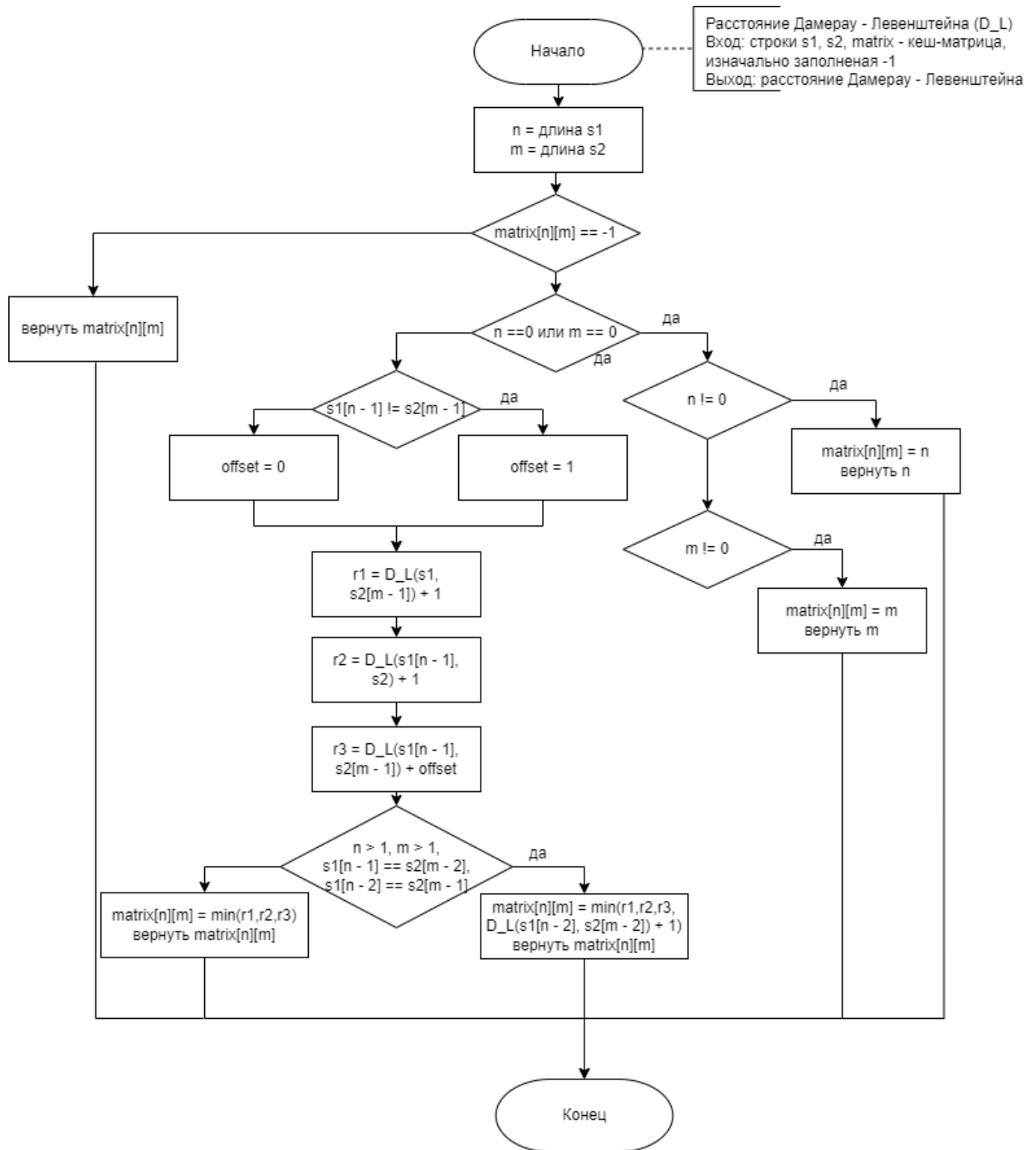


Рисунок 2.4 – Схема рекурсивного алгоритма поиска расстояния Дамерау-Левенштейна с кешем

Вывод

В данном разделе на основе теоретических данных были построены схемы требуемых алгоритмов, выбраны используемые типы данных, а также была проведена оценка затрачиваемого объёма памяти.

3 Технологическая часть

3.1 Требования к программному обеспечению

В программе должна присутствовать возможность:

- 1) ввода исходных слов, для которых будут находиться расстояния Левенштейна и Дамерау-Левенштейна;
- 2) поиска искомого расстояния для введенных строк с помощью одного из четырех рассматриваемых алгоритмов;
- 3) замера процессорного времени выполнения реализаций алгоритмов поиска расстояний Левенштейна и Дамерау-Левенштейна.

3.2 Выбор языка программирования

Для реализации алгоритмов поиска редакционных расстояний был выбран язык программирования C в силу наличия точных библиотек для замеров процессорного времени и быстродействия языка.

3.3 Выбор библиотеки и способа для замера времени

Для замера процессорного времени выполнения реализаций алгоритмов была выбрана не стандартная функция библиотеки `<time.h>` языка C — `clock()`, которая недостаточно четко работает при замерах небольших промежутков времени, а `QueryPerformanceCounter` - API-интерфейс, использующийся для получения меток времени с высоким разрешением или измерения интервалов времени.

Для облегчения работы с данным инструментом были самостоятельно написаны обертки-макросы, представленные на листинге 3.1.

Листинг 3.1 – Листинг макросов

```
1  #define TIMER_INIT \  
2      LARGE_INTEGER frequency; \  
3      LARGE_INTEGER t1,t2; \  
4      double elapsedTime; \  
5      QueryPerformanceFrequency(&frequency);  
6  
7  #define TIMER_START QueryPerformanceCounter(&t1);  
8  
9  #define TIMER_STOP \  
10     QueryPerformanceCounter(&t2); \  
11     elapsedTime=(float)(t2.QuadPart1.QuadPart)/frequency.  
12     QuadPart/COUNT*MICRO; \  
    printf("%lf", elapsedTime);
```

В силу существования явления вытеснения процессов из ядра, квантования процессорного времени все процессорное время не отдается какой-либо одной задаче, поэтому для получения точных результатов необходимо усреднить результаты вычислений: замерить совокупное время выполнения реализации алгоритма N раз и вычислить среднее время выполнения.

3.4 Реализации алгоритмов

В листингах 3.2, 3.3, 3.4, 3.5 приведены реализации алгоритмов поиска расстояний Левенштейна (итерационный), Дамерау-Левенштейна (итерационный), Дамерау-Левенштейна (рекурсивный), Дамерау-Левенштейна (рекурсивный с кешем) соответственно.

Листинг 3.2 – Листинг итерационного алгоритма поиска расстояния Левенштейна

```
1  int lev(char *str_1, char *str_2, int print_table_flag)  
2  {  
3      matrix_t *m = create_matrix(strlen(str_1) + 1, strlen(  
4          str_2) + 1);  
5      m->elements[0][0] = 0;  
6      for (size_t i = 1; i < m->rows; ++i)  
7          m->elements[i][0] = i;
```



```

8      for (size_t i = 1; i < m->cols; ++i)
9          m->elements[0][i] = i;
10
11     for (size_t i = 1; i < m->rows; ++i)
12         for (size_t j = 1; j < m->cols; ++j)
13     {
14         int offset = str_1[i - 1] == str_2[j - 1] ? 0 : 1;
15         m->elements[i][j] = min(3, m->elements[i][j - 1] +
16                                 1,
17                                 m->elements[i - 1][j] +
18                                 1,
19                                 m->elements[i - 1][j -
20                                     1] + offset);
21     }
22
23     int res = m->elements[m->rows - 1][m->cols - 1];
24     free_matrix(m);
25
26     return res;
27 }

```

Листинг 3.3 – Листинг итерационного алгоритма поиска расстояния
Дамерау-Левенштейна

```

1  int dameray_lev(char *str_1, char *str_2, int print_table_flag
2  )
3  {
4      matrix_t *m = create_matrix(strlen(str_1) + 1, strlen(
5          str_2) + 1);
6
7      m->elements[0][0] = 0;
8      for (size_t i = 1; i < m->rows; ++i)
9          m->elements[i][0] = i;
10     for (size_t i = 1; i < m->cols; ++i)
11         m->elements[0][i] = i;
12
13     for (size_t i = 1; i < m->rows; ++i)
14         for (size_t j = 1; j < m->cols; ++j) {
15             int offset = str_1[i - 1] == str_2[j - 1] ? 0 : 1;
16             m->elements[i][j] = min(3, m->elements[i][j - 1] +
17                                     1,
18                                     m->elements[i - 1][j] +
19                                     1,
20                                     m->elements[i - 1][j -
21                                         1] + offset);
22         }
23     }
24
25     return m->elements[m->rows - 1][m->cols - 1];
26 }

```

```

16         1,
           m->elements[i - 1][j -
17             1] + offset);
           if (j > 1 && i > 1 && str_1[i - 1] == str_2[j - 2]
18             && str_1[i - 2] == str_2[j - 1])
           m->elements[i][j] = min(2, m->elements[i][j],
19             m->elements[i - 2][j - 2] + 1);
20     }
21     int res = m->elements[m->rows - 1][m->cols - 1];
22     free_matrix(m);
23
24     return res;
25 }

```

Листинг 3.4 – Листинг рекурсивного алгоритма поиска расстояния
Дамерау-Левенштейна

```

1  int dameray_lev_rec(char *str_1, char *str_2, int len_1, int
   len_2)
2  {
3      if (len_1 == 0)
4          return len_2;
5      if (len_2 == 0)
6          return len_1;
7
8      int offset = str_1[len_1 - 1] == str_2[len_2 - 1] ? 0 : 1;
9      int res = min(3, dameray_lev_rec(str_1, str_2, len_1 - 1,
   len_2)
10                     + 1,
                     dameray_lev_rec(str_1, str_2, len_1,
11                     len_2 - 1) + 1,
                     dameray_lev_rec(str_1, str_2, len_1 - 1,
12                     len_2 - 1) + offset);
13      if (len_1 > 1 && len_2 > 1 && str_1[len_1 - 1] == str_2[
   len_2 - 2] && str_1[len_1 - 2] == str_2[len_2 - 1])
14          res = min(2, res, dameray_lev_rec(str_1, str_2, len_1
   - 2, len_2 - 2) + 1);
15
16      return res;
17  }

```

Листинг 3.5 – Листинг рекурсивного алгоритма поиска расстояния
Дамерау-Левенштейна с кешем

```
1  int dameray_lev_rec_cache_call(char *str_1, char *str_2, int
   len_1, int len_2, matrix_t *mat)
2  {
3      if (len_1 == 0)
4          return (mat->elements)[len_1][len_2] = len_2;
5      if (len_2 == 0)
6          return (mat->elements)[len_1][len_2] = len_1;
7
8      if (mat->elements[len_1 - 1][len_2] == -1)
9          dameray_lev_rec_hash_call(str_1, str_2, len_1 - 1,
   len_2, mat);
10     if (mat->elements[len_1][len_2 - 1] == -1)
11         dameray_lev_rec_hash_call(str_1, str_2, len_1, len_2 -
   1, mat);
12     if (mat->elements[len_1 - 1][len_2 - 1] == -1)
13         dameray_lev_rec_hash_call(str_1, str_2, len_1 - 1,
   len_2 - 1, mat);
14
15     int offset = str_1[len_1 - 1] == str_2[len_2 - 1] ? 0 : 1;
16     mat->elements[len_1][len_2] = min(3, mat->elements[len_1][
   len_2 - 1] + 1,
17                                     mat->elements[len_1 - 1][len_2] + 1,
18                                     mat->elements[len_1 -
   1][len_2 - 1] +
   offset);
19     if (len_1 > 1 && len_2 > 1 && str_1[len_1 - 1] == str_2[
   len_2 - 2] && str_1[len_1 - 2] == str_2[len_2 - 1]) {
20         if (mat->elements[len_1 - 2][len_2 - 2] == -1)
21             dameray_lev_rec_hash_call(str_1, str_2, len_1 - 2,
   len_2 - 2, mat);
22         mat->elements[len_1][len_2] = min(2, mat->elements[
   len_1][len_2], mat->elements[len_1 - 2][len_2 - 2]
   + 1);
23     }
24
25     return mat->elements[len_1][len_2];
26 }
27
28 int dameray_lev_rec_cache(char *str_1, char *str_2)
```

```

29  {
30      matrix_t *m = create_matrix(strlen(str_1) + 1, strlen(
31          str_2) + 1);
32      clear_matrix(m);
33
34      int res = damerau_lev_rec_hash_call(str_1, str_2, strlen(
35          str_1), strlen(str_2), m);
36      free_matrix(m);
37
38      return res;
39  }

```

3.5 Тестирование алгоритмов

В таблице 3.1 приведены проведенные тесты для функций, реализующих алгоритмы поиска расстояний Левенштейна и Дамерау-Левенштейна.

Таблица 3.1 – Тестирование функций

Первое слово	Второе слово	Расстояние Лев.	Расстояние Дам.-Лев.
λ	λ	0	0
<i>cat</i>	λ	3	3
<i>cat</i>	<i>cat</i>	0	0
<i>cat</i>	<i>catdog</i>	3	3
<i>hotdog</i>	<i>hodtog</i>	2	1

Вывод

В данном разделе были разработаны алгоритмы поиска расстояний Левенштейна и Дамерау-Левенштейна, а также проведено тестирование.

4 Экспериментальная часть

4.1 Технические характеристики

Ниже приведены технические характеристики устройства, на котором было проведено тестирование программного обеспечения:

- 1) операционная система Windows-10, 64-bit;
- 2) оперативная память 8 ГБ;
- 3) процессор Intel(R) Core(TM) i3-7020U CPU @ 2.30GHz, 2304 МГц, ядер 2, логических процессоров 4.

4.2 Замеры времени

В таблице 4.1 приведены результаты замеров в микросекундах времени алгоритмов для входных строк разной длины.

Таблица 4.1 – Таблица замера времени выполнения алгоритмов на строках, имеющих разные длины

Длина строки	Лев.	Дам.-Лев.	Дам.-Лев.(рек.)	Дам.-Лев.(рек. с кэшем)
1	1.15	1.12	0.11	0.96
2	1.57	1.08	0.24	1.96
3	2.71	2.84	1.36	1.61
4	1.56	2.56	6.30	3.94
5	2.85	3.71	44.91	4.11
6	5.35	6.73	161.43	5.47
7	6.60	7.40	1061.13	6.12
8	5.43	6.27	5355.47	4.91
9	3.70	4.47	30525.51	8.21
10	5.24	6.52	180535.97	6.95

Зависимость времени работы алгоритмов поиска расстояний Левенштейна (итерационный), Дамерау-Левенштейна (итерационный), Дамерау-

Левенштейна (рекурсивный с кешем) от длины входных строк представлена на рисунке 4.1.

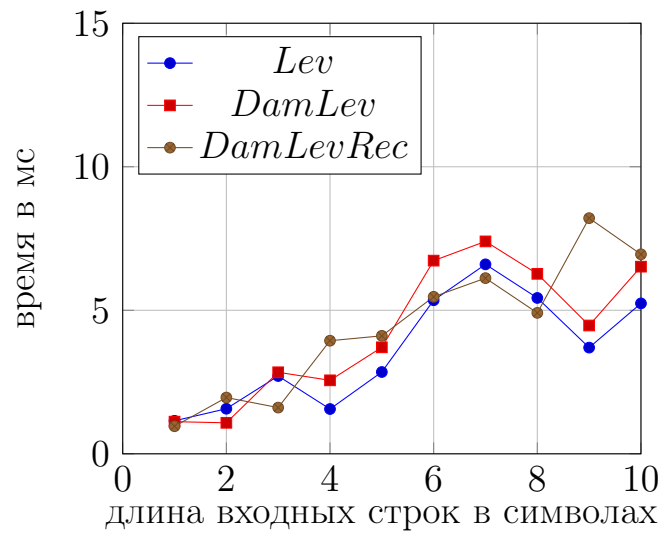


Рисунок 4.1 – Зависимость времени от длины входных строк

Зависимость времени работы алгоритма поиска расстояний Дамерау-Левенштейна (рекурсивный) от длины входных строк представлена на рисунке 4.2.

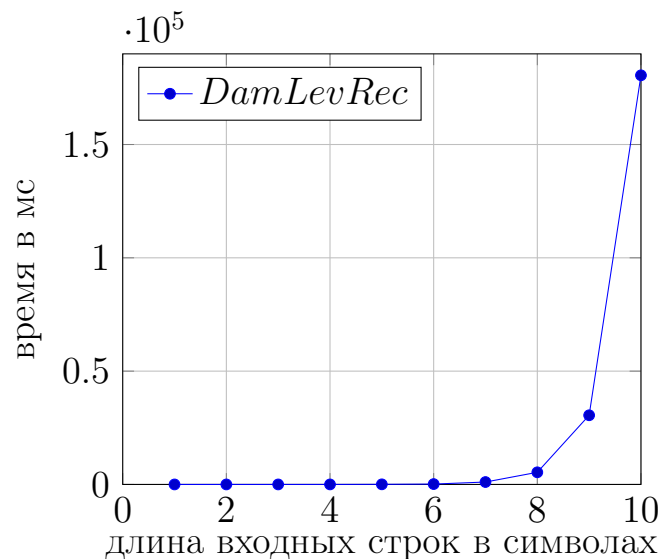


Рисунок 4.2 – Зависимость времени от длины входных строк

Вывод

Результаты замеров показали, что рекурсивный алгоритм поиска расстояния Дameraу-Левенштейна работает дольше всего. При этом его оптимизация с кешем работает в разы быстрее. Итерационные алгоритмы нахождения расстояний Левенштейна и Дameraу-Левенштейна оказались наиболее быстрыми, причем на длине строки от 1 до 10 элементов итерационный алгоритм Левенштейна в среднем работает немного быстрее, нежели итерационный алгоритм нахождения расстояния Дamera-Левенштейна.

Заключение

Цель лабораторной работы достигнута – были получены навыки динамического программирования на примере разработки алгоритмов поиска редакционных расстояний. Все задачи решены:

- 1) изучены расстояния Левенштейна и Дamerau-Левенштейна;
- 2) разработаны указанные алгоритмы поиска расстояний;
- 3) реализованы разработанные алгоритмы;
- 4) была рассчитана их трудоемкость по памяти;
- 5) проведен сравнительный анализ реализаций алгоритмов по затраченному процессорному времени и памяти на основе экспериментальных данных;
- 6) описаны и обоснованы полученные результаты в отчете о выполненной лабораторной работе.

На основании проведенных экспериментов было определено, что время работы алгоритмов с увеличением длины строк увеличивается в геометрической прогрессии. Самым медленным по скорости выполнения, но наименее затратным по памяти оказался рекурсивный алгоритм определения расстояния Дamerau-Левенштейна. Самым быстрым же оказался итерационный алгоритм нахождения расстояния Левенштейна, который одновременно с этим, оказался одним из самых затратных по памяти, так как в его реализации дополнительно используется матрица, размером, соответствующим длине входных строк.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- [1] Погорелов, Д. А. Сравнительный анализ алгоритмов редакционного расстояния Левенштейна и Дамерау-Левенштейна / Д. А. Погорелов, А. М. Тарзанов. - Синергия Наук. - 2019. URL: <https://elibrary.ru/item.asp?id=36907767> (дата обращения: 25.10.2022).
- [2] Microsoft. QueryPerformanceCounter function (profileapi.h). URL: <https://learn.microsoft.com/en-us/windows/win32/api/profileapi/nf-profileapi-queryperformancecounter> (дата обращения: 25.10.2022).
- [3] The LaTeX project. Core Documentation. URL: <https://www.latex-project.org/help/documentation/> (дата обращения: 20.10.2022).