



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 «Анализы Алгоритмов»

на тему: «Динамическое программирование»

Студент группы ИУ7-56Б

(Подпись, дата)

Разин А. В.
(Фамилия И.О.)

Преподаватель

(Подпись, дата)

Волкова Л. Л.
(Фамилия И.О.)

Преподаватель

(Подпись, дата)

Строганов Ю. В..
(Фамилия И.О.)

Москва — 2023 г.

Содержание

Введение	3
1 Аналитическая часть	4
1.1 Расстояние Левенштейна	4
1.1.1 Нерекурсивный алгоритм нахождения расстояния Ле- венштейна	6
1.2 Расстояние Дамерау-Левенштейна	7
1.2.1 Рекурсивный алгоритм нахождения расстояния Дамерау- Левенштейна	8
1.2.2 Рекурсивный алгоритм нахождения расстояния Дамерау- Левенштейна с кешированием	9
1.2.3 Нерекурсивный алгоритм нахождения расстояния Дамерау- Левенштейна	9
2 Конструкторская часть	11
2.1 Требования к программному обеспечению	11
2.2 Разработка алгоритмов	11
2.3 Описание используемых типов данных	17
3 Технологическая часть	18
3.1 Средства реализации	18
3.2 Сведения о модулях программы	19
3.3 Реализация алгоритмов	19
3.4 Функциональные тесты	26
4 Исследовательская часть	27
4.1 Технические характеристики	27
4.2 Демонстрация работы программы	27
4.3 Временные характеристики	29
4.4 Характеристики по памяти	30
4.5 Вывод	36

Заключение	37
Список использованных источников	38

Введение

Часто при работе со словами, необходимо их сравнивать, причем необходима конкретная метрика, которая покажет, насколько посимвольно одно слово отличается от другого. Одной из таких метрик является Расстояние Левенштейна. Данное расстояние является метрикой, измеряющей по модулю разность между двумя последовательностями символов.

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

Расстояние Левенштейна активно используется и по сей день:

- исправление ошибок в слове(в поисковых системах, базах данных, при вводе текста, при автоматическом распознавании отсканированного текста или речи);
- сравнение текстовых файлов утилитой diff;
- для сравнения геномов, хромосом и белков в биоинформатике.

Данная метрика была модифицирована Фредриком Дамерау, путем введения операции перестановки соседних символов.

1 Аналитическая часть

Целью данной лабораторной работы является описание и исследование алгоритмов поиска расстояний Левенштейна и Дamerau-Левенштейна.

Для поставленной цели необходимо выполнить следующие задачи.

- 1) Исследовать расстояние Левенштейна.
- 2) Разработка алгоритмов поиска расстояний Левенштейна, Дamerau-Левенштейна.
- 3) Создать программное обеспечение, реализующее следующие алгоритмы.
- 4) Провести исследование, затрачиваемого процессорного времени и памяти при различных реализациях алгоритмов.
- 5) Провести сравнительный анализ алгоритмов.

1.1 Расстояние Левенштейна

Расстояние Левенштейна [1] (редакционное расстояние, дистанция редактирования) — метрика, измеряющая разность между двумя последовательностями символов.

Вводятся операции нескольких типов:

- 1 I (англ. insert) — операция вставки символа.
- 2 D (англ. delete) — операция удаления символа.
- 3 R (англ. replace) — операция замены символа.

Также введем символ λ , обозначающий пустой символ строки, не входящий ни в одну из рассматриваемых строк.

Будем считать стоимость каждой вышеизложенной операции равной 1:

— $D(a, b) = 1$, $a \neq b$, в противном случае замена не происходит;

- $D(\lambda, b) = 1;$
- $D(a, \lambda) = 1.$

Также обозначим совпадение символов как M (англ. match), таким образом $D(a, a) = 0$. Существует проблема взаимного выравнивания строк, в случае когда строки различной длины существует множество способов сопоставить соответствующие символы.

Решим проблему введением рекуррентной формулы, обозначим:

1. $L1$ — длина S_1 .
2. $L2$ — длина S_2 .
3. $S_1[1...i]$ — подстрока S_1 длиной i , начиная с первого символа.
4. $S_2[1...j]$ — подстрока S_2 длиной j , начиная с первого символа.

Расстояние Левенштейна между двумя строками S_1 и S_2 длиной M и N соответственно рассчитывается по рекуррентной формуле:

$$D(S_1[1...i], S_2[1...j]) = \begin{cases} 0, i = 0, j = 0 \\ i, j = 0, i > 0 \\ j, i = 0, j > 0 \\ \min \begin{cases} D(S_1[1...i], S_1[1...j-1]) + 1, \\ D(S_1[1...i-1], S_1[1...j]) + 1, \\ D(S_1[1...i-1], S_1[1...j-1]) + m(S_1[i], S_2[j]), \end{cases} & i > 0 \end{cases} \quad (1.1)$$

где сравнение символов строк S_1 и S_2 рассчитывается как:

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

В расстоянии Дамерау-Левенштейна вводится еще одна операция, обозначим ее как S (англ. swap) — данная операция применима, только если $S_1[i] = S_2[j-1]$ и $S_1[i-1] = S_2[j]$. Рекуррентная формула данной

метрики выглядит следующим образом:

$$D(i, j) = \begin{cases} 0, i = 0, j = 0, \\ i, j = 0, i > 0, \\ j, i = 0, j > 0, \\ \min \begin{cases} D(S_1[1...i], S_2[1...j]) + 1, \\ D(S_1[1...i - 1], S_2[1...j]) + 1, \\ D(S_1[1...i - 1], S_2[1...j - 1]) + m(S_1[i], S_2[j]), \\ D(S_1[1...i - 2], S_2[1...j - 2]) + s(S_1[i], S_2[j], S_1[i - 1], S_2[j - 1]), \end{cases} \end{cases} \quad (1.3)$$

где стоимость перестановки пар символов рассчитывается как:

$$s(a, b, c, d) = \begin{cases} 1 & \text{если } a = d, b = c \\ +\infty & \text{иначе.} \end{cases} \quad (1.4)$$

Идея в том, что после замены пары символов местами, полученные пары были равны друг другу.

1.1.1 Нерекурсивный алгоритм нахождения расстояния Левенштейна

Рекурсивная реализация алгоритма Левенштейна малоэффективна по времени при больших M и N , так как множество промежуточных значений. Для оптимизации можно использовать итерационную реализацию заполнения матрицы промежуточными значениями $D(i, j)$.

В качестве структуры данных для хранения промежуточных значений можно использовать матрицу, имеющую размеры:

$$(N + 1) \times (M + 1) \quad (1.5)$$

Значения в ячейке $[i, j]$ равно значению $D(S_1[1...i], S_2[1...j])$. Первый элемент матрицы заполнен нулем. Всю таблицу заполнять в соответствии с формулой (??).

Однако матричный алгоритм является малоэффективным по памяти по сравнению с рекурсивным при больших M и N , т.к. множество промежуточных значений $D(i, j)$ хранится в памяти после их использования. Для оптимизации по памяти рекурсивного алгоритма нахождения расстояния Левенштейна можно использовать кеш, т.е. пару строк, содержащую значения $D(i, j)$, вычисленные в предыдущей итерации, алгоритма и значения $D(i, j)$, вычисляемые в текущей итерации.

1.2 Расстояние Дамерау-Левенштейна

Расстояние Дамерау-Левенштейна, названное в честь ученых Фредерика Дамерау и Владимир 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 трем базовым операциям добавляется операция транспозиции T (от англ. transposition).

Расстояние Дамерау-Левенштейна может быть вычислено по рекуррентной формуле:

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0, \\ i, & j = 0, i > 0, \\ j, & i = 0, j > 0, \\ \min \begin{cases} D(i, j - 1) + 1, \\ D(S_1[1...i - 1], S_2[1...j]) + 1, \\ D(S_1[1...i - 1], S_2[1...j - 1]) + m(S_1[i], S_2[j]), \\ D(S_1[1...i - 2], S_2[1...j - 2]) + 1, \end{cases} & \begin{aligned} & \text{если } i > 1, j > 1, \\ & S_1[i] = S_2[j - 1] \\ & S_1[i - 1] = S_2[j] \end{aligned} \\ \min \begin{cases} D(i, j - 1) + 1, \\ D(S_1[1...i - 1], S_2[1...j]) + 1, \\ D(S_1[1...i - 1], S_2[1...j - 1]) + m(S_1[i], S_2[j]), \end{cases} & \text{иначе.} \end{cases} \quad (1.6)$$

1.2.1 Рекурсивный алгоритм нахождения расстояния Дameraу-Левенштейна

Рекурсивный алгоритм реализует формулу (1.6), функция D составлена таким образом, что верно следующее.

- 1) Для передачи из пустой строки в пустую требуется ноль операций.
- 2) Для перевода из пустой строки в строку a требуется $|a|$ операций.
- 3) Для перевода из строки a в пустую строку требуется $|a|$ операций.
- 4) Для перевода из строки a в строку b требуется выполнить последовательно некоторое количество операций удаления, вставки, замены, транспозиции в некоторой последовательности. Последовательность поведения любых двух операций можно поменять, порядок поведения операций не имеет никакого значения. Если полагать, что a' , b' – строки a и b без последнего символа соответственно, а a'' , b'' – строки a и b без двух последних символов, то цена преобразования из строки a в b выражается из элементов, представленных ниже:
 - сумма цены преобразования строки a' в b и цены проведения операции удаления, которая необходима для преобразования a' в a ;
 - сумма цены преобразования строки a в b' и цены проведения операции вставки, которая необходима для преобразования b' в b ;
 - сумма цены преобразования из a' в b' и операции замены, предполагая, что a и b оканчиваются на разные символы;
 - сумма цены преобразования из a'' в b'' и операции перестановки, предполагая, что длины a'' и b'' больше 1 и последние два символа a'' , поменянные местами, совпадут с двумя последними символами b'' ;
 - цена преобразования из a' в b' , предполагая, что a и b оканчиваются на один и тот же символ.

Минимальной стоимостью преобразования будет минимальное значение приведенных вариантов.

1.2.2 Рекурсивный алгоритм нахождения расстояния Дамерау-Левенштейна с кешированием

Рекурсивная реализация алгоритма Дамерау-Левенштейна малоэффективна по времени при больших M и N по причине проблемы повторных вычислений значений расстояний между подстроками. Для оптимизации алгоритма нахождения расстояния Левенштейна можно использовать матрицу в целях хранения соответствующих промежуточных значений. В таком случае алгоритм представляет собой рекурсивное заполнение матрицы $A_{|a|,|b|}$ промежуточными значениями $D(i, j)$, такое хранение промежуточных данных можно назвать кешем для рекурсивного алгоритма.

1.2.3 Нерекурсивный алгоритм нахождения расстояния Дамерау-Левенштейна

Рекурсивная реализация алгоритма Левенштейна с кешированием малоэффективна по времени при больших M и N . Для оптимизации можно использовать итерационную реализацию заполнения матрицы промежуточными значениями $D(i, j)$.

В качестве структуры данных для хранения промежуточных значений можно использовать *матрицу*, имеющую размеры:

$$(N + 1) \times (M + 1), \quad (1.7)$$

Значение в ячейке $[i, j]$ равно значению $D(S_1[1...i], S_2[1...j])$. Первый элемент заполнен нулем. Всю таблицу заполняем в соответствии с формулой (1.6).

Вывод

В данном разделе были рассмотрены алгоритмы динамического программирования — алгоритмы нахождения расстояний Левенштейна и Дамерау-Левенштейна, формулы которых задаются рекуррентно, а следовательно, данные алгоритмы могут быть реализованы рекурсивно и итеративно. На вход алгоритмам поступают две строки, которые могут содержать как русские, так и английские буквы, также будет предусмотрен ввод пустых строк.

2 Конструкторская часть

В данном разделе будут приведены схемы алгоритмов нахождения расстояний Левенштейна и Дамерау-Левенштейна, приведены описание используемых типов данных, оценки памяти, а также описана структура программного обеспечения.

2.1 Требования к программному обеспечению

К программе предъявлен ряд функциональных требований: входные данные — две строки, выходные данные — результат работы всех алгоритмов поиска расстояний, целое число.

К программе предъявлен ряд требований:

- наличие интерфейса для выбора действий;
- должна обрабатывать строки;
- возможность обработки строк, включающих буквы как на латинице, так и на кириллице;
- наличие функциональности замера процессорного времени работы реализаций алгоритмов поиска расстояний Левенштейна и Дамерау-Левенштейна.

2.2 Разработка алгоритмов

На вход алгоритмов подаются строки S_1 и S_2 .

На рисунке 2.1 представлена схема алгоритма поиска расстояния Левенштейна. На рисунках 2.2 – 2.5 представлены схемы алгоритмов поиска Дамерау-Левенштейна.

img/levmatr.png

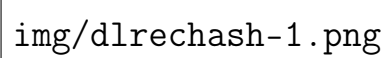
Рисунок 2.1 – Схема нерекурсивного алгоритма нахождения расстояния
Левенштейна

img/dlitter.png

Рисунок 2.2 – Схема нерекурсивного алгоритма нахождения расстояния
Дамерау-Левенштейна

img/dlrec.png

Рисунок 2.3 – Схема рекурсивного алгоритма нахождения расстояния
Дамерау-Левенштейна



img/dlrehash-1.png

Рисунок 2.4 – Схема рекурсивного алгоритма нахождения расстояния Дамерау-Левенштейна с кешированием

img/dlrehash-2.png

Рисунок 2.5 – Схема алгоритма рекурсивного заполнения матрицы путем поиска расстояния Дамерау-Левенштейна

2.3 Описание используемых типов данных

При реализации алгоритмов будут использованы следующие структуры данных:

- строка — массив типа *wchar* размером длины строки;
- длина строки — целое число типа *int*;
- матрица — двумерный массив значений типа *int*.

Вывод

В данном разделе на основе теоретических данных были построены схемы требуемых алгоритмов, выбраны используемые типы данных.

3 Технологическая часть

В данном разделе будут приведены требования к программному обеспечению, средства реализации, листинг кода и функциональные тесты.

3.1 Средства реализации

Для реализации данной лабораторной работы был выбран язык *C++* [2]. Данный выбор обусловлен наличием у языка встроенного модуля *ctime* измерения процессорного времени и типа данных, позволяющего хранить как кириллические символы, так и латинские — *std::wstring*, что позволит удовлетворить третьему требованию из п.2.1 соответствуют выдвинутым техническим требованиям.

Время работы было замерено с помощью функции *clock()* из библиотеки *ctime* [3]

3.2 Сведения о модулях программы

Данная программа разбита на следующие модули.

- `main.cpp` — файл, содержащий точку входа в программу, из которой происходит вызов алгоритмов по разработанному интерфейсу.
- `algorithms.cpp` — файл содержит функции поиска расстояния Левенштейна и Дамерау-Левенштейна.
- `allocate.cpp` — файл содержит функции динамического выделения и очищения памяти для матрицы.
- `print_mtr_lev.cpp` — файл содержит функцию вывода матрицы для итеративных алгоритмов поиска расстояния Левенштейна и Дамерау-Левенштейна, включая строки.
- `cru_time.cpp` — файл содержит функции, измеряющее процессорное время алгоритмов поиска расстояния Левенштейна и Дамерау-Левенштейна.
- `memory.cpp` — файл содержит функции, измеряющее память итеративного и рекурсивного алгоритмов поиска расстояния Левенштейна.

3.3 Реализация алгоритмов

В листингах 3.1 – 3.4 приведены реализации алгоритмов поиска расстояний Левенштейна (только нерекурсивный алгоритм) и Дамерау-Левенштейна (нерекурсивный, рекурсивный и рекурсивный с кешированием). В листингах 3.5 – 3.6 приведены реализации алгоритмов выделения памяти под матрицу и вывод матрицы.

Листинг 3.1 – Функция нахождения расстояния Левенштейна с использованием матрицы

```
1 int lev_mtr(wstring &str1, wstring &str2, bool print ) {
2     size_t n = str1.length();
3     size_t m = str2.length();
4     int **mtr = malloc_mtr(n + 1, m + 1);
5     int res = 0;
6     for (int i = 0; i <= n; i++)
7         for (int j = 0; j <= m; j++)
8             if (i == 0 && j == 0)
9                 mtr[i][j] = 0;
10            else if (i > 0 && j == 0)
11                mtr[i][j] = i;
12            else if (j > 0 && i == 0)
13                mtr[i][j] = j;
14        else {
15            int change = 0;
16            if (str1[i - 1] != str2[j - 1])
17                change = 1;
18
19            mtr[i][j] = std::min(mtr[i][j - 1] + 1,
20                               std::min(mtr[i - 1][j] + 1,
21                                         mtr[i - 1][j - 1] + change));
22        }
23        if (print)
24            print_mtr_lev(str1, str2, mtr, n, m);
25        res = mtr[n][m];
26        free_mtr(mtr, n);
27        return res;
28 }
```

Листинг 3.2 – Функция нахождения расстояния Дамерау-Левенштейна с использованием матрицы

```
1 int dameray_lev_mtr(wstring &str1, wstring &str2, bool print)
2 {
3     size_t n = str1.length();
4     size_t m = str2.length();
5     int **mtr = malloc_mtr(n + 1, m + 1);
6     int res = 0;
7
8     for (int i = 0; i <= n; i++)
9         for (int j = 0; j <= m; j++) {
10             if (i == 0 && j == 0)
11                 mtr[i][j] = 0;
12             else if (i > 0 && j == 0)
13                 mtr[i][j] = i;
14             else if (j > 0 && i == 0)
15                 mtr[i][j] = j;
16             else {
17                 int change = 0;
18                 if (str1[i - 1] != str2[j - 1])
19                     change = 1;
20
21                 mtr[i][j] = min(mtr[i][j - 1] + 1,
22                               min(mtr[i - 1][j] + 1,
23                                   mtr[i - 1][j - 1] + change));
24
25                 if (i > 1 && j > 1 &&
26                     str1[i - 1] == str2[j - 2] &&
27                     str1[i - 2] == str2[j - 1])
28                     mtr[i][j] = min(mtr[i][j], mtr[i - 2][j -
29                                     2] + 1);
30             }
31         }
32
33     if (print)
34         print_mtr_lev(str1, str2, mtr, n, m);
35     res = mtr[n][m];
36     free_mtr(mtr, n);
37
38     return res;
39 }
```

Листинг 3.3 – Функция нахождения расстояния Дameraу-Левенштейна рекурсивно

```
1 int damera_lev_rec_t(wstring &str1, wstring &str2, size_t n,
2   size_t m) {
3     if (n == 0)
4       return m;
5     if (m == 0)
6       return n;
7
8     int change = 0;
9     int res = 0;
10    if (str1[n - 1] != str2[m - 1])
11      change = 1;
12
13    res = min(damera_lev_rec_t(str1, str2, n, m - 1) + 1,
14             min(damera_lev_rec_t(str1, str2, n - 1, m) + 1,
15                 damera_lev_rec_t(str1, str2, n - 1, m - 1) +
16                   change));
17
18    if (n > 1 && m > 1 &&
19        str1[n - 1] == str2[m - 2] &&
20        str1[n - 2] == str2[m - 1])
21      res = std::min(res, damera_lev_rec_t(str1, str2, n -
22        2, m - 2) + 1);
23    return res;
24 }
25
26 int damera_lev_rec(wstring &str1, wstring &str2)
27 {
28   return damera_lev_rec_t(str1, str2, str1.length(),
29                             str2.length());
30 }
```

Листинг 3.4 – Функция нахождения расстояния Дамерау-Левенштейна рекурсивно с кешированием

```
1  int dameray_lev_rec_hash_t(wstring &str1, wstring &str2, int
   **mtr, size_t n, size_t m) {
2      if (n == 0)
3          return mtr[n][m] = m;
4      if (m == 0)
5          return mtr[n][m] = n;
6      int change = 0;
7      if (str1[n - 1] != str2[m - 1])
8          change = 1;
9      mtr[n][m] = min(dameray_lev_rec_hash_t(str1, str2, mtr, n,
   m - 1) + 1,
10                     min(dameray_lev_rec_hash_t(str1, str2, mtr, n -
   1, m) + 1,
11                        dameray_lev_rec_hash_t(str1, str2, mtr, n -
   1, m - 1) + change));
12      if (n > 1 && m > 1 &&
13          str1[n - 1] == str2[m - 2] &&
14          str1[n - 2] == str2[m - 1])
15          mtr[n][m] = min(mtr[n][m], dameray_lev_rec_hash_t(str1,
   str2, mtr, n - 2, m - 2) + 1);
16      return mtr[n][m];
17 }
18
19 int dameray_lev_rec_hash(wstring &str1, wstring &str2, bool
   print) {
20     size_t n = str1.length();
21     size_t m = str2.length();
22     int **mtr = malloc_mtr(n + 1, m + 1);
23     for (int i = 0; i <= n; i++)
24         for (int j = 0; j <= m; j++) {
25             mtr[i][j] = -1;
26         }
27     int res = dameray_lev_rec_hash_t(str1, str2, mtr, n, m);
28     if (print)
29         print_mtr_lev(str1, str2, mtr, n, m);
30     free_mtr(mtr, n);
31     return res;
32 }
```


Листинг 3.5 – Функции динамического выделения и очищения памяти под матрицу

```
1 void free_mtr(int **mtr, std::size_t n) {
2     if (mtr != nullptr)
3     {
4         for (std::size_t i = 0; i < n; i++)
5             if (mtr[i] != nullptr)
6                 free(mtr[i]);
7         free(mtr);
8     }
9 }
10
11 int **malloc_mtr(std::size_t n, std::size_t m)
12 {
13     if (n == 0)
14         return nullptr;
15
16     int **mtr = static_cast<int **>(malloc(n * sizeof(int *)));
17     if (mtr != nullptr)
18         for (std::size_t i = 0; mtr[i] != nullptr && i < n; i++) {
19             mtr[i] = static_cast<int *>(malloc(m * sizeof(int)));
20             if (mtr[i] == nullptr)
21                 free_mtr(mtr, n);
22         }
23
24     return mtr;
25 }
```

Листинг 3.6 – Функции вывода матрицы для алгоритмов поиска расстояния Левенштейна и Дамерау-Левенштейна

```
1 void print_mtr_lev(std::wstring str1, std::wstring str2,
2 int **mtr, std::size_t n, std::size_t m)
3 {
4     for(std::size_t i = 0; i <= n + 1; i++)
5     {
6         for(std::size_t j = 0; j <= m + 1; j++)
7         {
8             if (i == 0 && j == 0)
9                 std::wcout << " ";
10            else if (i == 0)
11                if (j == 1)
12                    std::wcout << "- ";
13            else
14                std::wcout << str2[j - 2] << " ";
15            else if (j == 0)
16                if (i == 1)
17                    std::wcout << "- ";
18            else
19                std::wcout << str1[i - 2] << " ";
20            else
21                std::wcout << mtr[i - 1][j - 1] << " ";
22        }
23        std::wcout << std::endl;
24    }
25 }
```

3.4 Функциональные тесты

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

Таблица 3.1 – Функциональные тесты

Входные данные		Расстояние и алгоритм			
Строка 1	Строка 2	Левенштейна	Дамерау-Левенштейна		
		Итеративный	Итеративный	Рекурсивный	
				Без кеша	С кешом
а	б	1	1	1	1
а	а	0	0	0	0
кот	скат	2	2	2	2
друзья	рдузия	3	2	2	2
вагон	гонки	4	4	4	4
бар	раб	2	2	2	2
слон	слоны	1	1	1	1

Вывод

Были реализованы алгоритмы поиска расстояния Левенштейна итеративно, а также поиска расстояния Дамерау—Левенштейна итеративно, рекурсивно и рекурсивного с кешированием. Проведено тестирование реализаций алгоритмов.

4 Исследовательская часть

4.1 Технические характеристики

Технические характеристики устройства, на котором выполнялись замеры по времени, представлены далее.

- Процессор: Intel(R) Core(TM) i5-10300H CPU 2.50 ГГц [4].
- Оперативная память: 16 ГБайт.
- Операционная система: Windows 10 Pro 64-разрядная система версии 21H2 [5].

При замерах времени ноутбук был включен в сеть электропитания и был нагружен только системными приложениями.

4.2 Демонстрация работы программы

На рисунке 4.1 представлена демонстрация работы разработанного программного обеспечения, а именно показаны результаты вычислений реализаций алгоритмов поиска расстояний Левенштейна и Дамерау-Левенштейна на примере двух строк «друзья» и «рдузия». При этом выводятся матрицы для алгоритмов, использующих матрицы для промежуточных результатов.

img/example.png

Рисунок 4.1 – Демонстрация работы программы при поиске расстояний
Левенштейна и Дамерау-Левенштейна

4.3 Временные характеристики

Результаты эксперимента замеров по времени приведены в таблице 4.1, в которой есть поля, обозначенные «-». Это обусловлено тем, что для рекурсивной реализации алгоритмов достаточно приведенных замеров для построения графика. По полученным замерам по времени для рекурсивной реализации понятно, что проведения замеров на длин строк больше 15 будет достаточно долгим, поэтому нет смысла проводить замеры по времени рекурсивных реализаций алгоритмов поиска расстояния.

Замеры проводились на одинаковых длин строк от 1 до 200 с различным шагом.

Таблица 4.1 – Замер по времени для строк, размер которых от 1 до 200

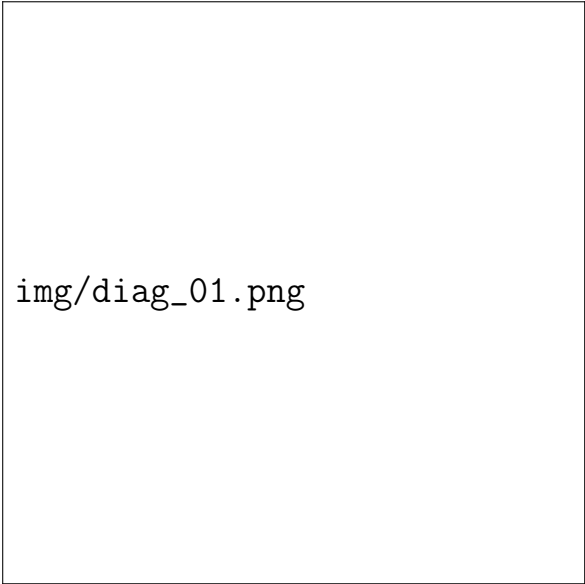
Длина (символ)	Время, нс			
	Левенштейн	Дамерау-Левенштейн		
	Итеративный	Итеративный	Рекурсивный	
			Без кеша	С кешом

Отдельно сравниваются итеративные алгоритмы поиска расстояний Левенштейна и Дамерау-Левенштейна. Сравнение будет производится на основе данных, представленных в таблице 4.1. Результат можно рассмотреть на рисунке 4.2.

При длинах строк менее 30 символов разница по времени между итеративными реализациями незначительна, однако при увеличении длины строки алгоритм поиска расстояния Левенштейна оказывается быстрее вплоть до полутора раз (при длинах строк равных 200). Это обосновывается тем, что у алгоритма поиска расстояния Дамерау-Левенштейна задействуется дополнительная операция, которая замедляет алгоритм.

Также сравним рекурсивную и итеративную реализации алгоритма поиска расстояния Дамерау-Левенштейна. Данные представлены в таблице 4.1 и отображены на рисунке 4.3.

На рисунке 4.3 продемонстрировано, что рекурсивный алгоритм становится менее эффективным по времени (вплоть до 21 раз при длине строк равной 7 элементов), чем итеративный.



img/diag_01.png

Рисунок 4.2 – Сравнение по времени нерекурсивных реализаций алгоритмов поиска расстояний Левенштейна и Дамерау-Левенштейна

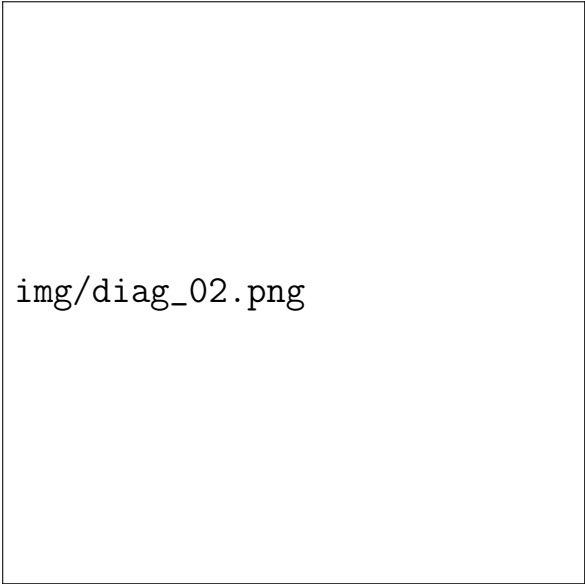
Кроме того, согласно данным, приведенным в таблице 4.1, рекурсивные алгоритмы при длинах строк более 10 элементов не пригодны к использованию в силу экспоненциально роста затрат процессорного времени, в то время, как затраты итеративных алгоритмов по времени линейны.

4.4 Характеристики по памяти

Введем следующие обозначения:

- n — длина строки S_1 ;
- m — длина строки S_2 ;
- $size()$ — функция вычисляющая размер в байтах;
- *string* — строковый тип;
- *int* — целочисленный тип;
- *size_t* — беззнаковый целочисленный тип.

Максимальная глубина стека вызовов при рекурсивной реализации нахождения расстояния Дамерау-Левенштейна равна сумме входящих строк,



img/diag_02.png

Рисунок 4.3 – Сравнение по времени алгоритмов поиска расстояния
Дамерау-Левенштейна

а на каждый вызов требуется 2 дополнительные переменные, соответственно, максимальный расход памяти равен:

$$(n + m) \cdot (2 \cdot \text{size}(\text{string}) + 3 \cdot \text{size}(\text{int}) + 2 \cdot \text{sizeof}(\text{size_t})), \quad (4.1)$$

где:

- $2 \cdot \text{size}(\text{string})$ — хранение двух строк;
- $2 \cdot \text{size}(\text{size_t})$ — хранение размеров строк;
- $2 \cdot \text{size}(\text{int})$ — дополнительные переменные;
- $\text{size}(\text{int})$ — адрес возврата.

Для рекурсивного алгоритма с кешированием поиска расстояния Дамерау-Левенштейна будет теоретически схож с расчетом в формуле (4.1), но также учитывается матрица, соответственно, максимальный расход памяти равен:

$$(n + m) \cdot (2 \cdot \text{size}(\text{string}) + 3 \cdot \text{size}(\text{int}) + 2 \cdot \text{size}(\text{size_t})) + \\ + (n + 1) \cdot (m + 1) \cdot \text{size}(\text{int}) \quad (4.2)$$

Использование памяти при итеративной реализации алгоритма поиска рас-

стояния Левенштейна теоретически равно:

$$(n + 1) \cdot (m + 1) \cdot \text{size}(\text{int}) + 2 \cdot \text{size}(\text{string}) + 2 \cdot \text{size}(\text{size_t}) + \text{size}(\text{int} **) + (n + 1) \cdot \text{size}(\text{int} *) + 2 \cdot \text{size}(\text{int}), \quad (4.3)$$

где

- $2 \cdot \text{size}(\text{string})$ — хранение двух строк;
- $2 \cdot \text{size}(\text{size_t})$ — хранение размеров матрицы;
- $(n + 1) \cdot (m + 1) \cdot \text{size}(\text{int})$ — хранение матрицы;
- $\text{size}(\text{int} **) + (n + 1) \cdot \text{size}(\text{int} *)$ — указатель на матрицу;
- $\text{size}(\text{int})$ — дополнительная переменная для хранения результата;
- $\text{size}(\text{int})$ — адрес возврата.

Использование памяти при итеративной реализации алгоритма поиска расстояния Дамерау-Левенштейна теоретически равно:

$$(n + 1) \cdot (m + 1) \cdot \text{size}(\text{int}) + 2 \cdot \text{size}(\text{string}) + 2 \cdot \text{size}(\text{size_t}) + \text{size}(\text{int} **) + (n + 1) \cdot \text{size}(\text{int} *) + 3 \cdot \text{size}(\text{int}), \quad (4.4)$$

где

- $2 \cdot \text{size}(\text{string})$ — хранение двух строк;
- $2 \cdot \text{size}(\text{size_t})$ — хранение размеров матрицы;
- $(n + 1) \cdot (m + 1) \cdot \text{size}(\text{int})$ — хранение матрицы;
- $\text{size}(\text{int} **) + (n + 1) \cdot \text{size}(\text{int} *)$ — указатель на матрицу;
- $2 \cdot \text{size}(\text{int})$ — дополнительные переменные;
- $\text{size}(\text{int})$ — адрес возврата.

По расходу памяти итеративные алгоритмы проигрывают рекурсивным: максимальный размер используемой памяти в итеративном растет как произведение длин строк, в то время как у рекурсивного алгоритма — как сумма длин строк.

По формулам 4.1 – 4.3 затрат по памяти в программе были написаны соответствующие функции для подсчета расходуемой памяти, результаты расчетов, которых представлены в таблице 4.2, где размеры строк находятся в диапазоне от 10 до 200 с шагом 10.

Таблица 4.2 – Замер памяти для строк, размером от 10 до 200

Длина (символ)	Размер в байтах			
	Левенштейн	Дамерау-Левенштейн		
	Итеративный	Итеративный	Рекурсивный	
			Без кеша	С кешом

Из данных, приведенных в таблице 4.2, понятно, что рекурсивные алгоритмы являются более эффективными по памяти, так как используется только память под локальные переменные, передаваемые аргументы и возвращаемое значение, в то время как итеративные алгоритмы затрачивают память линейно пропорционально длинам обрабатываемых строк.

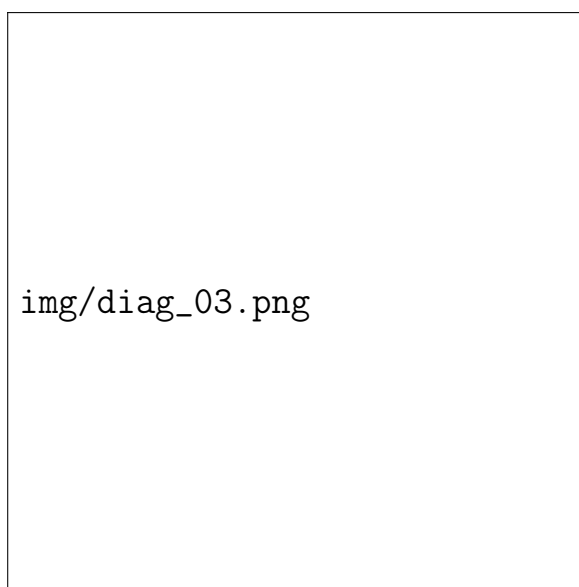
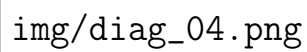


Рисунок 4.4 – Сравнение по памяти алгоритмов поиска расстояния Левенштейна и Дамерау-Левенштейна — итеративной и рекурсивной реализации

Из рисунка 4.5 понятно, что рекурсивная реализация алгоритма поиска расстояния Дамерау-Левенштейна эффективная по памяти, чем итеративная.



img/diag_04.png

Рисунок 4.5 – Сравнение по памяти алгоритмов поиска расстояния Дамерау-Левенштейна — итеративной и рекурсивной реализации

4.5 Вывод

В данном разделе было произведено сравнение количества затраченного времени и памяти алгоритмов поиска расстояний Левенштейна и Дамерау-Левенштейна. Наименее затратным по времени оказался итеративный алгоритм нахождения расстояния Левенштейна.

Приведенные характеристики показывают, что рекурсивная реализация алгоритма в 21 раз проигрывает по времени. В связи с этим, рекурсивные алгоритмы следует использовать лишь для малых размерностей строк (1 – 4 символа).

Так как во время печати очень часто возникают ошибки связанные с транспозицией букв [6], алгоритм поиска расстояния Дамерау-Левенштейна является наиболее предпочтительным, не смотря на то, что он проигрывает по времени и памяти алгоритму Левенштейна.

Рекурсивная реализация алгоритма поиска расстояния Дамерау-Левенштейна будет более затратным по времени по сравнению с итеративной реализацией алгоритма поиска расстояния Дамерау-Левенштейна, но менее затратным по памяти по отношению к итеративному алгоритму Дамерау-Левенштейна.

Заключение

В результате исследования было определено, что время алгоритмов нахождения расстояний Левенштейна и Дамерау-Левенштейна растет в геометрической прогрессии при увеличении длин строк. Лучшие показатели по времени дает матричная реализация алгоритма нахождения расстояния Дамерау-Левенштейна и его рекурсивная реализация с кешем, использование которых приводит к 21-кратному превосходству по времени работы уже на длине строки в 4 символа за счет сохранения необходимых промежуточных вычислений. При этом итеративная реализации с использованием матрицы занимают довольно много памяти при большой длине строк.

Цель данной лабораторной работы были достигнуты, а именно описание и исследование особенностей задач динамического программирования на алгоритмах Левенштейна и Дамерау-Левенштейна.

Для достижения поставленной целей были выполнены следующие задачи.

- 1) Описаны алгоритмы поиска расстояния Левенштейна и Дамерау-Левенштейна;
- 2) Создано программное обеспечение, реализующее следующие алгоритмы.
 - нерекурсивный метод поиска расстояния Левенштейна;
 - нерекурсивный метод поиска расстояния Дамерау-Левенштейна;
 - рекурсивный метод поиска расстояния Дамерау-Левенштейна;
 - рекурсивный с кешированием метод поиска расстояния Дамерау-Левенштейна.
- 3) Выбраны инструменты для замера процессорного времени выполнения реализаций алгоритмов.
- 4) Проведены анализ затрат работы программы по времени и по памяти, выяснить влияющие на них характеристики.

Список использованных источников

- 1 И. Левенштейн В. Двоичные коды с исправлением выпадений, вставок и замещений символов. — М.: Издательство «Наука», Доклады АН СССР, 1965. Т. 163.
- 2 Документация по Microsoft C++ [Электронный ресурс]. — Режим доступа: <https://learn.microsoft.com/ru-ru/cpp/?view=msvc-170&viewFallbackFrom=vs-2017> (дата обращения: 25.09.2022).
- 3 C library function clock() [Электронный ресурс]. — Режим доступа: https://www.tutorialspoint.com/c_standard_library/c_function_clock.htm (дата обращения: 25.09.2022).
- 4 Intel [Электронный ресурс]. — Режим доступа: <https://ark.intel.com/content/www/ru/ru/ark/products/201839/intel-core-i510300h-processor-8m-cache-up-to-4-50-ghz.html> (дата обращения: 25.09.2022).
- 5 Windows 10 Pro 2h21 64-bit [Электронный ресурс]. — Режим доступа: <https://www.microsoft.com/ru-ru/software-download/windows10> (дата обращения: 25.09.2022).
- 6 В. Ульянов М. Ресурсно-эффективные компьютерные алгоритмы: учебное пособие. — М.: Издательство «Наука», ФИЗМАТЛИ, 2007.