

Москва, 2020

# Оглавление

<b>Введение</b>	<b>2</b>
<b>1 Аналитическая часть</b>	<b>3</b>
<b>2 Конструкторская часть</b>	<b>5</b>
2.1 Разработка алгоритмов . . . . .	5
<b>3 Технологическая часть</b>	<b>9</b>
3.1 Требования к программному обеспечению . . . . .	9
3.2 Средства реализации . . . . .	9
3.3 Листинг кода алгоритмов . . . . .	10
<b>4 Экспериментальная часть</b>	<b>13</b>
4.1 Тестирование программы . . . . .	13
4.2 Постановка эксперимента по замеру времени и памяти . . . . .	14
4.3 Сравнительный анализ на материале экспериментальных данных . . . . .	15
<b>Список использованной литературы</b>	<b>21</b>

# Введение

**Расстояние Левенштейна** - минимальное количество редакторских операций, которое необходимо для превращения одной строки в другую (или другими словами - расстояние между строками).

Расстояние Левенштейна используется в следующих сферах:

- поисковики - строка поиска Яндекс: автоисправление, автозамена;
- бионформатика - в белке каждая молекула представляется буквой из ограниченного алфавита, и белки представляются строками, которые можно сравнивать.

Целью данной лабораторной работы является изучение метода динамического программирования на материале алгоритмов Левенштейна и Дамерау-Левенштейна.

Задачами данной лабораторной являются:

1. изучение алгоритмов Левенштейна и Дамерау-Левенштейна нахождения расстояния между строками;
2. применение метода динамического программирования для матричной реализации указанных алгоритмов;
3. получение практических навыков реализации указанных алгоритмов: двух алгоритмов в матричной версии и одного из алгоритмов в рекурсивной версии;
4. сравнительный анализ линейной и рекурсивной реализаций выбранного алгоритма определения расстояния между строками по затрачиваемым ресурсам (времени и памяти);
5. экспериментальное подтверждение различий во временной эффективности рекурсивной и нерекурсивной реализаций выбранного алгоритма определения расстояния между строками при помощи разработанного программного обеспечения на материале замеров процессорного времени выполнения реализации на варьирующихся длинах строк;
6. описание и обоснование полученных результатов в отчете о выполненной лабораторной работе, выполненного как расчётно-пояснительная записка к работе.

# 1. Аналитическая часть

Задача по нахождению расстояния Левенштейна заключается в поиске минимального количества операций вставки/удаления/замены для превращения одной строки в другую.

При нахождении расстояния Дамерау — Левенштейна добавляется операция транспозиции (перестановки соседних символов).

Действия обозначаются так:

1. D (delete) — удалить;
2. I (insert) — вставить;
3. R (replace) — заменить;
4. M(match) - совпадение.

Пусть  $S_1$  и  $S_2$  — две строки (длиной  $M$  и  $N$  соответственно) над некоторым алфавитом, тогда расстояние Левенштейна можно подсчитать по следующей рекуррентной формуле, описанной в книге [1]:

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & j = 0, i > 0 \\ j, & i = 0, j > 0 \\ \min( & \\ \quad D(S_1[i], S_2[j - 1]) + 1, & j > 0 \quad //I \\ \quad D(S_1[i - 1], S_2[j]) + 1, & i > 0 \quad //D \\ \quad D(S_1[i - 1], S_2[j - 1]) + & \\ \quad \begin{cases} 0, & \text{если } S_1[i] == S_2[j], \quad //M \\ 1, & \text{иначе} \quad //R \end{cases} & \\ ) & \end{cases}$$

где  $\min()$  возвращает наименьший из аргументов.

Расстояние Дамерау-Левенштейна вычисляется по следующей рекуррентной формуле, описанной в книге [1]:

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & j = 0, i > 0 \\ j, & i = 0, j > 0 \\ \min( & \\ D(S_1[i], S_2[j - 1]) + 1, & j > 0 \quad //I \\ D(S_1[i - 1], S_2[j]) + 1, & i > 0 \quad //D \\ D(S_1[i - 1], S_2[j - 1]) + & \\ \left[ \begin{array}{ll} 0, & \text{если } S_1[i] == S_2[j], \quad //M \\ 1, & \text{иначе} \quad //R \end{array} \right. & \\ D(S_1[i - 2], S_2[j - 2]) + 1 & \text{если } i = 1, j = 1, S_1[i] = S_2[j - 1], S_1[i - 1] = S_2[j] \\ ) & \end{cases}$$

## Вывод

Была рассмотрена теоретическая часть алгоритмов нахождения минимального редакционного расстояния: Левенштейна и Дамерау-Левенштейна. Отличие которых — наличие дополнительной операции во втором алгоритме.

## 2. Конструкторская часть

В данном разделе описываются схемы алгоритмов для нахождения минимального редакционного расстояния.

### 2.1 Разработка алгоритмов

На рисунках 1 - 4 представлены схемы алгоритмов, сделанные на основе теоретического материала, представленного в книге [2]

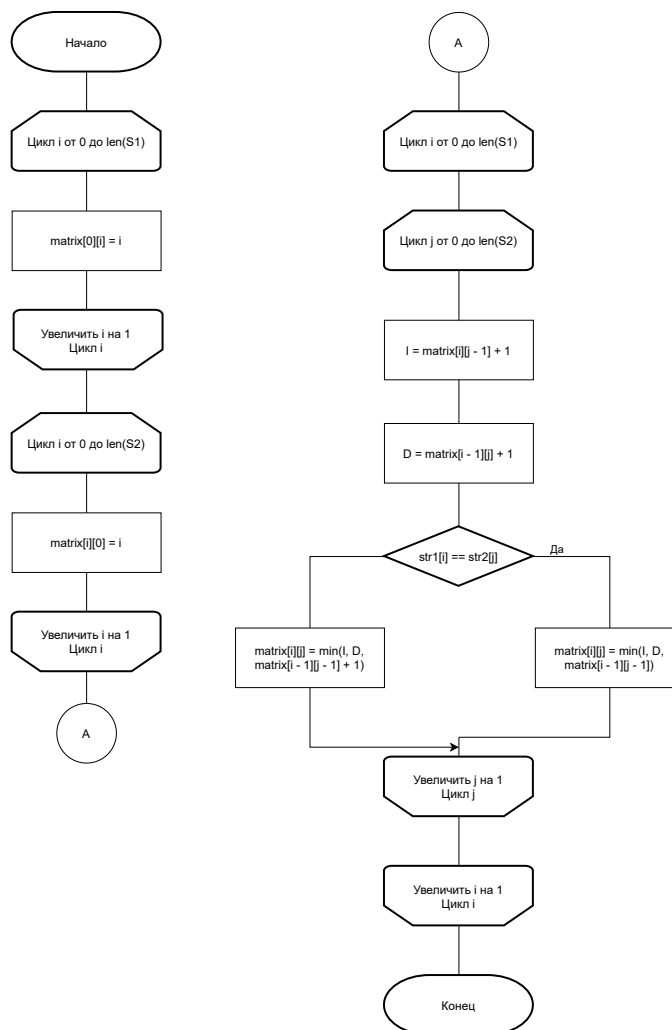


Рисунок 1: Схема матричного алгоритма Левенштейна

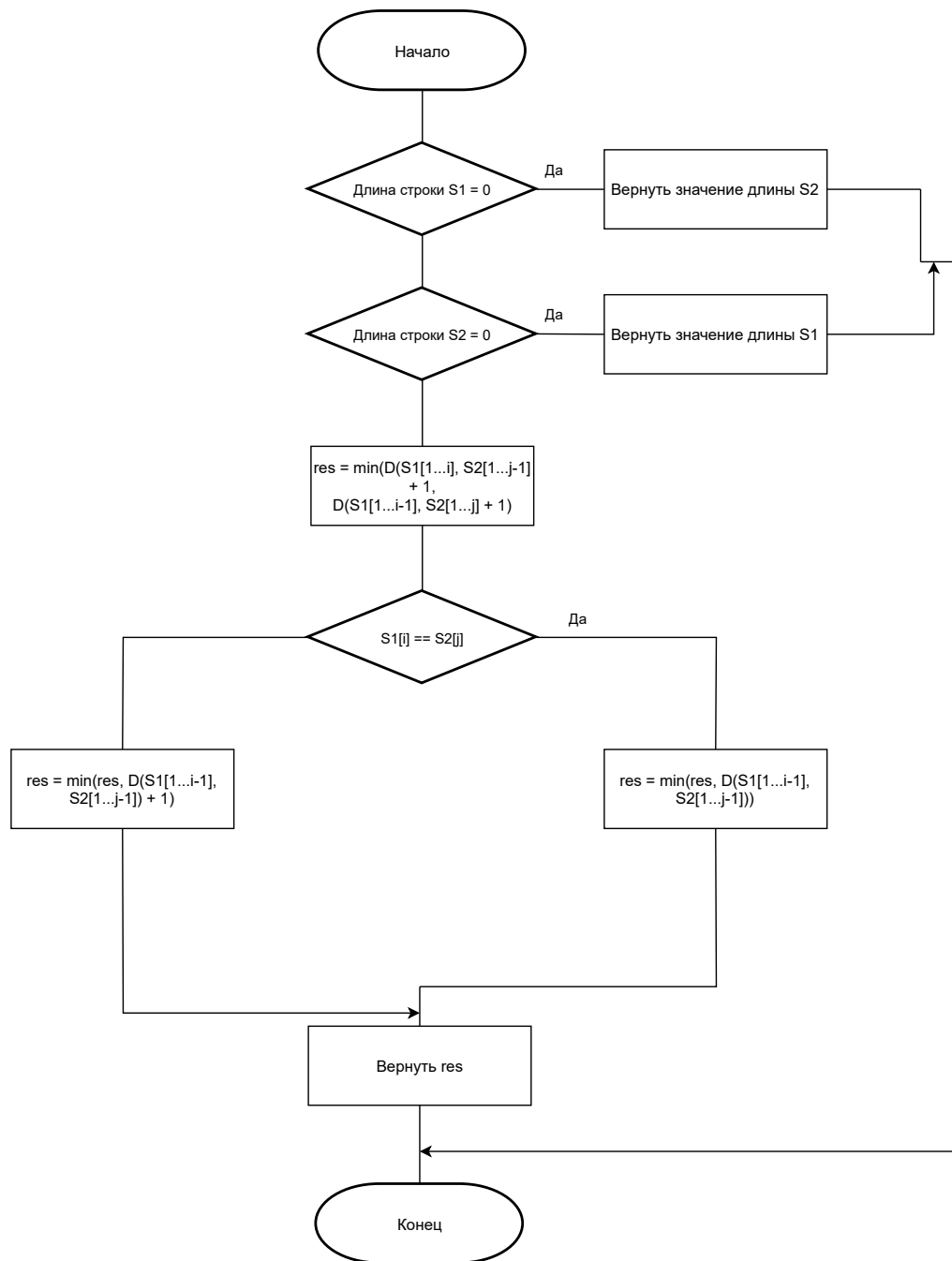


Рисунок 2: Схема рекурсивного алгоритма нахождения расстояния Левенштейна

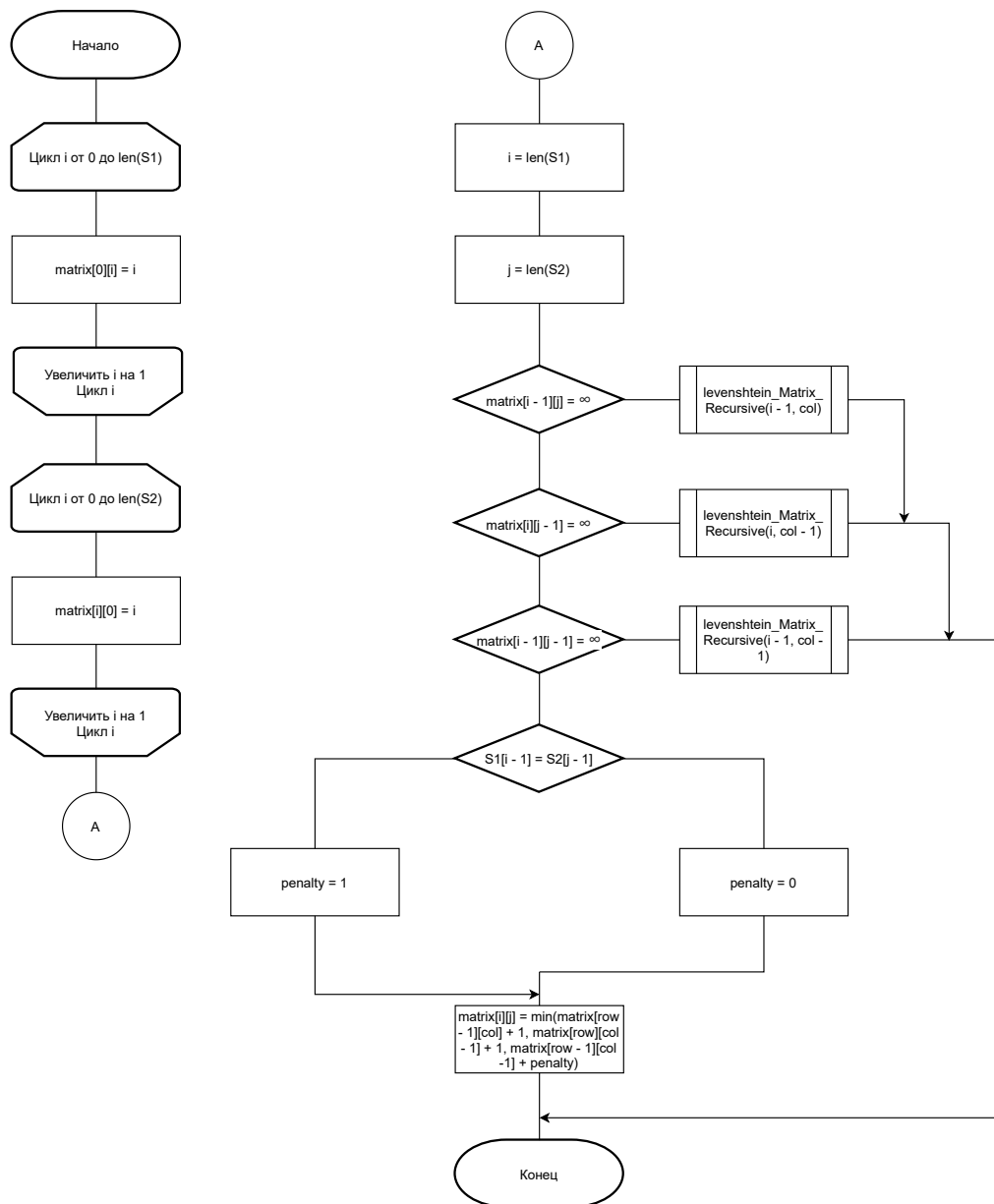


Рисунок 3: Схема рекурсивного алгоритма нахождения расстояния Левенштейна с заполнением матрицы и дополнительными проверками



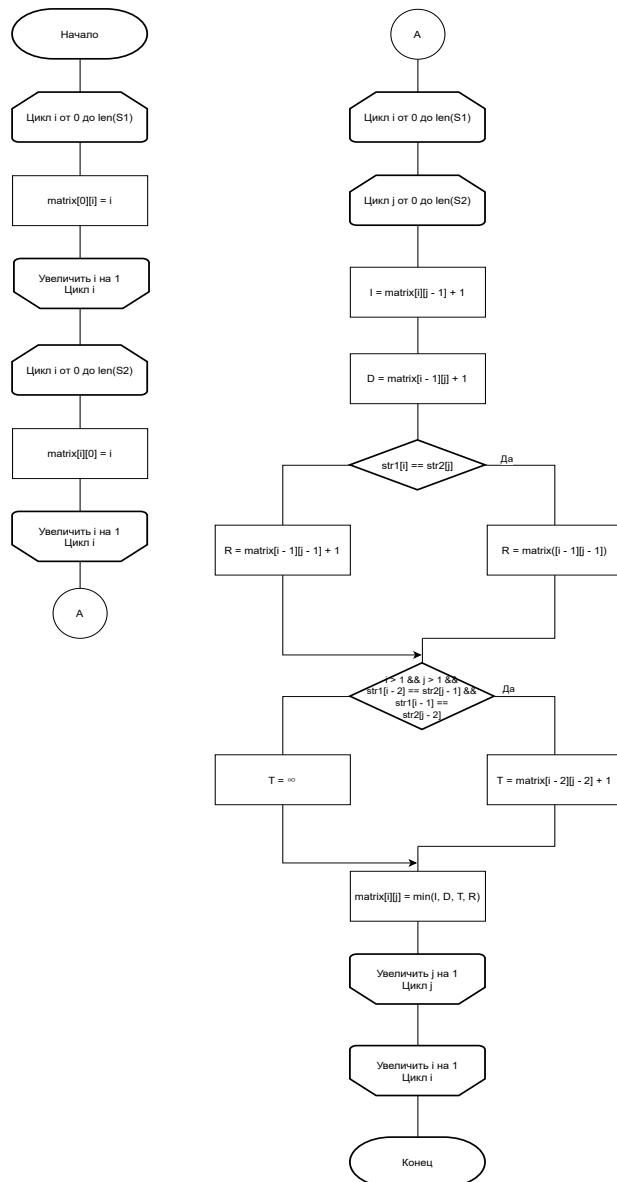


Рисунок 4: Схема матричного алгоритма нахождения расстояния Дamerau - Левенштейна

## Вывод

В данном разделе были проиллюстрированы схемы алгоритмов нахождения минимального редакционного расстояния.

## 3. Технологическая часть

В данном разделе описываются средства реализации, рассматриваются требования к программному обеспечению, а также приводится реализация алгоритмов на языке C++.

### 3.1 Требования к программному обеспечению

**Требования к вводу:**

1. на вход подаются две строки;
2. буквы в верхнем и нижнем регистрах считаются разными.

**Требования к программе:** Две пустые строки - корректный ввод, программа не должна завершаться аварийно.

### 3.2 Средства реализации

Для выполнения лабораторной работы мною был выбран язык C++, так как он является одним из самых высокопроизводительных языков программирования. Для замера времени использовалась встроенная библиотека `std::chrono`[4].

### 3.3 Листинг кода алгоритмов

На листинге 3.1 представлена реализация стандартного алгоритма Левенштейна:

Листинг 3.1: Функция нахождения расстояния Левенштейна с использованием матрицы

```
1 Matrix<size_t> Levenstein(string_view first , string_view second) {
2     auto rows = first.size() + 1;
3     auto cols = second.size() + 1;
4     Matrix<size_t> result(rows, cols);
5     for (size_t i = 0; i < rows; i++) {
6         for (size_t j = 0; j < cols; j++) {
7             if (i * j == 0) {
8                 result[i][j] = i | j;
9             } else if (first[i - 1] == second[j - 1]) {
10                result[i][j] = result[i - 1][j - 1];
11            } else {
12                result[i][j] = 1 + u_min(
13                    result[i - 1][j - 1],
14                    result[i][j - 1],
15                    result[i - 1][j]
16                );
17            }
18        }
19    }
20    return result;
21 }
```

На листинге 3.2 представлена реализация алгоритма Левенштейна с использованием рекурсии:

Листинг 3.2: Функция нахождения расстояния Левенштейна с использованием рекурсии

```
1 size_t Recursivempl(string_view first , string_view second ,
2                     const size_t i, const size_t j)
3 {
4     if (i * j == 0) {
5         return i | j;
6     } else if (first[i - 1] == second[j - 1]) {
7         return Recursivempl(first , second , i - 1, j - 1);
8     }
9     return 1 + u_min(
10         Recursivempl(first , second , i , j - 1),
11         Recursivempl(first , second , i - 1, j),
12         Recursivempl(first , second , i - 1, j - 1)
13     );
14 }
15
16 size_t RecursiveLevenstein(string_view first , string_view secnd) {
17     return Recursivempl(first , secnd , first.size() + 1, secnd.size() + 1);
18 }
```

На листинге 3.3 представлена реализация алгоритма Левенштейна с использованием рекурсии и заполнением матрицы:

Листинг 3.3: Функция нахождения расстояния Левенштейна с использованием рекурсии и заполнением матрицы

```
1 size_t RecursiveMatrImpl(Matrix<size_t>& matrix, string_view first ,
2     string_view second, const size_t i, const size_t j)
3 {
4     if (matrix[i][j] != -1) {
5         return matrix[i][j];
6     }
7     if (i * j == 0) {
8         return matrix[i][j] = i | j;
9     }
10    auto cost = (first[i - 1] != second[j - 1]);
11    return matrix[i][j] = cost + u_min(
12        RecursiveMatrImpl(matrix, first, second, i, j - 1),
13        RecursiveMatrImpl(matrix, first, second, i - 1, j),
14        RecursiveMatrImpl(matrix, first, second, i - 1, j - 1)
15    );
16 }
17
18 Matrix<size_t> RecursiveMatrLevenstein(string_view first ,
19     string_view second)
20 {
21     Matrix<size_t> result(first.size() + 1, second.size() + 1);
22     result.assign(-1);
23     RecursiveMatrImpl(result, first, second, first.size(), second.size());
24     return result;
25 }
```

На листинге 3.4 представлена реализация стандартного алгоритма Дамерау-Левенштейна:

Листинг 3.4: Функция нахождения расстояния Дамерау-Левенштейна с использованием матрицы

```
1 Matrix<size_t> DamerauLevenstein(string_view first, string_view second) {
2     auto rows = first.size() + 1;
3     auto cols = second.size() + 1;
4     Matrix<size_t> result(rows, cols);
5
6     for (size_t i = 0; i < rows; i++) {
7         for (size_t j = 0; j < cols; j++) {
8             if (i * j == 0) {
9                 result[i][j] = i | j;
10            } else if (first[i - 1] == second[j - 1]) {
11                result[i][j] = result[i - 1][j - 1];
12            } else {
13                result[i][j] = 1 + u_min(
14                    result[i - 1][j - 1],
15                    result[i][j - 1],
16                    result[i - 1][j]
17                );
18            }
19            if (i > 1 && j > 1 && first[i - 1] == second[j - 2]
20                && first[i - 2] == second[j - 1])
21            {
22                result[i][j] = u_min(
23                    result[i][j],
24                    result[i - 2][j - 1] + 1
25                );
26            }
27        }
28    }
29    return result;
30 }
```

## Вывод

В данном разделе были приведены реализации алгоритмов поиска минимального редакционного расстояния на языке программирования C++.

## 4. Экспериментальная часть

В данном разделе иллюстрируются примеры работы программы, результаты тестирования, а также описывается постановка эксперимента по замеру времени и памяти.

### 4.1 Тестирование программы

Примеры работы для каждой из операций:

```
Enter your choice: 1
Input first string:
some
Input second string:
som

Result:
  0   1   2   3
  1   0   1   2
  2   1   0   1
  3   2   1   0
  4   3   2   1
```

(1) удаление символов из строки

```
Input first string:
cat
Input second string:
cot

Result:
  0   1   2   3
  1   0   1   2
  2   1   1   2
  3   2   2   1
```

(2) замена

```
Input first string:
cat
Input second string:
ca

Result:
  0   1   2
  1   0   1
  2   1   0
  3   2   1
```

(3) добавление

```
Input first string:
kick
Input second string:
kick

Result:
  0   1   2   3   4
  1   0   1   2   2
  2   1   0   1   2
  3   2   1   0   1
  4   2   2   1   0
```

(4) ввод одинаковых строк

```
Input first string:
cat
Input second string:
cta

Result:
  0   1   2   3
  1   0   1   2
  2   1   1   1
  3   2   1   1
```

(5) перестановка

Рисунок 5: Результаты работы программы

Также было проведено тестирование методом черного ящика, результаты показаны в таблице 4.2:

№	Первое слово	Второе слово	Ожидаемый результат(Левенштейн, Дамерау- Левенштейн)	Полученный результат
1	kot	skat	2 2	2 2
2	kate	ktae	2 1	2 1
3	abacaba	aabcaab	4 2	4 2
4	sobaka	sboku	3 3	3 3
5	qwerty	queue	4 4	4 4
6	apple	aplpe	2 1	2 1
7	bmstu	utsmb	4 4	4 4

Таблица 4.2: Таблица результатов тестов

## 4.2 Постановка эксперимента по замеру времени и памяти

В рамках данной лабораторной работы были проведены следующие эксперименты:

1. теоретический замер потребляемой памяти всеми четырьмя алгоритмами на строках длины 10;
2. сравнение скорости работы матричных алгоритмов Левенштейна и Дамерау-Левенштейна на строках длиной от 1000 до 2000 включительно. Один эксперимент проводился 50 раз;
3. сравнение скорости работы рекурсивного алгоритма Левенштейна и рекурсивного алгоритма Левенштейна с матрицей на строках длиной от 1 до 12 включительно. Один эксперимент ставился 5 раз;
4. сравнение трех алгоритмов Левенштейна на строках длиной от 1 до 12 включительно. Один эксперимент проводился 5 раз.

## 4.3 Сравнительный анализ на материале экспериментальных данных

В таблице 4.3 представлен список основных сущностей, используемых в алгоритмах, и количество памяти, которое занимает каждая из них.

Сущность	Описание	Размер(64 битная архитектура) в байтах
Целое число	Машинно зависимая величина.	4
Массив рун	Руна - это псевдоним для целого беззнакового 32 битного числа	$\text{len}(\text{arr}) * 4$
Вектор / Строка	Данная сущность по сути представляет собой динамический массив.	24
Ссылка	Объект, указывающий на определенные данные, но не хранящий их.	4

Таблица 4.3: Память, занимаемая различными сущностями в C++

В таблице 4.4 видно, что потребляемая память у рекурсивных алгоритмов зависит от максимальной глубины рекурсии, которая равна максимальной длине из двух слов.

Алгоритм	Используемые сущности	Память, занимаемая каждой сущностью(в байтах)	Потребляемая память для всего алгоритма(64 битная архитектура)
Матричный Левенштейн	Две строки	$\text{len}(s1) * 2 + \text{len}(s2) * 2 + 48$	$2 * (\text{len}(s1) + \text{len}(s2)) + 88 + (\text{len}(s1) + 1) * (24 + (\text{len}(s2) + 1) * 4)$
	Вектор векторов для представления матрицы	$(\text{len}(s1) + 1) * (24 + (\text{len}(s2) + 1) * 4) + 24$	
	Четыре переменные типа <code>size_t</code> для итерации в цикле и хранения длины строк	16	
Рекурсивный Левенштейн	Две строки	$\text{len}(s1) * 2 + \text{len}(s2) * 2 + 48$	$(2 * (\text{len}(s1) + \text{len}(s2)) + 48 + 68 * (\text{макс.уровень рекурсии}))$
	Две переменные типа <code>size_t</code> для индексации строк	8	
	Указатель для адреса возврата	8	
	Переменная типа <code>size_t</code> для ответа	4	



Рекурсивный Левенштейн с Матрицей	Две строки	$\text{len}(s1) * 2 + \text{len}(s2) * 2 + 48$	$(2 * (\text{len}(s1) + \text{len}(s2)) + 68 * (\text{макс.уровень рекурсии}) + (\text{len}(s1) + 1) * (24 + (\text{len}(s2) + 1) * 4) + 72$
	Вектор век- торов для представления матрицы	$(\text{len}(s1) + 1) * (24 + (\text{len}(s2) + 1) * 4) + 24$	
	Две перемен- ные типа size_t для индексации строк	8	
	Указатель для адреса возврата	8	
	Переменная ти- па size_t для ответа	4	
Матричный Дамерау- Левенштейн	Две строки	$\text{len}(s1) * 2 + \text{len}(s2) * 2 + 48$	$2 * (\text{len}(s1) + \text{len}(s2)) + 88 + (\text{len}(s1) + 1) * (24 + (\text{len}(s2) + 1) * 4)$
	Вектор век- торов для представления матрицы	$(\text{len}(s1) + 1) * (24 + (\text{len}(s2) + 1) * 4) + 24$	
	Четыре пере- менные типа size_t для ите- рации в цикле и хранения длины строк	16	

Таблица 4.4: Теоретическая оценка памяти, потребляемой алгоритмами

В таблице 4.5 видно, что рекурсивный алгоритм Левенштейна наименее затратный по памяти. Самым затратным же алгоритмом является рекурсивный Левенштейн с заполнением матрицы, так как помимо потворного вызова рекурсивной функции, в самом начале также аллоцируется память под матрицу, которая в итоге и заполняется.

Алгоритм	Теоретическое потребле- ние памяти при $\text{len}(s1) = \text{len}(s2) = 10$
Матричный Левенштейн	876
Рекурсивный Левенштейн	768
Рекурсивный Левенштейн с матрицей	1540
Матричный Дамерау Ле- венштейн	876

Таблица 4.5: Теоретическое потребление памяти алгоритмами при строках длины 10

Взглянув на рисунок 6 можно сделать вывод, что матричный алгоритм Дамерау-Левенштейна работает медленнее матричного алгоритма Левенштейна, в силу того, что

в Дамерау на каждой итерации цикла нам нужно дополнительно проверять на истинность еще одно условие и в случае прохождения проверки производить дополнительные вычисления.

В силу того, что рекурсивный алгоритм работает достаточно медленно, график скорости работы двух рекурсивных алгоритмов был разбит на 2 рисунка. Сравнив рисунки 7 и 8 можно сделать вывод, что несмотря на то, что рекурсивный алгоритм Левенштейна с заполнением матрицы требует больше памяти, чем рекурсивный алгоритм без нее, благодаря использованию матрицы он значительно выигрывает в скорости, так как в ходе выполнения не происходит лишних вызовов на повторное вычисление значений, вместо этого они сразу берутся из матрицы.

Взглянув на рисунок 6 можно сделать вывод, что матричный алгоритм Дамерау-Левенштейна работает медленнее матричного алгоритма Левенштейна, в силу того, что в Дамерау на каждой итерации цикла нам нужно дополнительно проверять на истинность еще одно условие и в случае прохождения проверки производить дополнительные вычисления.

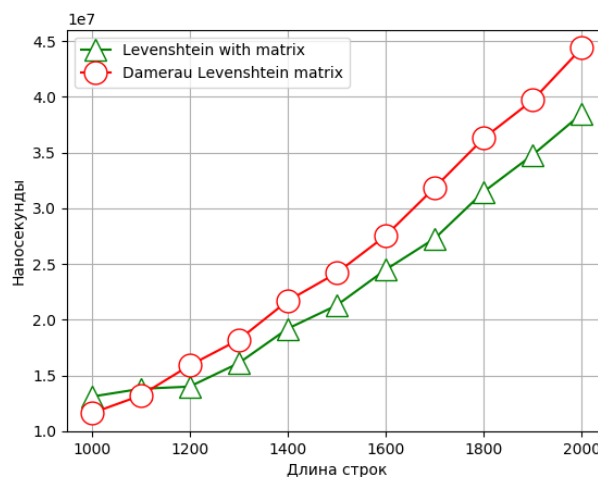


Рисунок 6: График работы матричных алгоритмов Левенштейна и Дамерау Левенштейна

В силу того, что рекурсивный алгоритм работает достаточно медленно, график скорости работы двух рекурсивных алгоритмов был разбит на 2 рисунка. Сравнив рисунки 7 и 8 можно сделать вывод, что несмотря на то, что рекурсивный алгоритм Левенштейна с заполнением матрицы требует больше памяти, чем рекурсивный алгоритм без нее, благодаря использованию матрицы он значительно выигрывает в скорости, так как в ходе выполнения не происходит лишних вызовов на повторное вычисление значений, вместо этого они сразу берутся из матрицы.

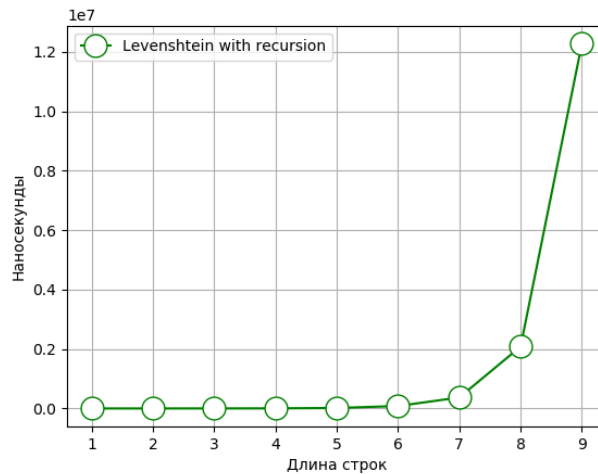


Рисунок 7: График работы рекурсивного алгоритма Левенштейна

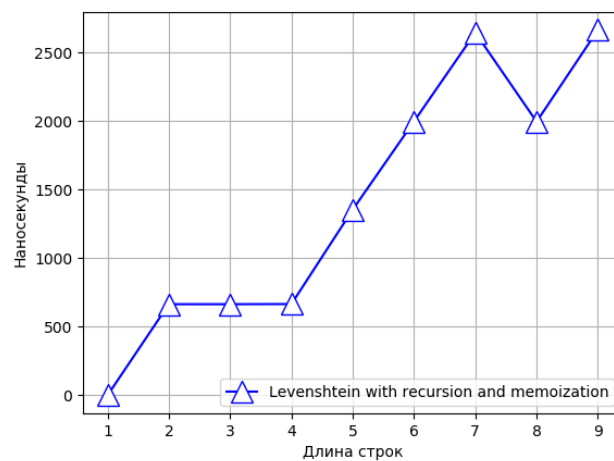


Рисунок 8: График работы рекурсивного алгоритма с заполнением матрицы Левенштейна

Сравнив рисунок 8 и рисунок 9 можно сделать вывод, что из трех алгоритмов Левенштейна самым медленным оказывается рекурсивный алгоритм, причем при увеличении длины слов появляется проблема того, что результат не может быть получен за приемлимое время. Однако как видно на рисунке 9 добавление матрицы существенно сокращает время выполнения.

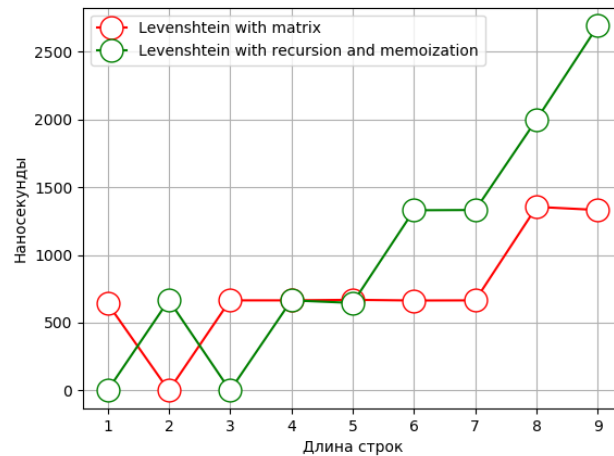


Рисунок 9: График работы алгоритмов Левенштейна

## Вывод

В данном разделе были продемонстрированы результаты экспериментов и сделаны выводы в соответствии с анализом их результатов. Было установлено, что рекурсивный алгоритм работает медленнее остальных, но требует меньше памяти. Самым быстрым алгоритмом оказался матричный Левенштейн.

# Заключение

В ходе выполнения лабораторной работы были выполнены следующие задачи:

1. изучил алгоритмы Левенштейна и Дameraу-Левенштейна нахождения расстояния между строками;
2. применил метод динамического программирования для матричной реализации указанных алгоритмов;
3. получил практические навыки реализации указанных алгоритмов: двух алгоритмов в матричной версии и одного из алгоритмов в рекурсивной версии;
4. провел сравнительный анализ линейной и рекурсивной реализаций выбранного алгоритма определения расстояния между строками по затрачиваемым ресурсам (времени и памяти);
5. экспериментально подтвердил различия во временной эффективности рекурсивной и нерекурсивной реализаций выбранного алгоритма определения расстояния между строками при помощи разработанного программного обеспечения на материале замеров процессорного времени выполнения реализации на варьирующихся длинах строк;
6. описал и обосновал полученные результаты в отчете о выполненной лабораторной работе, выполненного как расчётно-пояснительная записка к работе.

# Список использованной литературы

- [1] Скиена Стивен. Алгоритмы.Руководство к разработке. – М.:БХВ–Петербург, 2018. – 720 с.
- [2] Дэн Гасфилд. Строки, деревья и последовательности в алгоритмах. – М.:БХВ–Петербург, 2003. – 654 с.
- [3] Справка по C++ // cppreference URL: <https://en.cppreference.com/w/> (дата обращения: 06.10.2020).
- [4] Библиотека для замера времени std::chrono // cppreference URL: <https://en.cppreference.com/w/cpp/header/chrono> (дата обращения: 06.10.2020)