



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"

Тема Расстояние Левенштейна и Дамерау-Левенштейна

Студент Талышева О.Н.

Группа ИУ7-55Б

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

Москва — 2024 г.

Содержание

Введение	3
1 Аналитическая часть	4
1.1 Редакционное расстояние между двумя строками	4
1.2 Выравнивание строк	4
1.3 Расстояние Левенштейна	4
1.4 Расстояние Дамерау-Левенштейна	6
2 Конструкторская часть	8
3 Технологическая часть	21
3.1 Требования к программному обеспечению	21
3.2 Средства реализации	21
3.3 Реализации алгоритмов	21
3.4 Тесты	24
4 Исследовательская часть	26
4.1 Сравнение работы матричной, рекурсивной и рекурсивно-матричной ре- ализаций алгоритмов	27
4.2 Сравнение работы алгоритмов Левенштейна и Дамерау-Левенштейна (от- дельно каждый способ)	28
4.3 Сравнение работы матричных и рекурсивно-матричных алгоритмов Ле- венштейна и Дамерау-Левенштейна	30
Заключение	34
Список использованных источников	35

ВВЕДЕНИЕ

Цель лабораторной работы: исследовать алгоритмы вычисления расстояния Левенштейна и Дамерау-Левенштейна в матричной, рекурсивно-матричной и рекурсивной реализациях.

Для достижения этой цели были поставлены следующие задачи:

- рассмотреть алгоритм вычисления расстояния Левенштейна;
- рассмотреть алгоритм вычисления расстояния Дамерау-Левенштейна;
- применить метод динамического программирования для матричных реализаций алгоритмов;
- сравнить матричную, рекурсивно-матричную и рекурсивную реализации алгоритмов;
- сравнить алгоритмы вычисления расстояния Левенштейна и Дамерау-Левенштейна.

1. Аналитическая часть

1.1. Редакционное расстояние между двумя строками

Часто требуется измерить различие или расстояние между двумя строками (например, в эволюционных, структуральных или функциональных исследованиях биологических строк, в хранении текстовых баз данных, в методах проверки правописания). Есть несколько способов формализации понятия расстояния между строками. Одна общая и простая, формализация называется редакционным расстоянием; она основана на преобразовании (или редактировании) одной строки в другую серией операций редактирования, выполняемых над отдельными символами. Разрешенные операции редактирования — это вставка (I - insertion) символа в первую строку, удаление (D - deletion) символа из первой строки и подстановка или замена (substitution или, лучше, R - replace) символа из первой строки символом из второй строки. Обозначим M — “не-операцию” над правильной буквой (от match).

Строка над алфавитом Σ , D, R, M, которая описывает преобразование одной строки в другую, называется редакционным предписанием (предписанием) этих двух строк.

Редакционное расстояние между двумя строками определяется как минимальное число редакционных операций — вставок, удалений и подстановок, необходимое для преобразования первой строки во вторую.

Подчеркнем, что совпадения операциями не являются и не засчитываются. Редакционное расстояние иногда называют расстоянием Левенштейна по статье В. Левенштейна, где оно рассматривалось, вероятно, впервые.[4]

1.2. Выравнивание строк

Редакционное предписание — это способ представления конкретного преобразования одной строки в другую. Альтернативный (и часто предпочтительный) способ заключается в показе явного выравнивания (alignment) этих двух строк. (Глобальное) выравнивание двух строк, S1 и S2, получается вставкой пробелов в строки S1 и S2 (возможно, на их концах) и размещением двух получившихся строк друг над другом так, чтобы каждый символ или пробел одной строки оказался напротив одного символа или пробела другой строки.

Термин «глобальный» подчёркивает, что обе строки участвуют в выравнивании полностью.[3]

1.3. Расстояние Левенштейна

Расстояние Левенштейна, или редакционное расстояние, — метрика сходства между двумя строковыми последовательностями. Чем больше расстояние, тем более различны строки. По сути, это минимальное число односимвольных преобразований

(удаления, вставки или замены), необходимых, чтобы превратить одну последовательность в другую.

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

- $D(a, b)$ — цена замены символа a на символ b
- $D(\lambda, b)$ — цена вставки символа b
- $D(a, \lambda)$ — цена удаления символа a

Необходимо найти последовательность замен, минимизирующую суммарную цену. Расстояние Левенштейна является частным случаем этой задачи при ценах:

- $D(a, a) = 0$
- $D(a, b) = 1$, при $a \neq b$
- $D(\lambda, b) = 1$
- $D(a, \lambda) = 1$

Как частный случай, так и задачу для произвольных D , решает алгоритм Вагнера — Фишера, приведённый ниже. Здесь и ниже считается, что все D неотрицательны, и действует неравенство треугольника: замена двух последовательных операций одной не увеличит общую цену (например, замена символа x на y , а потом y на z не лучше, чем сразу x на z).

Например, $D(\text{'hello'}, \text{'hallo'}) = 1$, так как потребуется провести одну замену 'e' на 'a'.

Алгоритм реализуется по следующей формуле:

$$d(S_1, S_2) = D(M, N), \text{ где } D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & i > 0, j = 0 \\ j, & i = 0, j > 0 \\ \min \begin{cases} D(i-1, j) + 1 & \text{(удаление)} \\ D(i, j-1) + 1 & \text{(вставка)} \\ D(i-1, j-1) + 1_{\text{если } a_i \neq b_j} & \text{(замена)} \end{cases} & i > 0, j > 0 \end{cases} \quad (1)$$

Таким образом, требуется вычислить матрицу расстояний размерностью $\text{len}(str_1) * \text{len}(str_2)$, следовательно, объем требуемой памяти растет как $O(\text{len}(str_1) * \text{len}(str_2))$. Иными словами, для двух мегабайтных строк потребуются гигабайты памяти. Фактически

в кэше будет храниться почти все матрица редактирований, а она не нужна целиком. Искомая цель – правый нижний элемент.

		Л	А	Б	Р	А	Д	О	Р
	0	1	2	3	4	5	6	7	8
Г	1	1	2	3	4	5	6	7	8
И	2	2	2	3	4	5	6	7	8
Б	3	3	3	2	3	4	5	6	7
Р	4	4	4	3	2	3	4	5	6
А	5	5	4	4	3	2	3	4	5
Л	6	5	5	5	4	3	3	4	5
Т	7	6	6	6	5	4	4	4	5
А	8	7	6	7	6	5	5	5	5
Р	9	8	7	7	7	6	6	6	5

Рисунок 1 – Пример нахождения расстояния Левенштейна

Для его поиска можно обойтись лишь парой рядов: текущим и предыдущим. А остальные ряды не хранить в памяти. Так будет достигнут конец таблицы, и нижний правый угол и будет искомым значением.

Чтобы использовать еще меньше памяти, можно поменять местами строки, чтобы длина рядов была минимальна. Это существенно экономит память, если одна из строк длинная, а другая короткая.

1.4. Расстояние Дамерау-Левенштейна

Если к списку разрешённых операций добавить транспозицию (два соседних символа меняются местами), получается расстояние Дамерау — Левенштейна. Для неё также существует алгоритм, требующий $O(\text{len}(\text{str1}) * \text{len}(\text{str2}))$ операций. Дамерау показал, что 80% ошибок при наборе текста человеком являются транспозициями. Кроме того, расстояние Дамерау-Левенштейна используется и в биоинформатике.

Цена операции транспозиция также равна 1. При работе алгоритма Левенштейна эта операция реализовалась бы двумя заменами и стоила бы 2. Таким образом, расстояние Дамерау-Левенштейна в некоторых случаях даёт меньший результат, чем расстояние Левенштейна.

В формулу 1 добавляется следующая часть:

$$\begin{cases} i > 1 \\ j > 1 \\ \text{str}_1[i-1] = \text{str}_2[j] \\ \text{str}_1[i] = \text{str}_2[j-1] \end{cases} \quad (2)$$

В результате получается следующая формула для алгоритма Дамерау-Левенштейна:

$$d(S_1, S_2) = D(M, N), \text{ где } D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i, & i > 0, j = 0 \\ j, & i = 0, j > 0 \\ \min \begin{cases} D(i-1, j) + 1 & (\text{удаление}) \\ D(i, j-1) + 1 & (\text{вставка}) \\ D(i-1, j-1) + 1_{\text{если } a_i \neq b_j} & (\text{замена}) \\ \begin{cases} i > 1 \\ j > 1 \\ \text{str}_1[i-1] = \text{str}_2[j] \\ \text{str}_1[i] = \text{str}_2[j-1] \end{cases} & (\text{транспозиция}) \end{cases} & i > 0, j > 0 \end{cases} \quad (3)$$

2. Конструкторская часть

Описания алгоритмов

На основании теоретических измышлений были разработаны алгоритмы, вычисляющие расстояние Левенштейна и Дameraу-Левенштейна тремя способами: матричным, рекурсивным и рекурсивно-матричным. Блок-схемы этих алгоритмов приведены на рисунках 2, 3, 4, 5, 6, 7 (Левенштейн) и рисунках 8, 9, 10, 11, 12, 13 (Дameraу-Левенштейн).

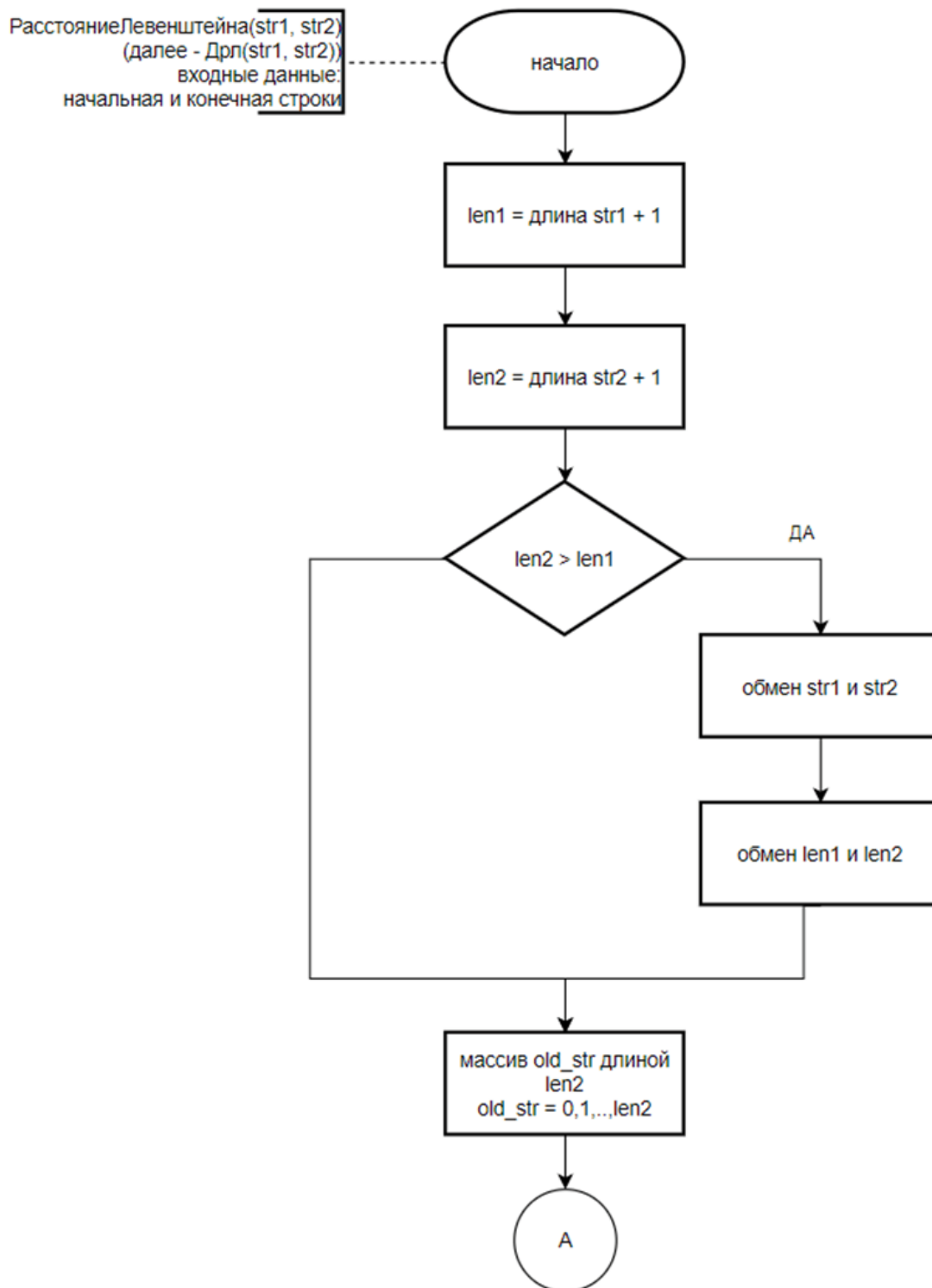


Рисунок 2 – Блок-схема алгоритма Левенштейна (матричная реализация)

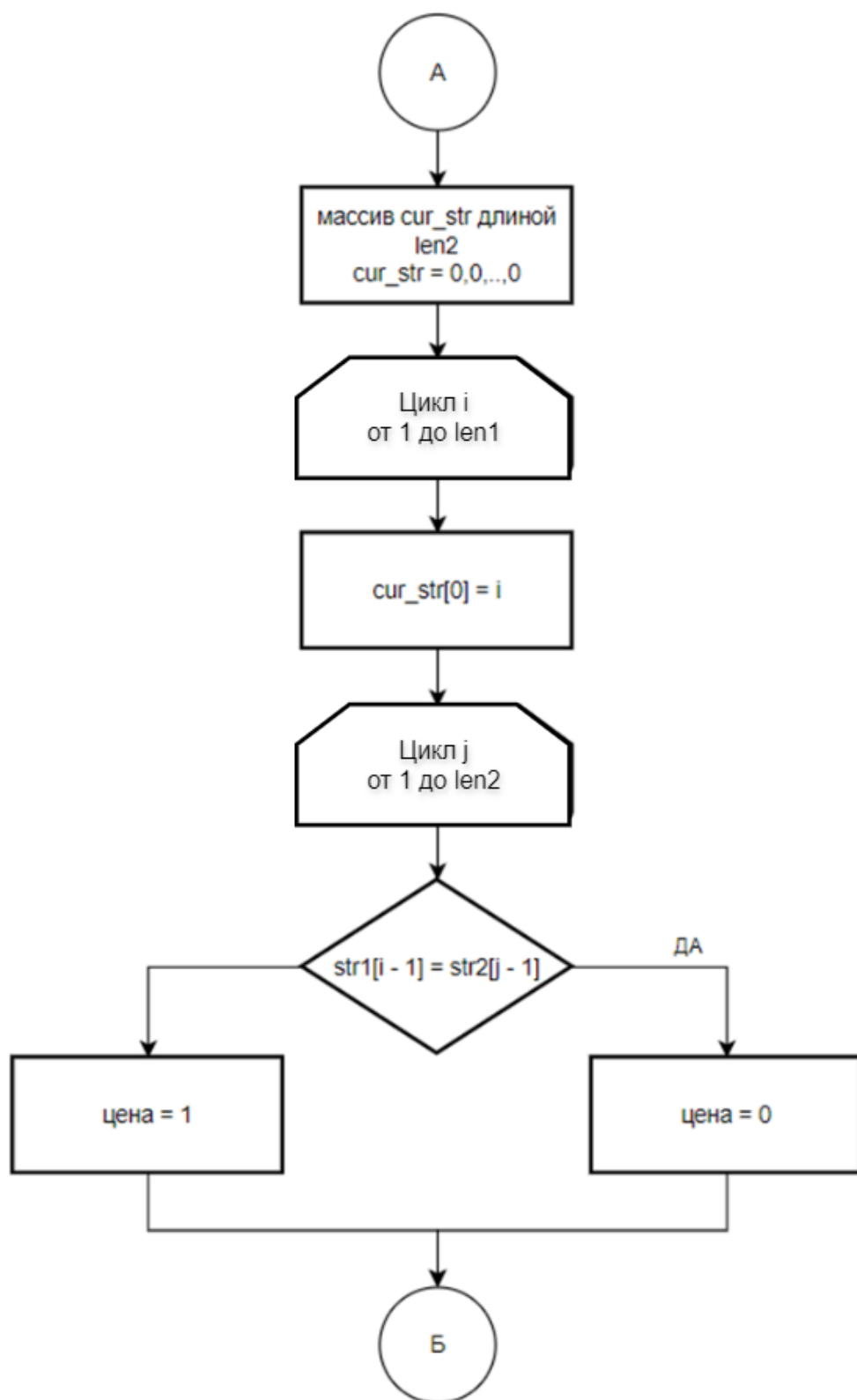


Рисунок 3 – Блок-схема алгоритма Левенштейна (матричная реализация)
(продолжение)

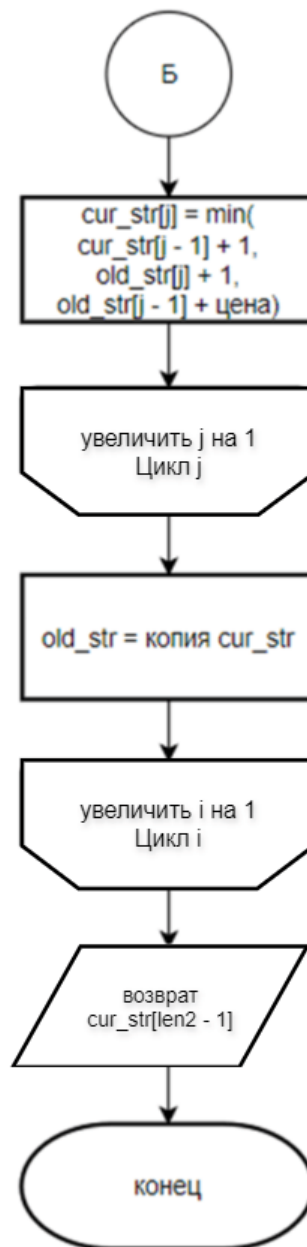


Рисунок 4 – Блок-схема алгоритма Левенштейна (матричная реализация)
(продолжение (2))

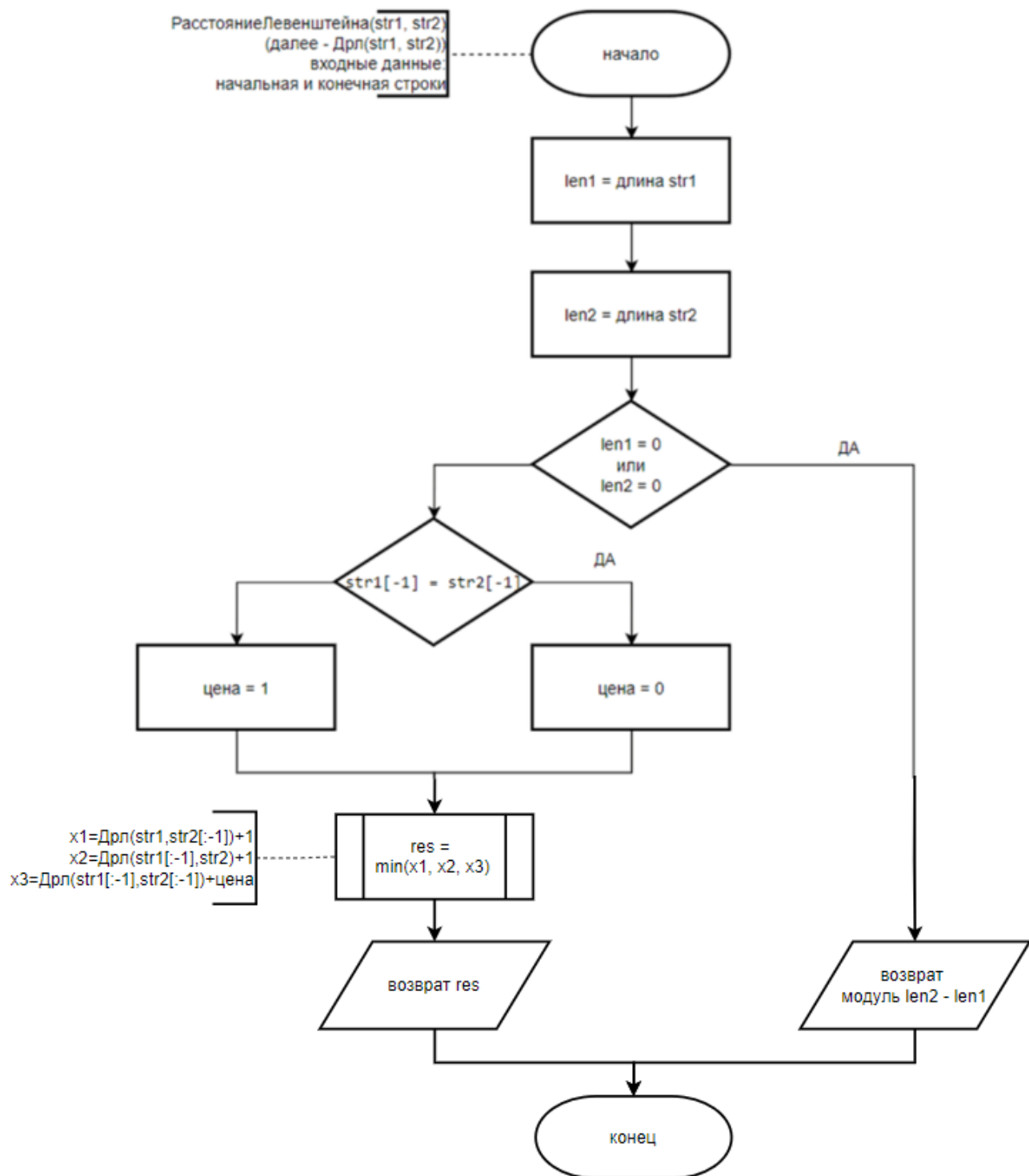


Рисунок 5 – Блок-схема алгоритма Левенштейна (рекурсивная реализация)

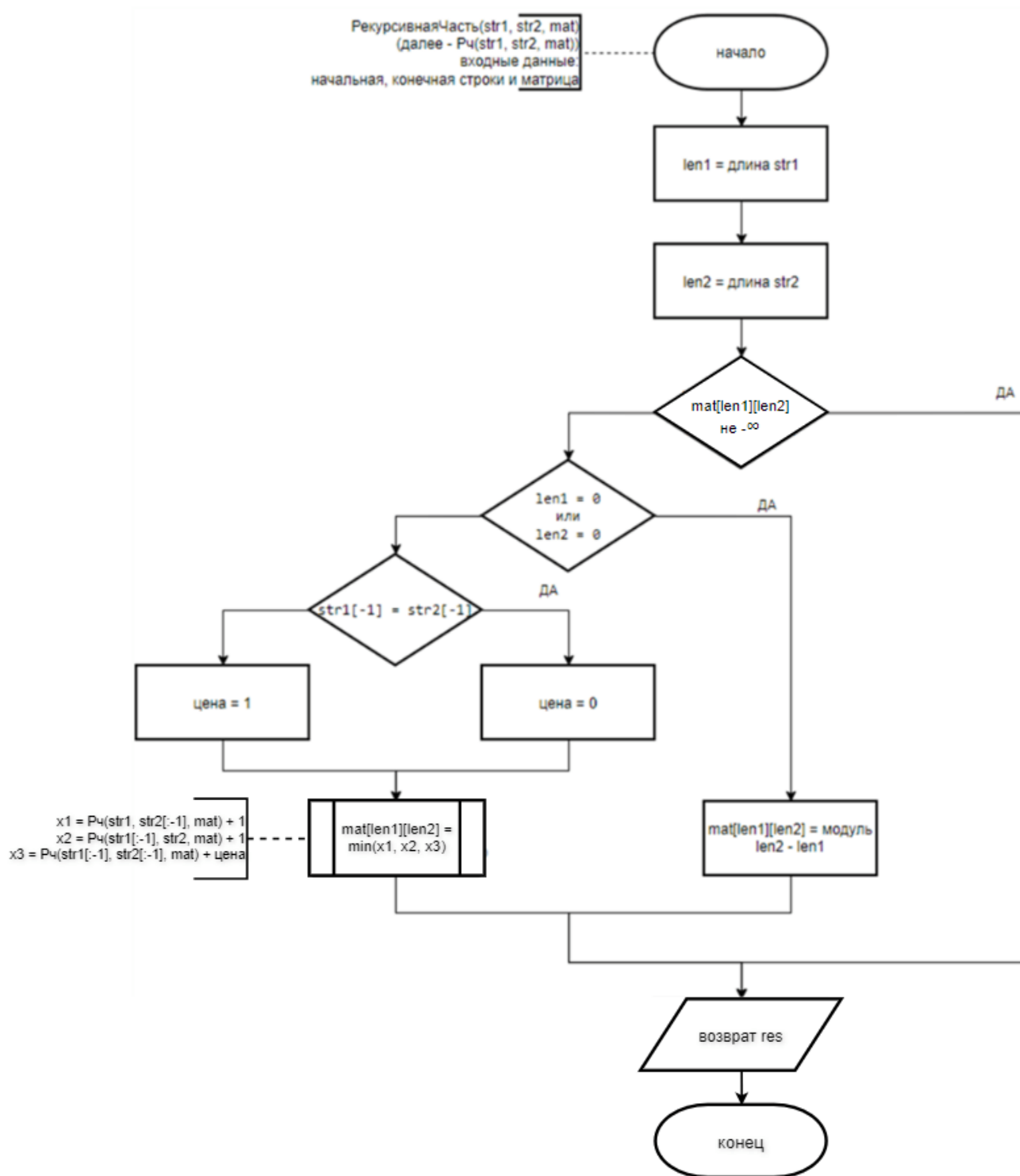


Рисунок 6 – Блок-схема алгоритма Левенштейна (рекурсивно-матричная реализация)

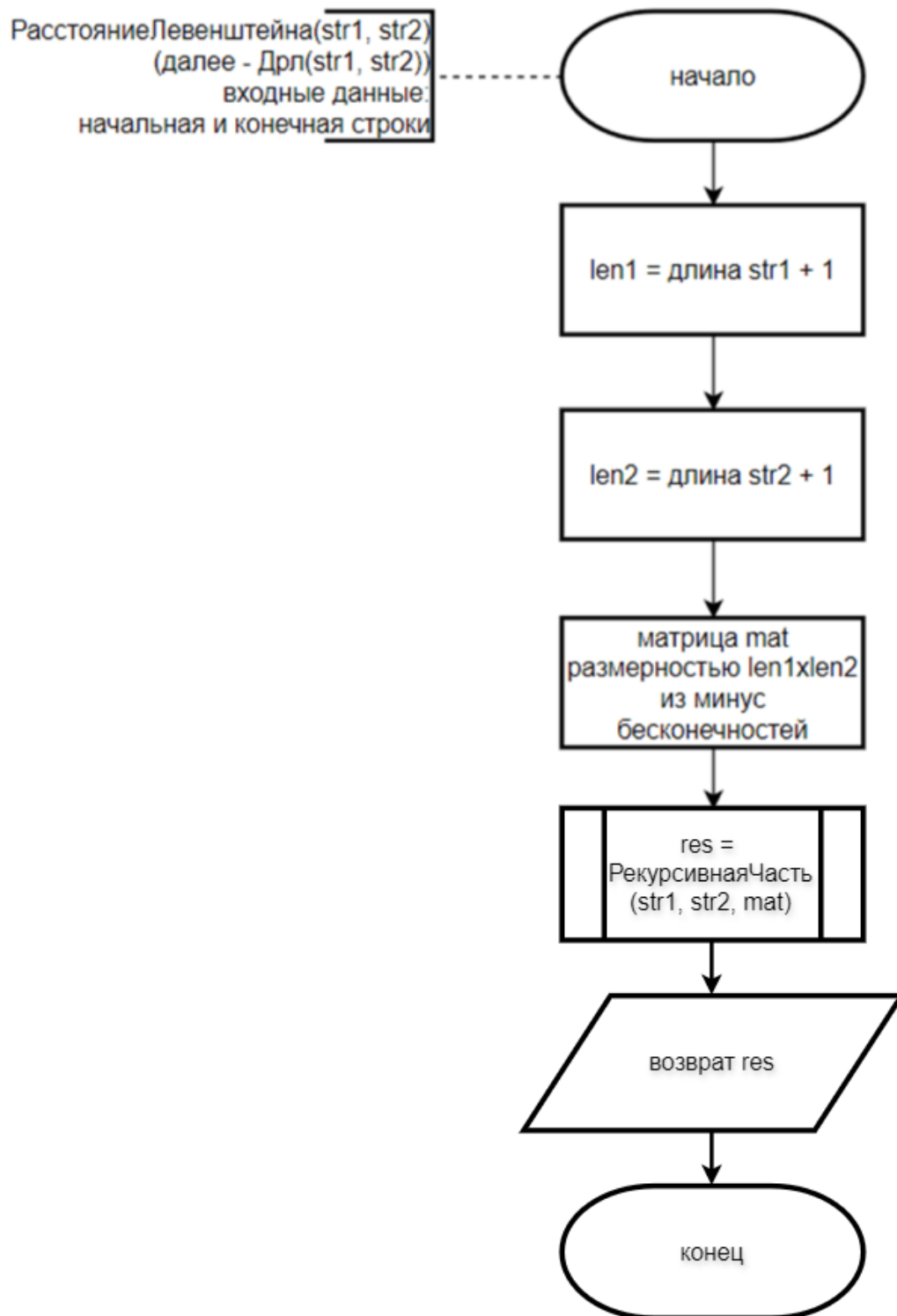


Рисунок 7 – Блок-схема алгоритма Левенштейна (рекурсивно-матричная реализация (продолжение))

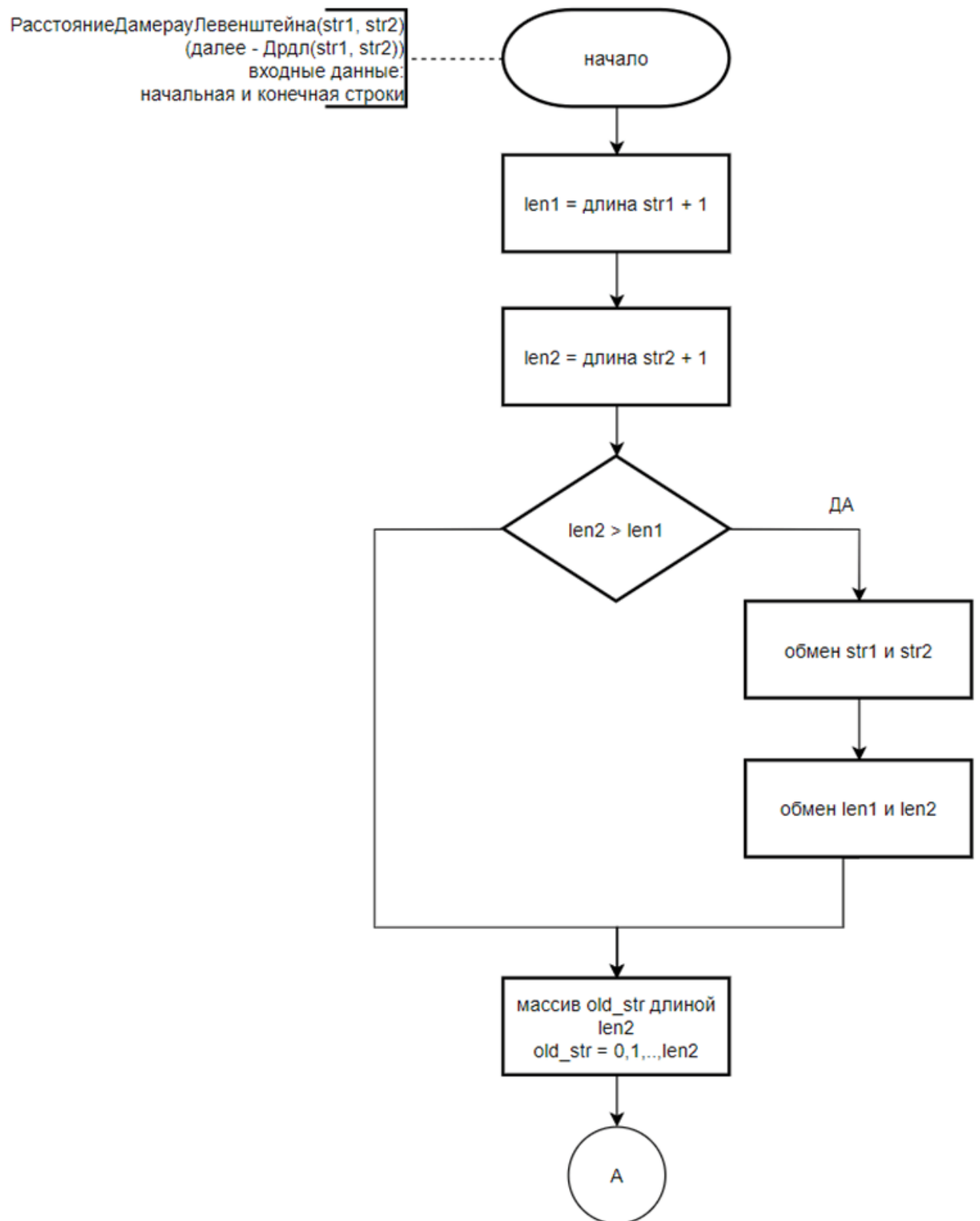


Рисунок 8 – Блок-схема алгоритма Дамерау-Левенштейна (матричная реализация)

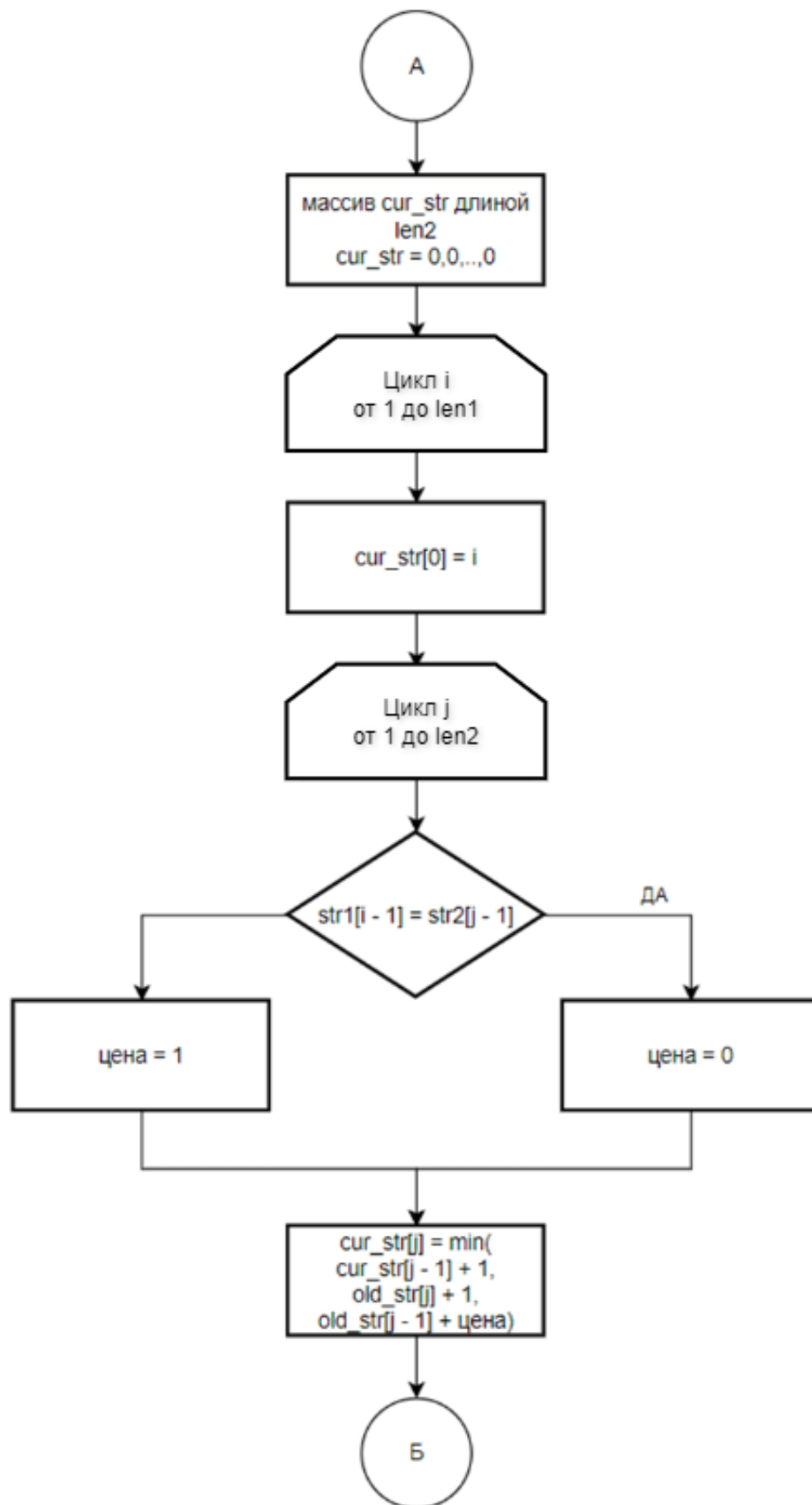


Рисунок 9 – Блок-схема алгоритма Дамерау-Левенштейна (матричная реализация)
(продолжение)

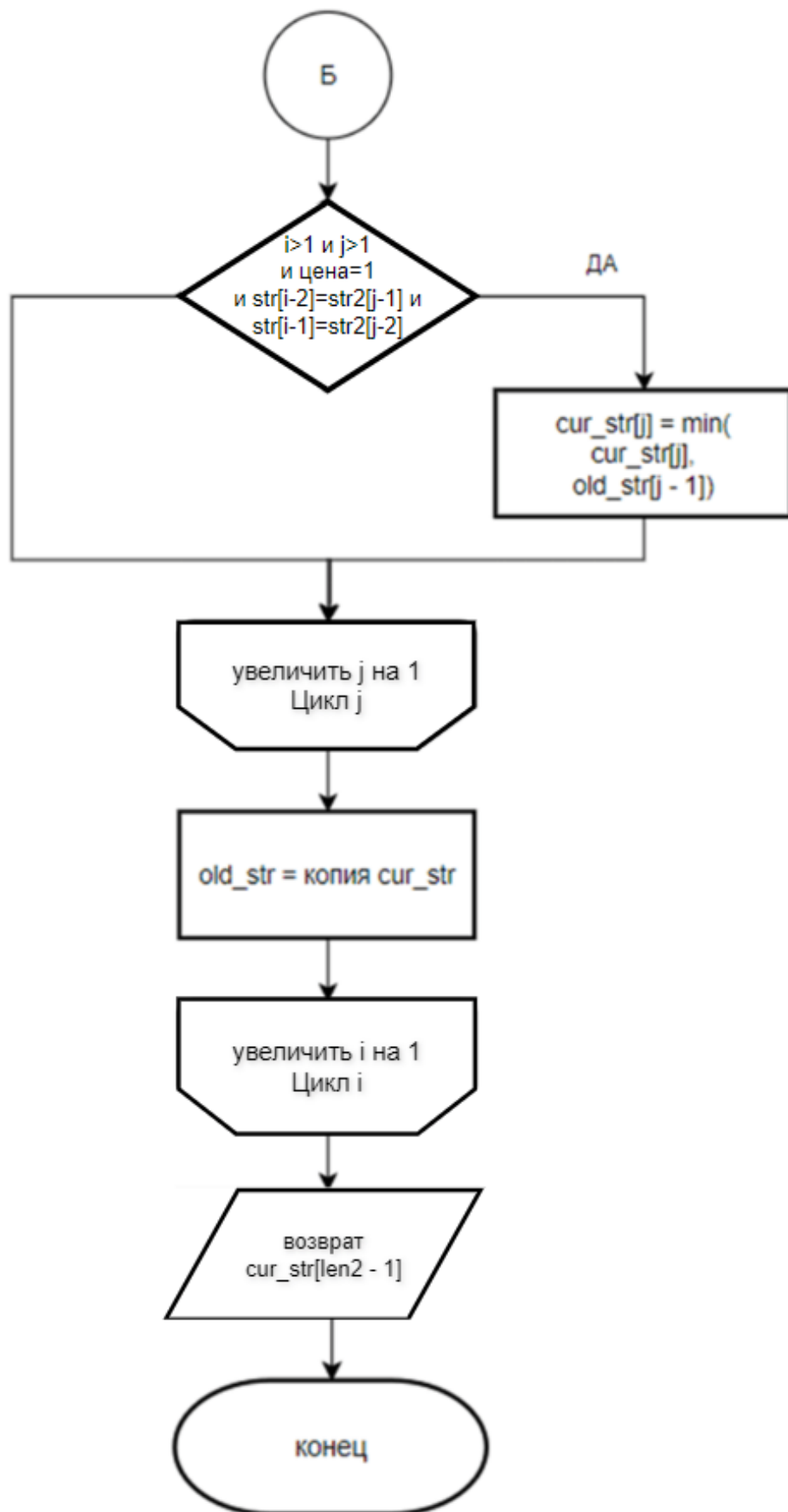


Рисунок 10 – Блок-схема алгоритма Дамерау-Левенштейна (матричная реализация)
(продолжение (2))

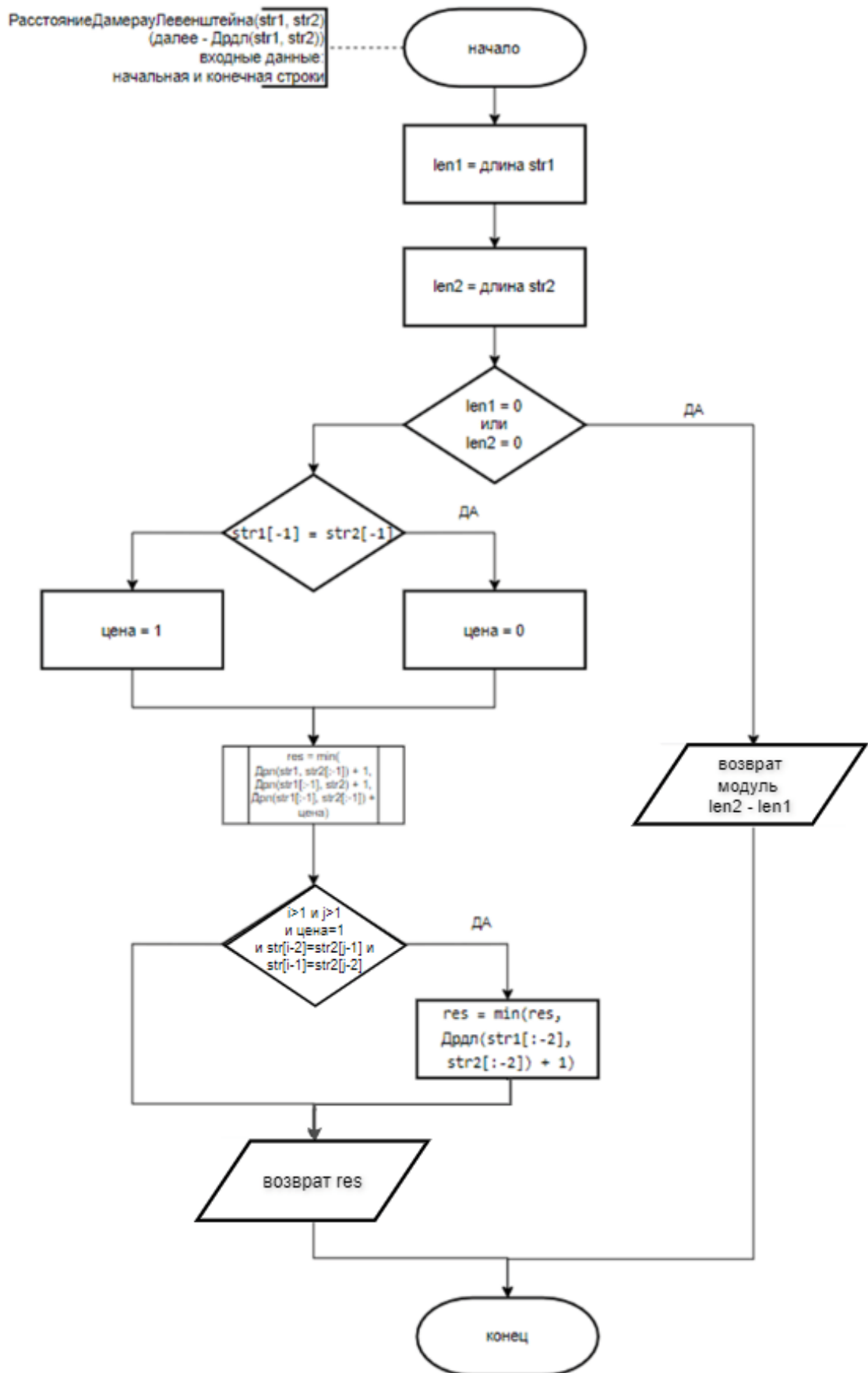


Рисунок 11 – Блок-схема алгоритма Дameraу-Левенштейна (рекурсивная реализация)

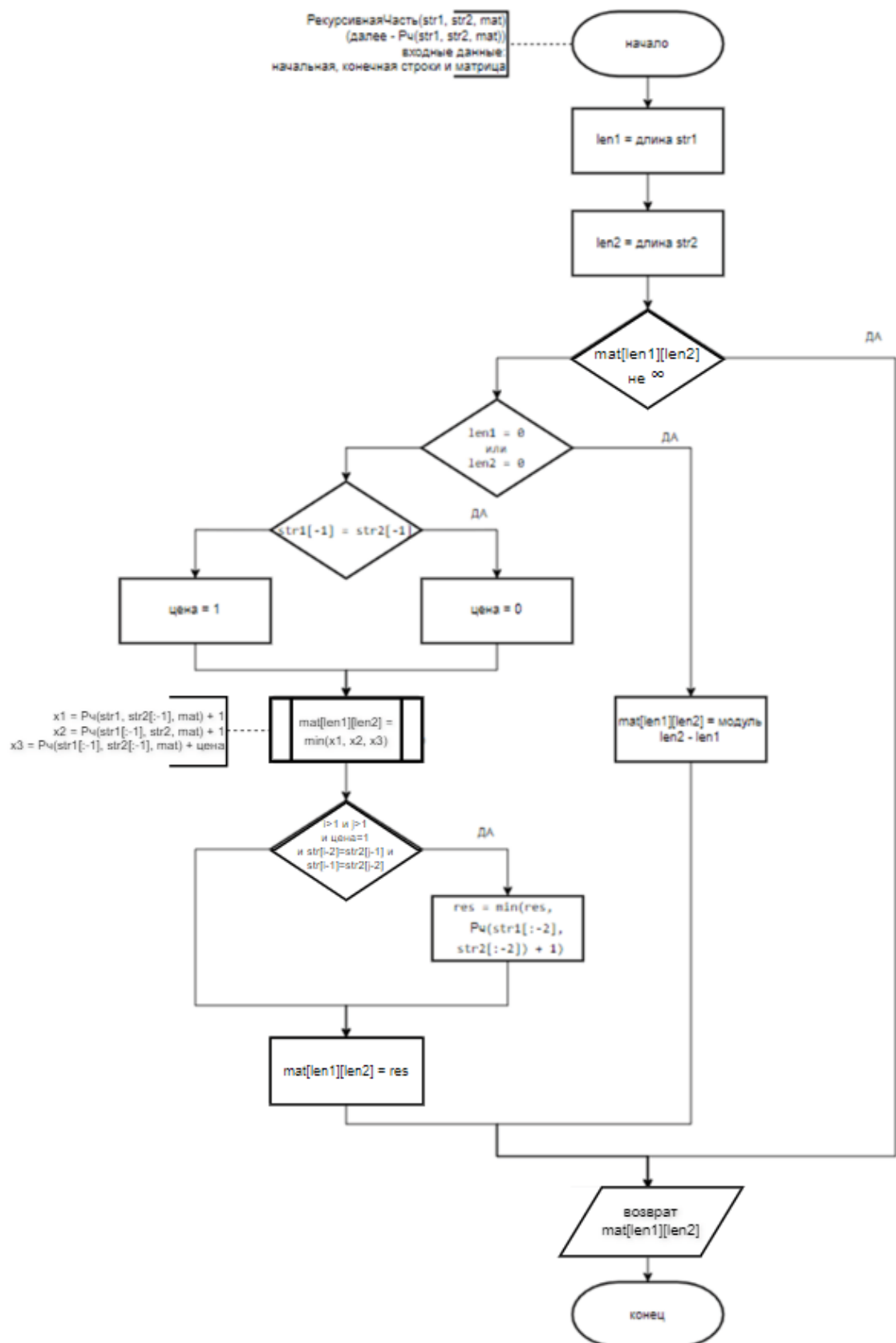


Рисунок 12 – Блок-схема алгоритма Дамерау-Левенштейна (рекурсивно-матричная реализация)

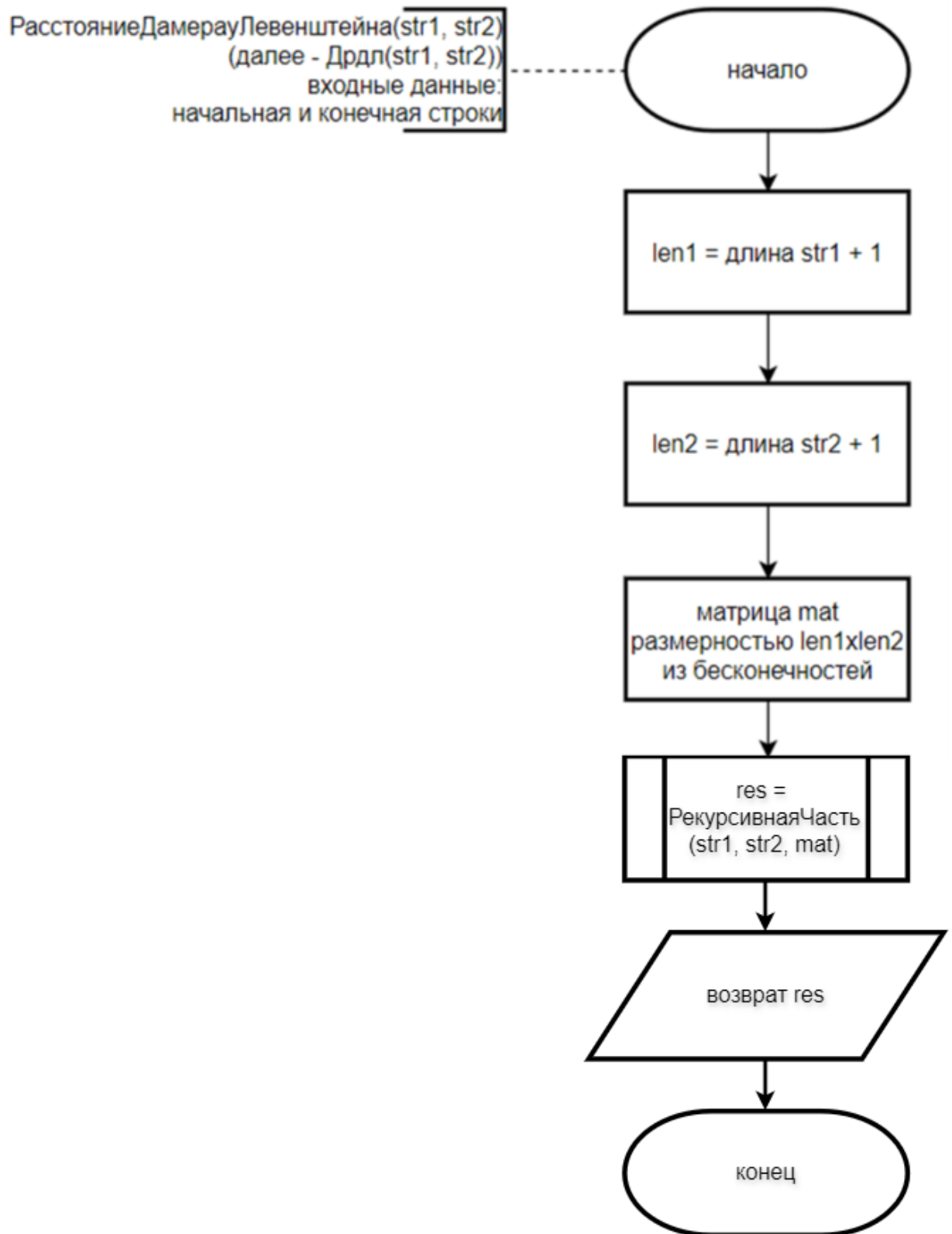


Рисунок 13 – Блок-схема алгоритма Дameraу-Левенштейна (рекурсивно-матричная реализация (продолжение))

3. Технологическая часть

3.1. Требования к программному обеспечению

На вход программе подаются 2 строки из символов, которые входят в таблицу Юникода (UTF-8).

На выход программа выдаёт число – расстояние между строками, вычисленное алгоритмом Левенштейна или Дамерау-Левенштейна матричной или рекурсивной реализацией. Для матричных реализаций также выводится матрица расстояний. Также в зависимости от выбранного пункта меню программа замеряет время работы алгоритмов и рисует получившиеся графики.

3.2. Средства реализации

Python эффективно обрабатывает строки, поэтому программа была реализована на этом языке программирования. Для замеров времени была использована функция `process_time()` из библиотеки `time`, вычисляющая процессорное время [1].

3.3. Реализации алгоритмов

Ниже приведены реализации алгоритмов поиска расстояния Левенштейна (листинги 1, 2 и 3) и Дамерау-Левенштейна (листинги 4, 5 и 6) матричным, рекурсивным и рекурсивно-матричным способами на Python.

```
1 def algo_Levenstein_matrix(str1: str, str2: str) -> int:
2     len1, len2 = len(str1) + 1, len(str2) + 1
3     if len2 > len1:
4         str1, str2 = str2, str1
5         len1, len2 = len2, len1
6     old_str = [i for i in range(len2)]
7     cur_str = [0 for _ in range(len2)]
8     for i in range(1, len1):
9         cur_str[0] = i
10        for j in range(1, len2):
11            cur_str[j] = min(cur_str[j - 1] + 1,
12                            old_str[j] + 1,
13                            old_str[j - 1] + (str1[i - 1] != str2[j - 1]))
14        old_str = cur_str.copy()
15    return cur_str[-1]
```

Листинг 1 – реализация матричного алгоритма Левенштейна

```
1 def algo_Levenstein_recursion(str1: str, str2: str) -> int:
2     len1, len2 = len(str1), len(str2)
3     if len1 * len2 == 0:
```

```

4         return abs(len2 - len1)
5     return min(algo_Levenstein_recursion(str1, str2[:-1]) + 1,
6               algo_Levenstein_recursion(str1[:-1], str2) + 1,
7               algo_Levenstein_recursion(str1[:-1], str2[:-1]) + (str1
8               [-1] != str2[-1]))

```

Листинг 2 – реализация рекурсивного алгоритма Левенштейна

```

1 def algo_Levenstein_recursion_matrix(str1: str, str2: str) -> int:
2     len1, len2 = len(str1) + 1, len(str2) + 1
3     mat = [[float("-inf") for i in range(len2)] for j in range(len1)]
4     # the recursive part itself
5     def recursion_part(str1: str, str2: str, mat: List[float] = []) -> int:
6         :
7         len1, len2 = len(str1), len(str2)
8         if mat[len1][len2] > float("-inf"):
9             pass
10        elif len1 * len2 == 0:
11            mat[len1][len2] = abs(len2 - len1)
12        else:
13            mat[len1][len2] = min(recursion_part(str1, str2[:-1], mat) +
14                                1,
15                                recursion_part(str1[:-1], str2, mat) + 1,
16                                recursion_part(str1[:-1], str2[:-1], mat) + (str1[-1]
17                                != str2[-1]))
18        return mat[len1][len2]
19    return recursion_part(str1, str2, mat)

```

Листинг 3 – реализация рекурсивно-матричного алгоритма Левенштейна

```

1 def algo_Damerau_Levenstein_matrix(str1: str, str2: str) -> int:
2     len1, len2 = len(str1) + 1, len(str2) + 1
3     if len2 > len1:
4         str1, str2 = str2, str1
5         len1, len2 = len2, len1
6     old_str = [i for i in range(len2)]
7     cur_str = [0 for _ in range(len2)]
8     for i in range(1, len1):
9         cur_str[0] = i
10        for j in range(1, len2):
11            m = str1[i - 1] != str2[j - 1]
12            cur_str[j] = min(cur_str[j - 1] + 1,
13                            old_str[j] + 1,
14                            old_str[j - 1] + m)
15            if (i > 1) and (j > 1) and m and (str1[i - 2] == str2[j - 1])
16                and (str1[i - 1] == str2[j - 2]):
17                cur_str[j] = min(cur_str[j], old_str[j - 1])

```

```

17     old_str = cur_str.copy()
18     return cur_str[-1]

```

Листинг 4 – реализация матричного алгоритма Дамерау-Левенштейна

```

1 def algo_Damerau_Levenstein_recursion(str1: str, str2: str) -> int:
2     len1, len2 = len(str1), len(str2)
3     if len1 * len2 == 0:
4         return abs(len2 - len1)
5     res = min(algo_Damerau_Levenstein_recursion(str1, str2[:-1]) + 1,
6               algo_Damerau_Levenstein_recursion(str1[:-1], str2) + 1,
7               algo_Damerau_Levenstein_recursion(str1[:-1], str2[:-1]) +
8               (str1[-1] != str2[-1]))
9     if ((len(str1) >= 2) and (len(str2) >= 2) and (str1[-1] == str2[-2])
10         and (str1[-2] == str2[-1])):
11         res = min(res, algo_Damerau_Levenstein_recursion(str1[:-2], str2
12                                                            [:-2]) + 1)
13     return res

```

Листинг 5 – реализация рекурсивного алгоритма Дамерау-Левенштейна

```

1 def algo_Damerau_Levenshtein_recursion_matrix(str1: str, str2: str) -> int
2 :
3     len1, len2 = len(str1) + 1, len(str2) + 1
4     mat = [[float("inf") for i in range(len2)] for j in range(len1)]
5     # the recursive part itself
6     def recursion_part(str1: str, str2: str, mat: List[float] = []) -> int
7     :
8         len1, len2 = len(str1), len(str2)
9         if mat[len1][len2] < float("inf"):
10             pass
11         elif len1 * len2 == 0:
12             mat[len1][len2] = abs(len2 - len1)
13         else:
14             mat[len1][len2] = min(recursion_part(str1, str2[:-1], mat) +
15                                   1,
16                                   recursion_part(str1[:-1], str2, mat) + 1,
17                                   recursion_part(str1[:-1], str2[:-1], mat) + (str1[-1]
18                                   != str2[-1]))
19             if ((len(str1) >= 2) and (len(str2) >= 2) and (str1[-1] ==
20                                                             str2[-2]) and (str1[-2] == str2[-1])):
21                 mat[len1][len2] = min(mat[len1][len2], recursion_part(str1
22                                                                          [:-2], str2[:-2], mat) + 1)
23         return mat[len1][len2]
24     return recursion_part(str1, str2, mat)

```

Листинг 6 – реализация рекурсивно-матричного алгоритма Дамерау-Левенштейна

3.4. Тесты

Для тестирования алгоритмов Левенштейна и Дамерау-Левенштейна были составлены таблицы с входными данными (2 строки, возможно, пустые - λ), ожидаемым результатом (расстоянием) и полученным результатом от всех трёх способов.

Строка 1	Строка 2	Ожидание	матричный	рекурсивный	рекурсивно-матричный
λ	λ	0	0	0	0
a	a	0	0	0	0
abc	abc	0	0	0	0
λ	a	1	1	1	1
a	λ	1	1	1	1
a	b	1	1	1	1
abc	abs	1	1	1	1
odc	abc	2	2	2	2
ods	abc	3	3	3	3
abcs	abc	1	1	1	1
bc	abc	1	1	1	1
bac	abc	2	2	2	2

Таблица 1 – Таблица тестов для алгоритмов Левенштейна

Строка 1	Строка 2	Ожидание	матричный	рекурсивный	рекурсивно-матричный
λ	λ	0	0	0	0
a	a	0	0	0	0
abc	abc	0	0	0	0
λ	a	1	1	1	1
a	λ	1	1	1	1
a	b	1	1	1	1
abc	abs	1	1	1	1
odc	abc	2	2	2	2
ods	abc	3	3	3	3
abcs	abc	1	1	1	1
bc	abc	1	1	1	1
bac	abc	1	1	1	1

Таблица 2 – Таблица тестов для алгоритмов Дамерау-Левенштейна

В ходе проведённого тестирования (с помощью pytest) ошибок в алгоритмах не выявлено:

```

1 pytest
2 ===== test session starts =====
3 platform win32 -- Python 3.11.9, pytest-8.2.2, pluggy-1.5.0
4 rootdir: L:\sem_5\Algorithm_analysis\lab_01
5 plugins: Faker-28.4.1

```



```
6 collected 72 items
7
8 test_algo.py ..... [100%]
9
10 ===== 72 passed in 0.69s =====
```

Листинг 7 – тестирование алгоритмов с помощью pytest

4. Исследовательская часть

Для сравнения времени реализаций алгоритмов Левенштейна и Дамерау-Левенштейна в их матричных, рекурсивных и рекурсивно-матричных реализациях программа была запущена на случайно сгенерированных строках длинами от 1 до 9 с шагом 2 по 50 замеров каждая строка, среднее значение было вынесено в таблицу и для наглядности изображено на графике (рисунок 14 и таблица 3)

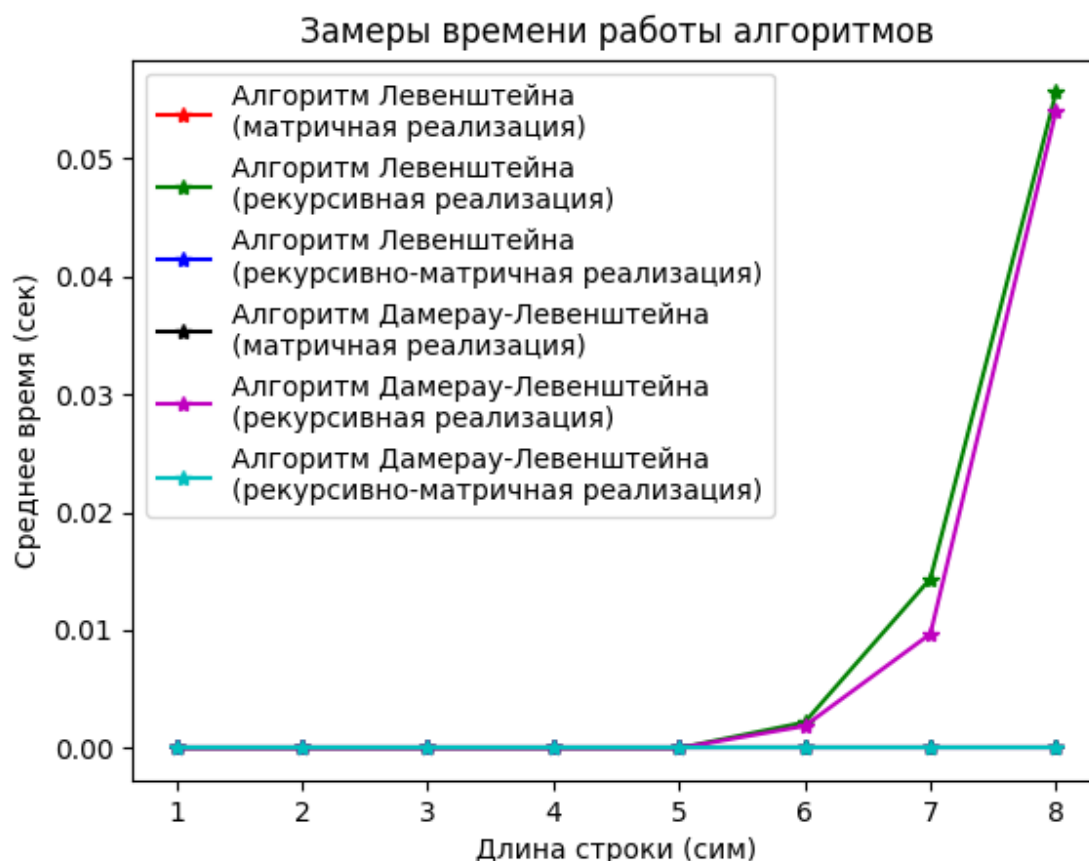


Рисунок 14 – График времени работы всех алгоритмов в зависимости от длин строк

Алгоритм	1	2	3	4	5	6	7	8
Левенштейна (матричный)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Левенштейна (рекурсивный)	0.0	0.0	0.0	0.0	0.0	0.0022	0.014	0.056
Левенштейна (рекурс-мат.)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Дамерау-Левенштейна (матричный)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Дамерау-Левенштейна (рекурсивный)	0.0	0.0	0.0	0.0	0.0	0.0019	0.0097	0.054
Дамерау-Левенштейна (рекурс-мат.)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Таблица 3 – Таблица времени (сек) работы всех алгоритмов в зависимости от длин строк (сим)

Кроме того, замеры времени работы всех алгоритмов были проведены на микроконтроллерах STM32. На графике 15 и в таблице 4 приведены замеры на длинах строк

от 1 символа до 6 с шагом 1 и запуском каждого по 10 раз.

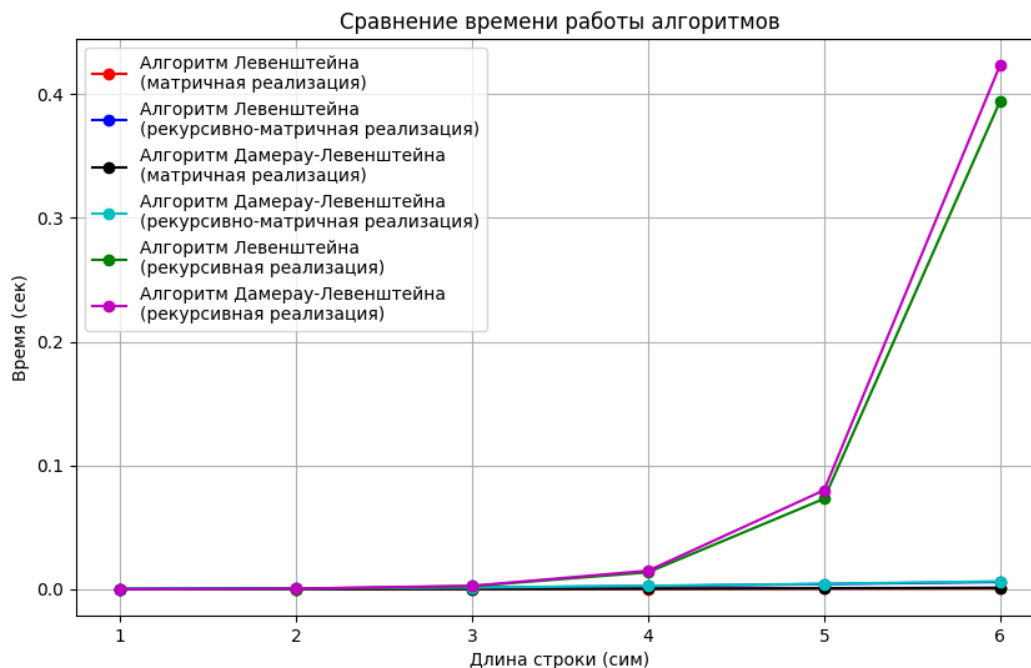


Рисунок 15 – График времени работы всех алгоритмов в зависимости от длин строк (замеры на микроконтроллерах)

Алгоритм	1	2	3	4	5	6
Левенштейна (матричный)	0.0001	0.0002	0.0003	0.0005	0.0007	0.0009
Левенштейна (рекурсивный)	0.0001	0.0005	0.0026	0.0138	0.0732	0.3941
Левенштейна (рекурс-мат.)	0.0003	0.0007	0.0015	0.0026	0.0043	0.0060
Дамерау-Левенштейна (матрич.)	0.0000	0.0002	0.0004	0.0006	0.0009	0.00013
Дамерау-Левенштейна (рекурс.)	0.0001	0.0006	0.0029	0.0151	0.0801	0.4233
Дамерау-Левенштейна (рек-мат.)	0.0003	0.0007	0.0015	0.0029	0.0045	0.0065

Таблица 4 – Таблица времени (сек) работы всех алгоритмов в зависимости от длин строк (сим) (замеры на микроконтроллерах)

4.1. Сравнение работы матричной, рекурсивной и рекурсивно-матричной реализаций алгоритмов

Из графиков, приведённых выше, очевидно, что матричная реализация обоих алгоритмов быстро становится эффективнее рекурсивной на много порядков. Это происходит из-за того, что при рекурсии даже на небольшой длине строк происходит много рекурсивных вызовов для подстрок, на что тратится большое количество времени и памяти. В то время как для матричной реализации данные, на основе которых вычисляются следующие значения, хранятся в двух массивах длиной в кратчайшую из двух строк, что экономит как время, так и память. При этом рекурсивно-матричная реализация оказалась почти столь же быстрой, как и матричная благодаря исключению

повторных вычислений идентичных веток рекурсии, что в разы сократило количество вычислений.

4.2. Сравнение работы алгоритмов Левенштейна и Дамерау-Левенштейна (отдельно каждый способ)

Отдельно было измерено время работы алгоритмов Левенштейна и Дамерау-Левенштейна в их матричных реализациях на большем диапазоне длин случайно сгенерированных строк (от 25 символов до 125 с шагом 25) по 50 замеров каждая строка, среднее значение было вынесено в таблицу и для наглядности изображено на графике (рисунок 16 и таблица 5).

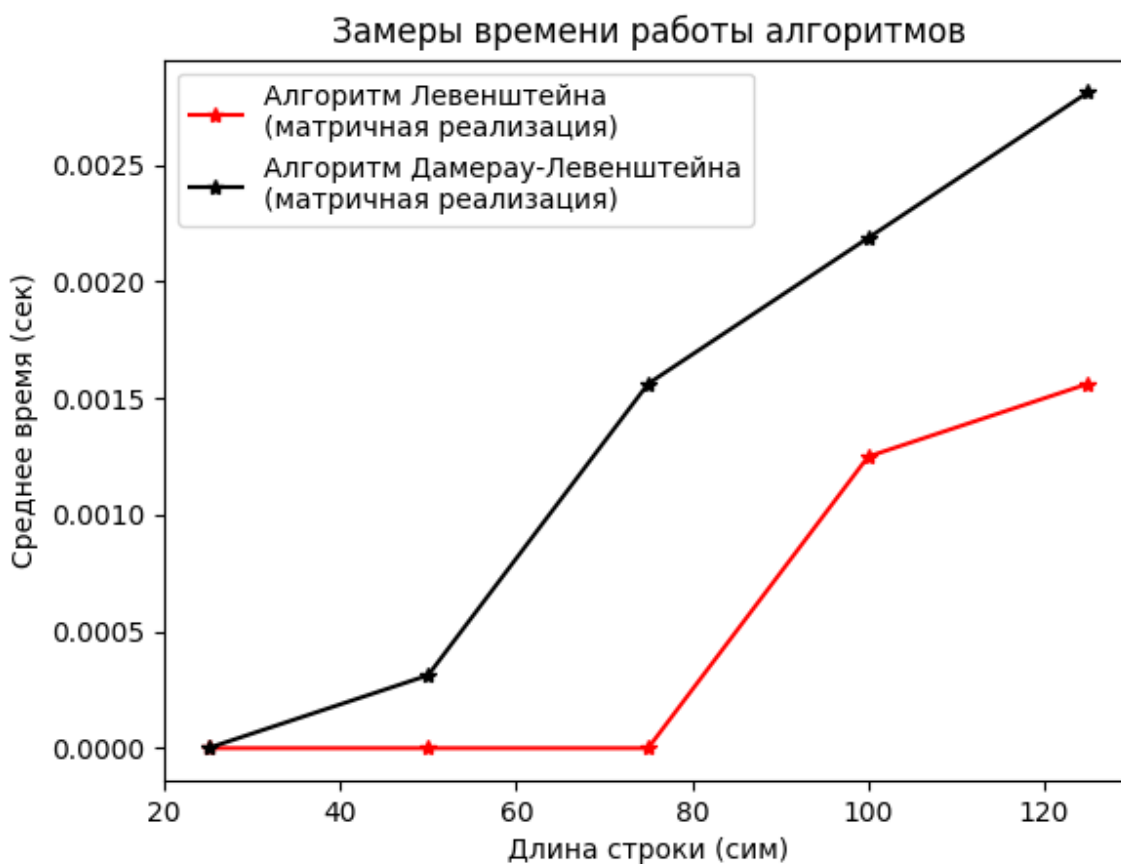


Рисунок 16 – График времени работы матричных реализаций алгоритмов в зависимости от длин строк

Алгоритм	25	50	75	100	125
Левенштейна (матричный)	0.0	0.0	0.0	0.0013	0.0016
Дамерау-Левенштейна (матричный)	0.0	0.00031	0.0016	0.0022	0.0028

Таблица 5 – Таблица времени (сек) работы матричных реализаций алгоритмов в зависимости от длин строк (сим)

Отдельно было измерено время работы алгоритмов Левенштейна и Дамерау-

Левенштейна в их рекурсивных реализациях на малом диапазоне длин случайно сгенерированных строк (от 1 символа до 9 с шагом 2) по 50 замеров каждая строка, среднее значение было вынесено в таблицу и для наглядности изображено на графике (рисунок 17 и таблица 6).

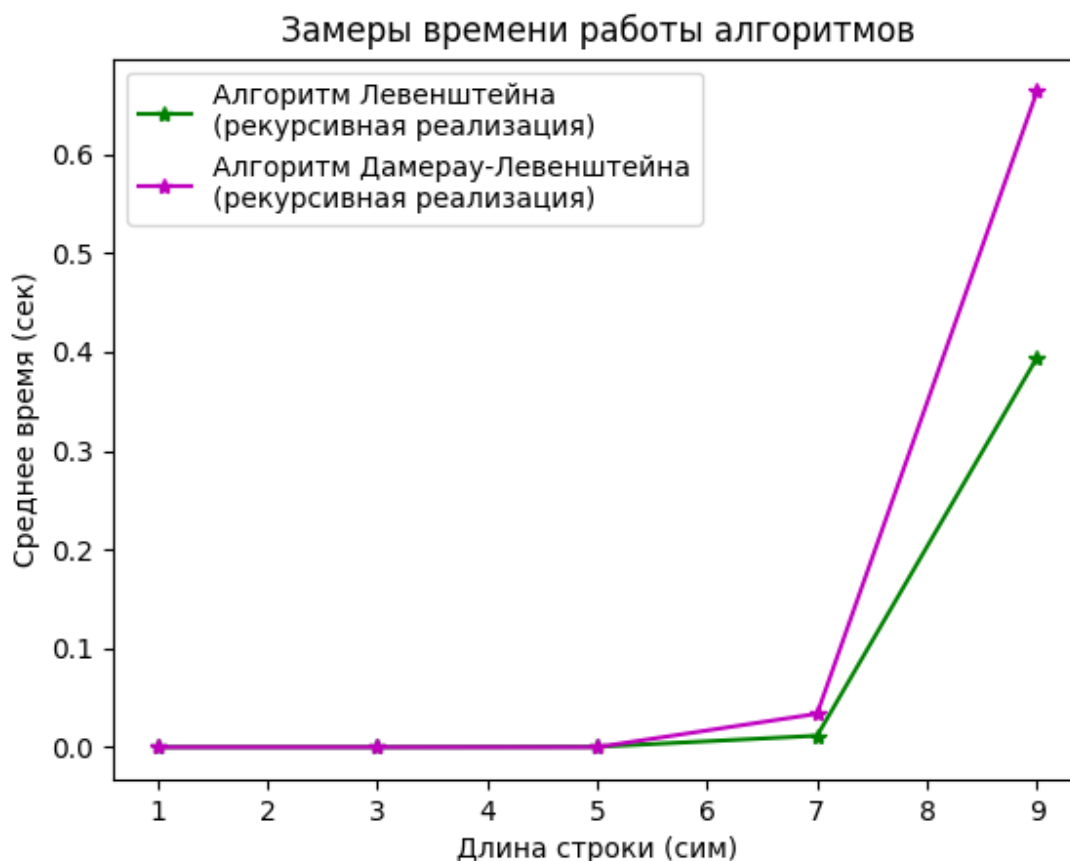


Рисунок 17 – График времени работы рекурсивных реализаций алгоритмов в зависимости от длин строк

Алгоритм	1	3	5	7	9
Левенштейна (рекурсивный)	0.0	0.0	0.00031	0.012	0.39
Дамерау-Левенштейна (рекурсивный)	0.0	0.0	0.0	0.034	0.66

Таблица 6 – Таблица времени (сек) работы рекурсивных реализаций алгоритмов в зависимости от длин строк (сим)

Отдельно было измерено время работы алгоритмов Левенштейна и Дамерау-Левенштейна в их рекурсивно-матричных реализациях на большем диапазоне длин случайно сгенерированных строк (от 25 символа до 125 с шагом 25) по 50 замеров каждая строка, среднее значение было вынесено в таблицу и для наглядности изображено на графике (рисунок 18 и таблица 7).

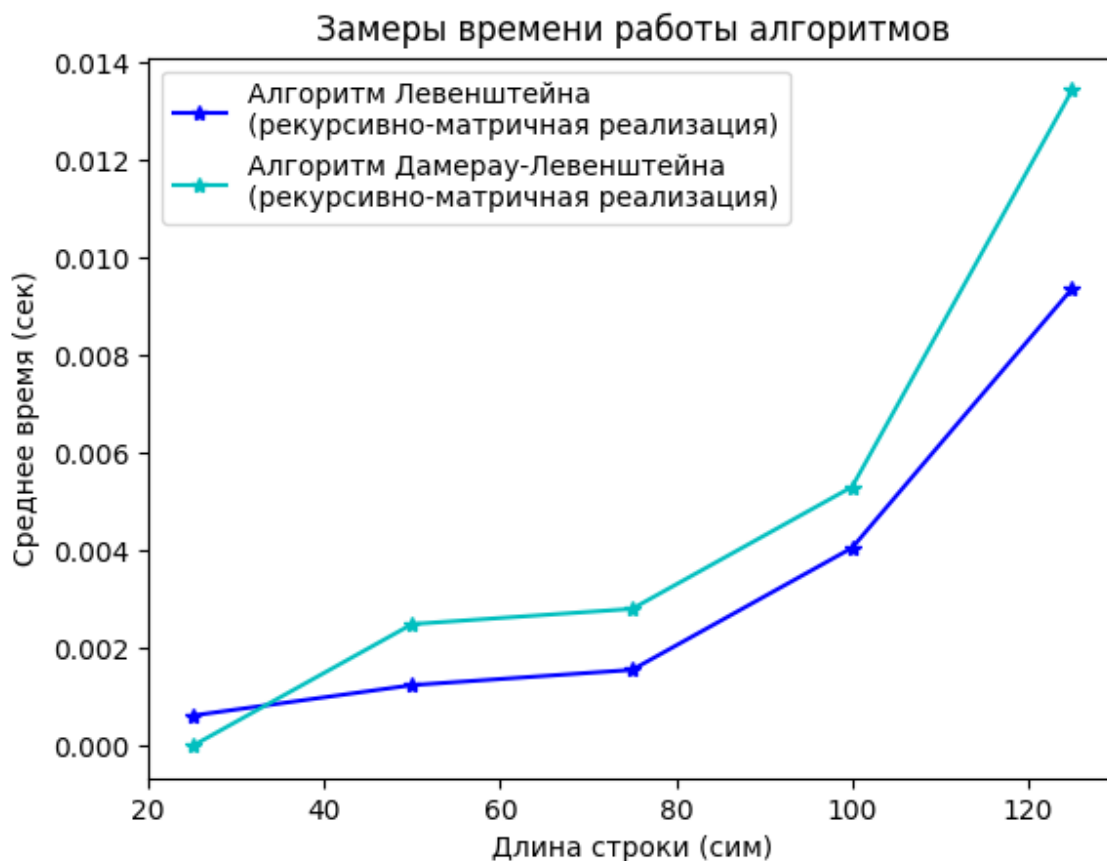


Рисунок 18 – График времени работы рекурсивно-матричных реализаций алгоритмов в зависимости от длин строк

Алгоритм	25	50	75	100	125
Левенштейна (рек-мат.)	0.00063	0.0013	0.0016	0.0041	0.0094
Дамерау-Левенштейна (рек-мат.)	0.0	0.0025	0.0028	0.0053	0.013

Таблица 7 – Таблица времени (сек) работы рекурсивно-матричных реализаций алгоритмов в зависимости от длин строк (сим)

Видно, что алгоритм Левенштейна оказался немного быстрее алгоритма Дамерау-Левенштейна из-за дополнительной проверки во втором, что компенсируется большим расстоянием в результате первого при наличии перестановок букв в строках.

4.3. Сравнение работы матричных и рекурсивно-матричных алгоритмов Левенштейна и Дамерау-Левенштейна

Так как на общем графике матричный и рекурсивно-матричный алгоритмы были очень близки по скорости, были проведены отдельные замеры (на 5-и точках с длиной строк от 25 до 125 символов с шагом 25 по 50 запусков), среднее значение было вынесено в таблицу и для наглядности изображено на графике (рисунок 19 и таблица 8).

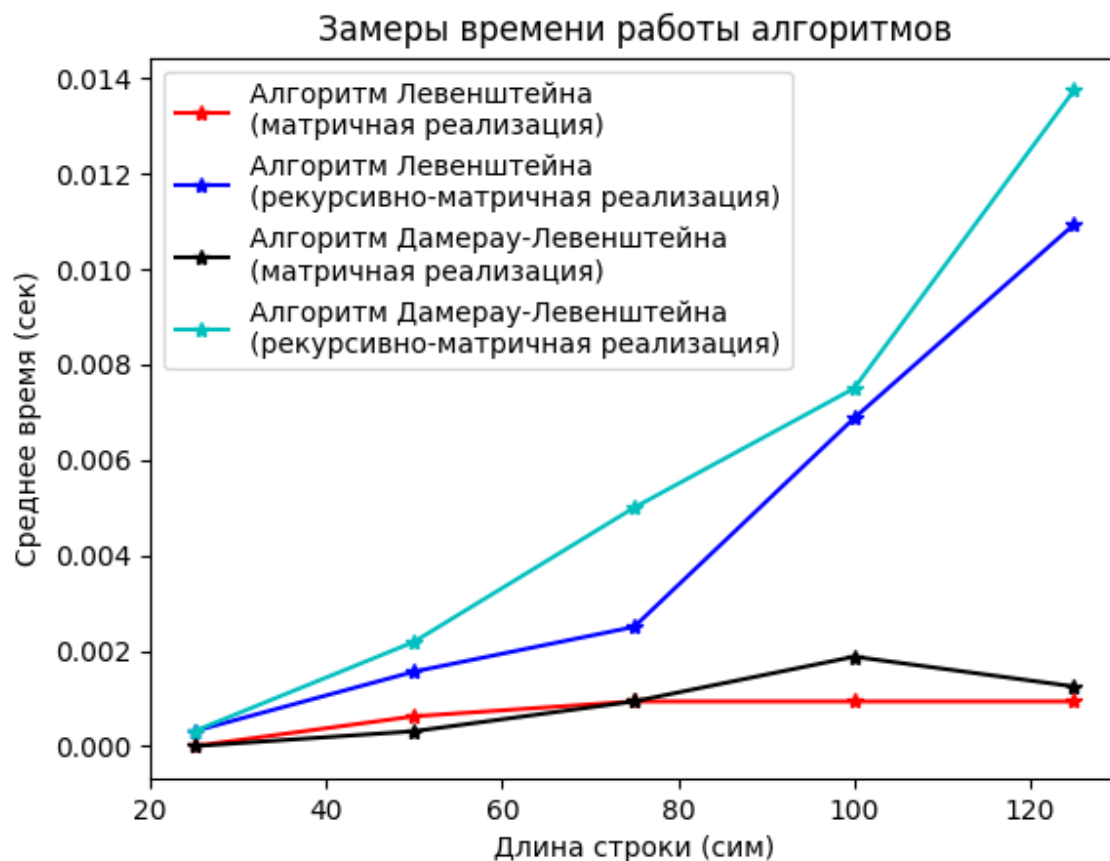


Рисунок 19 – График времени работы матричных и рекурсивно-матричных реализаций алгоритмов в зависимости от длин строк

Алгоритм	25	50	75	100	125
Левенштейна (матричный)	0.0	0.00063	0.00094	0.00094	0.00094
Левенштейна (рек-мат.)	0.00031	0.0016	0.0025	0.0069	0.011
Дамерау-Левенштейна (матричный)	0.0	0.00031	0.00094	0.0019	0.0013
Дамерау-Левенштейна (рек-мат.)	0.00031	0.0022	0.005	0.0075	0.014

Таблица 8 – Таблица времени (сек) работы матричных и рекурсивно-матричных реализаций алгоритмов в зависимости от длин строк (сим)

Кроме того, замеры времени работы матричных и рекурсивно-матричных алгоритмов были проведены на микроконтроллерах STM32. На графике 20 и в таблице 9 приведены замеры на длинах строк от 5 символов до 45 с шагом 5 и запуском каждого по 20 раз.

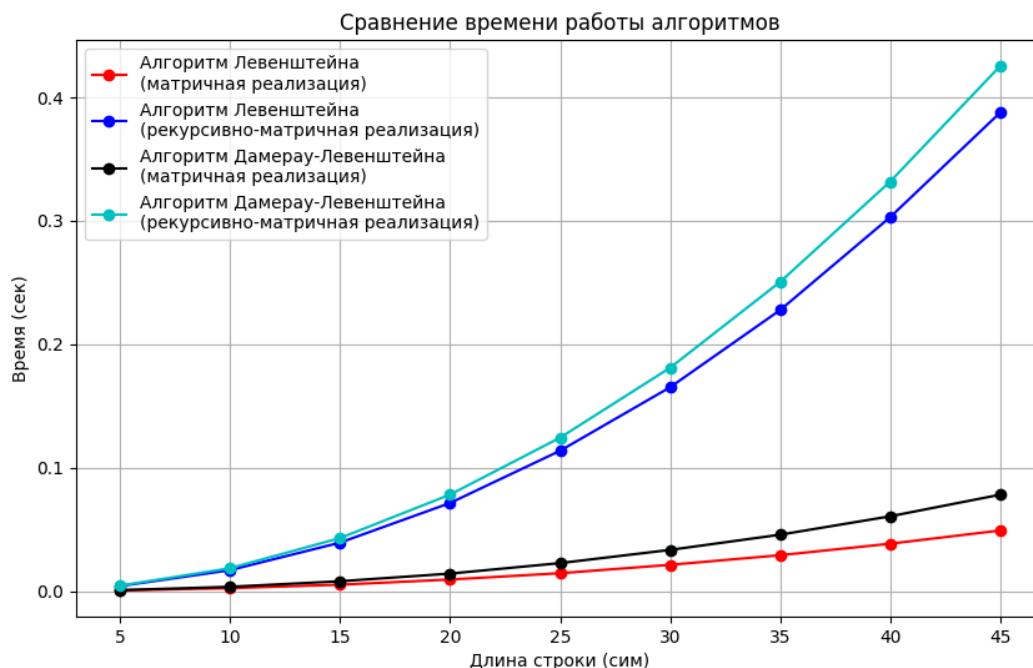


Рисунок 20 – График времени работы матричных и рекурсивно-матричных реализаций алгоритмов в зависимости от длин строк (замеры на микроконтроллерах)

Алгоритм	5	10	15	20	25	30	35	40	45
Левенш(м)	0.0007	0.0025	0.0053	0.0095	0.0146	0.0214	0.0292	0.0385	0.0492
Левенш(р-м)	0.0042	0.0171	0.0394	0.0715	0.1138	0.1652	0.2277	0.3029	0.3876
Дам-Лев(м)	0.0009	0.0036	0.0081	0.0143	0.0228	0.0335	0.0458	0.0607	0.0782
Дам-Лев(р-м)	0.0045	0.0187	0.0431	0.0783	0.1244	0.1810	0.2505	0.3317	0.4253

Таблица 9 – Таблица времени (сек) работы матричных и рекурсивно-матричных реализаций алгоритмов в зависимости от длин строк (сим) (замеры на микроконтроллерах)

По результатам приведённых графиков видно, что и в алгоритме Левенштейна и Дамерау-Левенштейна рекурсивно-матричный метод работает дольше матричного. Это объясняется затратами на вызов функции при рекурсии и на дополнительные проверки является ли искомое значение уже посчитанным. На небольших длинах строк разница в скорости работы алгоритмов отличается несущественно, однако с увеличением данных растёт и разница во времени.

Вывод

По проведённым исследованиям была выявлена большая скорость работы алгоритма Левенштейна над алгоритмом Дамерау-Левенштейна за счёт уменьшения числа проверок, что, однако, даёт иной результат при наличии возможности перестановок символов в строках. При этом матричный вариант выигрывает по скорости в обоих алгоритмах, на втором месте оказался рекурсивно-матричный метод, который делает

меньше рекурсивных вызовов, чем рекурсивный метод, и исключает повторные вычисления идентичных веток, так как при вызове каждой новой функции в этом методе передаётся в качестве аргумента ссылка на матрицу, которая хранит уже посчитанные значения, но на эту матрицу также необходима память, а проверки на уже вычисленные значения не всегда приносят положительный результат и занимают время.

ЗАКЛЮЧЕНИЕ

В результате выполнения лабораторной работы были исследованы алгоритмы вычисления расстояния Левенштейна и Дамерау-Левенштейна в матричной, рекурсивно-матричной и рекурсивной реализациях.

В частности:

- были рассмотрены алгоритмы вычисления расстояния Левенштейна и Дамерау-Левенштейна;
- применён метод динамического программирования для матричных реализаций алгоритмов;
- сравнены матричная, рекурсивно-матричная и рекурсивная реализации алгоритмов;
- сравнены алгоритмы вычисления расстояния Левенштейна и Дамерау-Левенштейна.

В ходе лабораторной работы были рассмотрены, спроектированы и запрограммированы алгоритмы нахождения расстояний Левенштейна и Дамерау-Левенштейна в их матричных, рекурсивных и рекурсивно-матричных реализациях.

Были разработаны тесты для всех алгоритмов, учитывающие крайние случаи, ожидаемых результатов которых достигли все реализации.

Сравнения полученных программ показали, что алгоритм Левенштейна работает быстрее алгоритма Дамерау-Левенштейна за счёт меньшего числа проверок, однако это приводит к другим результатам, если возможны перестановки символов. Матричный вариант оказался самым быстрым среди всех, на втором месте — рекурсивно-матричный метод, который снижает количество рекурсивных вызовов и избегает повторных вычислений, используя матрицу для хранения уже найденных значений. Однако этот метод требует дополнительной памяти, а проверки на уже вычисленные значения не всегда ускоряют процесс и также занимают время.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- [1] Python Documentation. `time.process_time()` - Документация по стандартной библиотеке Python. Дата обращения: 03 сентября 2024 г. [Электронный ресурс]. Доступно по адресу: <https://docs-python.ru/standart-library/modul-time-python/funktsija-process-time-modulja-time/>
- [2] Tirinox. Алгоритм Левенштейна на Python: реализация и объяснение. Дата обращения: 02 сентября 2024 г. [Электронный ресурс]. Доступно по адресу: <https://tirinox.ru/levenstein-python/>
- [3] Гасфилд Дэн. Строки, деревья и последовательности в алгоритмах: Информатика и вычислительная биология / Пер. с англ. И. В. Романовского. — СПб.: Невский Диалект; БХВ-Петербург, 2003. — 654 с: ил.
- [4] Левенштейн В. И. Двоичные коды с исправлением выпадений, вставок и замещений символов. Доклады Академий Наук СССР, 1965. 163.4:845-848.
- [5] Ниёзов Д. Л. Применение методов нечеткого сравнения строк в прикладных задачах: Выпускная квалификационная работа (Бакалаврская работа). — Тольятти: Тольяттинский государственный университет, 2020. — 45 стр.