



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 Редакционные расстояния

Студент Морозов Д.В.

Группа ИУ7-52Б

Оценка (баллы) _____

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

Москва — 2022 г.

Оглавление

| | |
|--|-----------|
| Введение | 3 |
| 1 Аналитическая часть | 5 |
| 1.1 Матричный алгоритм нахождения расстояния Левенштейна | 5 |
| 1.2 Рекурсивный алгоритм нахождения расстояния Дамерау- Левенштейна | 6 |
| 1.3 Рекурсивный алгоритм нахождения расстояния Дамерау- Левенштейна с использованием кеша | 7 |
| 1.4 Матричный алгоритм нахождения расстояния Дамерау-Левенштейна | 8 |
| 2 Конструкторская часть | 10 |
| 2.1 Матричный алгоритм поиска расстояния Левенштейна | 10 |
| 2.2 Рекурсивный алгоритм нахождения расстояния Дамерау- Левенштейна | 12 |
| 2.3 Рекурсивный алгоритм нахождения расстояния Дамерау- Левенштейна с использованием кеша | 13 |
| 2.4 Матричный алгоритм нахождения расстояния Дамерау- Левенштейна | 15 |
| 3 Технологическая часть | 16 |
| 3.1 Требования к ПО | 16 |
| 3.2 Средства реализации | 16 |
| 3.3 Листинг кода | 17 |
| 3.4 Функциональные тесты | 23 |
| 4 Исследовательская часть | 24 |
| 4.1 Технические характеристики | 24 |
| 4.2 Время выполнения алгоритмов | 24 |
| 4.3 Использование памяти | 27 |

| | |
|---|-----------|
| Заключение | 30 |
| Список использованных источников | 31 |

Введение

Целью данной лабораторной работы является получение практических навыков динамического программирования на примере реализации алгоритмов Левенштейна и Дамерау-Левенштейна.

Расстояние Левенштейна (редакционное расстояние, дистанция редактирования) — метрика, измеряющая разность между двумя последовательностями символов. Она определяется как минимальное количество односимвольных операций (вставки, удаления, замены), необходимых для превращения одной строки в другую [1].

Расстояние Левенштейна и его обобщения применяются:

- для исправления ошибок в слове (в поисковых системах, базах данных, при вводе текста, при автоматическом распознавании отсканированного текста или речи);
- для сравнения текстовых файлов утилитой diff и ей подобными (здесь роль «символов» играют строки, а роль «строк» — файлы);
- в биоинформатике [2].

Расстояние Дамерау-Левенштейна (названо в честь учёных Фредерика Дамерау и Владимира 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, так как к операциям вставки, удаления и замены символов, определённых в расстоянии Левенштейна добавлена операция транспозиции (перестановки) символов [3].

Задачами данной лабораторной являются:

- изучение и реализация алгоритмов Левенштейна и Дamerau-Левенштейна нахождения редакционного расстояния между строками;
- выполнение теоретической либо экспериментальной оценки затрат алгоритмов по памяти;
- выполнение экспериментальной оценки затрат алгоритмов по времени;
- сравнение алгоритмов по проведённым оценкам.

1 Аналитическая часть

В этом разделе будут представлены описания алгоритмов нахождения расстояний Левенштейна и Дamerau-Левенштейна и их практическое применение.

1.1 Матричный алгоритм нахождения расстояния Левенштейна

Расстояние Левенштейна между двумя строками — это минимальное количество операций вставки, удаления и замены, необходимых для превращения одной строки в другую. При этом каждая операция имеет свою цену (штраф).

Рассмотрим матрицу A размером $(length(s_1)+1) \cdot ((length(s_2)+1))$, где $length(S)$ — длина строки S . Пусть значение в ячейке $[i, j]$ матрицы равно расстоянию между префиксом s_1 длины i и префиксом s_2 длины j . У элементов первой строки значение равно индексу столбца, у элементов первого столбца — индексу строки.

Остальные ячейки заполняем в соответствии с формулой (1.1).

$$A[i][j] = \min \begin{cases} A[i-1][j] + 1; \\ A[i][j-1] + 1; \\ A[i-1][j-1] + m(s_1[i], s_2[j]). \end{cases} \quad (1.1)$$

Функция m определена как:

$$m(s_1[i], s_2[j]) = \begin{cases} 0, & \text{если } s_1[i-1] = s_2[j-1]; \\ 1, & \text{иначе.} \end{cases} \quad (1.2)$$

В результате расстоянием Левенштейна будет ячейка матрицы с индексами $i = \text{length}(s_1)$ и $j = \text{length}(s_2)$.

1.2 Рекурсивный алгоритм нахождения расстояния Дамерау-Левенштейна

Расстояние Дамерау-Левенштейна между двумя строками — это минимальное количество операций вставки, удаления, замены и транспозиции (перестановки двух соседних символов), необходимых для перевода одной строки в другую. Является модификацией расстояния Левенштейна.

Рекурсивный алгоритм считает редакционное расстояние для двух строк s_1 и s_2 по рекуррентной формуле (1.3), где i — длина префикса

s_1, j — длина префикса s_2 :

$$d(i, j) = \begin{cases} \max(i, j), \text{ если } \min(i, j) = 0; \\ \min\{ \\ \quad d(i, j - 1) + 1; \\ \quad d(i - 1, j) + 1; \\ \quad d(i - 1, j - 1) + m(a[i], b[j]); \\ \quad \left\{ \begin{array}{ll} d(i - 2, j - 2) + 1, & \text{если } i, j > 1, \\ s_1[i] = s_2[j - 1], \\ s_2[j] = s_1[i - 1]; \\ \infty, & \text{иначе.} \end{array} \right. \\ \quad \}, \text{ иначе.} \end{cases} \quad (1.3)$$

Функция f определена как:

$$m(i, j) = \begin{cases} 0, & \text{если } s_1[i - 1] = s_2[j - 1]; \\ 1, & \text{иначе.} \end{cases} \quad (1.4)$$

1.3 Рекурсивный алгоритм нахождения расстояния Дамерау-Левенштейна с использованием кеша

Рекурсивный алгоритм можно оптимизировать, если записывать найденные промежуточные расстояния в кеш-матрицу. Перед началом расчёта требуется инициализировать ячейки матрицы значени-

ем -1 . При рекурсивном вызове требуется проверить, было ли значение вычислено ранее, — проверить, находится ли в соответствующей ячейке матрицы -1 . Таким образом, на поиск уже найденных расстояний время не тратится — их значения берутся из матрицы.

1.4 Матричный алгоритм нахождения расстояния Дameraу-Левенштейна

При больших i, j прямая реализация формулы (1.3) может быть неэффективна по времени, так как некоторые значения $D(i, j)$ вычисляются несколько раз. Для оптимизации алгоритма можно хранить промежуточные значения в матрице размером $(length(s_1)+1) \times (length(s_2)+1)$, где $length(S)$ — длина строки S . Значение в ячейке $[i, j]$ матрицы A равно расстоянию между префиксом s_1 длины i и префиксом s_2 длины j . У элементов первой строки значение равно индексу столбца, у элементов первого столбца — индексу строки.

Остальные ячейки заполняем в соответствии с формулой (1.5).

$$A[i][j] = \min \begin{cases} A[i-1][j] + 1; \\ A[i][j-1] + 1; \\ A[i-1][j-1] + m(s_1[i], s_2[j]); \\ \begin{cases} A[i-2][j-2] + 1, & \text{если } i, j > 1, \\ s_1[i] = s_2[j-1], \\ s_2[j] = s_1[i-1]; \\ \infty, & \text{иначе.} \end{cases} \end{cases} \quad (1.5)$$

Функция m определена как

$$m(s_1[i], s_2[j]) = \begin{cases} 0, & \text{если } s_1[i-1] = s_2[j-1]; \\ 1, & \text{иначе.} \end{cases} \quad (1.6)$$

В результате расстоянием Дameraу-Левенштейна будет ячейка матрицы с индексами $i = \text{length}(s_1)$ и $j = \text{length}(s_2)$.

2 Конструкторская часть

В этом разделе будут приведены схемы алгоритмов нахождения расстояний Левенштейна и Дамерау-Левенштейна.

2.1 Матричный алгоритм поиска расстояния Левенштейна

На рисунке 2.1 приведена схема матричного алгоритма нахождения расстояния Левенштейна.

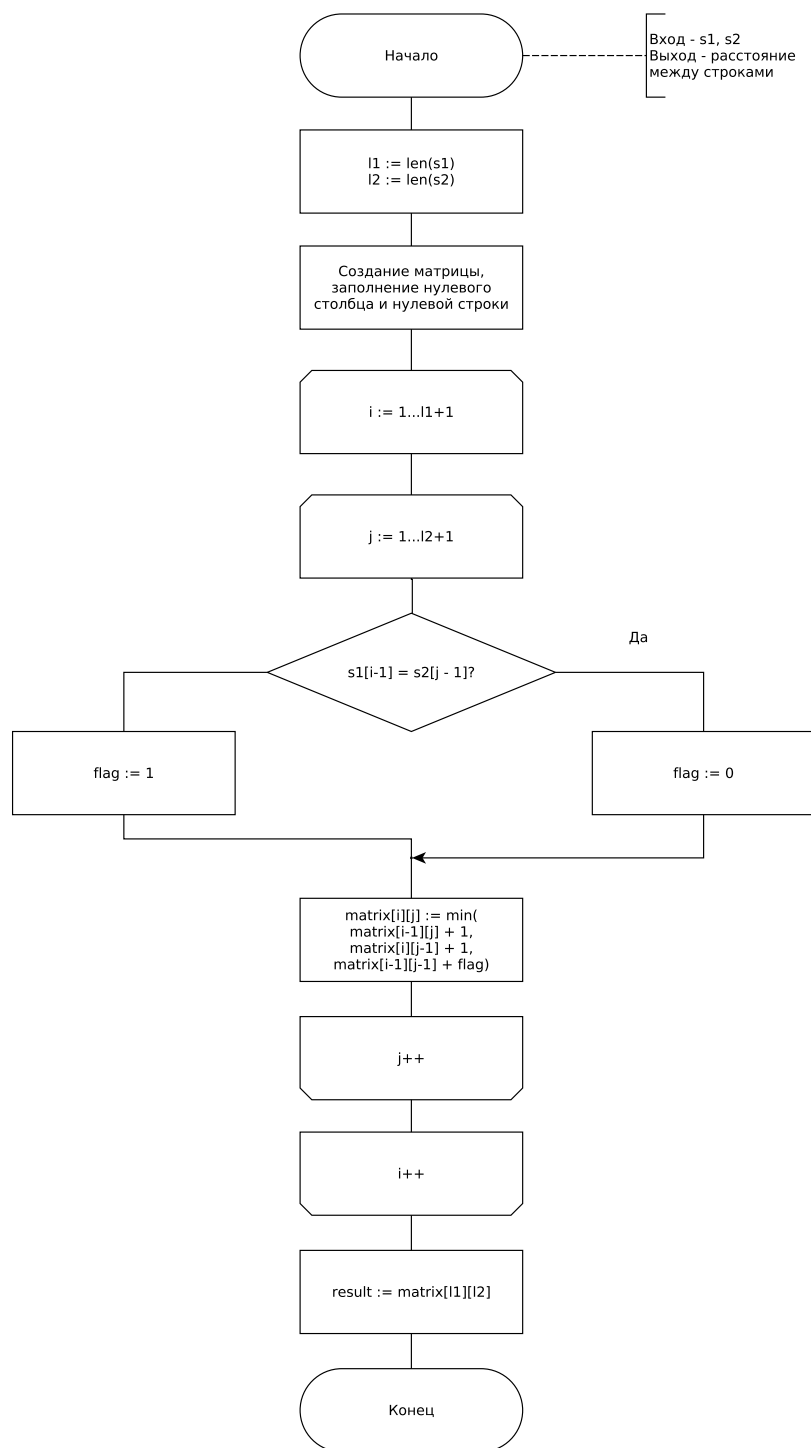


Рисунок 2.1 – Схема матричного алгоритма нахождения расстояния Левенштейна

2.2 Рекурсивный алгоритм нахождения расстояния Дамерау-Левенштейна

На рисунке 2.2 приведена схема рекурсивного алгоритма нахождения расстояния Дамерау-Левенштейна.

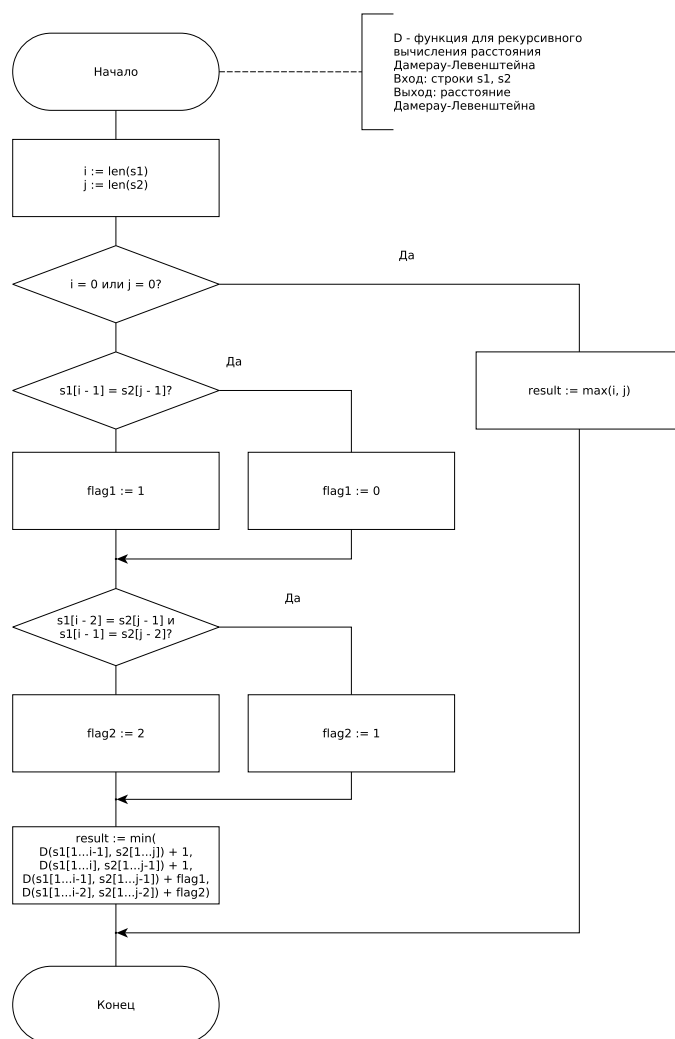


Рисунок 2.2 – Схема рекурсивного алгоритма нахождения расстояния Дамерау-Левенштейна

2.3 Рекурсивный алгоритм нахождения расстояния Дамерау-Левенштейна с использованием кеша

На рис. 2.3 приведена схема рекурсивного алгоритма поиска расстояния Дамерау-Левенштейна с использованием кеша — матрицы.

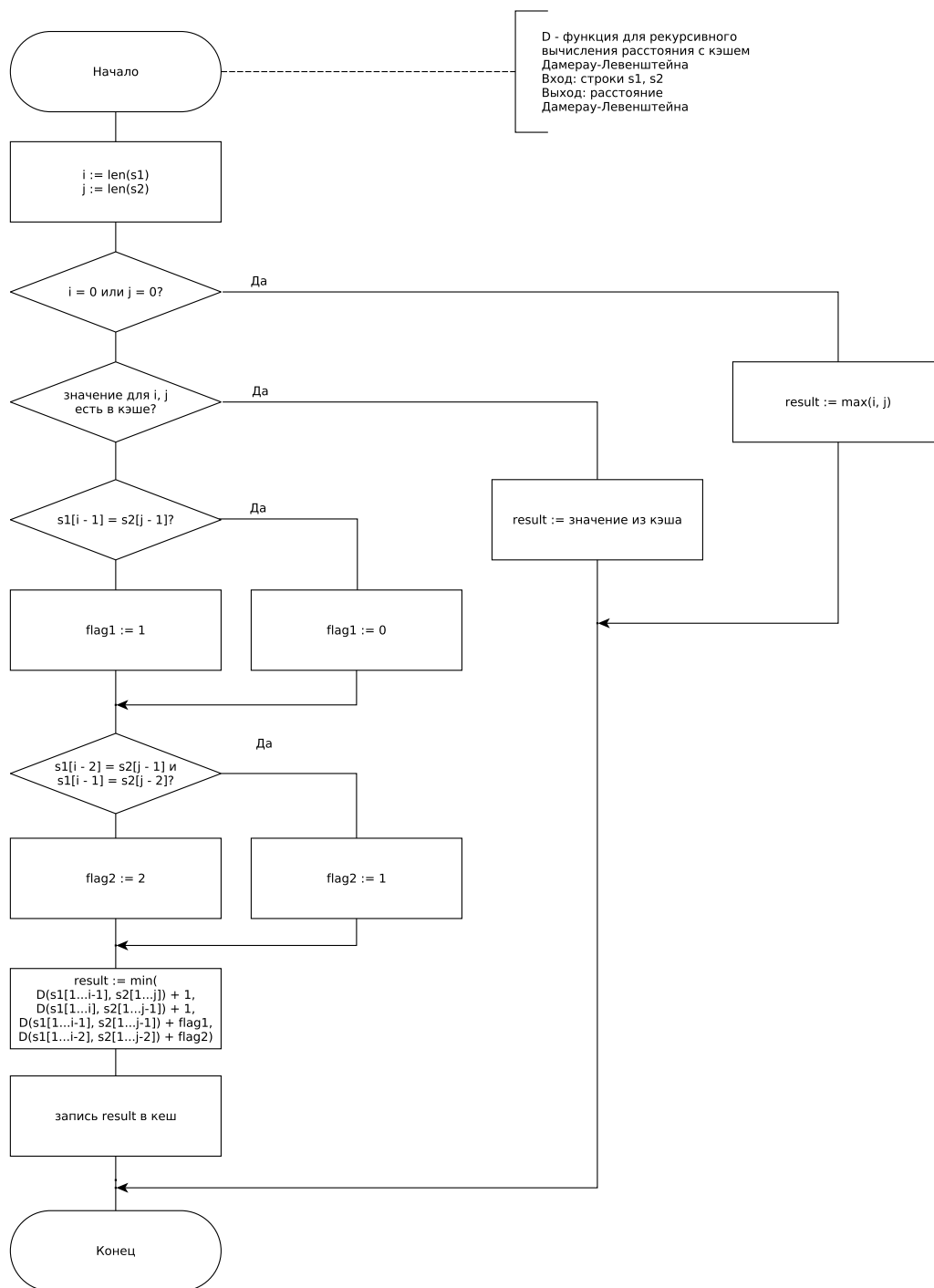


Рисунок 2.3 – Схема рекурсивного алгоритма нахождения расстояния Дамерау-Левенштейна с использованием кэша (матрицы)

2.4 Матричный алгоритм нахождения расстояния Дамерау-Левенштейна

На рисунке 2.4 приведена схема матричного алгоритма нахождения расстояния Дамерау-Левенштейна.

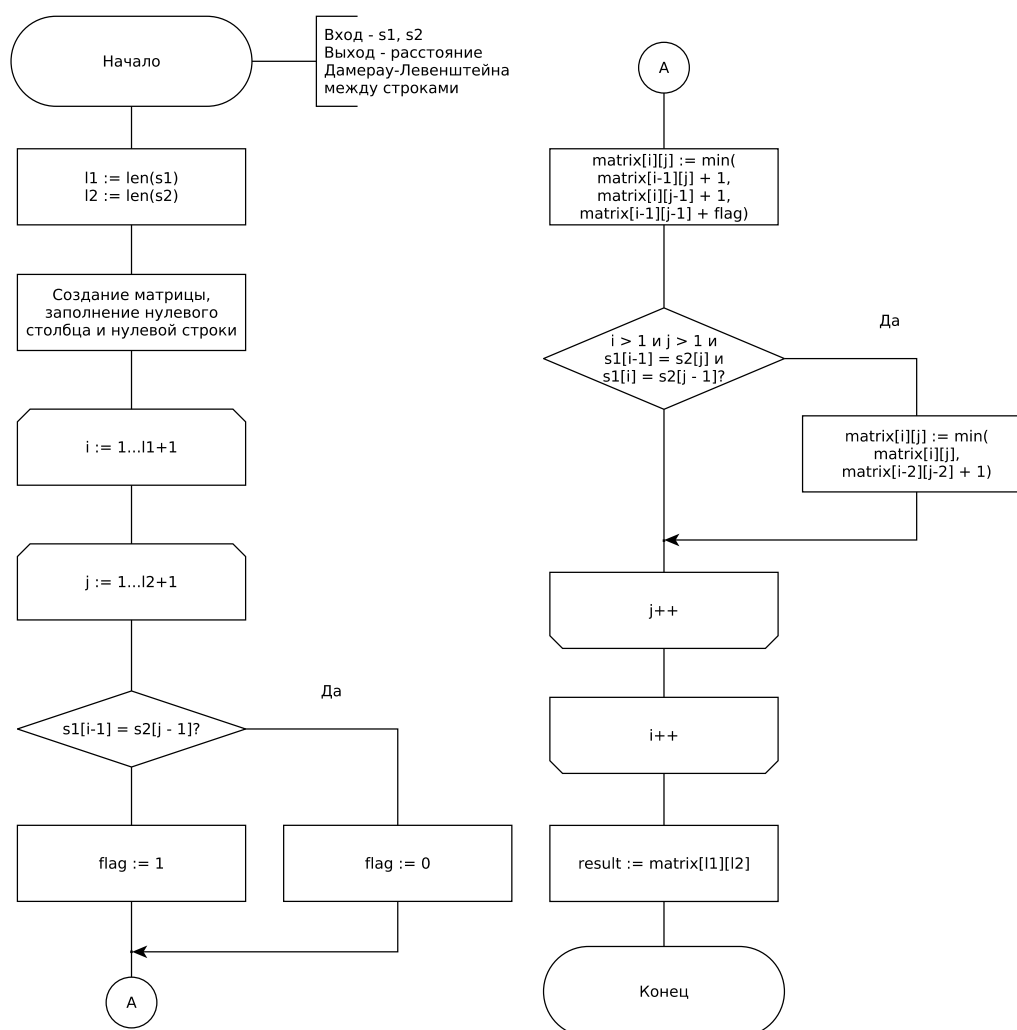


Рисунок 2.4 – Схема матричного алгоритма нахождения расстояния Дамерау-Левенштейна

3 Технологическая часть

В данном разделе приведены требования к ПО, средства реализации и листинги кода.

3.1 Требования к ПО

Используемое ПО должно предоставлять возможность измерения процессорного времени.

3.2 Средства реализации

Для реализации данной лабораторной работы был выбран язык программирования Golang [4] и среда разработки Goland, которая позволяет замерять процессорное время с помощью пакета `CGo` [5].

3.3 Листинг кода

В листингах 3.1–3.6 приведены реализации алгоритмов нахождения расстояния Левенштейна и Дамерау-Левенштейна.

Листинг 3.1 – Функция нахождения расстояния Левенштейна с заполнением матрицы (начало)

```
1 func LevenshteinMatrix(str1 , str2 string) int {  
2     n := len(str1)  
3     m := len(str2)  
4     if n == 0 {  
5         return m  
6     } else if m == 0 {  
7         return n  
8     }  
9  
10    b1 := make([]int , n+1)  
11    b2 := make([]int , n+1)  
12  
13    for i := 0; i < n+1; i++ {  
14        b2[i] = i  
15    }
```

Листинг 3.2 – Функция нахождения расстояния Левенштейна с
заполнением матрицы (окончание)

```
1  for i := 1; i < m+1; i++ {
2      swap(&b1, &b2)
3
4      b2[0] = i
5      for j := 1; j < n+1; j++ {
6          flag := 1
7          if str1[j-1] == str2[i-1] {
8              flag = 0
9          }
10
11         res := min(b2[j-1]+1, b1[j]+1, b1[j-1]+flag)
12
13         b2[j] = res
14     }
15 }
16
17 return b2[n]
18 }
```

Листинг 3.3 – Функция нахождения расстояния
Дамерау-Левенштейна без использования рекурсии (начало)

```
1 func DamerauLevenshteinMatrix(str1, str2 string) int {
2     n := len(str1)
3     m := len(str2)
4
5     if n == 0 {
6         return m
7     } else if m == 0 {
8         return n
9     }
```

Листинг 3.4 – Функция нахождения расстояния

Дамерау-Левенштейна без использования рекурсии (окончание)

```
1    b1 := make ([] int , n+1)
2    b2 := make ([] int , n+1)
3    b3 := make ([] int , n+1)
4
5    for i := 0; i < n+1; i++ {
6        b3[i] = i
7    }
8
9    for i := 1; i < m+1; i++ {
10       swapLeft(&b1, &b2, &b3)
11
12       b3[0] = i
13       for j := 1; j < n+1; j++ {
14           flag := 1
15           if str1[j-1] == str2[i-1] {
16               flag = 0
17           }
18
19           res := min(b3[j-1]+1, b2[j]+1, b2[j-1]+flag)
20
21           if i > 1 && j > 1 && str1[j-1] == str2[i-2] &&
               str1[j-2] == str2[i-1] {
22               res = min(res, b1[j-2] + 1)
23           }
24
25           b3[j] = res
26       }
27     }
28
29     return b3[n]
30 }
```

Листинг 3.5 – Функция нахождения расстояния
Дамерау-Левенштейна с использованием рекурсии.

```
1 func DamerauLevenshteinRec(str1, str2 string) int {
2     n := len(str1)
3     m := len(str2)
4
5     if n == 0 {
6         return m
7     } else if m == 0 {
8         return n
9     }
10
11     flag := 1
12     if str1[n-1] == str2[m-1] {
13         flag = 0
14     }
15
16     res := min(DamerauLevenshteinRec(str1[:n-1], str2)+1,
17     DamerauLevenshteinRec(str1, str2[:m-1])+1,
18     DamerauLevenshteinRec(str1[:n-1], str2[:m-1])+flag)
19
20     if n >= 2 && m >= 2 && str1[n-1] == str2[m-2] &&
21     str1[n-2] == str2[m-1] {
22         cur := DamerauLevenshteinRec(str1[:n-2], str2[:m-2])
23         + 1
24         if cur < res {
25             res = cur
26         }
27     }
28 }
```

Листинг 3.6 – Функция нахождения расстояния

Дамерау-Левенштейна с использованием рекурсии с кешем (начало)

```
1 func damerauLevenshteinRecCache(str1, str2 string, cache
  [][]int) int {
2   n := len(str1)
3   m := len(str2)
4
5   if n == 0 {
6     return m
7   } else if m == 0 {
8     return n
9   }
10
11  if cache[n-1][m-1] != -1 {
12    return cache[n-1][m-1]
13  }
14
15  res := damerauLevenshteinRecCache(str1[:n-1], str2,
    cache) + 1
16
17  cur := damerauLevenshteinRecCache(str1, str2[:m-1],
    cache) + 1
18  if cur < res {
19    res = cur
20  }
21
22  flag := 1
23  if str1[n-1] == str2[m-1] {
24    flag = 0
25  }
26
27  cur = damerauLevenshteinRecCache(str1[:n-1], str2[:m-1],
    cache) + flag
```

Листинг 3.7 – Функция нахождения расстояния
Дамерау-Левенштейна с использованием рекурсии с кешем
(окончание)

```
1  if cur < res {
2      res = cur
3  }
4
5  if n >= 2 && m >= 2 && str1[n-1] == str2[m-2] &&
   str1[n-2] == str2[m-1] {
6      cur = damerauLevenshteinRecCache(str1[:n-2],
   str2[:m-2], cache) + 1
7      if cur < res {
8          res = cur
9      }
10 }
11
12 cache[n-1][m-1] = res
13 return res
14 }
```

3.4 Функциональные тесты

В таблице 3.1 приведены функциональные тесты для алгоритмов вычисления расстояния Левенштейна (в таблице столбец подписан "Левенштейн") и Дамерау-Левенштейна (в таблице — "Дамерау-Л."). Все тесты пройдены успешно.

Таблица 3.1 – Функциональные тесты

| № | Входные данные | | Ожидаемый результат | |
|----|----------------|---------------|---------------------|------------|
| | Строка 1 | Строка 2 | Левенштейн | Дамерау-Л. |
| 1 | cat | cute | 2 | 2 |
| 2 | cute | cat | 2 | 2 |
| 3 | dog | dog | 0 | 0 |
| 4 | toook | t | 4 | 4 |
| 5 | пустая строка | пустая строка | 0 | 0 |
| 8 | пустая строка | 33 | 2 | 2 |
| 9 | 1234 | пустая строка | 4 | 4 |
| 10 | sample | samlpee | 2 | 2 |
| 11 | abcde | abced | 2 | 1 |
| 12 | abcde | based | 4 | 2 |

4 Исследовательская часть

В данном разделе произведено сравнение алгоритмов.

4.1 Технические характеристики

Технические характеристики устройства, на котором выполнялись замеры времени:

- операционная система — Ubuntu 22.04.1 Linux x86_64;
- оперативная память — 8 ГБ;
- процессор — AMD Ryzen 5 3550H [6].

Замеры проводились на ноутбуке, включенном в сеть электропитания. Во время замеров ноутбук не был нагружен сторонними приложениями.

4.2 Время выполнения алгоритмов

На рисунке 4.1 представлен график, иллюстрирующий зависимость времени работы от длины строк для матричного алгоритма поиска расстояния Левенштейна и матричного алгоритма поиска расстояния Дамерау-Левенштейна.

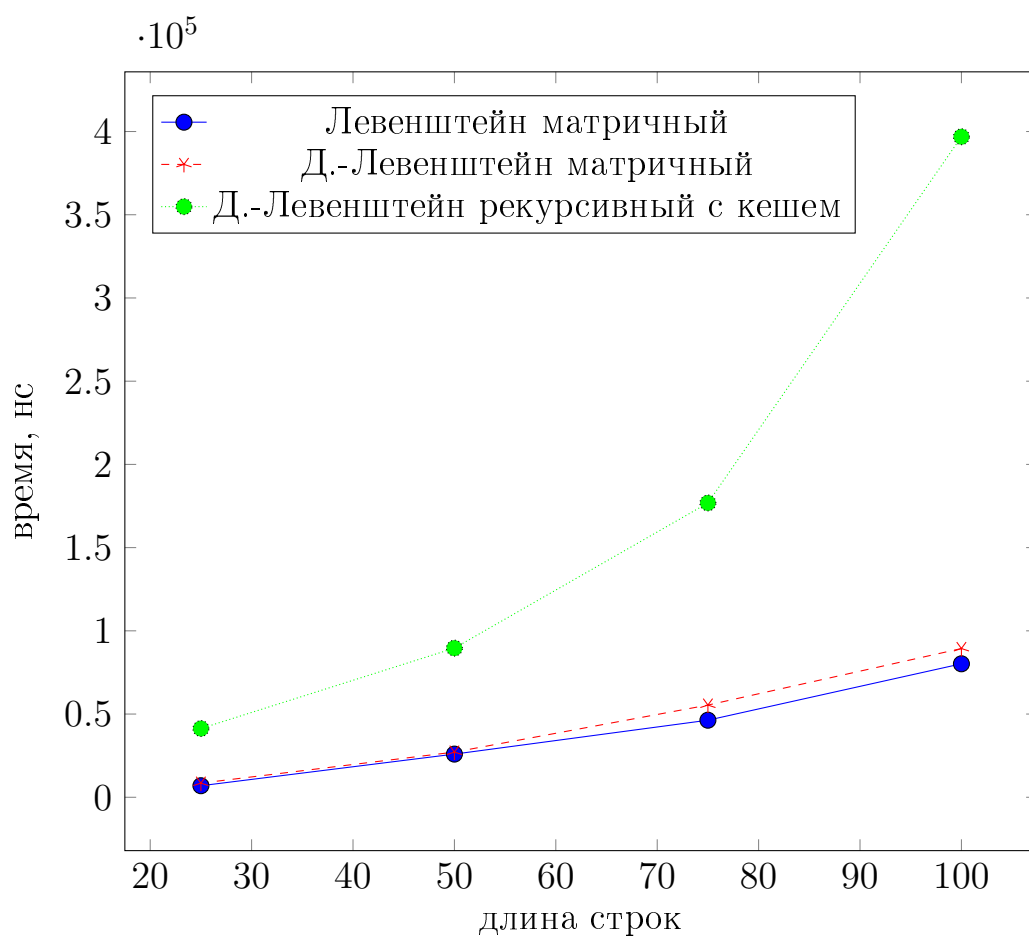


Рисунок 4.1 – Сравнение матричного алгоритма поиска расстояния Левенштейна и матричного алгоритма поиска расстояния Дамерау-Левенштейна

На рисунке 4.2 представлен график, иллюстрирующий зависимость времени работы от длины строк для рекурсивных алгоритмов поиска расстояния Дameraу-Левенштейна с использованием кеша и без.

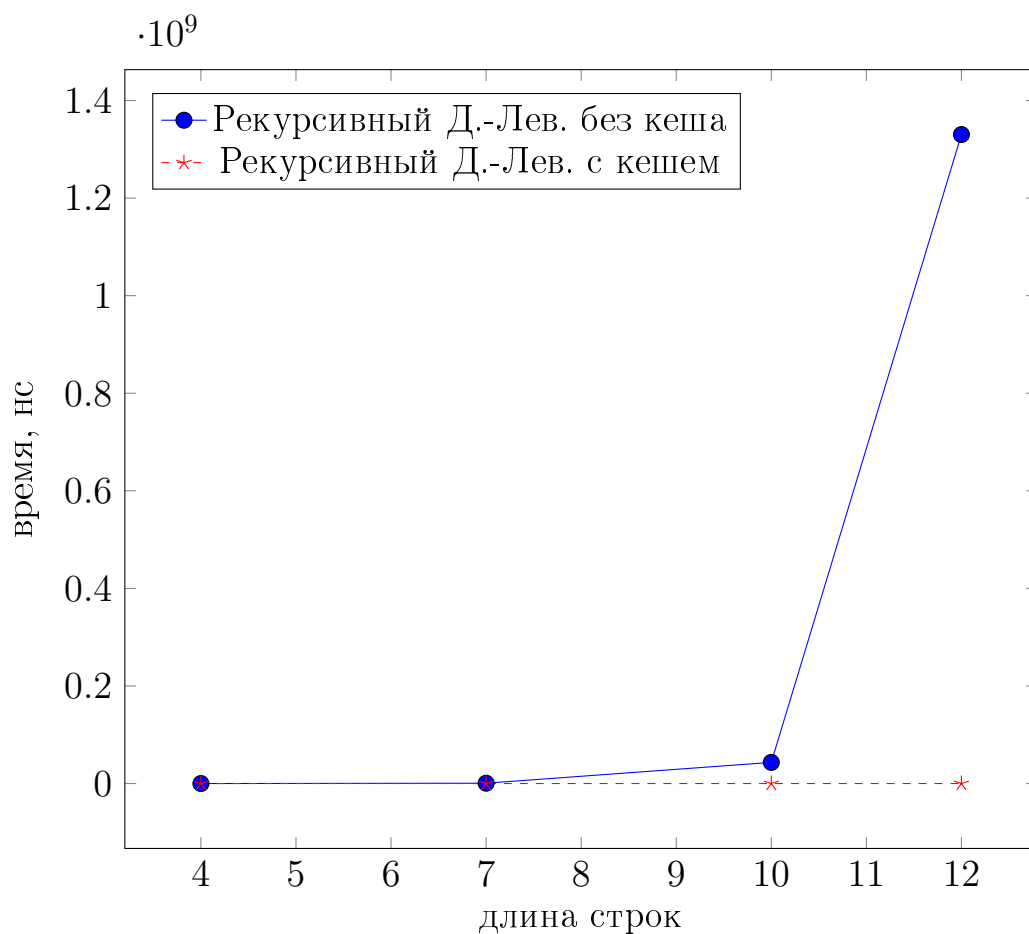


Рисунок 4.2 – Сравнение рекурсивных алгоритмов поиска расстояния Дameraу-Левенштейна с использованием кеша и без

4.3 Использование памяти

Пусть длина строки S_1 — n , длина строки S_2 — m . Максимальная глубина стека вызовов при рекурсивной реализации равна сумме длин входящий строк.

Обозначим: `char` — тип, используемый для хранения символа строки, `int` — тип, используемый для хранения чисел.

Рассчитаем затраты по памяти для матричных алгоритмов поиска расстояния Левенштейна и Дамерау-Левенштейна:

- строки S_1, S_2 — $(m + n) \cdot \text{sizeof}(\text{char})$;
- матрица — $(m + 1) \cdot (n + 1) \cdot \text{sizeof}(\text{int})$;
- длины строк — $2 \cdot \text{sizeof}(\text{int})$;
- вспомогательные переменные — $3 \cdot \text{sizeof}(\text{int})$;
- итого — $(m + n) \cdot \text{sizeof}(\text{char}) + (m + 1) \cdot (n + 1) \cdot \text{sizeof}(\text{int}) + 5 \cdot \text{sizeof}(\text{int})$.

Рассчитаем затраты по памяти для рекурсивного алгоритма поиска расстояния Дамерау-Левенштейна (для каждого вызова):

- строки S_1, S_2 — $(m + n) \cdot \text{sizeof}(\text{char})$;
- длины строк — $2 \cdot \text{sizeof}(\text{int})$;
- вспомогательные переменные — $3 \cdot \text{sizeof}(\text{int})$;
- адрес возврата — 8 байт;
- итого — $(m + n) \cdot \text{sizeof}(\text{char}) + 5 \cdot \text{sizeof}(\text{int}) + 8$ байт.

Высота дерева рекурсивных вызовов $\max(m, n) + 1$. Тогда максимальная глубина стека равна

$$M_{rec} = (\min(m, n) + 1) \cdot ((m + n) \cdot \text{sizeof}(\text{char}) + 5 \cdot \text{sizeof}(\text{int}) + 8). \quad (4.1)$$

Рассчитаем затраты по памяти для рекурсивного алгоритма поиска расстояния Дамерау-Левенштейна с использованием кеша (для каждого вызова):

- строки S_1, S_2 — $(m + n) \cdot \text{sizeof}(\text{char})$;
- длины строк — $2 \cdot \text{sizeof}(\text{int})$;
- вспомогательные переменные — $3 \cdot \text{sizeof}(\text{int})$;
- ссылка на матрицу — 8 байт;
- адрес возврата — 8 байт;
- итого — $(m + n) \cdot \text{sizeof}(\text{char}) + 5 \cdot \text{sizeof}(\text{int}) + 16$ байт.

Память для хранения матрицы (для всех вызовов общая)

$$M_{matr} = (m + 1) \cdot (n + 1) \cdot \text{sizeof}(\text{int}). \quad (4.2)$$

Максимальная глубина стека равна

$$M_{cash} = (\min(m, n) + 1) \cdot ((m + n) \cdot \text{sizeof}(\text{char}) + 5 \cdot \text{sizeof}(\text{int}) + 16) + (m + 1) \cdot (n + 1) \cdot \text{sizeof}(\text{int}). \quad (4.3)$$

Вывод

Алгоритм нахождения расстояния Дамерау-Левенштейна по времени выполнения незначительно отличается от алгоритма нахождения расстояния Левенштейна (для слов длиной 100 символов 89 мкс против 80 мкс).

По расходу памяти итеративный алгоритм проигрывают рекурсивному: максимальный размер используемой памяти в них растёт как произведение длин строк, в то время как у рекурсивного алгоритма — как сумма длин строк.

Рекурсивный алгоритм с заполнением матрицы превосходит по времени работы простой рекурсивный (для слов длиной 10 символов 16 мкс против 43353 мкс) и сравним с матричным алгоритмом.

Заключение

Цель достигнута. В ходе выполнения лабораторной работы были решены следующие задачи:

- изучены и реализованы алгоритмы нахождения расстояний Левенштейна и Дамерау-Левенштейна;
- выполнена теоретическая оценка затрат алгоритмов по памяти;
- выполнена экспериментальная оценка затрат алгоритмов по времени;
- проведено сравнение алгоритмов по проведённым оценкам.

Алгоритм нахождения расстояния Дамерау-Левенштейна по производительности схож с алгоритмом нахождения расстояния Левенштейна (для слов длиной 100 быстрее на 9%).

Рекурсивный алгоритм с заполнением матрицы эффективнее по времени работы, чем простой рекурсивный (для слов длиной 10 в 2700 раз быстрее) и незначительно отличается от матричной реализации (для слов длиной 100 в 4 раз медленнее). Однако по расходу памяти рекурсивные алгоритмы эффективнее матричных, так как максимальный размер используемой памяти в них растёт как произведение длин строк, в то время как у рекурсивного алгоритма — как сумма длин строк.

Список использованных источников

1. Левенштейн В. И. Двоичные коды с исправлением выпадений, вставок и замещений символов // Доклады АН СССР. – М.: Наука, 1965. Т. 163. С. 845–848.
2. А.С.Гуменюк Н.Н.Поздниченко И.Н.Родионов С.Н.Шпынов. О средствах формального анализа строя нуклеотидных цепей // Математическая биология и биоинформатика. Омский государственный технический университет, 2013. Т. 8. С. 373–397.
3. Черненький В. М. Гапанюк Ю. Е. Методика идентификации пассажира по установочным данным // Вестник МГТУ им. Н.Э. Баумана. Сер. “Приборостроение”. – М.: Издательство МГТУ им. Н.Э. Баумана., 2012. Т. 163. С. 30–34.
4. The Go Programming Language [Электронный ресурс]. Режим доступа: <https://golang.org/> (дата обращения: 14.09.2022).
5. C? Go? Cgo! – The Go Programming Language [Электронный ресурс]. Режим доступа: <https://go.dev/blog/cgo> (дата обращения: 14.09.2022).
6. AMD Ryzen™ 5 3550H [Электронный ресурс]. Режим доступа: <https://www.amd.com/en/products/apu/amd-ryzen-5-3550h> (дата обращения: 14.09.2022).