

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ ИМ. Н. Э. БАУМАНА

УДК _____

№ госрегистрации _____

Инв. № _____

УТВЕРЖДАЮ

Преподаватель

«_____» _____ 2019 г.

ДИСЦИПЛИНА АНАЛИЗ АЛГОРИТМОВ
ОТЧЁТ
ПО ЛАБОРАТОРНОЙ РАБОТЕ №1

Расстояние Левенштейна
(промежуточный)

Студент

_____ Ф.М. Набиев

Преподаватели

Л.Л. Волкова, Ю.В. Строганов

Москва, 2019

СОДЕРЖАНИЕ

Введение	3
1 Аналитический раздел	4
1.1 Описание алгоритмов	4
1.1.1 Алгоритм Левенштейна	4
1.1.2 Алгоритм Дамерау-Левенштейна	5
2 Конструкторский раздел	6
2.1 Модель	6
2.2 Разработка алгоритмов	6
2.2.1 Алгоритм Вагнера-Фишера	6
2.2.2 Матричный алгоритм Дамерау-Левенштейна	7
2.2.3 Рекурсивный алгоритм Дамерау-Левенштейна	10
2.3 Сравнительный анализ реализаций	12
2.3.1 Оценка сложности	13
2.3.2 Оценка памяти	13
2.3.3 Итог	14
3 Технологический раздел	15
3.1 Требования к программному обеспечению	15
3.2 Средства реализации	15
3.3 Листинги кода	16
3.4 Описание тестирования	18
4 Исследовательский раздел	20
Заключение	21

ВВЕДЕНИЕ

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

Данные алгоритмы решают проблему поиска редакционного расстояния между двумя строками. Редакционное расстояние определяется количеством некоторых операций, необходимых для превращения одного слова в другое, а так же стоимостью этих операций.

Для достижения поставленной цели необходимо решить следующие задачи:

- изучение алгоритмов Левенштейна и Дамерау-Левенштейна нахождения расстояния между строками;
- применение метода динамического программирования для матричной реализации указанных алгоритмов;
- получение практических навыков реализации указанных алгоритмов: двух алгоритмов в матричной версии и одного из алгоритмов в рекурсивной версии;
- сравнительный анализ линейной и рекурсивной реализаций выбранного алгоритма определения расстояния между строками по затрачиваемым ресурсам (времени и памяти);
- 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;
- описание и обоснование полученных результатов в отчете о выполненной лабораторной работе, выполненного как расчётно-пояснительная записка к работе.

1 Аналитический раздел

1.1 Описание алгоритмов

Алгоритм Дамерау-Левенштейна является модификацией алгоритма Левенштейна. Рассмотрим данные методы подробнее.

1.1.1 Алгоритм Левенштейна

Расстояние Левенштейна между двумя строками - это минимальная сумма произведений количества операций вставки, удаления и замены одного символа, необходимых для превращения одной строки в другую, на их стоимость.

Вышеописанные операции имеют следующие обозначения:

- I (insert) - вставка;
- D (delete) - удаление;
- R (replace) - замена;

При этом $\text{cost}(x)$ есть обозначение стоимости некоторой операции x . Будем считать, что символы в строках нумеруются с первого. Пусть S_1 и S_2 - две строки с длинами N и M соответственно. Тогда расстояние Левенштейна $D(M, N)$ вычисляется по формуле (1.1):

$$D(i, j) = \begin{cases} 0, & i = 0, j = 0 \\ i * \text{cost}(D), & j = 0, i > 0 \\ j * \text{cost}(I), & i = 0, j > 0 \\ D(i - 1, j - 1), & S_1[i] = S_2[j] \\ \min(& \\ D(i, j - 1) + \text{cost}(I), & \\ D(i - 1, j) + \text{cost}(D), & j > 0, i > 0, S_1[i] \neq S_2[j] \\ D(i - 1, j - 1) + \text{cost}(R) & \\), & \end{cases} \quad (1.1)$$

где $\min(a, b, c)$ возвращает наименьшее значение из a, b, c .

1.1.2 Алгоритм Дамерау-Левенштейна

Определение расстояния Дамерау-Левенштейна аналогично определению расстояния Левенштейна с учётом новой операции - перестановки соседних символов (транспозиции). Соответственно, обозначения операций:

- I (insert) - вставка;
- D (delete) - удаление;
- R (replace) - замена;
- T (transpose) - перестановка соседних символов.

При тех же обозначениях имеем формулы (1.2) и (1.3):

$$D(i,j) = \begin{cases} \min(A, D(i-2, j-2) + \text{cost}(T), & i > 1, j > 1, \\ & S_1[i] = S_2[j-1], \\ & S_1[i-1] = S_2[j] \\ A & \text{Иначе} \end{cases} \quad (1.2)$$

где A:

$$A = \begin{cases} 0, & i = 0, j = 0 \\ i * \text{cost}(D), & j = 0, i > 0 \\ j * \text{cost}(I), & i = 0, j > 0 \\ D(i-1, j-1), & S_1[i] = S_2[j] \\ \min(& \\ D(i, j-1) + \text{cost}(I), & \\ D(i-1, j) + \text{cost}(D), & j > 0, i > 0, S_1[i] \neq S_2[j] \\ D(i-1, j-1) + \text{cost}(R) & \\), & \end{cases} \quad (1.3)$$

2 Конструкторский раздел

2.1 Модель

IDEFØ модель задачи вычисления редакционного расстояния приведена на рисунке 2.1.

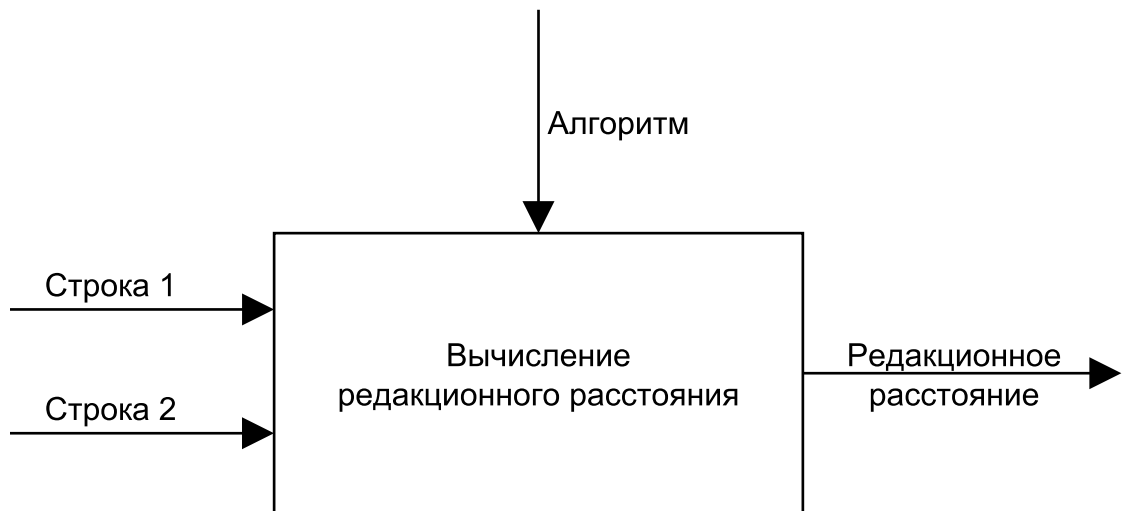


Рисунок 2.1 — IDEFØ модель

2.2 Разработка алгоритмов

Для непосредственной реализации вышеописанных алгоритмов важно иметь их некоторые упрощённые формальные представления, так как чтение таких представлений упрощает написание кода. Подходящим для этого вариантом визуализации являются схемы алгоритмов.

2.2.1 Алгоритм Вагнера-Фишера

Алгоритм Вагнера-Фишера является матричной реализацией поиска расстояния Левенштейна. Схема данного алгоритма приведена на рисунке 2.2.

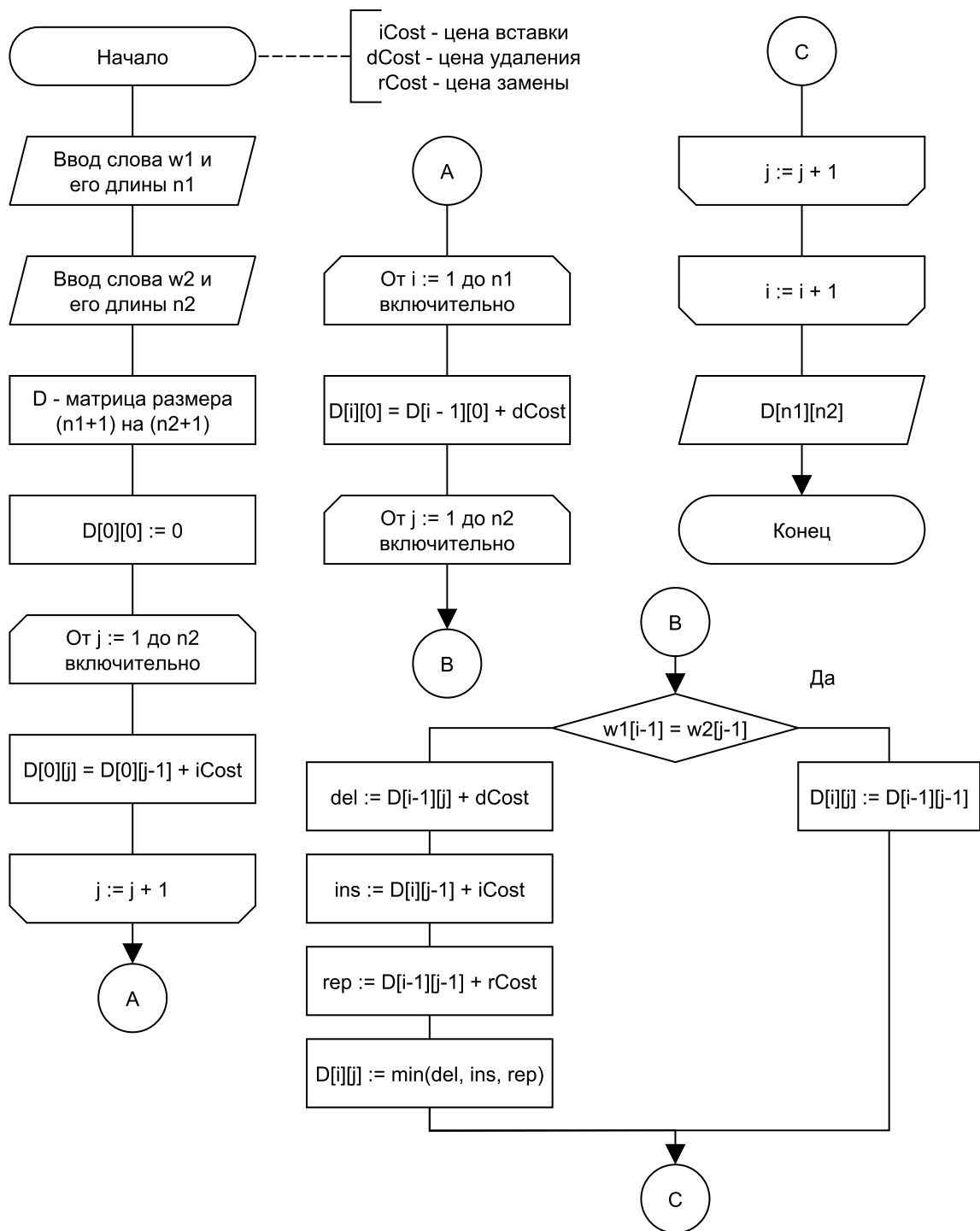


Рисунок 2.2 — Алгоритм Вагнера-Фишера

2.2.2 Матричный алгоритм Дамерау-Левенштейна

Матричный алгоритм Дамерау-Левенштейна представляет из себя модификацию алгоритма Вагнера-Фишера, в котором происходит дополнительная

проверка на возможность проведения операции транспозиции. Схема данного алгоритма приведена на рисунках 2.3 и 2.4

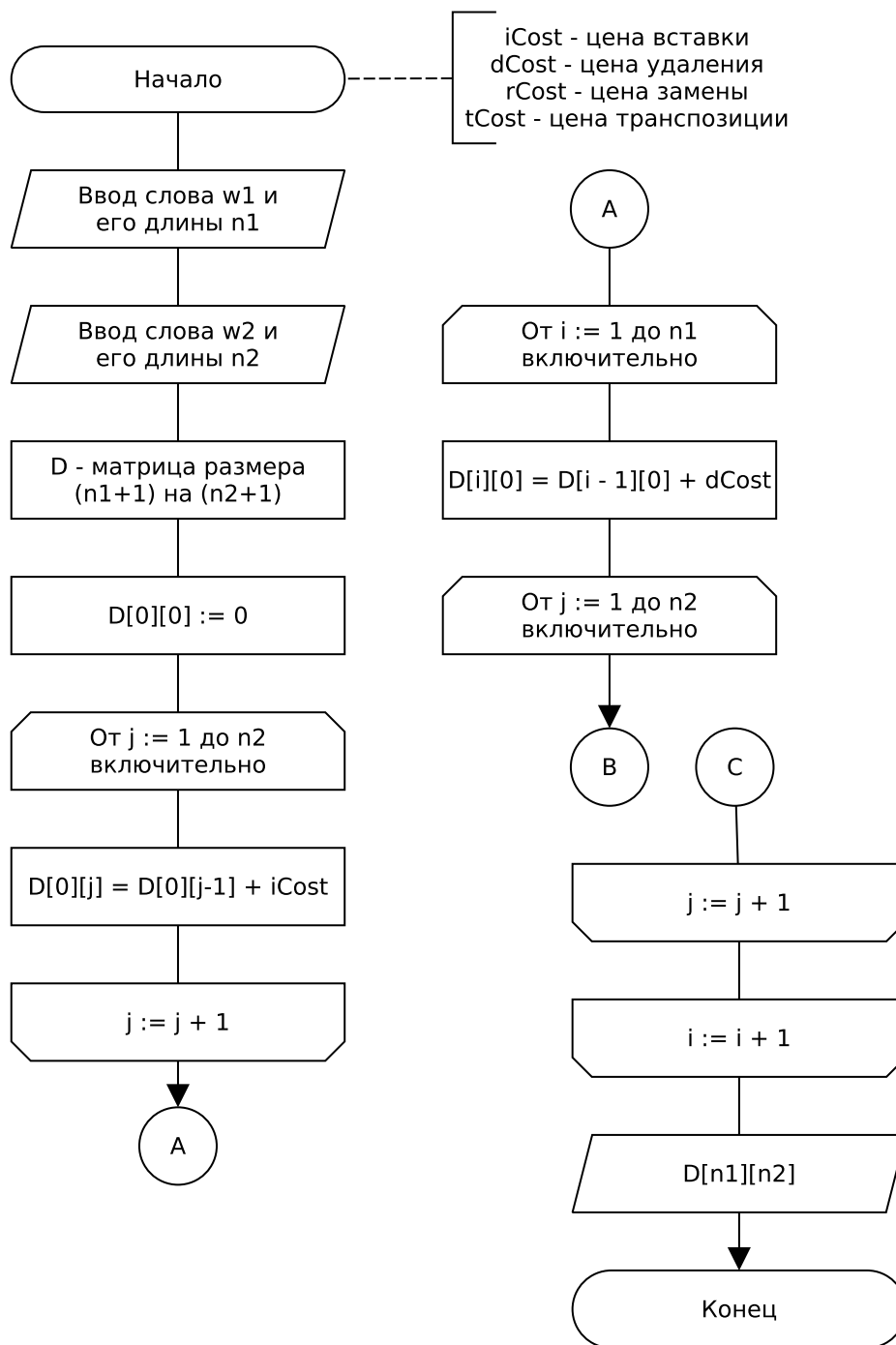


Рисунок 2.3 — Матричный алгоритм Дамерау-Левенштейна, часть 1

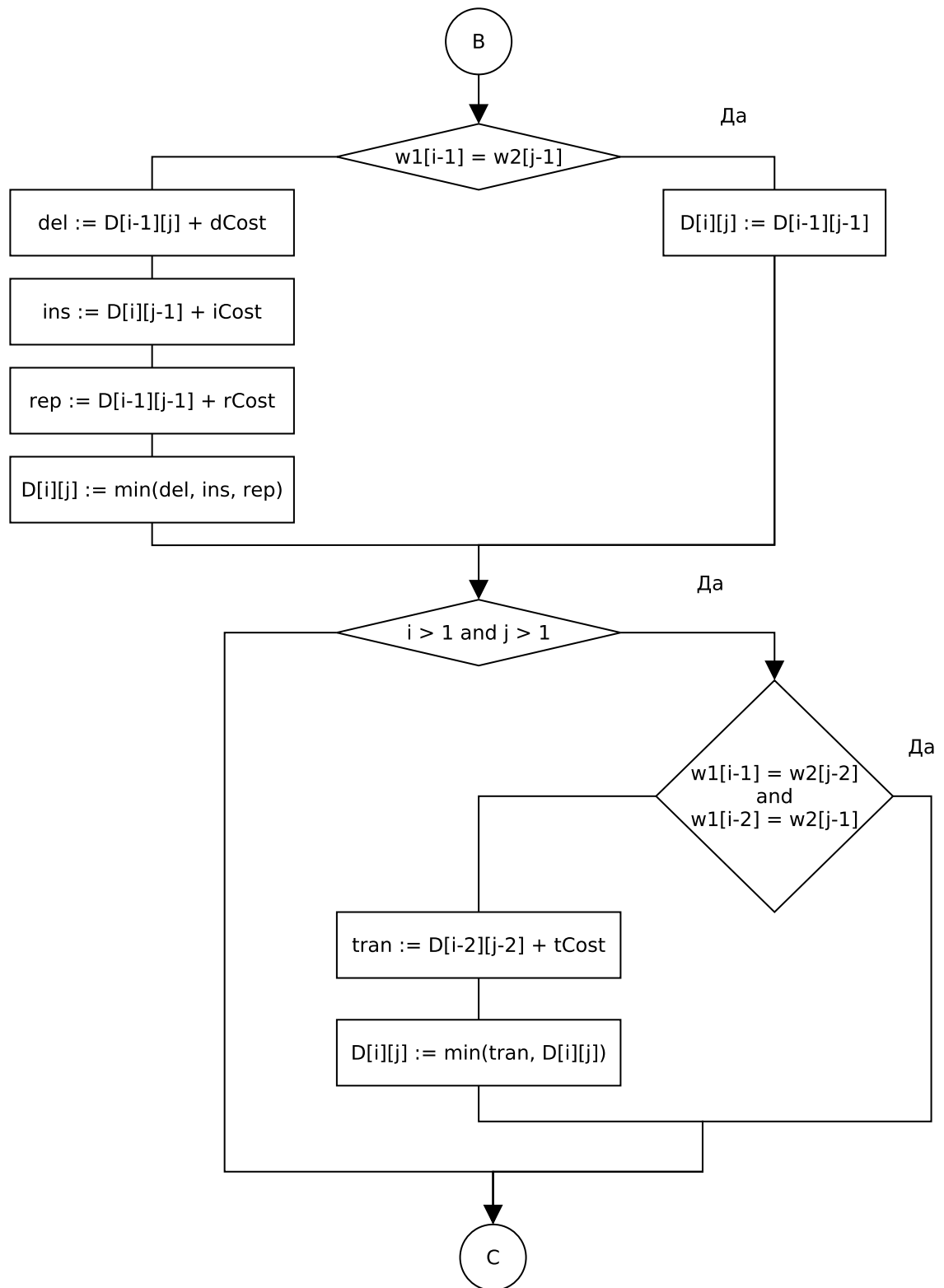


Рисунок 2.4 — Матричный алгоритм Дамерау-Левенштейна, часть 2

2.2.3 Рекурсивный алгоритм Дамерау-Левенштейна

Суть рекурсивного алгоритма Левенштейна состоит в сведении поиска редакционного расстояния до тривиального случая, когда длина хотя бы одного из слов равна 0. Отличие алгоритма Левенштейна от алгоритма Дамерау-Левенштейна было описано ранее. Схема данного алгоритма приведена на рисунках 2.5, 2.6 и 2.7.

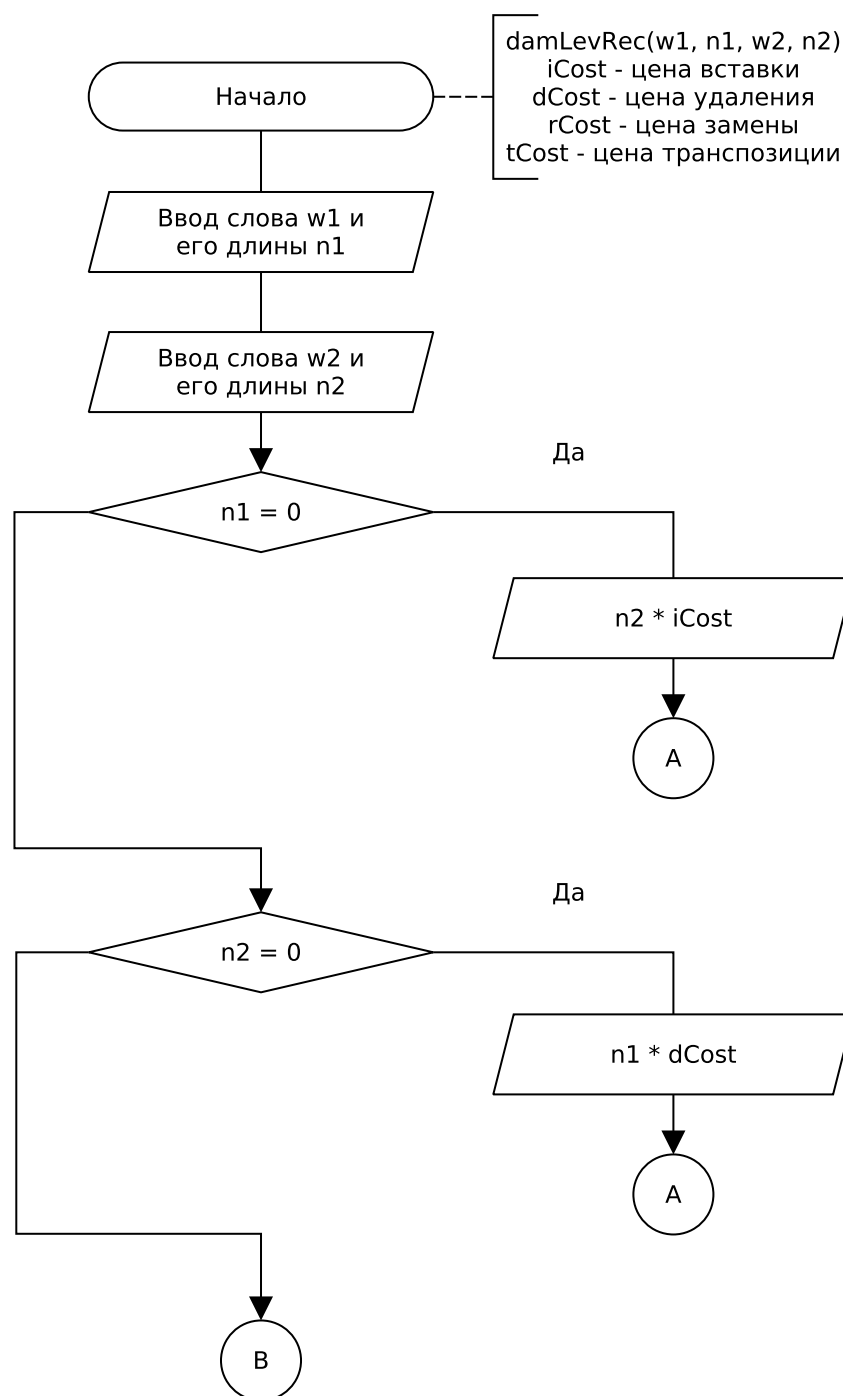


Рисунок 2.5 — Рекурсивный алгоритм Дамерау-Левенштейна, часть 1

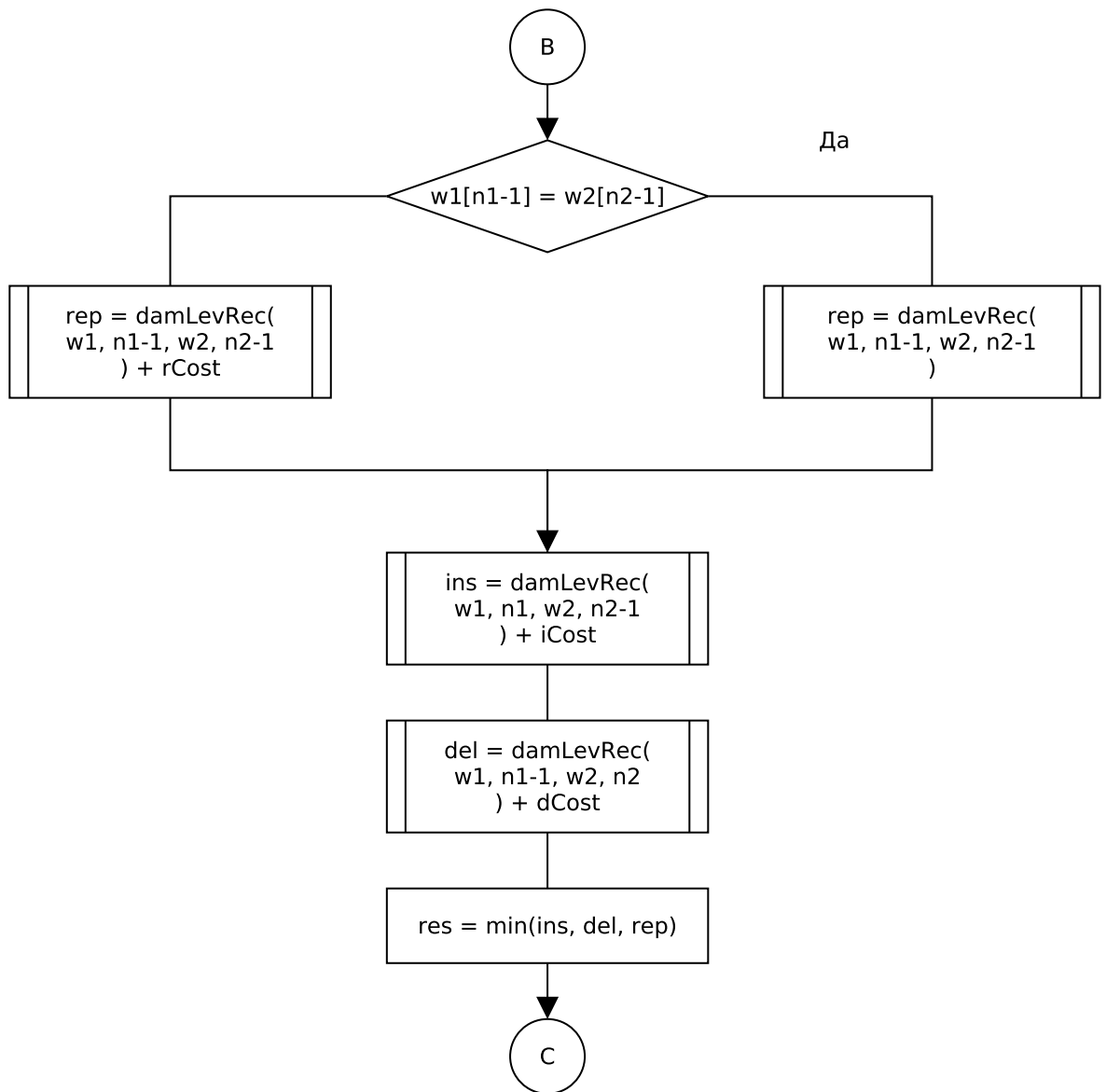


Рисунок 2.6 — Рекурсивный алгоритм Дамерау-Левенштейна, часть 2

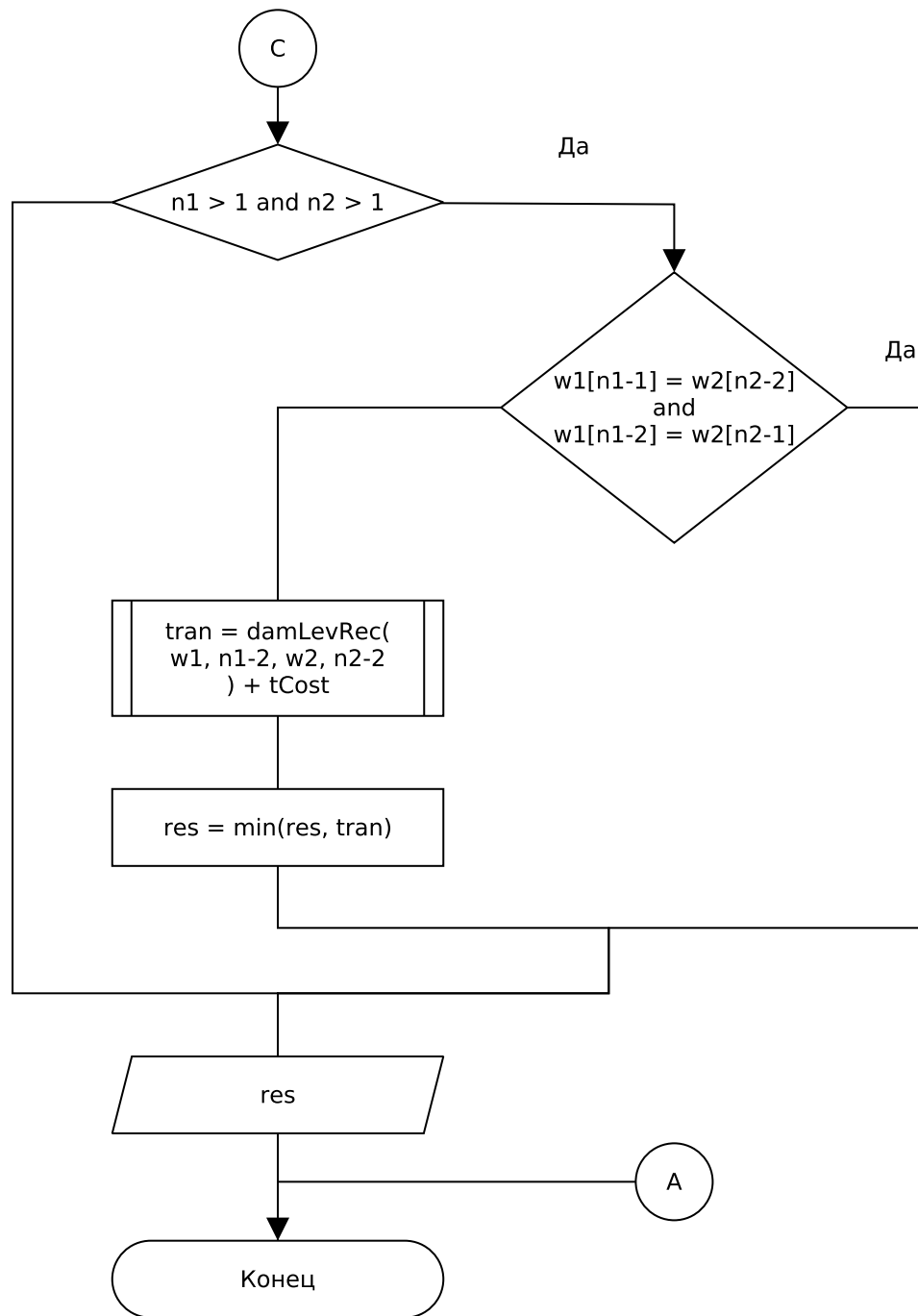


Рисунок 2.7 — Рекурсивный алгоритм Дамерау-Левенштейна, часть 3

2.3 Сравнительный анализ реализаций

Алгоритм Дамерау-Левенштейна был рассмотрен в двух вариациях: рекурсивной и матричной. Сравним эти реализации.

2.3.1 Оценка сложности

Произведем оценку общей сложности рекурсивного алгоритма. В рассмотренной реализации присутствуют 4 точки входа в рекурсию. Условием выхода из рекурсии является равенство длины хотя бы одной из строк нулю. Из этого можно сделать вывод, что рекурсивная вариация алгоритма Дамерау-Левенштейна имеет сложность $O(4^{\max(n_1, n_2)})$, где n_1 и n_2 - длины обрабатываемых строк.

Что касается матричной реализации, её задача сводится к полному обходу матрицы размера $(n_1 + 1) \cdot (n_2 + 1)$, где n_1 и n_2 - длины обрабатываемых строк. Следовательно, данный алгоритм имеет сложность $O(n_1 \cdot n_2)$.

2.3.2 Оценка памяти

Память, затрачиваемая на выполнение рассматриваемых алгоритмов зависит от используемых типов данных и соглашения о вызовах. В качестве примера, будем считать, что обрабатываемые строки передаются в функции по указателю, размер одного символа составляет 1 байт, целочисленного типа - 4 байта, указателя - 8 байт, используется cdecl (параметры функции передаются через стек, в который так же помещаются значения адреса возврата и указателя на верхушку текущего стекового кадра).

Как было показано ранее, в рекурсивной реализации происходит $4^{\max(n_1, n_2)}$ вызовов функции. Допустим, что в аргументах функции передаются два указателя на обрабатываемые строки и 2 целочисленные переменные, означающие длины этих строк. Без учета возможного использования локальных переменных, имеем следующие затраты памяти:

$$4^{\max(n_1, n_2)} \cdot (8 + 8 + 8 + 4 + 8 + 4) = 10 \cdot 4^{\max(n_1, n_2) + 1} \quad (2.1)$$

В случае матричной вариации алгоритма Дамерау-Левенштейна будем считать, что функция использует матрицу размера $(n_1 + 1) \cdot (n_2 + 1)$, указатель на начало этой матрицы и два целочисленных счётчика для циклов. Тогда

имеем следующие затраты памяти:

$$\begin{aligned} 40 + (n_1 + 1) \cdot (n_2 + 1) \cdot 4 + 8 + 4 + 4 = \\ 56 + 4 \cdot (n_1 + 1) \cdot (n_2 + 1) \end{aligned} \tag{2.2}$$

2.3.3 Итог

Таким образом, можно сделать вывод о том, что матричная реализация алгоритма Дамерау-Левенштейна работает быстрее и потребляет меньше памяти, чем рекурсивная.

3 Технологический раздел

3.1 Требования к программному обеспечению

Требования к вводу:

- На вход подаются две строки, разделённые переносом строки;
- Одна и та же буква в верхнем и нижнем регистрах считается как разные символы;
- Пустая строка может быть введена.

Требования к выводу:

- Редакционное расстояние;
- В случае матричной реализации алгоритма выводить матрицу, полученную в ходе вычисления расстояния.

Требования к программе:

- Выбор алгоритма происходит через аргументы командной строки.

3.2 Средства реализации

Для реализации программы вычисления редакционного расстояния мной был выбран язык программирования C++. В рамках текущей задачи данный язык программирования имеет ряд существенных преимуществ:

- Статическая типизация;
- Близость к низкоуровневому C при наличии многих возможностей высокоуровневых языков;
- Встроенная библиотека `std::chrono`, позволяющая измерять процессорное время.

3.3 Листинги кода

Реализации алгоритмов Вагнера-Фишера и матричной и рекурсивной вариаций Дамерау-Левенштейна приведены в листингах 3.1, 3.2, 3.3 соответственно.

Листинг 3.1 — Расстояние Левенштейна (матричная реализация)

```
1  template<typename _Word_t>
2  int  WagnerFischer<_Word_t>::distance(_Word_t w1, int n1,
3                                         _Word_t w2, int n2) {
4      int **D = Util::createMatrix<int>(n1 + 1, n2 + 1);
5      D[0][0] = 0;
6
7      for (int j = 1; j <= n2; j++) {
8          D[0][j] = D[0][j - 1] + insertCost;
9      }
10
11     for (int i = 1; i <= n1; i++) {
12         D[i][0] = D[i - 1][0] + deleteCost;
13
14         for (int j = 1; j <= n2; j++) {
15             if (w1[i - 1] == w2[j - 1]) {
16                 D[i][j] = D[i - 1][j - 1];
17             }
18             else {
19                 D[i][j] = std::min(std::min(D[i - 1][j] + deleteCost,
20                                             D[i][j - 1] + insertCost),
21                                   D[i - 1][j - 1] + replaceCost);
22             }
23         }
24     }
25
26     int res = D[n1][n2];
27     delete [] D;
28
29     return res;
30 }
```


Листинг 3.2 — Расстояние Дамерау-Левенштейна (матричная реализация)

```

1  template<typename _Word_t>
2  int  DamerauLevenshtein<_Word_t>::distance(_Word_t w1, int n1,
3                                             _Word_t w2, int n2) {
4      int **D = Util::createMatrix<int>(n1 + 1, n2 + 1);
5      D[0][0] = 0;
6
7      for (int j = 1; j <= n2; j++) {
8          D[0][j] = D[0][j - 1] + insertCost;
9      }
10
11     for (int i = 1; i <= n1; i++) {
12         D[i][0] = D[i - 1][0] + deleteCost;
13
14         for (int j = 1; j <= n2; j++) {
15             if (w1[i - 1] == w2[j - 1]) {
16                 D[i][j] = D[i - 1][j - 1];
17             }
18             else {
19                 D[i][j] = std::min(std::min(D[i - 1][j] + deleteCost,
20                                             D[i][j - 1] + insertCost),
21                                   D[i - 1][j - 1] + replaceCost);
22
23                 if (i > 1 && j > 1) {
24                     if (w1[i - 1] == w2[j - 2] && w1[i - 2] == w2[j - 1]) {
25                         D[i][j] = std::min(D[i - 2][j - 2] + transposeCost,
26                                             D[i][j]);
27                     }
28                 }
29             }
30         }
31     }
32
33     int res = D[n1][n2];
34     delete [] D;
35
36     return res;
37 }
```

Листинг 3.3 — Расстояние Дameraу-Левенштейна (рекурсивная реализация)

```
1  template<typename _Word_t>
2  int  DamerauLevenshteinRecursive<_Word_t>::distance(_Word_t w1, int n1,
3                                                     _Word_t w2, int n2) {
4      if (n1 == 0) {
5          return n2 * insertCost;
6      }
7      if (n2 == 0) {
8          return n1 * deleteCost;
9      }
10
11     bool isSame = w1[n1 - 1] == w2[n2 - 1];
12
13     int res = std::min(std::min(distance(w1, n1 - 1, w2, n2) + deleteCost,
14                                     distance(w1, n1, w2, n2 - 1) + insertCost),
15                       distance(w1, n1 - 1, w2, n2 - 1) + replaceCost *
16                               !isSame);
17
18     if (n1 > 1 && n2 > 1) {
19         if (w1[n1 - 1] == w2[n2 - 2] && w1[n1 - 2] == w2[n2 - 1]) {
20             res = std::min(distance(w1, n1 - 2, w2, n2 - 2) +
21                             transposeCost, res);
22         }
23     }
24
25     return res;
26 }
```

3.4 Описание тестирования

Для тестирования программы были подготовлены данные, представленные в таблице 3.1.

Таблица 3.1 — Тестовые данные

№	Строка 1	Строка 2	Ожидаемое расстояние Левенштейна	Ожидаемое расстояние Дамерау-Левенштейна
1	some	any	4	4
2	man	person	6	6
3		nothing	7	7
4			0	0
5	bashrc	bashcr	2	1
6	goose	duck	5	5
7	bus	BuS	2	2
8	electricity	city	7	7
9	powerfull	powerless	4	4
10	grow	flow	2	2
11	rise	rice	1	1
12	fake	lake	1	1
13	same	same	0	0

4 Исследовательский раздел

ЗАКЛЮЧЕНИЕ