

Московский авиационный институт  
(национальный исследовательский университет)

Факультет информационных технологий и прикладной математики

Кафедра вычислительной математики и программирования

**Лабораторная работа №5 по курсу «Дискретный  
анализ»**

Студент: П.А. Земсков

Преподаватель: Н.К. Макаров

Группа: М8О-301Б

Дата: \_\_\_\_\_

Оценка: \_\_\_\_\_

Подпись: \_\_\_\_\_

Москва, 2025

# Содержание

1	Лабораторная работа №5	2
2	Описание	2
3	Исходный код	2
4	Консоль	6
5	Тест производительности	6
6	Выводы	7
7	Список литературы	7

# 1 Лабораторная работа №5

## Задача

Требуется разработать программу для поиска образцов в заранее известном тексте с использованием суффиксного массива.

**Вариант алгоритма:** Построение суффиксного массива с последующим поиском образцов методом бинарного поиска.

## 2 Описание

Суффиксным массивом (англ. suffix array, суфмасс) строки  $s$  называется перестановка индексов начал её суффиксов, которая задаёт порядок их лексикографической сортировки. Иными словами, чтобы его построить, нужно выполнить сортировку всех суффиксов заданной строки. Далее ниже приведен исходный код программы

## 3 Исходный код

**Общая идея:** Программа реализует поиск подстроки в тексте с использованием суффиксного массива. Сначала функция `buildSuffixArray` строит массив суффиксов строки, где каждый элемент представляет индекс начала суффикса. Алгоритм работает за  $O(n \log^2 n)$ : сначала сортируются суффиксы по первым символам, затем итеративно происходит удвоение длины сравниваемых префиксов. После построения массива суффиксов поиск подстроки выполняется бинарным поиском по этому массиву. Функции `lowerBound` и `upperBound` находят диапазон суффиксов, начинающихся с заданного шаблона. Разность этих индексов даёт количество вхождений подстроки, а сами позиции определяются как `sa[i] + 1`. Таким образом, поиск одной подстроки имеет сложность  $O(m \log n)$ . Главная идея алгоритма заключается в том, чтобы один раз построить индекс для текста, после чего искать любые подстроки значительно быстрее, чем при использовании линейного поиска.

Листинг 1: Исходный код

```

1 #include <iostream>
2 #include <string>
3 #include <vector>
4 #include <algorithm>
5
6 std::vector<int> buildSuffixArray(const std::string &s)
7 {
8     int n = s.length();
9     std::vector<int> sa(n), rank(n);
10    for (int i = 0; i < n; ++i)
11    {
12        sa[i] = i;
13        rank[i] = static_cast<unsigned char>(s[i]);
14    }
15    for (int k = 1; k < n; k <= 1)
16    {
17        std::stable_sort(sa.begin(), sa.end(), [&](int a, int b)
18            {
19                if (rank[a] != rank[b]) return rank[a] < rank[b];
20                int ra = (a + k < n) ? rank[a + k] : -1;
21                int rb = (b + k < n) ? rank[b + k] : -1;
22                return ra < rb; });
23        std::vector<int> new_rank(n);
24        new_rank[sa[0]] = 0;
25        for (int i = 1; i < n; ++i)
26        {
27            int a = sa[i - 1], b = sa[i];
28            int ra0 = rank[a], rb0 = rank[b];
29            int ra1 = (a + k < n) ? rank[a + k] : -1;
30            int rb1 = (b + k < n) ? rank[b + k] : -1;
31            new_rank[b] = new_rank[a] + (ra0 != rb0 || ra1 != rb1);
32        }
33        rank.swap(new_rank);
34        if (rank[sa[n - 1]] == n - 1)
35            break;
36    }
37    return sa;
38 }
39
40 int lowerBound(const std::string &text, const std::string &pattern,
41               const std::vector<int> &sa)
42 {
43     int n = (int)text.size();
44     int L = 0, R = n;
45     while (L < R)
46     {
47         int M = (L + R) >> 1;
48         int cmp = text.compare(sa[M], pattern.size(), pattern);
49         if (cmp < 0)
50             L = M + 1;

```

```

50         else
51             R = M;
52     }
53     return L;
54 }
55
56 int upperBound(const std::string &text, const std::string &pattern,
57               const std::vector<int> &sa)
58 {
59     int n = (int)text.size();
60     int L = 0, R = n;
61     while (L < R)
62     {
63         int M = (L + R) >> 1;
64         int cmp = text.compare(sa[M], pattern.size(), pattern);
65         if (cmp <= 0)
66             L = M + 1;
67         else
68             R = M;
69     }
70     return L;
71 }
72
73 int main()
74 {
75     std::ios::sync_with_stdio(false);
76     std::cin.tie(nullptr);
77
78     std::string text;
79     if (!std::getline(std::cin, text))
80         return 0;
81
82     std::vector<int> sa = buildSuffixArray(text);
83
84     std::string pattern;
85     int qid = 0;
86     while (std::getline(std::cin, pattern))
87     {
88         ++qid;
89
90         int lb = lowerBound(text, pattern, sa);
91         int ub = upperBound(text, pattern, sa);
92
93         if (lb >= ub)
94             continue;
95
96         std::vector<int> pos;
97         pos.reserve(ub - lb);
98         for (int i = lb; i < ub; ++i)
99             pos.push_back(sa[i] + 1);
100         std::sort(pos.begin(), pos.end());

```

```
100
101     std::cout << qid << ":\u";
102     for (size_t i = 0; i < pos.size(); ++i)
103     {
104         if (i)
105             std::cout << ",\u";
106         std::cout << pos[i];
107     }
108     std::cout << "\n";
109 }
110 return 0;
111 }
```

## 4 Консоль

Пример компиляции и демонстрация работы программы:

```
C:\Users\jocke\Documents\diskran\lab5> g++ --std=c++20 main.cpp -o main
C:\Users\jocke\Documents\diskran\lab5> ./main
abcdabc
abcd
bcd
bc
1: 1
2: 2
3: 2, 6
```

## 5 Тест производительности

При сравнении производительности с `std::string::find` результаты показывают, что для одного поиска стандартный метод быстрее, так как имеет меньшие константы и линейную сложность  $O(n \cdot m)$ .

```
C:\Users\jocke\Documents\diskran> g++ -std=c++20 lab5.cpp -o main
C:\Users\jocke\Documents\diskran> g++ -std=c++20 benchmark.cpp -o benchmark
C:\Users\jocke\Documents\diskran> .\benchmark
Search in suffix array time: 237 ms
std::find time: 4783 ms
```

### Результат:

Однако при множественных запросах суффиксный массив выигрывает за счёт того, что индекс строится один раз, а каждый последующий поиск выполняется в логарифмическое время. Например, при тексте длиной 10 символов и 1000 запросов `std::find` требует около 5 секунд, тогда как поиск с использованием суффиксного массива занимает около 0.2 секунды.

## 6 Выводы

Суффиксный массив эффективен при большом количестве запросов к одному и тому же тексту, в то время как для одиночного поиска выгоднее применять `std::find`.

## 7 Список литературы

1. Кормен Т. Х., Лейзерсон Ч. И., Ривест Р. Л., Штайн К. *Алгоритмы: построение и анализ. 3-е изд.* — М.: Издательский дом «Вильямс», 2010.
2. Knuth D. E. *The Art of Computer Programming. Vol. 3: Sorting and Searching. 2nd ed.* — Addison-Wesley, 1998.
3. *Суффиксный массив — Алгоритмика.* <https://ru.algorithmica.org/cs/string-structure/suffix-array>.