

Федеральное государственное бюджетное образовательное учреждение  
высшего образования «Сибирский государственный университет телекоммуникаций  
и информатики»

Факультет ИВТ

Кафедра вычислительных систем

**Курсовая работа**  
на тему «Полнотекстовый поиск по шаблону»  
Вариант 3.3 «Алгоритм Кнута-Морриса-Пратта»

Выполнил:  
студент гр. ИВ-221  
Анциферов Я.Г.

Проверил:  
Старший преподаватель Кафедры ВС  
Фульман В.О.

Новосибирск, 2023

## Тема курсовой работы: Алгоритм Кнута-Морриса-Практа

Тема выбранного раздела: Полнотекстовый поиск по шаблону

### Задание на курсовую работу

Реализовать программу **kmpmatcher** (Knuth–Morris–Pratt string MATCHER) полнотекстового поиска по шаблону. Шаблон и имя файла (директории), в которой осуществляется поиск, передаются через аргументы командной строки в следующем порядке:

\$ <b>kmpmatcher</b> "g*.le" ~/mydir	#Анализ всех файлов, расположенных в ~/mydir.
\$ <b>kmpmatcher -r</b> "g*.le" ~/mydir	#Рекурсивный поиск во всех директориях, расположенных ниже ~/mydir.

### Критерии оценки

- **Оценка «отлично»:** разработанная программа обеспечивает поиск текста по шаблону рекурсивно в заданной директории. Под рекурсивным поиском понимается анализ всех текстовых файлов в текущей директории, а также во всех вложенных директориях.
- **Оценка «хорошо»:** разработанная программа не предусматривает поиск по шаблону ИЛИ не способна выполнять рекурсивный поиск в дереве каталогов (поиск только в одном файле).
- **Оценка «удовлетворительно»:** реализован только алгоритм Кнута-Морриса-Практа.

### Указание к выполнению задания

Алгоритм Кнута-Морриса-Практа (КМП) основан на применении префикс-функции  $\pi$ , подробно описанной в общей информации к данному разделу. В листинге 5 приведен псевдокод алгоритма КМП.

#### Листинг 5. Псевдокод алгоритма КМП

```
KMP_MATCHER(T, P)
  m ← len(P)
  n ← len(T)
   $\pi P \leftarrow \text{COMPUTE\_PREFIX\_F}(P)$ 
  q ← 0
  for i ← 1 to n do
    while q > 0 и  $P[q + 1] \neq T[i]$  do
      q ←  $\pi P[q]$ 
    if  $P[q + 1] = T[i]$  then
      q ← q + 1
    if q = m then
      print "Образец обнаружен при сдвиге" i – m
    q ←  $\pi P[q]$ 
```

## Анализ задачи

1. Для реализации алгоритма Кнута-Морриса-Пратта необходимо реализовать префикс-функцию, которая является основой этого алгоритма. Ее предназначение заключается в вычислении наибольшей длины собственного суффикса строки (не совпадает со всей строкой), который совпадает с префиксом строки для каждого символа в ней.

### Пример:

Дана строка “abaabab”. После прохода префикс-функции результат будет таким: 0 0 1 1 2 3 2

### Принцип работы:

Для  $i = 0$

“abaabab”

Для нулевого элемента результат всегда будет равен 0, так как он не является собственным суффиксом.

Для  $i = 1$

“abaabab”

Результат 0, так как “a” != “b”

Для  $i = 2$

“abaabab”

Результат 1, так как “a” == “a”, но “ab” != “ba”

Для  $i = 3$

“abaabab”

Результат 1, так как “a” == “a”, но “ab” != “aa” и “aba” != “baa”

Для  $i = 4$

“abaabab”

Результат 2, так как “ab” == “ab”, но “aba” != “aab” и “abaa” != “baab”

Для  $i = 5$

“abaabab”

Результат 3, так как “aba” == “aba”, но “abaa” != “aaba” и т.д.

Для  $i = 6$

“abaabab”

Результат 2, так как “ab” == “ab”, далее суффиксы с префиксами не совпадают.

2. После реализации префикс-функции нужно составить алгоритм Кнута-Морриса-Пратта. Для его работы необходимы: шаблон и строка, где будет происходить поиск этого шаблона, также нужно создать массив для хранения данных от префикс-функции. Сам алгоритм выводит на экран индекс вхождения в строке. Преимущество этого алгоритма перед обычным поиском состоит в том, что он перемещается по строке намного быстрее (обычный поиск начнет работу с прошлой позиции плюс один, в то время как алгоритм может переместиться в случае неудачи максимум на длину шаблона и минимум на 1 от предыдущей позиции).

### Пример:

Дана строка  $S = \text{“abbcabbbcbba”}$ , шаблон  $P = \text{“bbcb”}$ , ожидаемый результат 5.

### Принцип работы:

Сначала находим префикс-функцию от  $P = 0\ 1\ 0\ 1\ 2$   
Затем начинаем поиск шаблона в строке

abbcabbcbba

bbcbb

Символы не совпадают, ищем дальше. Индекс = 0

abbcabbcbba

bbcbb

Символы совпадают, идем дальше по шаблону и строке. Индекс = 1

abbcabbcbba

bbcbb

Символы совпадают, идем дальше по шаблону и строке. Индекс = 2

abbcabbcbba

bbcbb

Символы совпадают, идем дальше по шаблону и строке. Индекс = 3

abbcabbcbbba

bbcbb

Символы не совпадают, поэтому мы берем значение индекса и уменьшаем его на один, так как в массиве префикс-функции в элементе [Индекс - 1] хранится 0, следовательно начинаем поиск заново с индекса 4.

abbcabbcbbba

bbcbb

Символы не совпадают, перемещаемся вперед на единицу. Индекс = 4

abbcabbcbbba

bbcbb

Символы совпадают, идем дальше по шаблону и строке. Индекс = 5

abbcabbcbba

bbcbb

Символы совпадают, идем дальше по шаблону и строке. Индекс = 6

abbcabbcbba

bbcbb

Символы совпадают, идем дальше по шаблону и строке. Индекс = 7

abbcabbcbba

bbcbb

Символы совпадают, идем дальше по шаблону и строке. Индекс = 8

abbcabbcbba

bbcbb

Шаблон найден, возвращаем индекс вхождения, он равен

$$9 \text{ (текущая позиция в строке)} - 5 \text{ (длина шаблона)} + 1 = 5$$

### Тестовые данные

Необходимо провести тестирование приложения. Предполагается, что при корректных данных будет проведен поиск шаблона в тексте и при некорректных выведена ошибка и информация о ней.

Тест на вхождение подстроки “ABC” в строку “ABCADBABCDACBABCAA”:

```
Entry test for 'ABCADBABCDACBABCAA':  
'ABC' in position 0  
'ABC' in position 6  
'ABC' in position 13  
No more matches found
```

Тест на вхождение подстроки “WASHINGTON” в строку  
“LONDONISTHECAPITALOFGREATBRITAIN”:

```
Entry test for 'LONDONISTHECAPITALOFGREATBRITAIN':  
No more matches found
```

Тест на вхождение подстроки “MOSCOW” в строку “MOSCOWISTHETHIRDROME”:

```
Entry test for 'MOSCOWISTHETHIRDROME':  
'MOSCOW' in position 0  
No more matches found
```

Тест на вхождение подстроки “NUMBERS” в строку “323353243411422131231”

```
Entry test for '323353243411422131231':  
No more matches found
```

## Листинг программы

### main.c

```
1 #include <stdio.h>
2
3 void algorithm_KMP (char* massive, char* pod_massive, int* array, int size_m, int size_p_m) {
4     printf("Entry test for '%s':\n", massive);
5
6     int i = 0;
7     int j = 0;
8
9     while ((i < size_m - 1) || (j < size_p_m - 1)) {
10         if (massive[i] == pod_massive[j]) {
11             i++;
12             j++;
13         }
14         else {
15             if (j == 0)
16                 i++;
17             else {
18                 j--;
19                 i -= array[j];
20                 j = 0;
21             }
22         }
23         if (j == size_p_m - 1)
24             printf("'%' in position %d\n", pod_massive, i - (size_p_m - 1));
25         if (i == size_m - 1) {
26             printf("No more matches found\n");
27             break;
28         }
29     }
30 }
31
32
33 void prefix_find (char* pod_massive, int* array, int size) {
34     *array = 0;
35     for (int i = 1; i < size; i++) {
36         int j = array[i-1];
37         while ((j > 0) && (pod_massive[i] != pod_massive[j]))
38             j = array[j-1];
39         if (pod_massive[i] == pod_massive[j])
40             j++;
41         array[i] = j;
42     }
43 }
44
45 int main () {
46     char massive[] = "323353243411422131231";
47     char pod_massive[] = "NUMBERS";
48     int array[sizeof(pod_massive)-1];
49
50     prefix_find(pod_massive, array, sizeof(pod_massive)-1);
51     algorithm_KMP(massive, pod_massive, array, sizeof(massive)-1, sizeof(pod_massive)-1);
52
53     return 0;
54 }
```