

# ИНФОРМАТИКА, ВЫЧИСЛИТЕЛЬНАЯ ТЕХНИКА И УПРАВЛЕНИЕ INFORMATION TECHNOLOGY, COMPUTER SCIENCE, AND MANAGEMENT



УДК 519.689

<https://doi.org/10.23947/1992-5980-2019-19-3-290-300>

## Ускоренный препроцессинг в задаче поиска подстрок в строке\*

А. В. Мазуренко<sup>1</sup>, Н. В. Болдырихин<sup>2\*\*</sup>

<sup>1</sup> ООО «ДДОС-Гвард», г. Ростов-на-Дону, Российская Федерация

<sup>2</sup> Донской государственный технический университет, г. Ростов-на-Дону, Российская Федерация

## Accelerated preprocessing in task of searching substrings in a string\*\*\*

A. V. Mazurenko<sup>1</sup>, N. V. Boldyrikhin<sup>2\*\*</sup>

<sup>1</sup> DDoS-GUARD LLC, Rostov-on-Don, Russian Federation

<sup>2</sup> Don State Technical University, Rostov-on-Don, Russian Federation

*Введение.* Бурное развитие таких систем, как Yandex, Google и пр., предопределило актуальность задачи поиска подстрок в строке. На сегодняшний день активно исследуются подходы к ее решению. Эта задача используется при создании систем управления базами данных, поддерживающих ассоциативный поиск. Кроме того, она применима при решении вопросов информационной безопасности, создании антивирусных программ. Алгоритмы поиска подстрок в строке используются в задачах обнаружения, основанного на сигнатурах.

*Материалы и методы.* Решение задачи базируется на алгоритме Ахо — Корасик, который представляет собой классический способ осуществления поиска подстрок в строке. Вместе с тем применен новый подход в части, касающейся предварительной обработки.

*Результаты исследования.* Показана возможность построения функции перехода и суффиксных ссылок при помощи суффиксных массивов и специальных отображений. Исследована взаимосвязь между префиксным деревом и суффиксными массивами. Это дало возможность разработать принципиально новый способ построения функций перехода и ошибок.

Полученные результаты позволяют существенно сократить время, затрачиваемое на предвыборную обработку множества строк образцов при использовании целочисленного алфавита.

В статье приведено восемь алгоритмов. Оценены разработанные алгоритмы. Полученные результаты сопоставлены с ранее известными. Доказаны две теоремы и восемь лемм. Приведены два примера, иллюстрирующие особенности практического применения разработанной процедуры препроцессинга.

*Обсуждение и заключения.* Предложенная в данной статье процедура препроцессинга основывается на связи между суффиксным массивом, созданным на основе множества строк образцов, и построением функций перехода и ошибок на начальных этапах работы алгоритма Ахо — Корасик. Такой подход отличен от традиционного и требует

*Introduction.* A rapid development of the systems such as Yandex, Google, etc., has predetermined the relevance of the task of searching substrings in a string, and approaches to its solution are actively investigated. This task is used to create database management systems that support associative search. Besides, it is applicable in solving information security issues and creating antivirus programs. Algorithms of searching substring in a string are used in signature-based discovery tasks.

*Materials and Methods.* The solution to the problem is based on the Aho-Corasick algorithm which is a typical technique of searching substrings in a string. At the same time, a new approach regarding preprocessing is employed.

*Research Results.* The possibility of constructing the transition function and suffix references through suffix arrays and special mappings, is shown. The relationship between the prefix tree and suffix arrays was investigated, which provided the development of a fundamentally new method of constructing the transition and error functions. The results obtained enable to substantially shorten the time intervals spent on the pre-election processing of a set of pattern strings when using an integer alphabet. The paper lists eight algorithms. The developed algorithms are evaluated. The results obtained are compared to the formerly known. Two theorems and eight lemmas are proved. Two examples illustrating features of the practical application of the developed preprocessing procedure are given.

*Discussion and Conclusions.* The preprocessing procedure proposed in this paper is based on the communication between the suffix array built on the ground of a set of pattern strings and the construction of transition and error functions at the initial stages of the Aho-Corasick algorithm. This approach differs from the traditional one and requires the use of algo-

\* Работа выполнена в рамках инициативной НИР.

\*\* E-mail: mazurencoal@gmail.com, boldyrikhin@mail.ru

\*\*\* The research is done within the frame of the independent R&D.



использования алгоритмов, позволяющих построить суффиксный массив за линейное время. Таким образом, описаны алгоритмы, позволяющие существенно сократить время на предварительную обработку множества строк образцов при условии использования определенного типа алфавита по сравнению с известным подходом, предложенным А. Ахо и М. Корасик.

Результаты исследований, приведенные в статье, могут быть применены в антивирусных программах, использующих поиск сигнатур вредоносных информационных объектов в памяти вычислительной системы. Кроме того, данный подход к решению задачи поиска подстроки в строке позволяет значительно ускорить работу систем управления баз данных, применяющих ассоциативный поиск.

**Ключевые слова:** поиск подстроки, алгоритм Ахо — Корасик, префиксное дерево, суффиксный массив, поиск информации, функция ошибок, функция перехода.

**Образец для цитирования:** Мазуренко, А.В. Ускоренный препроцессинг в задаче поиска подстрок в строке / А. В. Мазуренко Н. В. Болдырихин // Вестник Дон. гос. техн. ун-та. — 2019. — Т. 19, № 3. — С. 290–300. <https://doi.org/10.23947/1992-5980-2019-19-3-290-300>

rithms providing a suffix array in linear time. Thus, the algorithms that enable to significantly reduce the time for preprocessing of a set of pattern strings under the condition of using a certain type of alphabet in comparison to the known approach proposed in the Aho- Corasick algorithm are described. The research results presented in the paper can be used in antivirus programs that apply searching for signatures of malicious data objects in the memory of a computer system. In addition, this approach to solving the problem on searching substrings in a string will significantly speed up the operation of database management systems using associative search.

**Keywords:** string searching, Aho-Corasick algorithm, prefix tree, suffix array, information search, error function, transition function

**For citation:** A.V. Mazurenko, N.V. Boldyrikhin. Accelerated preprocessing in task of searching substrings in a string. Vestnik of DSTU, 2019, vol. 19, no. 3, pp. 290–300. <https://doi.org/10.23947/1992-5980-2019-19-3-290-300>

**Введение.** В настоящее время существенно возрастает важность обеспечения кибербезопасности распределенных информационных систем и отдельных вычислительных средств [1]. Диапазон таких задач достаточно широк [1–10]. Особое место занимают вопросы создания эффективного антивирусного программного обеспечения (ПО). Одной из важнейших задач, решаемых таким ПО, является поиск подстрок в строке [1, 5, 6, 10–13].

**Материалы и методы.** Задача поиска подстрок заключается в нахождении всех строк в тексте  $T$  общей длиной  $m$ , совпадающих с каким-либо образцом из заданного множества образцов  $P$ . Положим, сумма длин всех элементов  $P$ , состоящих из символов алфавита  $I$ , равна  $n$ . Решение этой задачи было предложено А. Ахо и М. Корасик [6, 10]. В их алгоритме время предвыборной обработки составляет  $O(n|I|)$ , а время поиска —  $O(m|I| + k)$ . Здесь  $k$  — количество совпадений, найденных в тексте со строками, принадлежащими множеству образцов.

В настоящее время задача поиска подстроки в строке активно исследуется по двум причинам: — поисковые системы активно развиваются [11]; — процесс обнаружения в антивирусных программных продуктах основывается на сигнатурах [1].

В этой связи были созданы алгоритмы, которые приходится выбирать в зависимости от определенных потребностей пользователя. Новейшие результаты, полученные при решении задачи поиска множества подстрок, описаны в работе [13].

Представленные в данной статье результаты основываются на связи между суффиксным массивом, созданным на основе множества строк образцов, и построением функций перехода и ошибок на начальных этапах работы алгоритма Ахо — Корасик. Такой подход отличен от традиционного и требует использования алгоритмов, позволяющих построить суффиксный массив за линейное время. Итак, в статье описаны алгоритмы, с помощью которых время выполнения предвыборной обработки сокращается до  $O(n)$ .

Пусть задан алфавит  $I$ , множество образцов  $P = \{P_1, P_2, \dots, P_k\}$ , где  $P_i \in I^*$ ,  $i = \overline{1, k}$ . Обозначим через  $n = \sum_{i=1}^k |P_i|$ . Будем считать, что алфавит  $I$  представляет собой ограниченный промежуток целых чисел. Граница может зависеть от длины рассматриваемой строки  $s \in I^*$  либо представлять собой промежуток  $[0, c]$ , где  $c$  — натуральное число:  $c \geq |s|$ . Пусть  $\epsilon \in I$  — пустая строка.

Пусть *goto* — функция перехода, а *failure* — функция ошибок. Приводимые модификации касаются методов построения упомянутых функций, используемых в алгоритме Ахо — Корасик [6, 10].

Положим,  $SuffArr(s)$  — некоторый алгоритм построения суффиксного массива для строки  $s \in I^*$  за линейное время. Описание таких алгоритмов можно найти, например, в [12–15].

Пусть  $x, y \in I^*$ . Тогда,  $lcp(x, y)$  — наибольший общий префикс строк  $x$  и  $y$ .

Рассмотрим строку  $s \in I^*$ ,  $s = s[s[0]s[1]...s[n-1]]$ . Пусть  $s[s[i]s[i+1]...s[j]]$  — подстрока  $s$ , включающая символы от  $i$  до  $j$ , где  $i \leq j$ ,  $i, j = \overline{0, n-1}$ . Обозначим через  $p_s$  суффиксный массив, соответствующий строке  $s$ . Пусть

$$p_s = p_s[p_s[0]p_s[1]...p_s[n-1]],$$

то есть  $s[s[p_s[0]]...s[p_s[n-1]]] < s[s[p_s[1]]...s[p_s[n-1]]] < \dots < s[s[p_s[n-1]]...s[p_s[n-1]]]$ .

Для построения суффиксного массива будет использован алгоритм, описанный в [15].

Положим,  $\alpha_i \notin I$ ,  $\alpha_i \neq \alpha_j$ ,  $1 \leq i < j \leq k+1$ ,  $\alpha_1 < \alpha_2 < \dots < \alpha_{k+1}$ . Пусть  $\forall b \in I \alpha_i < b$ , где  $1 \leq i \leq k+1$ .

Допустим, что  $P \neq \emptyset$ ,  $alpha = \{\alpha_1, \alpha_2, \dots, \alpha_k, \alpha_{k+1}\}$ .

**Алгоритм обработки суффиксного массива  $p_s$**

Здесь  $s \in I^*$ :  $s = \alpha_1 P_1 \alpha_2 P_2 \dots \alpha_k P_k \alpha_{k+1}$ ,  $P_i \in I^*$ ,  $1 \leq i \leq k$ .

*Adaptation*( $s, p_s, alpha$ )

1.  $new\_array \leftarrow \varepsilon$
2. *for* ( $i \leftarrow |alpha|; i < |s|; i++$ ) {
3.  $j \leftarrow 0$
4. *while* ( $s[s[p_s[i]]...s[|s|-1]][j] \notin alpha$ ) {
5.  $new\_array[i][j] \leftarrow s[s[p_s[i]]...s[|s|-1]][j]$
6.  $j \leftarrow j + 1$
7. }
8. }
9.  $ordered\_list[0] \leftarrow new\_array[0]$
10. *for* ( $i \leftarrow 1; i < |s| - |alpha|; i++$ ) {
11.  $j \leftarrow 0$
12. *if* ( $new\_array[i] \neq new\_array[i-1]$ ) {
13.  $ordered\_list[i] \leftarrow new\_array[i]$
14.  $j \leftarrow j + 1$
15. }
16. }
17. *return*  $ordered\_list$

Лемма 1. Пусть  $P = \{P_1, P_2, \dots, P_k\}$ ,  $s = \alpha_1 P_1 \alpha_2 P_2 \dots \alpha_k P_k \alpha_{k+1}$ . Тогда алгоритм *Adaptation* строит массив лексикографически упорядоченных суффиксов образов, принадлежащих  $P$ , за время  $O(|s| - |alpha|)$ .

Доказательство. В цикле 2–8 выполняется построение массива  $new\_array$ ,  $i$ -й элемент которого представляет собой такой префикс соответствующего суффикса  $s$ , который включает в себя все символы данного суффикса, начиная с нулевой позиции до его первого элемента, принадлежащего множеству  $alpha$ . При этом при помощи суффиксного массива  $p_s$  перебираются все суффиксы  $s$  согласно их лексикографическому порядку. Таким образом, массив  $new\_array$  состоит из всех суффиксов образов, принадлежащих  $P$  согласно их лексикографическому упорядочиванию, причем возможно повторение некоторых суффиксов. Заметим, что из рассмотрения исключаются все строки, начинающиеся с символов, принадлежащих массиву  $alpha$ , то есть первые  $|alpha|$  суффиксы. Далее в цикле 10–16 при помощи массива  $new\_array$  строится массив  $ordered\_list$  путем исключения повторяющихся строк. Для этого в силу лексикографического порядка следования строк достаточно проверить, совпадает ли рассматриваемая строка с предыдущей.

Цикл 2–8 выполняется за время  $O(|s| - |alpha|)$ , поскольку из рассмотрения исключаются все строки, начинающиеся с символов, принадлежащих массиву  $alpha$ . В цикле 10–16 происходит  $|s| - |alpha|$  сравнений строк. Таким образом, получаем асимптотическую оценку времени работы алгоритма  $O(|s| - |alpha|)$ . Лемма доказана.

**Алгоритм разбиения согласно лексикографическому упорядочению**

Здесь  $s$  — массив лексикографически упорядоченных строк.

$DandC(s)$

1.  $sub[0] \leftarrow 0$
2.  $j \leftarrow 0$
3. *for* ( $i \leftarrow 0; i < |s| - 1; i++$ ) {
4. *if* ( $s[i] \neq s[i+1]$ ) {
5.  $sub[j] \leftarrow i + 1$
6.  $j \leftarrow j + 1$
7. }
8.  $sub[j] \leftarrow |s|$
9. *return*  $sub$

Лемма 2. Алгоритм  $DandC$  на основе массива лексикографически упорядоченных строк  $s$  строит массив  $sub$ , состоящий из целых положительных чисел, указывающих индексы, соответствующие первым строкам среди строк, у которых равны первые символы, за время  $O(|s|)$ .

Доказательство. Граница, соответствующая первому символу, начинается с 0, что соответствует присваиванию, выполняемому на шаге 1. В цикле 3–7 последовательно сравниваются первые символы  $i$ -й и  $(i + 1)$ -й строк, где  $i = \overline{0, |s| - 2}$ . Если символы не равны, то в массив  $sub$  записывается начало границы, соответствующей следующему символу. В противном случае выполнение цикла продолжается. Правая граница последнего символа соответствует количеству строк в массиве  $s$  (шаг 8).

Сравнение на шаге 4 происходит за время  $O(1)$ , как и запись на шаге 5 и инкрементирование на шаге 6. Таким образом, цикл 3–7 выполняется за время  $O(|s|)$ . Лемма доказана.

**Алгоритм построения первой связи**

Здесь  $tree$  — дерево,  $lex\_words \in I^*$ ,  $link\_num$  — номер некоторого символа строки  $lex\_words$ ,  $v$  — порядковый номер нового узла, который присоединяется к узлу с порядковым номером  $node\_number$ .

$BuildFirstLink(tree\&, lex\_words\&, v\&, link\_num, node\_number)$

1.  $new\ tree.\ node[v]$
2.  $tree.\ node[v].state \leftarrow lex\_words[lex\_words[0]..lex\_words[link]]$
3.  $new\ tree.\ node[node\_number].link \leftarrow tree.\ node[v]$
4.  $tree.\ node[node\_number].link.symbol \leftarrow lex\_words[link\_num]$
5.  $v \leftarrow v + 1$

Лемма 3. Алгоритм  $BuildFirstLink$  осуществляет построение нового узла с порядковым номером  $v$  и дуги, ведущей от узла  $node\_number$  к новому узлу  $v$ , в дереве  $tree$  за время  $O(1)$ .

**Алгоритм построения связи с подстрокой**

Здесь  $tree$  — дерево,  $lex\_words \in I^*$ ,  $v$  — порядковый номер нового узла, который присоединяется к узлу с порядковым номером  $start$ .

$BuildSubstringLink(tree\&, lex\_words\&, v\&, start)$

1. *for* ( $k \leftarrow start; k < |lex\_words|; k++$ ) {

2.  $new\ tree.node[v]$
3.  $tree.node[v].state \leftarrow lex\_words[lex\_words[0]..lex\_words[k]]$
4.  $new\ tree.node[v-1].link \leftarrow tree.node[v]$
5.  $tree.node[v-1].link.symbol \leftarrow lex\_words[k]$
6.  $v \leftarrow v+1$
7. }

Лемма 4. Алгоритм *BuildSubstringLink* осуществляет построение новых узлов в дереве *tree*, соответствующих всем префиксам строки *lex\_words*, начиная с префикса  $lex\_words[lex\_words[0]..lex\_words[start]]$  за время  $O(|lex\_words| - start)$ .

#### Алгоритм построения последней связи

Здесь *tree* — дерево,  $lex\_words \in I^*$ ,  $v$  — порядковый номер новой дуги,  $I$  — алфавит.

*BuildLastLink*(*tree*&, *lex\_words*,  $v, I$ )

1.  $new\ tree.node[0].link[v] \leftarrow tree.node[0]$
2.  $symbols \leftarrow \emptyset$
3.  $for\ (i \leftarrow 0; i < |lex\_words|; i++)\ \{$
4.  $symbols[i] \leftarrow lex\_words[i][0]$
5.  $j \leftarrow 0$
6.  $for\ (i \leftarrow 0; i < |I|; i++)\ \{$
7.  $if\ (I[i] \notin symbols)\ \{$
8.  $tree.node[0].link[v].symbol[j] \leftarrow I[i]$
9.  $j \leftarrow j+1$
10. }
11. }

Лемма 5. Алгоритм *BuildLastLink* строит при корневом узле петлю. Ее пометка соответствует множеству символов, по которым невозможно перейти в другие узлы дерева *tree* из корневого узла за время  $O(|lex\_words| + |I|)$ .

#### Алгоритм построения переходов

Здесь *lex\_words* — массив лексикографически упорядоченных строк.

*CreateLink*(*lex\_words*)

1.  $str \leftarrow \emptyset$
2.  $sub \leftarrow DandC(lex\_words)$
3.  $v \leftarrow 1$
4.  $tree \leftarrow \emptyset$
5.  $tree.node[0].state \leftarrow \varepsilon$
6.  $for\ (i \leftarrow 0; i < |sub|-1; i++)\ \{$
7.  $BuildFirstLink(tree, lex\_words[sub[i]], v, 0, 0)$
8.  $BuildSubstringLink(tree, lex\_words[sub[i]], v, 1)$
9.  $for\ (j \leftarrow sub[i]+1; j < sub[i+1]; j++)\ \{$
10.  $temp \leftarrow |lcp(lex\_words[j-1], lex\_words[j])| + 1$
11.  $z \leftarrow tree.getStateNumber(lcp(lex\_words[j-1], lex\_words[j]))$
12.  $BuildFirstLink(tree, lex\_words[j], v, temp, z)$

13.  $BuildSubstringLink(tree, lex\_words[j], v, temp)$
14. }
15. }
16.  $BuildLastLink(tree, lex\_words, v, lex\_words)$
17.  $return tree$

Лемма 6. Алгоритм  $CreateLink$  осуществляет построение префиксного дерева  $tree$  с петлей при корневом узле за время  $O\left(\sum_{i=0}^{|lex\_words|-1} |lex\_words[i]|\right)$ .

Доказательство. На шаге 2 выполняется алгоритм  $DandC$  (см. лемму 2), после чего на шаге 5 создается корневой узел с порядковым номером 0 дерева  $tree$ . Его состояние полагается равным пустой строке  $\varepsilon$ . Рассмотрим цикл 6–15 на  $i$ -м шаге работы.

На шаге 7 при помощи алгоритма  $BuildFirstLink$  создается узел, состояние которого соответствует первому символу строки  $lex\_words[sub[i]]$ . Учитывая построение массива  $sub$ , можно утверждать, что такой символ прежде не встречался среди первых символов предыдущих строк. Тогда на шаге 8 реализация алгоритма  $BuildSubstringLink$  последовательно создает узлы, состояние которых соответствует всем префиксам строки  $lex\_words[sub[i]]$ , исключая префикс, построенный в предшествующей строке.

В цикле 9–14 при помощи алгоритмов  $BuildFirstLink$  и  $BuildSubstringLink$  выполняем те же действия со строками, лежащими в целочисленном промежутке  $[sub[i]+1, sub[i+1]-1]$ . Поскольку у каждой такой строки есть общий ненулевой префикс с предыдущей строкой, то алгоритм немедленно переходит в состояние, соответствующее наибольшему общему префиксу, начиная с которого необходимо строить новые узлы. На шаге 16 при помощи алгоритма  $BuildLastLink$  создается петля при корневом узле.

Шаги 12 и 13 выполняются за время

$$\begin{aligned} &O(1) + O\left(|lex\_words[j]| - |lcp(lex\_words[j-1], lex\_words[j])| - 1\right) = \\ &= O\left(|lex\_words[j]| - |lcp(lex\_words[j-1], lex\_words[j])|\right). \end{aligned}$$

Таким образом, из лемм 2, 3 и 4 следует, что цикл 9–14 выполняется за время

$$O\left(\sum_{j=sub[i]+1}^{sub[i+1]-1} |lex\_words[j]| - |lcp(lex\_words[j-1], lex\_words[j])|\right).$$

Цикл 6–14 выполняется за время

$$\begin{aligned} &O\left(\sum_{i=0}^{|sub|-2} |lex\_words[sub[i]]|\right) + \\ &+ O\left(\sum_{i=0}^{|sub|-2} \sum_{j=sub[i]+1}^{sub[i+1]-1} |lex\_words[j]| - |lcp(lex\_words[j-1], lex\_words[j])|\right) = \\ &= O\left(\sum_{i=0}^{|sub|-1} |lex\_words[sub[i]]| - |lex\_words|\right). \end{aligned}$$

Из леммы 5 следует, что шаг 16 выполняется за время  $O\left(|lex\_words| + \sum_{i=0}^{|lex\_words|-1} |lex\_words[i]|\right)$ .

Таким образом, получаем асимптотическую оценку времени работы алгоритма

$$\begin{aligned} &O(|lex\_words|) + O\left(\sum_{i=0}^{|sub|-1} |lex\_words[sub[i]]| - |lex\_words|\right) + \\ &+ O\left(|lex\_words| + \sum_{i=0}^{|lex\_words|-1} |lex\_words[i]|\right) = O\left(\sum_{i=0}^{|lex\_words|-1} |lex\_words[i]|\right). \end{aligned}$$

Лемма доказана.

**Алгоритм построения функции goto**

Здесь  $P$  — множество строк образцов.

$ConstructGoto(P)$

1.  $s \leftarrow \alpha_1 P_1 \alpha_2 P_2 \dots \alpha_k P_k \alpha_{k+1}$

2.  $p_s \leftarrow \text{SuffArr}(s)$
3.  $\alpha \leftarrow \{\alpha_1, \alpha_2, \dots, \alpha_k, \alpha_{k+1}\}$
4.  $\text{ordered\_list} \leftarrow \text{Adaptation}(s, p_s, \alpha)$
5.  $j \leftarrow 0$
6.  $\text{lex\_words} \leftarrow \emptyset$
7.  $P\_length \leftarrow \{|P_1|, |P_2|, \dots, |P_k|\}$
8. for ( $i \leftarrow 0; i < |\text{ordered\_list}|; i++$ )
9. if ( $(|\text{ordered\_list}[i]| \in P\_length)$  and ( $\text{ordered\_list}[i] \in P$ ))
10.  $\text{lex\_words}[j++] \leftarrow \text{ordered\_list}[i]$
11.  $\text{goto} \leftarrow \text{CreateLink}(\text{lex\_words})$
12. return  $\text{goto}$

Напомним, что  $n = \sum_{i=1}^k |P_i|$  для множества образцов  $P = \{P_1, P_2, \dots, P_k\}$ .

**Теорема 1.** Алгоритм *ConstructGoto* строит функцию *goto* за время  $O(n)$ .

Доказательство. На шаге 2 строится суффиксный массив  $p_s$  для строки  $s$ . На шаге 4 при помощи алгоритма *Adaptation* в массив *ordered\_list* записываются все суффиксы строк, принадлежащих множеству образцов  $P$ . При этом не исключены повторения. В цикле 8–10 строится массив *lex\_words*, содержащий суффиксы, принадлежащие  $P$ , расположенные в лексикографическом порядке без повторений. На шаге 11 происходит построение префиксного дерева с петлей при корневом узле на основе строк, которые содержатся в массиве *lex\_words*. Структура данных, возвращаемая алгоритмом *CreateLink*, в точности задает функцию *goto*.

Шаг 2 выполняется за время  $O(n+k+1)$  [12]. Из леммы 1 следует, что шаг 4 выполняется за время  $O(n+k+1-k-1) = O(n)$ . В цикле 8–10 рассматриваются только строки, длина которых равна длине какого-либо образца.

Таким образом, потребуется не более  $O(n)$  проверок для нахождения образцов  $P$ . Из леммы 6 следует, что шаг 11 выполняется за время  $O\left(\sum_{i=0}^{|\text{lex\_words}|-1} |\text{lex\_words}[i]|\right) = O(n)$ . Так как  $k \leq n$ , получаем асимптотическую оценку времени работы алгоритма  $O(n) + O(n+k+1) = O(n)$ . Теорема доказана.

### Результаты исследования

Пример 1.

Пусть  $P = \{one, on, once, cell, lull, eye, near\}$ . Тогда

$$s = \alpha_1 one \alpha_2 on \alpha_3 once \alpha_4 cell \alpha_5 lull \alpha_6 eye \alpha_7 near \alpha_8. \quad (1)$$

В табл. 1 представлен результат работы алгоритма построения функции *goto* при поступлении на вход строки  $s$  (1).

Таблица 1

Структура префиксного дерева

Номер <i>node</i>	Состояние <i>node</i>	Состояния, в которые ведут ветви <i>link</i> от <i>node</i>	Символы на ветвях <i>link</i> от <i>node</i>
0	$\epsilon$	1. $c$ ; 2. $e$ ; 3. $l$ ; 4. $n$ ; 5. $o$	1. $c$ ; 2. $e$ ; 3. $l$ ; 4. $n$ ; 5. $o$
1	$c$	1. $ce$	1. $e$
2	$ce$	1. $cel$	1. $l$
3	$cel$	1. $cell$	1. $l$
4	$cell$	–	–
5	$e$	1. $ey$	1. $y$
6	$ey$	1. $eye$	1. $e$
7	$eye$	–	–
8	$l$	1. $lu$	1. $u$

Номер <i>node</i>	Состояние <i>node</i>	Состояния, в которые ведут ветви <i>link</i> от <i>node</i>	Символы на ветвях <i>link</i> от <i>node</i>
9	<i>lu</i>	1. <i>lul</i>	1. <i>l</i>
10	<i>lul</i>	1. <i>lull</i>	1. <i>l</i>
11	<i>lull</i>	–	–
12	<i>n</i>	1. <i>ne</i>	1. <i>e</i>
13	<i>ne</i>	1. <i>nea</i>	1. <i>a</i>
14	<i>nea</i>	1. <i>near</i>	1. <i>r</i>
15	<i>near</i>	–	–
16	<i>o</i>	1. <i>on</i>	1. <i>n</i>
17	<i>on</i>	1. <i>onc</i> ; 2. <i>one</i>	1. <i>c</i> ; 2. <i>e</i>
18	<i>onc</i>	1. <i>once</i>	1. <i>e</i>
19	<i>once</i>	–	–
20	<i>one</i>	–	–

Пусть  $\tilde{s} = \alpha_{k+1}\tilde{P}_k\alpha_k\dots\alpha_2\tilde{P}_1\alpha_1\tilde{P}_0\alpha_0$  — зеркальное отражение строки  $s$ .

**Алгоритм построения функции ошибок *failure***

Здесь  $P$  — множество строк образцов.

$FalseSuff(P)$

1.  $\tilde{s} \leftarrow \alpha_{k+1}\tilde{P}_k\alpha_k\dots\alpha_2\tilde{P}_1\alpha_1\tilde{P}_0\alpha_0$
2.  $p_{\tilde{s}} \leftarrow SuffArr(\tilde{s})$
3.  $alpha \leftarrow \{\alpha_1, \alpha_2, \dots, \alpha_k, \alpha_{k+1}\}$
4.  $ordered\_list \leftarrow Adaptation(\tilde{s}, p_{\tilde{s}}, alpha)$
5.  $link \leftarrow \emptyset$
6. **for** ( $i \leftarrow 0; i < |\tilde{s}|; i++$ )
7.  $inLink[i] \leftarrow \varepsilon$
8.  $sub \leftarrow DandC(ordered\_list)$
9.  $str \leftarrow \emptyset$
10. **for** ( $i \leftarrow 0; i < |sub| - 1; i++$ )
11. **for** ( $j \leftarrow sub[i]; j < sub[i+1] - 1; j++$ )
12.  $str[j] \leftarrow |lcp(ordered\_list[j], ordered\_list[j+1])|$
13. **for** ( $i \leftarrow 0; i < |sub| - 1; i++$ )
14. **for** ( $k \leftarrow sub[i+1] - 1; k > sub[i]; k--$ ) {
15. **for** ( $j \leftarrow sub[i]; j < k; j++$ )
16.  $min\_element \leftarrow \min(str[k-1], str[k-2], \dots, str[j])$
17. **if** ( $min\_element = |ordered\_list[j]|$ )
18.  $min\_temp[j - sub[i]] \leftarrow min\_element$
19. }
20.  $max\_element \leftarrow \max(min\_temp[0], min\_temp[1], \dots, min\_temp[w])$
21. найти  $max\_index$ :  $min\_temp[max\_index] = max\_element$
22.  $inLink[k] \leftarrow ordered\_list[max\_index + sub[i]]$
23. }
24. **for** ( $i \leftarrow 0; i < |inLink|; i++$ ) {

25.  $link[i][0] \leftarrow ordered\_list[i]^\sim$ ; //зеркальное отражение строки  
 26.  $link[i][1] \leftarrow inLink[i]^\sim$ ; //зеркальное отражение строки  
 27. }  
 28. return link  
 Замечание. В строке 20  $w \leq sub[i+1] - sub[i] - 1$ .

**Теорема 2.** Алгоритм *FalseSuff* строит функцию ошибок *failure* за время  $O(n)$ . Доказательство. На шаге 1 строим массив символов, содержащий зеркальные отображения строк, принадлежащих множеству образцов  $P$ , и некоторые уникальные символы. На шаге 2 строим суффиксный массив  $p_{\tilde{s}}$  для строки  $\tilde{s}$ . На шаге 4 при помощи алгоритма *Adaptation* в массив *ordered\_list* записываются все суффиксы строк, принадлежащих множеству образцов  $\tilde{P}$  (множество образцов, состоящее из зеркально отраженных строк  $P$ ), при этом не исключены повторения. На шаге 8 выполняется алгоритм *DandC* (см. лемму 2), после чего в цикле 10–12 находим длину наибольшего общего префикса между строками, у которых совпадает первый символ. Записываем результат в массив *str*. Заметим, что эта величина равна нулю у строк, для которых данное условие не выполняется. В цикле 13–23 строится специальное отображение между строками, у которых совпадает первый символ. Опишем это отображение. Зафиксируем некоторую строку

$$s \in ordered\_list[ordered\_list[sub[i]], \dots, ordered\_list[sub[i+1]-1]].$$

Рассмотрим множество, состоящее из строк, принадлежащих  $ordered\_list[ordered\_list[sub[i]], \dots, ordered\_list[sub[i+1]-1]]$ . Их длина равна длине наибольшего общего префикса с  $s$ , исключая само  $s$ . Из этого множества найдем строку  $s'$ , которая имеет наибольшую длину, и поставим ее в соответствие  $s$ . Очевидно, что построенное отображение является биекцией при условии  $s' \neq \epsilon$ . Полученный результат записывается в массив *inLink*. В цикле 24–27 при помощи массива *inLink* явно указываем построенное отображение, при этом зеркально отражая каждую из строк. Таким образом, ставим узлу  $\tilde{s}$  префиксного дерева, построенного на основе массива образцов  $P$ , в соответствие узел  $\tilde{s}'$ . Его состояние равно наибольшему собственному суффиксу  $\tilde{s}$ , который встречается среди множества состояний рассматриваемого префиксного дерева. Но согласно определению функции ошибок *failure* это и есть искомым результат.

Пусть  $n = \sum_{i=1}^k |P_i|$ . Шаг 2 выполняется за время  $O(n+k+1)$  [12]. Из леммы 1 следует, что шаг 4 выполняется за время  $O(n+k+1-k-1) = O(n)$ . Цикл 6–7 выполняется за время  $O(|\tilde{s}|) = O(n)$ . Из леммы 2 следует, что шаг 8 выполняется за время  $O(n)$ . Цикл 10–12 выполняется за время

$$O\left(\sum_{i=0}^{|sub|-2} \sum_{j=sub[i]+1}^{sub[i+1]-1} \gamma_j\right) = O(k), \quad \forall j \quad \gamma_j = 1. \quad \text{Цикл 14–23 выполняется за время}$$

$$O\left(\sum_{k=sub[i]+1}^{sub[i+1]-1} \sum_{j=sub[i]}^{k-1} \gamma_j\right) = O\left(\sum_{j=0}^{sub[i+1]-sub[i]-2} \gamma_j\right) = O(sub[i+1]-sub[i]-1), \quad \forall j \quad \gamma_j = 1.$$

Тогда цикл 13–23 выполняется за время

$$O\left(\sum_{i=0}^{|sub|-2} (sub[i+1]-sub[i]-1)\right) = O(sub[|sub|-1]) = O(k).$$

Поскольку  $|inLink| < n$ , то цикл 24–27 выполняется за время  $O(|inLink|) = O(n)$ . Итак, учитывая, что  $k \leq n$ , получаем асимптотическую оценку времени работы алгоритма  $O(n) + O(n+k+1) + O(k) = O(n)$ . Теорема доказана.

Пример 2.

Пусть, как и в примере 1,  $P = \{one, on, once, cell, lull, eye, near\}$ . Тогда

$$\tilde{s} = \alpha_8 ra \epsilon \alpha_7 \epsilon \epsilon \alpha_6 llul \alpha_5 lle \alpha_4 \epsilon \epsilon \alpha_3 n \alpha_2 \epsilon \alpha_1. \quad (2)$$

В табл. 2 представлен результат работы алгоритма построения функции *failure* при поступлении на вход строки  $\tilde{s}$  (2).

Таблица 2

Ложные связи между узлами

Массив <i>inLink</i>	Массив <i>link</i>	Массив <i>inLink</i>	Массив <i>link</i>
0. $\epsilon$	0. <i>nea</i> ; 1. $\epsilon$	10. <i>l</i>	10. 0. <i>cel</i> ; 1. <i>l</i>
1. $\epsilon$	1. 0. <i>c</i> ; 1. $\epsilon$	11. <i>l</i>	11. 0. <i>cell</i> ; 1. <i>l</i>
2. <i>c</i>	2. 0. <i>onc</i> ; 1. <i>c</i>	12. <i>l</i>	12. 0. <i>lull</i> ; 1. <i>l</i>
3. $\epsilon$	3. 0. <i>e</i> ; 1. $\epsilon$	13. <i>l</i>	13. 0. <i>lu</i> ; 1. <i>l</i>
4. <i>e</i>	4. 0. <i>ce</i> ; 1. <i>e</i>	14. $\epsilon$	14. 0. <i>n</i> ; 1. $\epsilon$
5. <i>ec</i>	5. 0. <i>once</i> ; 1. <i>ce</i>	15. <i>n</i>	15. 0. <i>on</i> ; 1. <i>n</i>
6. <i>e</i>	6. 0. <i>ne</i> ; 1. <i>e</i>	16. $\epsilon$	16. 0. <i>o</i> ; 1. $\epsilon$
7. <i>en</i>	7. 0. <i>one</i> ; 1. <i>ne</i>	17. $\epsilon$	17. 0. <i>near</i> ; 1. $\epsilon$
8. <i>e</i>	8. 0. <i>eye</i> ; 1. <i>e</i>	18. $\epsilon$	18. 0. <i>lu</i> ; 1. $\epsilon$
9. $\epsilon$	9. 0. <i>l</i> ; 1. $\epsilon$	19. $\epsilon$	19. 0. <i>ey</i> ; 1. $\epsilon$

Для всех узлов, для которых на рис. 1 не указаны ложные связи, полагаем, что ложная связь ведет к корневому узлу.

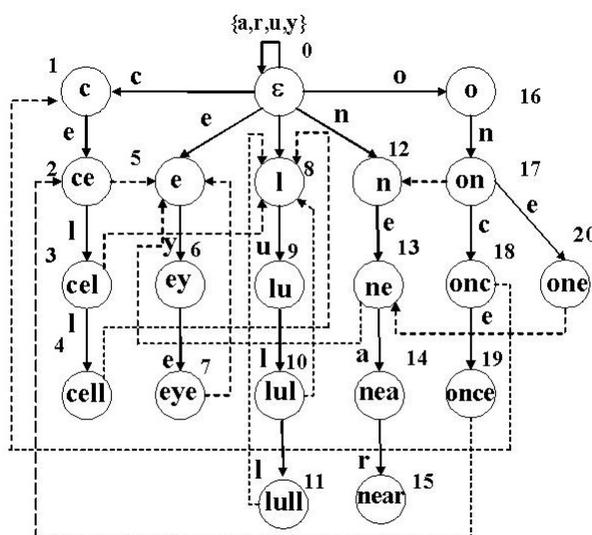


Рис. 1. Префиксное дерево с ложными связями

**Обсуждение и заключения.** Описана новая процедура препроцессинга в алгоритме Ахо — Корасик. Она выполняется за линейное время  $O(n)$ . Исследована связь между суффиксными массивами и префиксным деревом, что позволило предложить иной способ построения функций перехода и ошибок. Полученные результаты позволяют сократить время на предвыборную обработку множества строк образцов при использовании целочисленного алфавита.

**Библиографический список**

1. Stallings, W. Computer security: principles and practice / W. Stallings. — Boston : Pearson, 2012. — 182 p.
2. Исследование возможности применения генетических алгоритмов для реализации криптоанализа блочных криптосистем / Ю. О. Чернышев [и др.] // Вестник Дон. гос. техн. ун-та. — 2015. — Т. 15, № 3 (82). — С. 65–72.
3. Садовой, Н. Н. Программные утилиты для контроля и предотвращения сетевых атак на уровне доступа к сети / Н. Н. Садовой, Ю. В. Косолапов, В. В. Мкртчян // Вестник Дон. гос. техн. ун-та. — 2005. — Т. 5, № 2 (24). — С. 173–178.
4. Могилевская, Н. С. Пороговое разделение файлов на основе битовых масок: идея и возможное применение / Н. С. Могилевская, Р. В. Кульбикаян, Л. А. Журавлев // Вестник Дон. гос. техн. ун-та. — 2011 — Т. 11, № 10. — С. 1749–1755.

5. Шелудько, А. А. Поиск информационных объектов в памяти компьютера при решении задач обеспечения кибербезопасности / А. А. Шелудько, Н. В. Болдырихин // Молодой исследователь Дона. — 2018. — № 6 (15). — С. 81–86.
6. Мазуренко, А. В. Обнаружение, основанное на сигнатурах, с использованием алгоритма Ахо — Корасик / А. В. Мазуренко, Н. В. Болдырихин // Тр. Сев.-Кав. ф-ла Моск. техн. ун-та связи и информатики. — 2016. — № 1. — С. 339–344.
7. Мазуренко, А. В. Алгоритм проверки подлинности пользователя, основанный на графических ключах / А. В. Мазуренко, Н. С. Архангельская, Н. В. Болдырихин // Молодой исследователь Дона. — 2016. — № 3 (3). — С. 92–95.
8. Мазуренко, А. В. Геометрическая реализация метода проведения электронных выборов, основанного на пороговом разделении секрета / А. В. Мазуренко, В. А. Стукопин // Вестник Дон. гос. техн. ун-та. — 2018. — Т. 18, № 2. — С. 246–255.
9. Алгоритмическая оценка сложности системы кодирования и защиты информации, основанной на пороговом разделении секрета, на примере системы электронного голосования / Л. В. Черкесова [и др.] // Вестник Дон. гос. техн. ун-та. — 2017. — Т. 17, № 3. — С. 145–155.
10. Антонов, Е. С. Как найти миллион. Сравнение алгоритмов поиска множества подстрок / Е. С. Антонов // RSDN Magazine. — 2011. — № 1. — С. 60–67.
11. Tarakeswar, M. K. Search Engines: A Study / M. K. Tarakeswar, M. D. Kavitha // Journal of Computer Applications. — 2011. — Vol. 4, № 1. — P. 29–33.
12. Karkkainen, J. Linear work suffix array construction / J. Karkkainen, P. Sanders, S. Burkhardt // Journal of the ACM. — 2006. — Vol. 53, № 6. — P. 918–936.
13. Баклановский, М. В. Поведенческая идентификация программ / М. В. Баклановский, А. Р. Ханов // Моделирование и анализ информационных систем. — 2014. — Т. 21, № 6. — С. 120–130.
14. Efficient repeat finding in sets of strings via suffix arrays / V. Becher [et al.] // Discrete Mathematics and Theoretical Computer Science. — 2013. — Vol. 15, № 2. — P. 59–70.
15. Shrestha, A. M. S. A bioinformatician's guide to the forefront of suffix array construction algorithms / A. M. S. Shrestha, M. C. Frith, P. Horton // Briefings in Bioinformatics. — 2014. — Vol. 15, № 2. — P. 138–154.

Сдана в редакцию 22.01.2019

Принята к публикации 12.04.2019

#### **Об авторах:**

##### **Мазуренко Александр Вадимович,**

математик-программист ООО «ДЦОС-Гвард» (РФ, 344002, г. Ростов-на-Дону, пр. Буденновский, 62/2),  
ORCID: <http://orcid.org/0000-0001-9541-3374>  
[mazurencal@gmail.com](mailto:mazurencal@gmail.com)

##### **Болдырихин Николай Вячеславович,**

доцент кафедры «Кибербезопасность информационных систем» Донского государственного технического университета (РФ, 344000, г. Ростов-на-Дону, пл. Гагарина, 1), кандидат технических наук,  
ORCID: <http://orcid.org/0000-0002-9896-9543>,  
[boldyrikhin@mail.ru](mailto:boldyrikhin@mail.ru)