

On-Line Construction of Compact Directed Acyclic Word Graphs

Shunsuke Inenaga^{a,b} Hiromasa Hoshino^a Ayumi Shinohara^{a,b}
Masayuki Takeda^{a,b} Setsuo Arikawa^a Giancarlo Mauri^c
Giulio Pavesi^c

^a*Department of Informatics, Kyushu University, Japan*

^b*PRESTO, Japan Science and Technology Agency (JST)*

^c*Department of Computer Science, Systems and Communication, University of
Milan-Bicocca, Italy*

Abstract

Many different index structures, providing efficient solutions to problems related to pattern matching, have been introduced so far. Examples of these structures are *suffix trees* and *directed acyclic word graphs (DAWGs)*, which can be efficiently constructed in linear time and space. *Compact directed acyclic word graphs (CDAWGs)* are an index structure preserving some features of both suffix trees and DAWGs, and require less space than both of them. An algorithm which directly constructs CDAWGs in linear time and space was first introduced by Crochemore and V erin, based on McCreight's algorithm for constructing suffix trees. In this work, we present a novel *on-line* linear-time algorithm that builds the CDAWG for a single string as well as for a set of strings, inspired by Ukkonen's on-line algorithm for constructing suffix trees.

Key words: pattern matching on strings, compact directed acyclic word graphs, directed acyclic word graphs, suffix trees, on-line and linear-time algorithms

Email addresses: s-ine@i.kyushu-u.ac.jp (Shunsuke Inenaga),
hoshino@i.kyushu-u.ac.jp (Hiromasa Hoshino), ayumi@i.kyushu-u.ac.jp
(Ayumi Shinohara), takeda@i.kyushu-u.ac.jp (Masayuki Takeda),
arikawa@i.kyushu-u.ac.jp (Setsuo Arikawa), mauri@disco.uminib.it
(Giancarlo Mauri), pavesi@disco.uminib.it (Giulio Pavesi).

1 Introduction

Several different string problems, like those deriving from the analysis of biological sequences, can be efficiently solved by means of a suitable index structure [11,12,16,10]. The most widely known and studied structure of this kind seems to be *suffix trees* [36,30,6,34,27,33] perhaps because of their myriad applications [1]. For any string w the suffix tree of w requires only $O(n)$ space, and can be built in $O(n)$ time for a fixed alphabet, where n is the length of w . Although its theoretical space complexity is linear, much attention has been devoted to reduction of the practical space requirement of the structure. This has led to the introduction of more space-economical index structures like *suffix arrays* [14,29], *suffix cacti* [25], *compact suffix arrays* [28], *compressed suffix arrays* [15,31], and so on.

Blumer et al. [4] introduced *directed acyclic word graphs (DAWG)*. In [9] Crochemore pointed out that the DAWG of a string w is the smallest finite state automaton to recognize all suffixes of w . The smallest automaton accepting all factors of a string w , which is a variant of the DAWG for w , was also introduced in [4] and [9]. DAWGs are involved in several combinatorial algorithms on strings [5,17,3,35], since they serve as indexes of the string, as well as other index structures such as suffix tries and suffix trees. Some relationship between suffix trees and DAWGs can be found in [2].

In this work, we focus our attention on *compact directed acyclic word graphs (CDAWGs)* first introduced by Blumer et al. [5]. Crochemore and V erin displayed an overview relationship among suffix tries, suffix trees, DAWGs, and CDAWGs [13]. Suffix trees (resp. DAWGs) are the compacted (resp. minimized) version of suffix tries, as shown in Fig. 1. Similarly, CDAWGs can be obtained by either compacting DAWGs or minimizing suffix trees.

Not only in theory as stated above, but also in practice, do CDAWGs provide significant reductions of the memory space required by suffix trees and DAWGs, as experimental results have shown [5,13]. In Bioinformatics a considerable amount of DNA sequences has to be processed efficiently, both in space and in time. Therefore, from a practical viewpoint, CDAWGs could also play an important role in Bioinformatics.

The first algorithm to construct the CDAWG for a given string w was presented by Blumer et al. [4]. It once builds the DAWG of w , then removes every node of out-degree one and modifies its edges accordingly, so that the resulting structure becomes the CDAWG for w . It runs in linear time, but its main drawback is the construction of the DAWG as an intermediate structure, which takes larger space. A solution to this matter was provided by Crochemore and V erin [13]: a linear-time algorithm to construct the CDAWG for a string

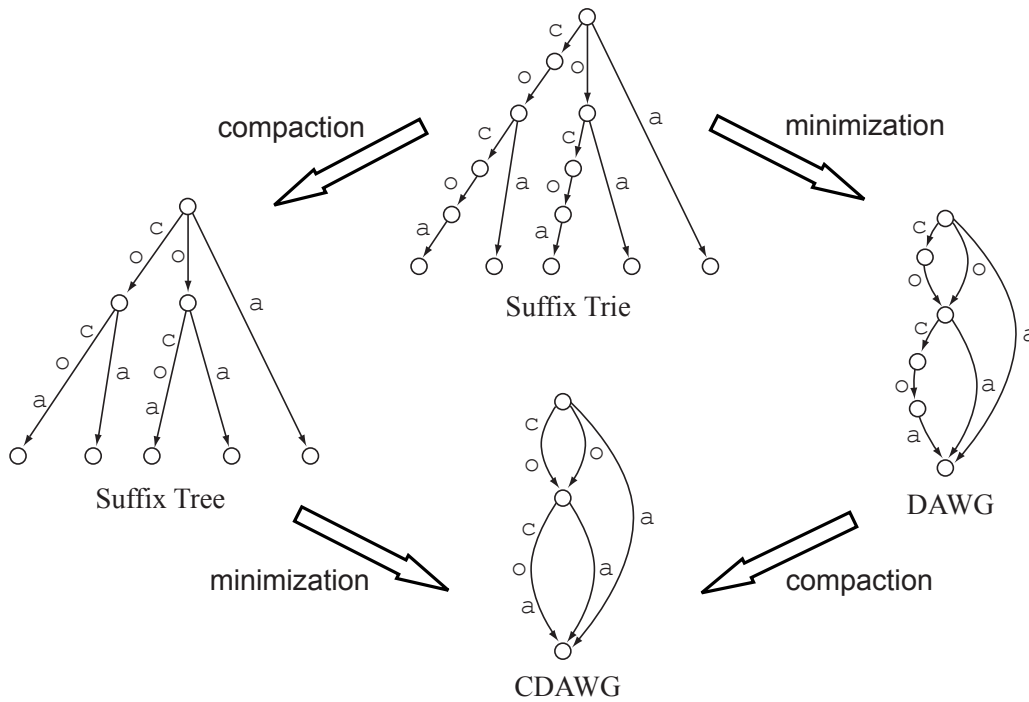


Fig. 1. Relationship among the suffix trie, the suffix tree, the DAWG, and the CDAWG for string *cocoa*.

directly. Their algorithm is based on McCreight’s suffix tree construction algorithm [30]. Both algorithms are *off-line*, that is, the whole input string has to be known beforehand. Thus, the structure (suffix tree or CDAWG) has to be rebuilt from scratch, when a new character is added to the input string. Table 1 summarizes some properties of typical algorithms to construct index structures. As seen there, a missing piece, which we have been looking for, is an *on-line algorithm for constructing CDAWGs*.

In this paper, we present a new linear-time algorithm to directly construct the CDAWG for a given string, which is based on Ukkonen’s suffix tree construction algorithm [34]. Our algorithm is *on-line*: it processes the characters of the input string from left to right, one by one, with no need to know the whole string beforehand. Our algorithm would be more efficient than the one in [13], in the sense that our algorithm allows us to update the input string. Furthermore, we show that the algorithm can be easily applied to building the CDAWG *for a set of strings*. Actually, the CDAWG for a set of strings can be constructed by the algorithm given in [5] which compacts the DAWG for the set. However, the drawback of this approach is that, when a new string is added to the set, the DAWG has to be built from scratch. Instead, our algorithm permits us the addition of a new string to the set. A previous version of this work appeared in [23].

The rest of the paper is organized as follows. Section 2 is devoted to introduction of basic notation and equivalence relations on strings. In Section 4 we give

<i>Index</i>	<i>Algorithm</i>	<i>linear time</i>	<i>on-line</i>	<i>multi strings</i>
suffix tries	Ukkonen [34]		✓	✓
suffix trees	Weiner [36]	✓		
	McCreight [30]	✓		
	Ukkonen [34]	✓	✓	✓
DAWGs	Blumer et al. [4]	✓	✓	
	Crochemore [9]	✓	✓	
	Blumer et al. [5]	✓	✓	✓
CDAWGs	Blumer et al. [5]	✓		✓
	Crochemore and V�erin [13]	✓		

Table 1

Properties of algorithms to construct index structures.

formal definitions of suffix tries, suffix trees, DAWGs, and CDAWGs, based on the equivalence relations on strings. The definitions give us a good ‘unified’ view for those index structures, revealing their relationship on the basis of the equivalence relations. We dedicate Section 5 to recalling the on-line algorithm to construct suffix tries given by Ukkonen in [34], which reminds us what an on-line algorithm is like. Ukkonen’s on-line suffix tree construction algorithm is also revisited in comparison to the suffix trie algorithm in Section 5. Section 6 introduces our new on-line algorithm that constructs the CDAWG for a given string. It is followed by Section 7 where it is shown that the CDAWG for a set of strings can also be built by the algorithm of Section 6. We finally conclude in Section 8.

2 Preliminaries

2.1 Notation

Let Σ be a finite alphabet. An element of Σ^* is called a *string*. Let x be a string such that $x = a_1a_2 \cdots a_n$ where $n \geq 1$ and $a_i \in \Sigma$ for $1 \leq i \leq n$. The *length* of x is n and denoted by $|x|$, that is, $|x| = n$. If $n = 0$, x is said to be the *empty string*. It is denoted by ε , that is, $|\varepsilon| = 0$. Let $\Sigma^+ = \Sigma^* - \{\varepsilon\}$.

Strings x , y , and z are said to be a *prefix*, *factor*, and *suffix* of string $w = xyz$, respectively. The sets of the prefixes, factors, and suffixes of a string w are denoted by $Prefix(w)$, $Factor(w)$, and $Suffix(w)$, respectively.

Let w be a string and $|w| = n$. The i -th character of w is denoted by $w[i]$ for $1 \leq i \leq n$, and the factor of w that begins at position i and ends at position j is denoted by $w[i : j]$ for $1 \leq i \leq j \leq n$. For convenience, let $w[i : j] = \varepsilon$ for

$j < i$.

For a set S of strings w_1, w_2, \dots, w_ℓ , let $|S|$ denote the cardinality of S , namely, $|S| = \ell$. We denote by $\|S\|$ the total length of strings in S .

The sets of prefixes, factors, and suffixes of the strings in S are denoted by $Prefix(S)$, $Factor(S)$, and $Suffix(S)$, respectively.

Definition 1 Let $S = \{w_1, \dots, w_k\}$ where $w_i \in \Sigma^*$ for $1 \leq i \leq k$ and $k \geq 1$. We say that S has the prefix property iff $w_i \notin Prefix(w_j)$ for any $1 \leq i \neq j \leq k$.

3 Equivalence Relations on Strings

Let $S \subseteq \Sigma^*$. For any string $x \in \Sigma^*$, let $Sx^{-1} = \{u \mid ux \in S\}$ and $x^{-1}S = \{u \mid xu \in S\}$.

Definition 2 Let $w \in \Sigma^*$. The equivalence relations \equiv_w^L and \equiv_w^R on Σ^* are defined by

$$\begin{aligned} x \equiv_w^L y &\Leftrightarrow Prefix(w)x^{-1} = Prefix(w)y^{-1}, \\ x \equiv_w^R y &\Leftrightarrow x^{-1}Suffix(w) = y^{-1}Suffix(w). \end{aligned}$$

The equivalence class of a string $x \in \Sigma^*$ with respect to \equiv_w^L (resp. \equiv_w^R) is denoted by $[x]_w^L$ (resp. $[x]_w^R$).

Note that all strings that are not in $Factor(w)$ form one equivalence class under \equiv_w^L . This equivalence class is called the *degenerate* class. All other classes are called *non-degenerate*. Similar arguments hold for \equiv_w^R .

Proposition 1 (Blumer et al. [5]) Let $w \in \Sigma^*$ and $x, y \in Factor(w)$. If $x \equiv_w^L y$, then either x is a prefix of y , or vice versa. If $x \equiv_w^R y$, then either x is a suffix of y , or vice versa.

Definition 3 For any string $x \in Factor(w)$, \overrightarrow{x}^w (resp. \overleftarrow{x}^w) denotes the longest member of $[x]_w^L$ (resp. $[x]_w^R$). We call \overrightarrow{x}^w (resp. \overleftarrow{x}^w) the representative of $[x]_w^L$ (resp. $[x]_w^R$).

What \overrightarrow{x}^w (resp. \overleftarrow{x}^w) means intuitively is that \overrightarrow{x}^w (resp. \overleftarrow{x}^w) is the string obtained by extending x in $[x]_w^L$ (resp. $[x]_w^R$) as long as possible. Each equivalence class in \equiv_w^L (\equiv_w^R) other than the degenerate class has a unique longest member. More formally:

Proposition 2 (Inenaga et al. [21], Inenaga [18]) Let $w \in \Sigma^*$. For any

string $x \in \text{Factor}(w)$, there uniquely exist strings $\alpha, \beta \in \Sigma^*$ such that $\overrightarrow{x} = x\alpha$ and $\overleftarrow{x} = \beta x$.

Example 1 Let $w = \text{coco}$. Then $\overrightarrow{\varepsilon} = \varepsilon$, $\overrightarrow{c} = \overrightarrow{c\delta} = \text{co}$, $\overrightarrow{co\delta} = \overrightarrow{c\delta} = \text{coco}$, $\overrightarrow{o} = \text{o}$, and $\overrightarrow{o\delta} = \overrightarrow{\delta o\delta} = \text{oco}$.

Example 2 Let $w = \text{coco}$. Then $\overleftarrow{\varepsilon} = \varepsilon$, $\overleftarrow{c} = \text{c}$, $\overleftarrow{o} = \overleftarrow{\delta o} = \text{co}$, $\overleftarrow{\delta c} = \overleftarrow{\delta o\delta} = \text{coc}$, and $\overleftarrow{\delta c\delta} = \overleftarrow{\delta o\delta} = \text{coco}$.

Proposition 3 (Inenaga [18]) Let $w \in \Sigma^*$ and $x \in \text{Factor}(w)$. Assume $\overrightarrow{x} = x$. Then, for any $y \in \text{Suffix}(x)$, $\overrightarrow{y} = y$.

We are now introducing a new equivalence class derived from $\overrightarrow{(\cdot)}$ and $\overleftarrow{(\cdot)}$.

Definition 4 For any string $x \in \text{Factor}(w)$, let \overleftarrow{x} be the string $\beta x \alpha$ ($\alpha, \beta \in \Sigma^*$) such that $\overrightarrow{x} = x\alpha$ and $\overleftarrow{x} = \beta x$.

What $\overleftarrow{x} = \beta x \alpha$ implies is that:

- (1) every time x occurs in w , it is preceded by β and followed by α .
- (2) α and β are the longest strings satisfying (1).

Definition 5 Let $x, y \in \Sigma^*$. We write $x \equiv_w y$ if,

- (1) $x, y \in \text{Factor}(w)$ and $\overleftarrow{x} = \overleftarrow{y}$, or
- (2) $x \notin \text{Factor}(w)$ and $y \notin \text{Factor}(w)$.

The equivalence class of a string $x \in \Sigma^*$ with respect to \equiv_w is denoted by $[x]_w$.

For any string $x \in \text{Factor}(w)$, \overleftarrow{x} is the unique longest member of $[x]_w$, and is called the *representative* of $[x]_w$.

Example 3 Let $w = \text{coco}$. Then $\overleftarrow{\varepsilon} = \varepsilon$, $\overleftarrow{c} = \overleftarrow{c\delta} = \overleftarrow{o} = \text{co}$, and $\overleftarrow{co\delta} = \overleftarrow{\delta o\delta} = \overleftarrow{\delta c} = \overleftarrow{\delta c\delta} = \text{coco}$.

Lemma 1 (Blumer et al. [5]) The equivalence relation \equiv_w is the transitive closure of the relation $\equiv_w^L \cup \equiv_w^R$.

It follows from the lemma above that:

Corollary 1 For any string $x \in \text{Factor}(w)$,

$$\overleftarrow{x} = \overrightarrow{\overleftarrow{x}} = \overleftarrow{\overrightarrow{x}}.$$

Note that, for a string $w \in \Sigma^*$, $|Factor(w)| = O(|w|^2)$. For example, consider string $a^n b^n$. However, considering set $S = \{x \mid x \in Factor(w) \text{ and } x = \overrightarrow{x}^w\}$, we have $|S| = O(|w|)$ for any $w \in \Sigma^*$. Similar arguments hold with \overleftarrow{x}^w and \overleftarrow{x}^w . The following lemma gives tighter upper-bounds.

Lemma 2 (Blumer et al. [4,5]) *Assume that $|w| > 1$. The number of non-degenerate equivalence classes in \equiv_w^L (or \equiv_w^R) is at most $2|w| - 1$. The number of non-degenerate equivalence classes in \equiv_w is at most $|w| + 1$.*

4 Index Structures for Text Strings

In this section, we recall four index structures, the suffix trie, the suffix tree, the directed acyclic word graph (DAWG), and the compact directed acyclic word graph (CDAWG) for a string $w \in \Sigma^*$ denoted by $STrie(w)$, $STree(w)$, $DAWG(w)$, and $CDAWG(w)$, respectively. All these structures represent every string $x \in Factor(w)$. We define them as edge-labeled graphs (V, E) with $E \subseteq V \times \Sigma^+ \times V$ where the second component of each edge represents its label. We also define the *suffix links* of each index structure. Suffix links are kinds of failure function often utilized for time-efficient construction of the index structures [36,30,34,4,5,13].

4.1 Suffix Tries

Definition 6 $STrie(w)$ is the tree (V, E) such that

$$\begin{aligned} V &= \{x \mid x \in Factor(w)\}, \\ E &= \{(x, a, xa) \mid x, xa \in Factor(w) \text{ and } a \in \Sigma\}, \end{aligned}$$

and its suffix links are elements of the set

$$F = \{(ax, x) \mid x, ax \in Factor(w) \text{ and } a \in \Sigma\}.$$

Each string $x \in Factor(w)$ has a one-to-one correspondence to a certain node in $STrie(w)$. $STrie(\text{coco})$ and $STrie(\text{cocoa})$ are displayed in Fig. 2 together with their suffix links.

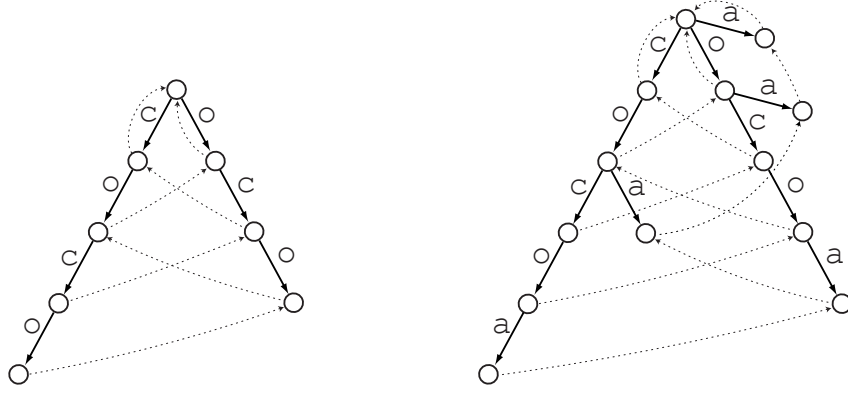


Fig. 2. $STrie(coco)$ on the left, and $STrie(cocoa)$ on the right. The solid arrows represent the edges, while the dotted arrows denote the suffix links.

4.2 Suffix Trees

Definition 7 $STree(w)$ is the tree (V, E) such that

$$V = \{ \overrightarrow{x} \mid x \in Factor(w) \},$$

$$E = \left\{ \left(\overrightarrow{x}, a\beta, \overrightarrow{x\hat{a}} \right) \mid \begin{array}{l} x, xa \in Factor(w), a \in \Sigma, \beta \in \Sigma^*, \\ \overrightarrow{x\hat{a}} = xa\beta, \text{ and } \overrightarrow{x} \neq \overrightarrow{x\hat{a}} \end{array} \right\},$$

and its suffix links are elements of the set

$$F = \{ (\overrightarrow{ax}, \overrightarrow{x}) \mid x, ax \in Factor(w), a \in \Sigma, \text{ and } \overrightarrow{ax} = a \cdot \overrightarrow{x} \}.$$

The node set of $STree(w)$ is a subset of that of $STrie(w)$, as seen in the definitions. It means that a string in $Factor(w)$ might be represented on an edge in $STree(w)$. In this case, we say that the string is represented by an *implicit* node. Conversely, every string in the node set V of $STree(w)$ is said to be represented by an *explicit* node. For example, in $STree(coco)$ of Fig. 3, string c is represented by an implicit node, while string co is on an explicit node.

$STree(w)$ can be seen as the compacted version of $STrie(w)$ with “ $\overrightarrow{(\cdot)}$ operation”. See $STrie(cocoa)$ in Fig. 2 and $STree(cocoa)$ in Fig. 3.

Theorem 1 (McCreight [30]) *Let $w \in \Sigma^*$ with $|w| > 1$. Let $STree(w) = (V, E)$. Then $|V| \leq 2|w| - 1$ and $|E| \leq 2|w| - 2$.*

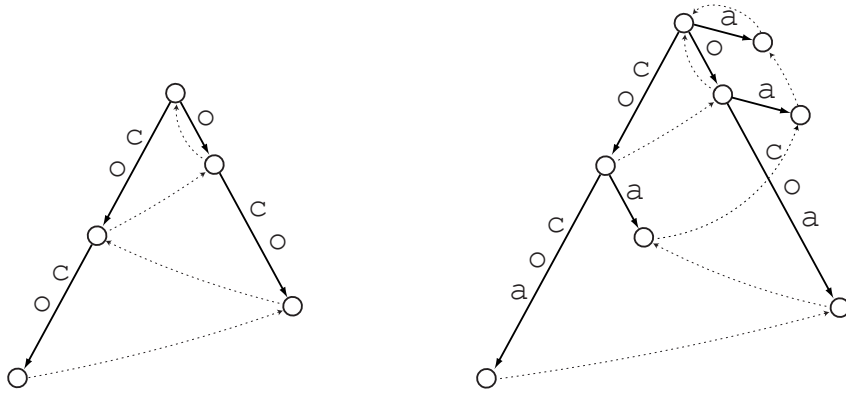


Fig. 3. $STree(coco)$ on the left, and $STree(cocoa)$ on the right. The solid arrows represent the edges, while the dotted arrows denote the suffix links.

4.3 DAWGs

Definition 8 $DAWG(w)$ is the directed acyclic graph (V, E) such that

$$V = \{[x]_w^R \mid x \in Factor(w)\},$$

$$E = \{([x]_w^R, a, [xa]_w^R) \mid x, xa \in Factor(w) \text{ and } a \in \Sigma\},$$

and its suffix links are elements of the set

$$F = \{([ax]_w^R, [x]_w^R) \mid x, ax \in Factor(w), a \in \Sigma, \text{ and } [ax]_w^R \neq [x]_w^R\}.$$

The node $[\varepsilon]_w^R$ is called the *source* node of $DAWG(w)$. A node of out-degree zero is called a *sink* node of $DAWG(w)$. $DAWG(w)$ has exactly one sink node for any $w \in \Sigma^*$. We define the *length* of a node $[x]_w^R$ by $|\overleftarrow{x}|$.

As seen in the definition, each node of $DAWG(w)$ is a non-degenerate equivalence class with respect to \equiv_w^R . In Fig. 4, one can see that nodes of $STrie(cocoa)$ are ‘merged’ in $DAWG(cocoa)$ by the equivalence class under \equiv_w^R . In this sense, $DAWG(w)$ can be seen as the minimized version of $STrie(w)$ with “[(\cdot)] $_w^R$ operation”.

Theorem 2 (Blumer et al. [4]) *Let $w \in \Sigma^*$ with $|w| > 1$. Let $DAWG(w) = (V, E)$. Then $|V| \leq 2|w| - 1$ and $|E| \leq 3|w| - 3$.*

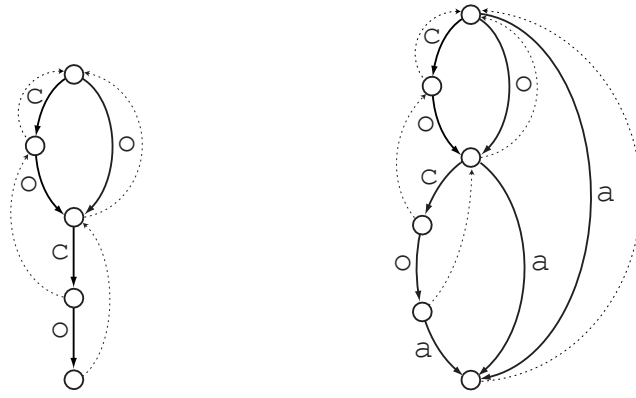


Fig. 4. $DAWG(\text{coco})$ on the left, and $DAWG(\text{cocoa})$ on the right. The solid arrows represent the edges, while the dotted arrows denote the suffix links.

4.4 CDAWGs

Definition 9 $CDAWG(w)$ is the directed acyclic graph (V, E) such that

$$V = \{[\overrightarrow{x}]_w^R \mid x \in \text{Factor}(w)\},$$

$$E = \left\{ ([\overrightarrow{x}]_w^R, a\beta, [\overrightarrow{x}\hat{a}]_w^R) \mid \begin{array}{l} x, xa \in \text{Factor}(w), a \in \Sigma, \beta \in \Sigma^*, \\ \overrightarrow{x}\hat{a} = xa\beta, \text{ and } \overrightarrow{x} \neq \overrightarrow{x}\hat{a} \end{array} \right\},$$

and its suffix links are elements of the set

$$F = \left\{ ([\overrightarrow{ax}]_w^R, [\overrightarrow{x}]_w^R) \mid \begin{array}{l} x, ax \in \text{Factor}(w), a \in \Sigma, \\ \overrightarrow{ax} = a \cdot \overrightarrow{x}, \text{ and } [\overrightarrow{x}]_w^R \neq [\overrightarrow{ax}]_w^R \end{array} \right\}.$$

The node $[\overrightarrow{\varepsilon}]_w^R$ is called the *source* node of $CDAWG(w)$. A node of out-degree zero is called a *sink* node of $CDAWG(w)$. $CDAWG(w)$ has exactly one sink

node for any $w \in \Sigma^*$. We define the *length* of a node $[\overrightarrow{x}]_w^R$ by $|\overleftarrow{(\overrightarrow{x})}| = |\overleftarrow{x}|$.

It follows from the definition that $CDAWG(w)$ is the minimization of $STree(w)$ with “ $[(\cdot)]_w^R$ operation”. In fact, $CDAWG(\text{cocoa})$ in Fig. 5 can be obtained by ‘merging’ the isomorphic subtrees in $STree(\text{cocoa})$. $CDAWG(w)$ can also be seen as the compaction of $DAWG(w)$ with “ $\overrightarrow{(\cdot)}$ operation”, as seen in $DAWG(\text{cocoa})$ in Fig. 4 and $CDAWG(\text{cocoa})$.

Theorem 3 (Blumer et al. [5], Crochemore and V erin [13])

Let $w \in \Sigma^*$ with $|w| > 1$. Let $CDAWG(w) = (V, E)$. Then $|V| \leq |w| + 1$ and $|E| \leq 2|w| - 2$.

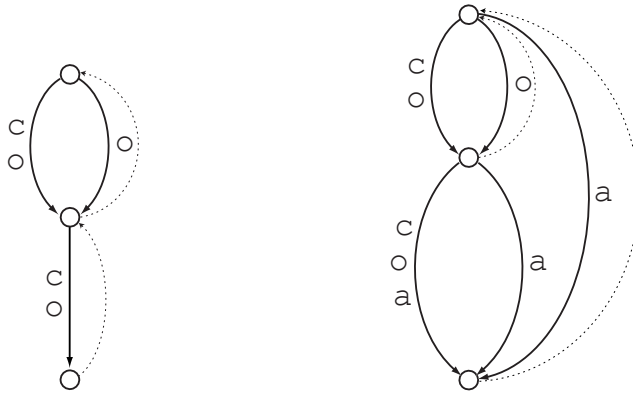


Fig. 5. $CDAG(coco)$ on the left, and $CDAG(cocoa)$ on the right. The solid arrows represent the edges, while the dotted arrows denote the suffix links.

5 General Idea of On-Line Construction Algorithms

5.1 On-Line Construction of Suffix Tries

In this section we recall the on-line algorithm to construct suffix tries given by Ukkonen [34], which is so intuitive that one can understand the basic concept of on-line algorithms.

For a string $x \in Factor(w)$, let $suf(x)$ denote the node reachable via the suffix link of the node x . It derives from Definition 6 that $suf(x) = y$ for some $y \in Factor(w)$ such that $x = ay$ for some character $a \in \Sigma$. In case that $x = \varepsilon$, let $suf(\varepsilon) = \perp$ where \perp is an auxiliary node called the *bottom* node. We suppose that there exists edge (\perp, a, ε) for any character $a \in \Sigma$. Assuming that the bottom node \perp corresponds to the inverse a^{-1} for any character $a \in \Sigma$, edge (\perp, a, ε) is consistently defined as well as other edges since $a^{-1} \cdot a = \varepsilon$. The auxiliary node \perp allows us to formalize the algorithm avoiding the distinction between the empty suffix and other non-empty suffixes (in other words, between the root node and other nodes). We leave $suf(\perp)$ undefined.

Suppose we have $STrie(w)$ with some $w \in \Sigma^*$. We now consider updating $STrie(w)$ to $STrie(wa)$ with an arbitrary character $a \in \Sigma$. What is necessary here is to insert suffixes of wa into $STrie(w)$.

Definition 10 For an arbitrary string $w \in \Sigma^*$ and character $a \in \Sigma$, the longest repeated suffix (LRS) of wa is the longest element of the set $Factor(w) \cap Suffix(wa)$.

The LRS of string w is denoted by $LRS(w)$.

Example 4 The LRS of string $coco$ is co . The LRS of string $cocoa$ is ε .

It is guaranteed that the LRS always exists for any string $w \in \Sigma^*$ since the empty string ε belongs to the set $Factor(w) \cap Suffix(wa)$ for any character $a \in \Sigma$.

Let $w \in \Sigma^*$ with $|w| = n$, and $a \in \Sigma$. Let $u_1, u_2, \dots, u_n, u_{n+1}, u_{n+2}$ be the suffixes of string wa sorted in decreasing order of their lengths, that is, $u_1 = wa$ and $u_{n+2} = \varepsilon$. These suffixes are divided into the following two groups as follows, by $LRS(wa)$.

- (1) u_1, \dots, u_ℓ where $u_{\ell+1} = LRS(wa)$.
- (2) $u_{\ell+1}, \dots, u_{n+2}$.

It follows from the definition of $LRS(wa)$ that the suffixes in the group (2) are already represented in $STrie(wa)$. Therefore, there is no need to newly insert any suffixes of the group (2) into $STrie(w)$.

Let v_1, \dots, v_{n+1} be the suffixes of w sorted in decreasing order of their lengths, that is, $v_1 = w$ and $v_{n+1} = \varepsilon$. Notice that $v_k a = u_k$ for any k with $1 \leq k \leq n+1$. The on-line algorithm starts with inserting into $STrie(w)$ the longest suffix u_1 of wa , which is wa itself. To do so, the algorithm creates a new edge labeled with a from the node $v_1 = w$ to a new node, which turns out to represent $wa = u_1$. It then focuses on the node that represents v_2 . Note that it is feasible in constant time to move from node v_1 to node v_2 by using the suffix link $suf(v_1)$. We then create a new edge (v_2, a, u_2) , and the suffix link of node u_1 is connected to u_2 , namely, $suf(u_1) = u_2$. We continue this operation until encountering $LRS(wa) = u_{\ell+1}$, and the resulting structure is $STrie(wa)$. The node with respect to $LRS(wa) = u_{\ell+1}$ is called the *end point* of $STrie(wa)$. This way suffix tries can be constructed in on-line fashion. However, since the space requirement of $STrie(w)$ is quadratic in $|w|$, this on-line algorithm does not run in linear time.

Theorem 4 (Ukkonen [34]) *Assume Σ is a fixed alphabet. For any string $w \in \Sigma^*$, $STrie(w)$ can be constructed on-line and in $O(|w|^2)$ time, using $O(|w|^2)$ space.*

On-line construction of $STrie(\text{cocoa})$ is illustrated in Fig. 6.

5.2 On-Line Construction of Suffix Trees

This section is devoted to recalling Ukkonen's on-line algorithm to construct suffix trees [34].

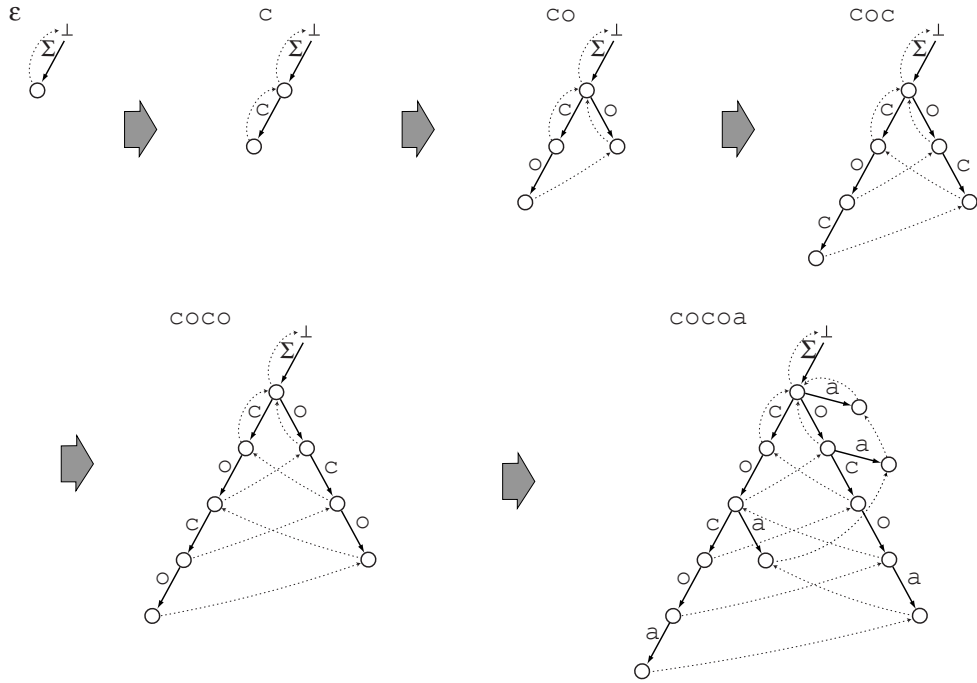


Fig. 6. On-line construction of $STrie(w)$ with $w = cocoa$.

5.2.1 Suffix Trees Redefined

For any string w , let $STree'(w)$ denote the tree obtained by removing all internal nodes of out-degree one from $STree(w)$. $STree'(coco)$ and $STree'(cocoa)$ are shown in Fig. 7. Actually, Ukkonen's suffix tree construction algorithm [34] builds $STree'(w)$, not $STree(w)$, since $STree'(w)$ is suitable in on-line string-processing scheme, as to be shown in the sequel.

To give a formal definition of $STree'(w)$, we introduce a relation X_w over Σ^* such that

$$X_w = \{(x, xa) \mid x \in Factor(w) \text{ and } a \in \Sigma \text{ is unique such that } xa \in Factor(w)\},$$

and let \equiv_w^L be the equivalence closure of X_w , i.e., the smallest superset of X_w that is symmetric, reflexive, and transitive.

Proposition 4 (Inenaga [18]) *For any string $w \in \Sigma^*$, \equiv_w^L is a refinement of \equiv_w^L , namely, every equivalence class under \equiv_w^L is a union of one or more equivalence classes in \equiv_w^L .*

For a string $x \in Factor(w)$, let \overrightarrow{x}^w denote the longest string in the equivalence class to which x belongs under the equivalence relation \equiv_w^L .

Example 5 *Let $w = coco$. Then $\overrightarrow{\varepsilon}^w = \varepsilon$, $\overrightarrow{c}^w = \overrightarrow{co}^w = \overrightarrow{coc}^w = \overrightarrow{coco}^w = coco$, and $\overrightarrow{o}^w = \overrightarrow{oc}^w = \overrightarrow{oco}^w = oco$.*

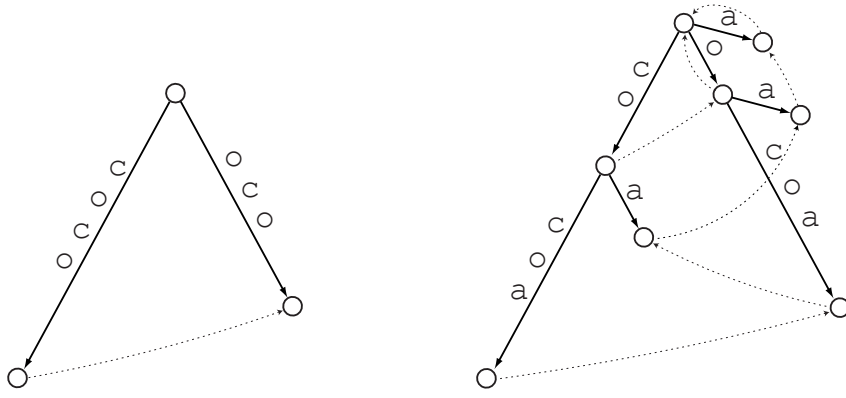


Fig. 7. $STree'(coco)$ on the left, and $STree'(cocoa)$ on the right. The solid arrows represent the edges, while the dotted arrows denote the suffix links.

Compare the above example with Example 1 for $\overrightarrow{(\cdot)}$.

Proposition 5 (Inenaga [18]) *For any string $x \in \Sigma^*$, \overrightarrow{x} is a prefix of \overrightarrow{x} . If $\overrightarrow{x} \neq \overrightarrow{x}$, then $\overrightarrow{x} \in Suffix(w)$.*

Proposition 6 (Inenaga [18]) *If set $Suffix(w) - \{\varepsilon\}$ satisfies the prefix property, $\overrightarrow{x} = \overrightarrow{x}$ for any string $x \in Factor(w)$.*

We are now ready to define $STree'(w)$.

Definition 11 $STree'(w)$ is the tree (V, E) such that

$$V = \{ \overrightarrow{x} \mid x \in Factor(w) \},$$

$$E = \left\{ \left(\overrightarrow{x}, a\beta, \overrightarrow{x\beta} \right) \left| \begin{array}{l} x, x\beta \in Factor(w), a \in \Sigma, \beta \in \Sigma^*, \\ \overrightarrow{x\beta} = xa\beta, \text{ and } \overrightarrow{x} \neq \overrightarrow{x\beta} \end{array} \right. \right\},$$

and its suffix links are elements of the set

$$F = \{ (\overrightarrow{ax}, \overrightarrow{x}) \mid x, ax \in Factor(w), a \in \Sigma, \text{ and } \overrightarrow{ax} = a \cdot \overrightarrow{x} \}.$$

This definition is the same as the one obtained by replacing “ $\overrightarrow{(\cdot)}$ operation” with “ $\overrightarrow{(\cdot)}$ operation” in Definition 7. Therefore, it follows from Proposition 6 that:

Corollary 2 *If $Suffix(w) - \{\varepsilon\}$ has the prefix property, $STree'(w) = STree(w)$.*

Even in case that the set $Suffix(w) - \{\varepsilon\}$ does not have the prefix property,

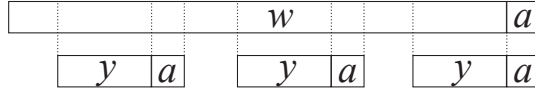


Fig. 8. Illustration for the proof of Lemma 7.

$STree'(w\$)$ is identical to $STree(w\$)$ where $\$$ is a special character that occurs nowhere in w .

As previously stated, Ukkonen's algorithm constructs $STree'(w)$ since $STree'(w)$ is suitable in on-line manner. The reasons are as follows.

Proposition 7 *Let $w \in \Sigma^*$ and $a \in \Sigma$. Let $z = LRS(w)$. Let $y \in Suffix(z)$. Assume $\overrightarrow{y}^w = y$, and there uniquely exists a character $a \in \Sigma$ such that $ya \in Factor(w)$. Then, $\overrightarrow{y}^{wa} = ya$.*

Proof. Since $y \in Suffix(z)$, $ya \in Suffix(za)$. Because a is the only character such that $ya \in Factor(w)$, it is also the only character such that $ya \in Factor(wa)$. Consequently we have $Prefix(wa)(y)^{-1} = Prefix(wa)(ya)^{-1}$, which implies that $y \equiv_{wa}^L ya$. It is clear that ya is the longest element of $[ya]_{wa}^L = [y]_{wa}^L$, meaning that $\overrightarrow{ya}^{wa} = \overrightarrow{y}^w = ya$. (See Fig. 8). \square

This proposition implies that an explicit node of $STree(w)$ may become implicit in $STree(wa)$. The maintenance of converting all of such explicit nodes into implicit nodes would take quadratic time in aggregate. This causes difficulty in updating $STree(w)$ into $STree(wa)$ in amortized constant time. However, we have the following proposition for $STree'(w)$.

Proposition 8 *Let $w \in \Sigma^*$ and $a \in \Sigma$. Let $z = LRS(w)$. Let $y \in Suffix(z)$. Assume $\overrightarrow{y}^w = y$. Then $\overrightarrow{y}^{wa} = y$ for any $a \in \Sigma$.*

Proof. From the assumption that $\overrightarrow{y}^w = y$, there exist at least two distinct characters b, c such that $yb, yc \in Factor(w)$. It is obvious that $yb, yc \in Factor(wa)$ for any $a \in \Sigma$, and thus we have $\overrightarrow{y}^{wa} = y$. \square

According to this proposition, in converting $STree'(w)$ into $STree'(wa)$ we do not need the maintenance that is required in converting $STree(w)$ into $STree(wa)$ due to Proposition 7. This is why $STree'(w)$ is suitable for on-line algorithm if we want it to run in linear time.

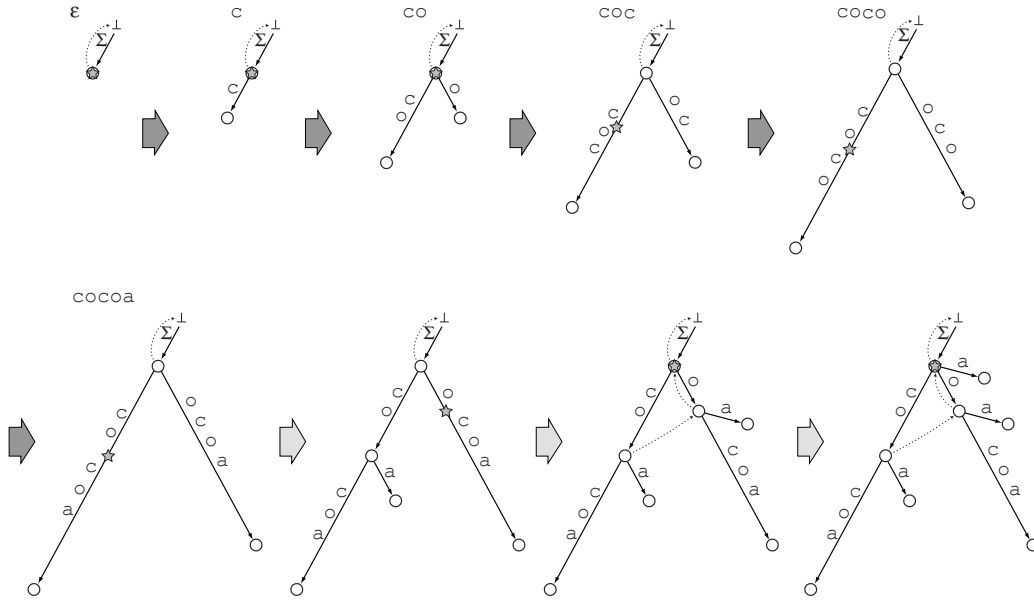


Fig. 9. On-line construction of $STree'(w)$ with $w = \text{cocoa}$. The star represents the active point for each step.

5.2.2 Ukkonen's Algorithm

Informal Description of the Algorithm.

Unlike the case of suffix tries mentioned in Section 5.1, Ukkonen's on-line suffix tree construction algorithm runs in linear time.

Theorem 5 (Ukkonen [34]) *Assume Σ is a fixed alphabet. For any string $w \in \Sigma^*$, $STree'(w)$ can be constructed on-line and in $O(|w|)$ time, using $O(|w|)$ space.*

We here summarize Ukkonen's suffix tree construction algorithm in the comparison with the previous suffix trie algorithm. Fig. 9 illustrates on-line construction of $STree'(\text{cocoa})$.

Focus on the update of $STree'(\text{co})$ to $STree'(\text{coc})$. Differently from that of $STrie(\text{co})$ to $STrie(\text{coc})$, the edges leading to the leaf nodes are *automatically* extended with the new character c in $STree'(\text{coc})$. This is feasible by the idea of so-called *open edges*.

See the first and second steps of the update of $STree'(\text{coco})$ to $STree'(\text{cocoa})$. The gray star mark indicates the *active point* from which a new edge is created in each step. After the new edge $(\text{co}, a, \text{coa})$ is inserted, the active point moves to the implicit node for string o . In case of the suffix trie, it is possible to move there by traversing the suffix link of node co . However, there is yet to be the suffix link of node co in the suffix tree. Thereof, Ukkonen's algorithm *simulates* the traversal of the suffix link as follows: First, it goes up to the explicit parent

node ε of node co which already has its suffix link. After that, it moves to the bottom node \perp via the suffix link of the root node, and then advances along the path spelling out co . (Note that the string co corresponds to the label of the edge the active point went up backward.) This way, in Ukkonen's algorithm the active point moves via 'implicit' suffix links. Since the suffix links of leaf nodes are never utilized in Ukkonen's algorithm, none of them are created.

Formal Description of the Algorithm.

Ukkonen's on-line suffix tree construction algorithm is based on the on-line algorithm to build suffix tries recalled in Section 5.1. As stated in Definition 11, an edge of $STree'(w)$ is labeled by a string $\alpha \in Factor(w)$. The key idea to achieve a linear-space implementation of the suffix tree is to label each edge $(\overrightarrow{x}, \alpha, \overrightarrow{x\alpha})$ in $STree'(w)$ by (k, p) , such that $w[k : p] = \alpha$.

An implicit node $y \in Factor(w)$ can be represented by an ordered pair $(\overrightarrow{x}, \alpha)$ of an explicit node \overrightarrow{x} and a string $\alpha \in Factor(w)$ such that $y = \overrightarrow{x} \cdot \alpha$. The ordered pair $(\overrightarrow{x}, \alpha)$ is called a *reference pair* for the implicit node y . Note that explicit nodes can also be represented by reference pairs. There can be one or more reference pairs for a node y . The reference pair $(\overrightarrow{x}, \alpha)$ for y in which $|\alpha|$ is minimized is called the *canonical* reference pair for y . The reference pair can also be represented using the integers h, p , as $(\overrightarrow{x}, (h, p))$ such that $w[h : p] = \alpha$.

Suppose we have $STree'(w)$ with some $w \in \Sigma^*$. We now consider updating $STree'(w)$ to $STree'(wa)$ with an arbitrary character $a \in \Sigma$. What is necessary here is to insert suffixes of wa into $STree'(w)$. The group (1) of the suffixes of wa , mentioned in the previous section, can moreover be divided into two sub-groups as follows, by $LRS(w)$.

- (1-a) u_1, \dots, u_k where $u_{k+1} = LRS(w) \cdot a$.
- (1-b) u_{k+1}, \dots, u_ℓ .

We remark that all the suffixes of the group (1-a) are those represented by the leaf nodes in $STree'(w)$. Note that, for any i with $1 \leq i \leq k$, we have $\overrightarrow{v_i} = v_i$ and $\overrightarrow{u_i} = u_i$. Moreover, $v_i a = u_i$ (see Section 5.1). This means that, intuitively, every leaf node of $STree'(w)$ is also a leaf node in $STree'(wa)$. This fact is crucial to Ukkonen's algorithm in order that it is able to *automatically* insert the suffixes in the group (1-a) into $STree'(wa)$, by means of *open edges* which will be recalled below. Suppose that $(\overrightarrow{x}, \alpha, \overrightarrow{x\alpha})$ is an edge of $STree'(w)$

where $\overrightarrow{x\alpha}$ is a leaf node. Let p be the integer such that $w[p : |w|] = \alpha$. Then we actually implement the edge label by (p, ∞) , where its implication is “a leaf node is always a leaf node.”. This way we need no explicit insertion of the suffixes of w in the group (1-a).

For a suffix u_j of group (1-b) where $k + 1 \leq j \leq \ell$, the location in $STree'(w)$ from which u_j is inserted is called the *active point* of u_j in $STree'(w)$. Note that the active point for u_{k+1} is the end point of $STree'(w)$. As well as the case of on-line construction of suffix tries, we insert the suffixes u_{k+1}, \dots, u_ℓ of wa into $STree'(w)$ in decreasing order of their lengths. To do so, we focus on the locations of the suffixes v_{k+1}, \dots, v_ℓ of w from which a new edge labeled by the new character a is created. Assume we are now inserting suffix u_j into $STree'(w)$ for some $k + 1 \leq j \leq \ell$. There are the two following cases regarding the active point.

(Case 1) The active point is on an explicit node v_j . In this case,

$$\overrightarrow{v_j} = \overrightarrow{v_j} = v_j.$$

Then a new edge $(\overrightarrow{v_j}, \alpha, \overrightarrow{v_j\alpha})$ is created, where α is a at the moment. Note $\overrightarrow{v_j\alpha} = \overrightarrow{v_j}a = v_ja = u_j$. In order for the implementation of an open edge, the edge is actually labeled by $(n + 1, \infty)$, where $|wa| = n + 1$ and $wa[n + 1] = a$.

After that, the active point moves to the explicit node $suf(\overrightarrow{v_j}) = \overrightarrow{v_{j+1}} = v_{j+1}$, in order to insert the next suffix u_{j+1} of wa .

(Case 2) The active point is on an implicit node v_j . In this case,

$$\overrightarrow{v_j} \neq v_j \text{ but } \overrightarrow{v_j} = v_j.$$

Let $(\overrightarrow{x}, \alpha)$ be the canonical reference pair for the active point, namely, $\overrightarrow{x} \cdot \alpha = v_j$. Focus on the edge $(\overrightarrow{x}, \alpha\beta, \overrightarrow{x\alpha\beta})$ with some non-empty string β . The implicit node representing v_j is located on this edge. The edge is then replaced by (separated into) the edges $(\overrightarrow{x}, \alpha, \overrightarrow{x\alpha})$ and $(\overrightarrow{x\alpha}, \beta, \overrightarrow{x\alpha\beta})$ where $\overrightarrow{x\alpha} = v_j$ is a new explicit node. Then a new edge $(\overrightarrow{x\alpha}, \gamma, \overrightarrow{x\alpha\gamma})$ is created, where $\gamma = a$ at the moment. Note $\overrightarrow{x\alpha\gamma} = \overrightarrow{v_j}a = v_ja = u_j$. In order for the implementation of an open edge, the edge is actually labeled by $(n + 1, \infty)$, where $|wa| = n + 1$ and $wa[n + 1] = a$.

After that, we need to move to the (implicit or explicit) node corresponding to v_{j+1} , the next active point, but the table suf is yet to be computed for the node $\overrightarrow{x\alpha}$ that has been created just now. Thus we once move to its parent node \overrightarrow{x} for which $suf(\overrightarrow{x})$ must have already been computed. Let

$suf(\overrightarrow{x}) = \overrightarrow{y}$. Then \overrightarrow{y} is explicit. Note that $\overrightarrow{y} \cdot \alpha = v_{j+1}$. We go down from the node \overrightarrow{y} while spelling out α , to obtain the canonical reference pair for the active point v_{j+1} . The node v_{j+1} either is already, or will in this step become, explicit. $suf(\overrightarrow{x\alpha}) = suf(\overrightarrow{v_j})$ is connected to v_{j+1} . This way the algorithm ‘simulates’ the suffix-link-traversal of suffix tries.

A pseudo-code for Ukkonen’s algorithm is shown in Fig. 10. Function *canonize* is a routine to canonize a given reference pair. Function *check_end_point* is one that returns true if a given reference pair is the end point, and false otherwise. Function *split_edge* splits an edge into two, by creating a new explicit node at the position to which the given reference pair corresponds.

6 On-Line Construction of CDAWGs

In this section we give an on-line algorithm to construct CDAWGs. In order to achieve a linear-time algorithm, we shall need CDAWGs that are defined on the basis of $\overrightarrow{(\cdot)}$ rather than $\overrightarrow{(\cdot)}$, like in case of suffix trees recalled in Section 5.2.

6.1 CDAWGs Redefined

Definition 12 *CDAWG'(w) is the directed acyclic graph (V, E) such that*

$$V = \{[\overrightarrow{x}]_w^R \mid x \in Factor(w)\},$$

$$E = \left\{ ([\overrightarrow{x}]_w^R, a\beta, [\overrightarrow{x\alpha}]_w^R) \mid \begin{array}{l} x, x\alpha \in Factor(w), a \in \Sigma, \beta \in \Sigma^*, \\ \overrightarrow{x\alpha} = xa\beta, \text{ and } \overrightarrow{x} \neq \overrightarrow{x\alpha} \end{array} \right\},$$

and its suffix links are elements of the set

$$F = \left\{ ([\overrightarrow{ax}]_w^R, [\overrightarrow{x}]_w^R) \mid \begin{array}{l} x, ax \in Factor(w), a \in \Sigma, \\ \overrightarrow{ax} = a \cdot \overrightarrow{x}, \text{ and } [\overrightarrow{x}]_w^R \neq [\overrightarrow{ax}]_w^R \end{array} \right\}.$$

In Fig. 7 and 11, one can see that nodes of *STree'(cocoa)* are ‘merged’ in *CDAWG'(cocoa)* by the equivalence class under \equiv_w^R . In this sense *CDAWG'(w)* can be seen as the minimized version of *STree'(w)* with “[(\cdot)] $_w^R$ operation”.

```

Algorithm for on-line construction of  $STree'(w\$)$ 
in alphabet  $\Sigma = \{w[-1], w[-2], \dots, w[-m]\}$ .
/* $ is the end-marker appearing nowhere in w. */
1 create nodes root and  $\perp$ ;
2 for  $j := 1$  to  $m$  do create edge  $(\perp, (-j, -j), root)$ ;
3  $suf(root) := \perp$ ;
4  $(s, k) := (root, 1)$ ;  $i := 0$ ;
5 repeat
6    $i := i + 1$ ;
7    $(s, k) := update(s, (k, i))$ ;
8 until  $w[i] = \$$ ;

function  $update(s, (k, p))$ : pair of integers;
/*  $(s, (k, p - 1))$  is the canonical reference pair for the active point. */
1  $c := w[p]$ ;  $olldr := \mathbf{nil}$ ;
2 while not  $check\_end\_point(s, (k, p - 1), c)$  do
3   if  $k \leq p - 1$  then  $r := split\_edge(s, (k, p - 1))$ ; /* implicit case. */
4   else  $r := s$ ; /* explicit case. */
5   create node  $r'$ ; create edge  $(r, (p, \infty), r')$ ;
6   if  $olldr \neq \mathbf{nil}$  then  $suf(olldr) := r$ ;
7    $olldr := r$ ;
8    $(s, k) := canonize(suf(s), (k, p - 1))$ ;
9 if  $olldr \neq \mathbf{nil}$  then  $suf(olldr) := s$ ;
10 return  $canonize(s, (k, p))$ ;

function  $check\_end\_point(s, (k, p), c)$ : boolean;
1 if  $k \leq p$  then /* implicit case. */
2   let  $(s, (k', p'), s')$  be the  $w[k]$ -edge from  $s$ ;
3   return  $(c = w[k' + p - k + 1])$ ;
4 else return (there is a  $c$ -edge from  $s$ );

function  $canonize(s, (k, p))$ : pair of node and integers;
1 if  $k > p$  then return  $(s, k)$ ; /* explicit case. */
2 find the  $w[k]$ -edge  $(s, (k', p'), s')$  from  $s$ ;
3 while  $p' - k' \leq p - k$  do
4    $k := k + p' - k' + 1$ ;  $s := s'$ ;
5   if  $k \leq p$  then find the  $w[k]$ -edge  $(s, (k', p'), s')$  from  $s$ ;
6 return  $(s, k)$ ;

function  $split\_edge(s, (k, p))$ : node;
1 let  $(s, (k', p'), s')$  be the  $w[k]$ -edge from  $s$ ;
2 create node  $r$ ;
3 replace the edge by  $(s, (k', k' + p - k), r)$  and  $(r, (k' + p - k + 1, p'), s')$ ;
4 return  $r$ ;

```

Fig. 10. Ukkonen's on-line algorithm for constructing suffix trees.

The node $\overrightarrow{[\varepsilon]}_w^R$ is called the *source* node of $CDAWG'(w)$. A node of out-degree zero is called a *sink* node of $CDAWG'(w)$. $CDAWG'(w)$ has exactly one sink node for any $w \in \Sigma^*$, which is $\overrightarrow{[w]}_w^R$. We define the *length* of a node $\overrightarrow{[x]}_w^R$ by $\left| \overleftarrow{[x]}_w \right|$, namely, by the length of the representative of $\overrightarrow{[x]}_w^R$. For examples, the

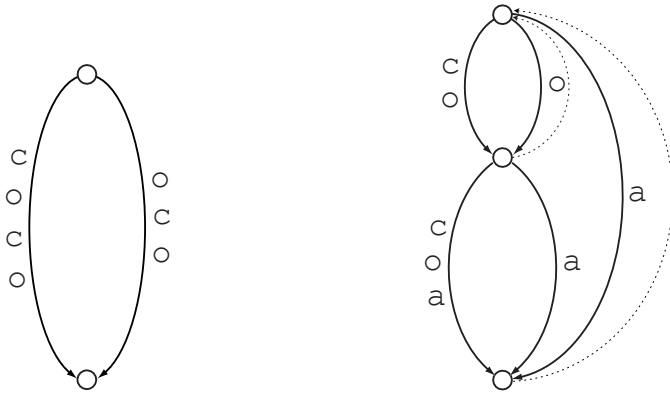


Fig. 11. $CDAWG'(coco)$ on the left, and $CDAWG'(cocoa)$ on the right. The solid arrows represent the edges, while the dotted arrows denote the suffix links.

length of the source node is 0 since $\left| \begin{smallmatrix} \leftarrow w \\ \overrightarrow{\varepsilon} \end{smallmatrix} \right| = \left| \leftarrow \varepsilon \right| = |\varepsilon| = 0$, and the length of the sink node is $|w|$ since $\left| \begin{smallmatrix} \leftarrow w \\ \overrightarrow{w} \end{smallmatrix} \right| = \left| \leftarrow w \right| = |w|$.

Definition 12 equals the one obtained by replacing “ $\begin{smallmatrix} \leftarrow \\ \overrightarrow{\cdot} \end{smallmatrix}$ operation” with “ $\begin{smallmatrix} \leftarrow \\ \overrightarrow{w} \end{smallmatrix}$ operation” in Definition 9. Therefore, it follows from Proposition 6 that:

Corollary 3 *If set $Suffix(w) - \{\varepsilon\}$ has the prefix property, $CDAWG'(w) = CDAWG(w)$.*

Even if $Suffix(w) - \{\varepsilon\}$ does not have the prefix property, $CDAWG'(w\$)$ is identical to $CDAWG(w\$)$ where $\$$ is a special character that occurs nowhere in w .

As well as the case with suffix trees, according to Propositions 8 and 7, $CDAWG'(w)$ is more suitable to treat in on-line manner. This is why we gave Definition 12 above, which our on-line algorithm actually constructs.

6.2 Our Algorithm

Informal Description of the Algorithm.

Before delving into the technical detail of the algorithm for on-line construction of CDAWGs, we informally describe how a CDAWG is built on-line. See Fig. 12 that shows on-line construction of $CDAWG'(cocoa)$, in comparison with Fig. 9 displaying on-line construction of $STree'(cocoa)$. Compare $CDAWG'(co)$ and $STree'(co)$. While strings co and o are separately represented in $STree'(co)$, they are in the same node in $CDAWG'(co)$. The destination of any open edge of a CDAWG is all the same, the sink node. Open

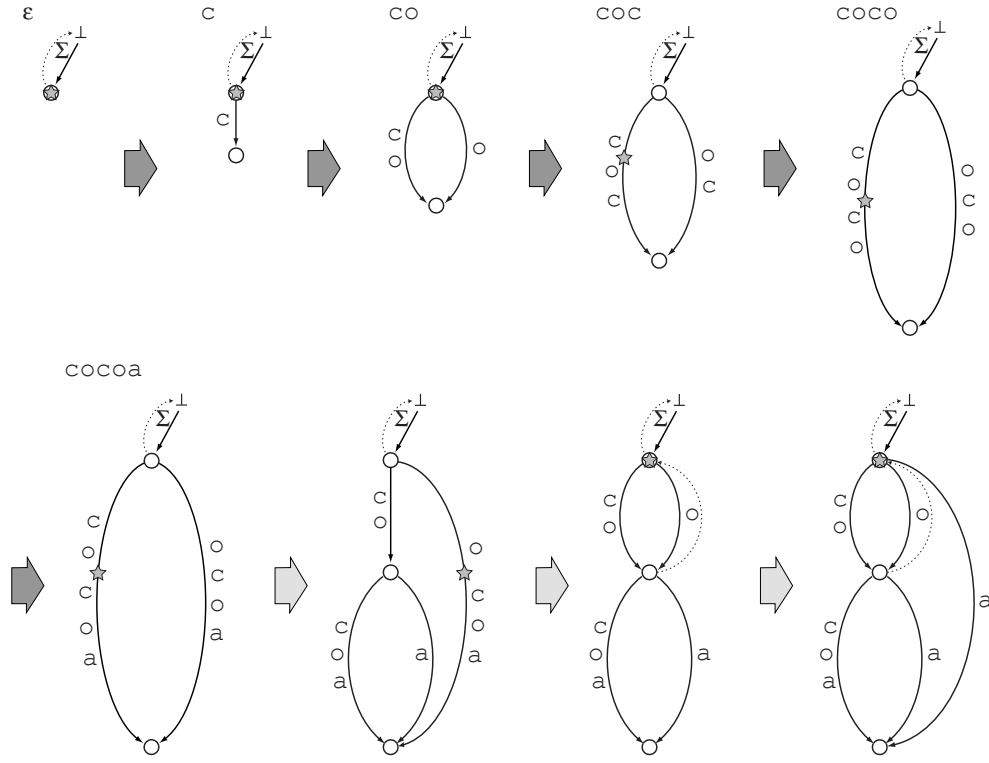


Fig. 12. On-line construction of $CDAWG'(w)$ with $w = cocoa$. The star mark represents the active point for each step.

edges of a CDAWG are also *automatically* extended, as well as those of a suffix tree (see $CDAWG'(coc)$ and $CDAWG'(coco)$).

Focus on the first step of the update of $CDAWG'(coco)$ to $CDAWG'(cocoa)$. String co there gets to be explicitly represented, and the active point is on implicit node o . In case of the construction of $STree'(cocoa)$, edge $(\varepsilon, oco, ocoa)$ is split into two edges (ε, o, o) and $(o, coa, ocoa)$, and then an open edge (o, a, oa) is newly created (see Fig. 9). However, in case of the CDAWG, edge $(\varepsilon, ocoa, ocoa)$ is *redirected* to node co , and the label is simultaneously modified. Since strings co and o are equivalent under the equivalence relation \equiv_{cocoa}^R , they are *merged* into a single node in $CDAWG'(cocoa)$.

Formal Description of the Algorithm.

The algorithm presented in this section for on-line construction of CDAWGs behaves similarly to Ukkonen's on-line suffix tree construction algorithm. Let $w \in \Sigma^*$ with $|w| = n$, and $a \in \Sigma$. Remember the set of the suffixes of w are divided into three sub-groups as mentioned in Sections 5.1 and 5.2. Let $u_1, u_2, \dots, u_n, u_{n+1}, u_{n+2}$ be the suffixes of string wa sorted in decreasing order of their lengths, and let $v_1, v_2, \dots, v_n, v_{n+1}$ be the suffixes of string w sorted in decreasing order of their lengths. Namely, $v_i a = u_i$ for any i with $1 \leq i \leq n+1$.

Assume $u_{\ell+1} = LRS(wa)$ and $v_{k+1} = LRS(w)$. The difference between on-line suffix tree construction and on-line CDAWG construction is summarized as follows:

- All the suffixes in the group (1) are equivalent under \equiv_{wa}^R . Thus all of them are represented by the same node $[\overset{wa}{\overrightarrow{u_1}}]_{wa}^R = [\overset{wa}{\overrightarrow{u_\ell}}]_{wa}^R$, which is in fact the sink node of $CDAWG'(wa)$. Namely, the destinations of the open edges are all the same. Moreover, $\overset{w}{\overrightarrow{v_i}} = v_i$ and $\overset{w}{\overrightarrow{u_i}} = u_i$ for any i with $1 \leq i \leq k$. According to these properties, we can generalize the idea of open edges as follows. For any open edge $(s, (p, \infty), t)$ of $CDAWG'(w)$ where t denotes the sink node, we actually implement it as $(s, (p, e), t)$ where e is a global variable that indicates the length of the current string, which is now $|w| = n$. When a new character a is added after w , we can extend all open edges only with increasing the value of e by 1. Thus $e = n + 1$ now. Obviously, it only takes $O(1)$ time.
- Recall **(Case 2)** regarding the position of the active point, mentioned in Section 5.2. Assume the active point is now on the *implicit* node representing a suffix v_j of w in $CDAWG'(w)$. Then, there can be a suffix v_h such that $h \neq j$ and $v_h \equiv_{wa}^R v_j$. In such case, they become *merged* into a single explicit node $[\overset{wa}{\overrightarrow{v_j}}]_{wa}^R$, during the update of $CDAWG'(w)$ to $CDAWG'(wa)$. The equivalence test is performed on the basis of Lemma 3 to be given in the sequel.
- Consider two distinct strings $x, y \in Factor(w)$ such that $\overset{w}{\overrightarrow{x}} = x$ and $\overset{w}{\overrightarrow{y}} = y$. Assume that $x \equiv_w^R y$, that is, they are represented by the same explicit node $[x]_w^R = [y]_w^R$ in $CDAWG'(w)$. Note that, however, x and y might not be equivalent under \equiv_{wa}^R . This means that when $CDAWG'(w)$ is updated to $CDAWG'(wa)$, then the node has to be *separated* into two nodes $[x]_{wa}^R$ and $[y]_{wa}^R$ if $x \not\equiv_{wa}^R y$. Here we can assume $|x| > |y|$ without loss of generality. Since this node separation happens only when $x \notin Suffix(wa)$ and $y \in Suffix(wa)$, we can do this procedure after we find the end point. The condition of the node separation will be given later on, in Lemma 4.

Merging Implicit Nodes.

As mentioned above, two or more nodes implicit in $CDAWG'(w)$ can be merged into one explicit node in $CDAWG'(wa)$. For a concrete example, we show in Fig. 13 $CDAWG'(w)$ and $CDAWG'(wa)$ with $w = \text{abcabcab}$ and $a = \text{a}$. It can be observed that the implicit nodes for abcab , bcab , and cab are merged into a single explicit node, and the implicit nodes for ab and b are also merged into another single explicit node in $CDAWG'(wa)$. The examination of whether to merge implicit nodes can be done by testing the equivalence of two nodes under the equivalence relation \equiv_{wa}^R . Here, we will often use the

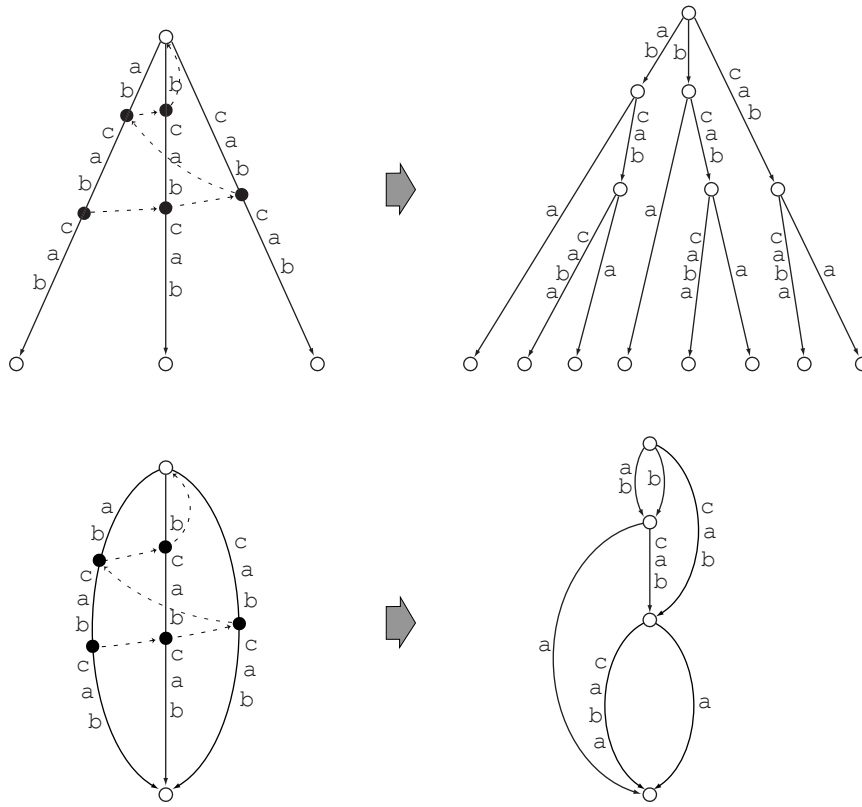


Fig. 13. Comparison of conversions. One is from $STree'(w)$ to $STree'(wa)$, while the other is from $CDAWG'(w)$ to $CDAWG'(wa)$ for $w = abcabcab$ and $a = a$. The black circles represent implicit nodes to be merged in the next step, connected by implicit suffix links corresponding to the traversal by the active point.

notation of $\overrightarrow{(\cdot)}^w$ rather than $\overleftarrow{(\cdot)}^w$ in order to refer to such a node which is implicit in $CDAWG'(w)$ but becomes explicit in $CDAWG'(wa)$. The equivalence test can be performed on the basis of the following proposition and lemma.

Proposition 9 *Let $w \in \Sigma^*$. For any string $x \in Factor(w)$, let $z = \overleftarrow{x}^w$. Then, x occurs within z exactly once.*

Proof. By Corollary 1 we have $\overleftarrow{x}^w = \overleftarrow{(\overrightarrow{x}^w)}^w$. Note that there exist strings $\alpha, \beta \in \Sigma^*$ such that $\overleftarrow{(\overrightarrow{x}^w)}^w = \overleftarrow{x}^w \alpha = \beta x \alpha$. According to Proposition 2, strings α, β have to be unique, which means that x can appear in w exactly once. \square

Lemma 3 *Let $w \in \Sigma^*$. For any strings $x, y \in Factor(w)$ with $y \in Suffix(x)$,*

$$x \equiv_w^R y \Leftrightarrow [\overrightarrow{x}^w]_w^R = [\overrightarrow{y}^w]_w^R.$$

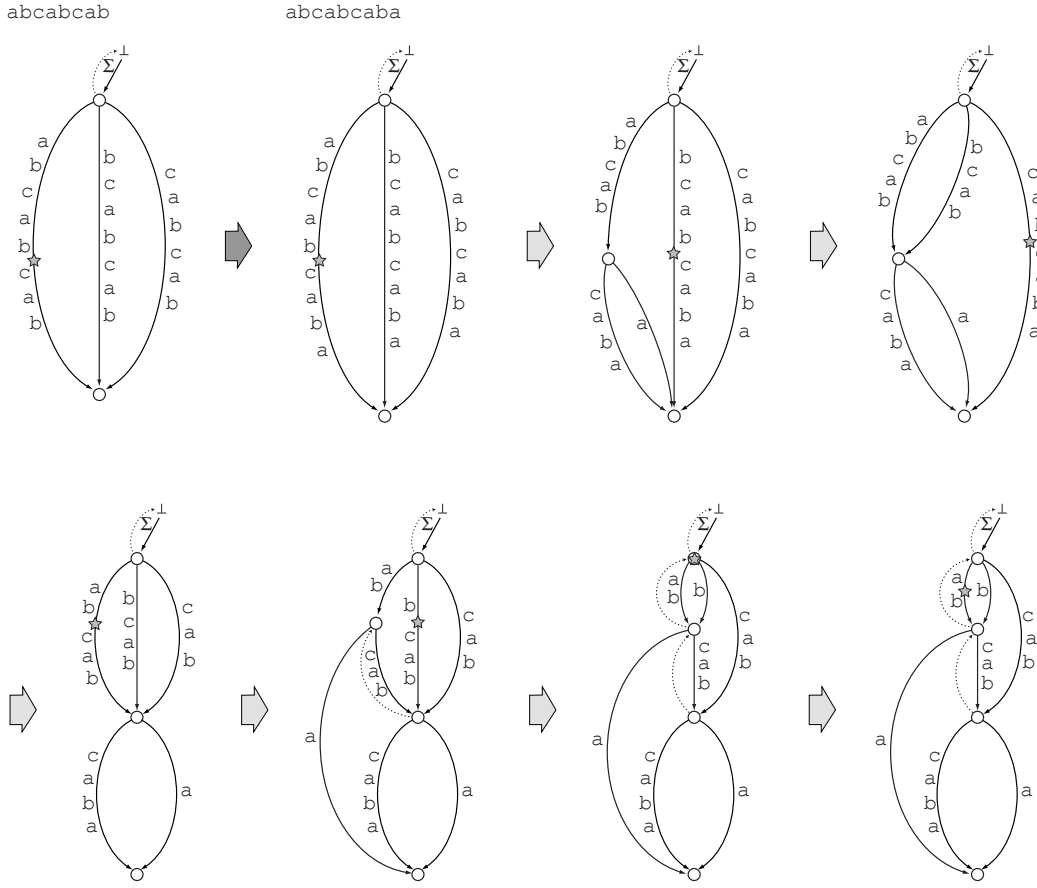


Fig. 14. Detailed conversion from $CDAWG'(w)$ to $CDAWG'(wa)$ for $w = abcabcab$ and $a = a$.

Proof. If $x \equiv_w^R y$, we have $\overleftarrow{x}^w = \overleftarrow{y}^w$ by Definition 3. By Corollary 1, we know $\overrightarrow{w} \overleftarrow{x}^w = \overleftarrow{w} \overrightarrow{x}^w$ and $\overrightarrow{w} \overleftarrow{y}^w = \overleftarrow{w} \overrightarrow{y}^w$, which yield $\overrightarrow{w} \overleftarrow{x}^w = \overrightarrow{w} \overleftarrow{y}^w$. Again by Definition 3, we have $[\overrightarrow{x}^w]^R_w = [\overrightarrow{y}^w]^R_w$.

Conversely, suppose $[\overrightarrow{x}^w]^R_w = [\overrightarrow{y}^w]^R_w$. Recall that $\overleftarrow{x}^w = \overleftarrow{w} \overrightarrow{x}^w$ by Corollary 1 and $\overleftarrow{w} \overrightarrow{x}^w$ is the unique longest member of $[\overrightarrow{x}^w]^R_w$. Similarly, \overleftarrow{y}^w is the unique longest member of $[\overrightarrow{y}^w]^R_w$. Thus we have $\overleftarrow{x}^w = \overleftarrow{y}^w$. Let $z = \overleftarrow{x}^w = \overleftarrow{y}^w$. Then $z = \alpha x \beta$ for some strings α and β . Since y is a suffix of x , there exists a string δ such that $x = \delta y$. We thus have $z = \alpha \delta y \beta$. This occurrence of y in z must be the only one due to Proposition 9. Since $\overleftarrow{y}^w = \alpha \delta y \beta$, we conclude that every occurrence of y within w must be preceded by δ . Thus we have $x \equiv_w^R y$. \square

For any string $x \in Factor(w)$, the equivalence class $[\overrightarrow{x}^w]^R_w$ is the closest explicit child of x in $CDAWG(w)$. Thus we can test the equivalence of two suffixes x, y of w with Lemma 3.

The matter is that, for a string $x \in \text{Suffix}(w)$, \overrightarrow{x}^w might not be explicit in $CDAWG'(w)$. Namely, on the equivalence test, we might refer to the node $[\overrightarrow{x}^w]_w^R$ instead of $[\overrightarrow{x}^w]_w^R$. Nevertheless, it does not actually happen in our on-line manner in which suffixes are processed in decreasing order of their length. See $CDAWG'(w)$ shown on the right of Fig. 13, where $w = \text{abcabcab}$. The black points are the implicit nodes the active point traverses in the next step via ‘implicit’ suffix links. In $CDAWG'(w)$, $[\overrightarrow{\text{cab}}]_w^R = [\overrightarrow{\text{ab}}]_w^R = [\overrightarrow{w}]_w^R$. However, in $CDAWG'(wa)$, $\text{cab} \not\equiv_{wa}^R \text{ab}$ where $a = \text{a}$. See Fig. 14 in which the detail of the update of $CDAWG'(w)$ to $CDAWG'(wa)$ is displayed. Notice that there is no trouble on merging the implicit nodes.

Separating Explicit Nodes.

When $CDAWG'(w)$ is updated to $CDAWG'(wa)$, an explicit node $[\overrightarrow{x}^w]_w^R$ with $x \in \text{Factor}(w)$ might become separated into two explicit nodes $[\overrightarrow{x}^w]_{wa}^R$ and $[\overrightarrow{y}^w]_{wa}^R$ if $x \notin \text{Suffix}(wa)$, $y \in \text{Suffix}(x)$, and $y \in \text{Suffix}(wa)$. It is inherently the same ‘phenomenon’ as the node separation occurring in the on-line construction of DAWGs [4]. Hereby we briefly recall the essence of the node separation of DAWGs. For $w \in \Sigma^*$ and $a \in \Sigma$, \equiv_{wa}^R is a refinement of \equiv_w^R . Furthermore, we have the following lemma.

Lemma 4 (Blumer et al. [4]) *Let $w \in \Sigma^*$ and $a \in \Sigma$. Let $z = \text{LRS}(wa)$. For a string $x \in \text{Factor}(w)$ assume $x = \overleftarrow{x}^w$, that is, x is the representative of $[x]_w^R$. Then,*

$$[x]_w^R = \begin{cases} [x]_{wa}^R \cup [z]_{wa}^R, & \text{if } z \in [x]_w^R \text{ and } x \neq z; \\ [x]_{wa}^R, & \text{otherwise.} \end{cases}$$

As stated in the above lemma, we need only to care about the node $[x]_w^R$ where $z \in [x]_w^R$ and $z = \text{LRS}(wa)$. Namely, at most one node can be separated when a DAWG is updated with a new character. If $z \neq \overleftarrow{x}^w$, it is separated into two nodes $[x]_{wa}^R$ and $[z]_{wa}^R$ when $DAWG(w)$ is updated to $DAWG(wa)$ (the former case). If $z = \overleftarrow{x}^w$, the node is not separated (the latter case). We examine whether $z = \overleftarrow{x}^w$ or not by checking the length of \overleftarrow{x}^w and z , as follows. Let $y \in \text{Factor}(w)$ be the string such that $\overleftarrow{y}^w \cdot a = z$. Note that there exists an

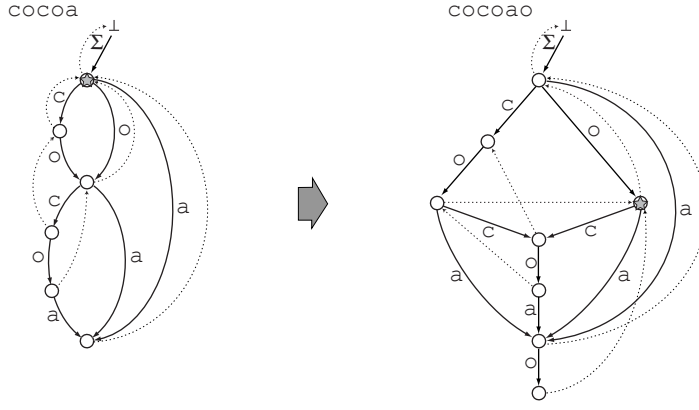


Fig. 15. Update of $DAWG(u)$ to $DAWG(ua)$, where $u = \text{cocoa}$ and $a = \text{o}$.

edge $([y]_w^R, a, [x]_w^R)$. Then,

$$z = \overset{w}{\leftarrow} x \Leftrightarrow \text{length}([y]_w^R) + |a| = \text{length}([x]_w^R), \text{ and}$$

$$z \neq \overset{w}{\leftarrow} x \Leftrightarrow \text{length}([y]_w^R) + |a| < \text{length}([x]_w^R).$$

If we define the length of the bottom node \perp by -1 , no contradiction occurs even in case that $z = \varepsilon$.

Fig. 15 shows the conversion from $DAWG(w)$ to $DAWG(wa)$ with $w = \text{cocoa}$ and $a = \text{o}$. The LRS of the string cocoa is o , therefore we focus on edge $([\varepsilon]_w^R, \text{o}, [\text{o}]_w^R)$. Since $\text{length}([\varepsilon]_w^R) + |\text{o}| = 1 < \text{length}([\text{o}]_w^R) = 2$, node $[\text{o}]_w^R$ is separated into two nodes $[\text{co}]_{wa}^R$ and $[\text{o}]_{wa}^R$, as seen in $DAWG(wa)$ of Fig. 15.

Now we go back to the update of $CDAWG'(w)$ to $CDAWG'(wa)$. The test of whether to separate a node when a CDAWG is updated can also be done on the basis of Lemma 4 in a very similar way. Since only explicit nodes can be separated, we only need to care about the case that $z = \overset{wa}{\leftarrow} z$ where $z = \text{LRS}(wa)$.

Lemma 5 *Let $w \in \Sigma^*$. Let $z = \text{LRS}(w)$. Then, if $z = \overset{w}{\leftarrow} z$, $\overset{w}{\leftarrow} x = \overset{w}{\leftarrow} x$ for any string $x \in \text{Factor}(w)$.*

Proof.

(1) When $x \in \text{Suffix}(z)$. Then we have $x \in \text{Suffix}(w)$ because $z = \text{LRS}(w)$.

From the definition of $(\overset{w}{\leftarrow} \cdot)$, therefore, it is obvious that $\overset{w}{\leftarrow} x = x$. Since $x \in$

$\text{Suffix}(w)$, $\overset{w}{\leftarrow} x = x$ from Proposition 3. Hence we have $\overset{w}{\leftarrow} x = \overset{w}{\leftarrow} x$.

(2) When $x \notin \text{Suffix}(z)$. We here have two sub-cases to consider:

(2-a) When x is a longer suffix of w than z . Since x occurs in w exactly once and $x \equiv_w^L x\alpha$ for any non-empty string α , we have $\overset{w}{\leftarrow} x = x$. Moreover, x

is not followed by any character in w , which means that $\overrightarrow{x}^w = x$. Thus we have $\overrightarrow{x}^w = \overrightarrow{x}^w$.

(2-b) When x is not a suffix of w . If there exist two distinct characters a, b such that $xa, xb \in Factor(w)$, then $\overrightarrow{x}^w = x$ from the definition. Moreover, $x \not\equiv_w^L xa$ since some occurrences of x in w are followed by b . Therefore we have $\overrightarrow{x}^w = x$, and thus $\overrightarrow{x}^w = \overrightarrow{x}^w$.

Now assume there exists a unique non-empty string β such that $\overrightarrow{x}^w = x\beta$. If $x\beta \in Suffix(w)$, then it falls into either Case (1) or (2-a). We now consider the case that $x\beta \notin Suffix(w)$. Since $\overrightarrow{x}^w = x\beta$, there exist two distinct characters a, b such that $x\beta a, x\beta b \in Factor(w)$, and it falls into the former part of Case (2-b).

Consequently, we have $\overrightarrow{x}^w = \overrightarrow{x}^w$ in any cases. \square

This lemma guarantees that the representative of $[\overrightarrow{x}^w]^R$ is equal to \overleftarrow{x}^w if the preconditions in the lemma are satisfied. We can therefore execute the node separation test as follows: If $z \neq \overleftarrow{x}^w$, it is separated into two nodes $[x]_{wa}^R$ and $[z]_{wa}^R$ when $CDAWG'(w)$ is updated to $CDAWG'(wa)$ (the former case). If $z = \overleftarrow{x}^w$, the node $[x]_w^R$ is not separated (the latter case). We examine if $z = \overleftarrow{x}^w$ or not by the length of \overleftarrow{x}^w and z in the following way. Let $y \in Factor(w)$ be the string such that $\overleftarrow{y}^w \cdot \alpha = z$ for some string $\alpha \in Factor(w)$. Note that there exists an edge $([y]_w^R, \alpha, [x]_w^R)$. Then,

$$\begin{aligned} z = \overleftarrow{x}^w &\Leftrightarrow length([y]_w^R) + |\alpha| = length([x]_w^R), \text{ and} \\ z \neq \overleftarrow{x}^w &\Leftrightarrow length([y]_w^R) + |\alpha| < length([x]_w^R). \end{aligned}$$

Fig. 16 shows the update of $CDAWG'(w)$ to $CDAWG'(wa)$, where $w = cocoa$ and $a = o$. The LRS of the string $cocoa$ is o , therefore we focus on edge $([\overrightarrow{\varepsilon}]_w^R, o, [\overrightarrow{o}]_w^R)$. Since $length([\overrightarrow{\varepsilon}]_w^R) + |o| = 1 < length([\overrightarrow{o}]_w^R) = 2$, node $[\overrightarrow{o}]_w^R$ is separated into two nodes $[\overrightarrow{co}]_{wa}^R$ and $[\overrightarrow{o}]_{wa}^R$, as seen in $CDAWG'(wa)$ of Fig. 16.

Pseudo-Code.

Our on-line algorithm to construct CDAWGs is described in Fig. 17 and Fig. 18. Function *extension* returns the explicit child node of a given node (implicit or explicit). Function *redirect_edge* redirects a given edge to a given node, with modifying the label of the edge accordingly. Function *split_edge* is

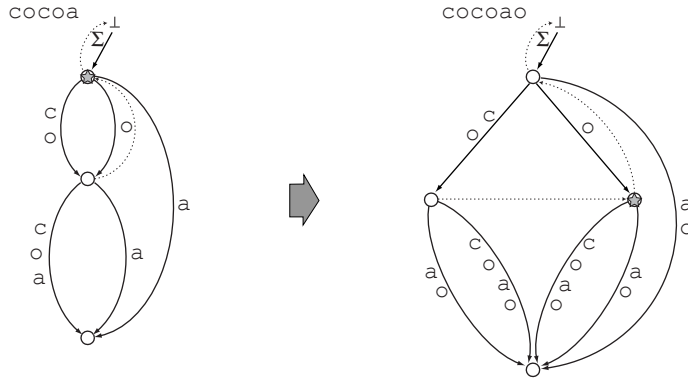


Fig. 16. Update of $CDAWG'(u)$ to $CDAWG'(ua)$, where $u = \text{cocoa}$ and $a = \text{o}$.

the same as the one used in Ukkonen's algorithm of Fig. 10, except that it also computes the length of nodes. Function *separate_node* separates a given node into two, if necessary. It is essentially the same as the separation procedure for $DAWG(w)$ given by Blumer et al. [4], except that implicit nodes are also treated.

Complexity of the Algorithm.

Theorem 6 *Assume Σ is a fixed alphabet. For any string $w \in \Sigma^*$, the proposed on-line algorithm to construct $CDAWG'(w)$ runs in $O(|w|)$ time.*

Proof. The linearity proof is in a sense the combination of that for the on-line algorithm for DAWGs [4] and that for the on-line algorithm for suffix trees [34]. We divide the time requirement into two components, both turn out to be linear. The first component consists of the total computation time by *canonize*. The second component consists of the rest.

Let $x \in \text{Factor}(w)$. We define the *suffix chain* started at x on w , denoted by $SC_w(x)$, to be the sequence of (possibly implicit) nodes reachable via suffix links from the (possibly implicit) node associated with x to the source node in $CDAWG'(w)$, as in [4]. We define its length by the number of nodes contained in the chain, and let $|SC_w(x)|$ denote it. Let k_1 be the number of iterations of the while loop of *update* and let k_2 be the number of iterations in the repeat-until loop in *separate_node*, when $CDAWG(w)$ is updated to $CDAWG(wa)$. By a similar argument as in [4], it can be derived that $|SC_{wa}(wa)| \leq |SC_w(w)| - (k_1 + k_2) + 2$. Initially $|SC_w(w)| = 1$ because $w = \varepsilon$, and then it grows at most two (possibly implicit) nodes longer in each call of *update*. Since $|SC_w(w)|$ decreases by an amount proportional to the sum of the number of iterations in the while loop and in the repeat-until loop on each call of *update*, the second time component is linear in the length of the input string.

For the analysis of the first time component we have only to consider the

```

Algorithm for on-line construction of  $CDAWG'(w\$)$ 
in alphabet  $\Sigma = \{w[-1], w[-2], \dots, w[-m]\}$ .
/* $ is the end-marker appearing nowhere in w. */
1 create nodes source, sink, and  $\perp$ ;
2 for  $j := 1$  to  $m$  do create a new edge  $(\perp, (-j, -j), source)$ ;
3  $suf(source) := \perp$ ;
4  $length(source) := 0$ ;  $length(\perp) := -1$ ;
5  $e := 0$ ;  $length(sink) := e$ ;
6  $(s, k) := (source, 1)$ ;  $i := 0$ ;
7 repeat
8    $i := i + 1$ ;  $e := i$ ; /* e is a global variable. */
9    $(s, k) := update(s, (k, i))$ ;
10 until  $w[i] = \$$ ;

function  $update(s, (k, p))$ : pair of node and integers;
/*  $(s, (k, p - 1))$  is the canonical reference pair for the active point. */
1  $c := w[p]$ ;  $oldr := \mathbf{nil}$ ;
2 while not  $check\_end\_point(s, (k, p - 1), c)$  do
3   if  $k \leq p - 1$  then /* implicit case. */
4     if  $s' = extension(s, (k, p - 1))$  then
5        $redirect\_edge(s, (k, p - 1), r)$ ;
6        $(s, k) := canonize(suf(s), (k, p - 1))$ ;
7       continue;
8     else
9        $s' := extension(s, (k, p - 1))$ ;
10       $r := split\_edge(s, (k, p - 1))$ ;
11    else /* explicit case. */
12       $r := s$ ;
13    create edge  $(r, (p, e), sink)$ ;
14    if  $oldr \neq \mathbf{nil}$  then  $suf(oldr) := r$ ;
15     $oldr := r$ ;
16     $(s, k) := canonize(suf(s), (k, p - 1))$ ;
17 if  $oldr \neq \mathbf{nil}$  then  $suf(oldr) := s$ ;
18 return  $separate\_node(s, (k, p))$ ;

```

Fig. 17. Main routine, **function** $update$, and **function** $check_end_point$ of the on-line algorithm to construct CDAWGs.

number of iterations in the while loop in $canonize$. By considering the calls of $canonize$ executed in the while loop in $update$, it results in that the total number of the iterations is linear (by the same argument as in [34]). Thus we shall consider the number of iterations of the while loop in $canonize$ called in $separate_node$. There are two cases to consider:

- (1) When the end point is on an implicit node. Then the computation in $canonize$ takes only constant time.
- (2) When the end point is on an explicit node. Let z be the LRS of w , which corresponds to the end point. Consider the last edge in the path spelling out z from the source node to the explicit node, and let the length of its label be k (≥ 1). The total number of iterations of the while loop of $canonize$ in the call of $separate_node$ is at most k . Since the value of k increases at most by 1 each time a new character is scanned, the time

```

function extension(s, (k, p)): node;
/* (s, (k, p)) is a canonical reference pair. */
1 if k > p then return s; /* explicit case. */
2 find the  $w[k]$ -edge (s, (k', p'), s') from s;
3 return s';

function redirect_edge(s, (k, p), r);
1 let (s, (k', p'), s') be the  $w[k]$ -edge from s;
2 replace the edge by edge (s, (k', k' + p - k), r);

function split_edge(s, (k, p)): node;
1 let (s, (k', p'), s') be the  $w[k]$ -edge from s;
2 create node r;
3 replace the edge by (s, (k', k' + p - k), r) and (r, (k' + p - k + 1, p'), s');
4  $length(r) := length(s) + (p - k + 1)$ ;
5 return r;

function separate_node(s, (k, p)): pair of node and integer;
1 (s', k') := canonicalize(s, (k, p));
2 if k' ≤ p then return (s', k'); /* implicit case. */
3 /* explicit case. */
4 if  $length(s') = length(s) + (p - k + 1)$  then return (s', k');
5 create node r' as a duplication of s' with the out-going edges;
6  $suf(r') := suf(s')$ ;  $suf(s') := r'$ ;
7  $length(r') := length(s) + (p - k + 1)$ ;
8 repeat
9   replace the  $w[k]$ -edge from s to s' by edge (s, (k, p), r');
10  (s, k) := canonicalize(suf(s), (k, p - 1));
11 until (s', k') ≠ canonicalize(s, (k, p));
12 return (r', p + 1);

```

Fig. 18. Other functions for the on-line algorithm to construct CDAWGs. Since **function** *check_end_point* and **function** *canonicalize* used here are identical to those shown in Fig. 10, they are omitted.

requirement of the while loop of *canonicalize* in *separate_node* is bounded by the total length of the input string.

As a result of the above discussion, we can finally conclude that the first and second components take overall linear time. \square

7 Construction of CDAWGs for a Set of Strings

In the previous chapters we discussed on-line construction of index structures for a single string $w \in \Sigma^*$. On the other hand, we now consider such case that we are given a set S of strings as an input. The suffix trie, suffix tree, DAWG, and CDAWG for S can all be well-defined. Any index structure for S must represent all strings in $Factor(S)$.

Blumer et al. [5] introduced DAWGs for a set of strings, and presented an algo-

rithm that builds $DAWG(S)$ in $O(\|S\|)$ time. They also introduced CDAWGs for a set of strings. Their algorithm for construction of $CDAWG(S)$ runs in linear time in the input size (namely, in $O(\|S\|)$ time), but not in time linear in the output size because it first builds $DAWG(S)$ and then converts it to $CDAWG(S)$ by deleting internal nodes of out-degree one and concatenating their edges accordingly. Kosaraju [26] introduced the suffix tree for a set S of strings, often referred to the *generalized suffix tree* for S . A slight modification of Ukkonen’s algorithm recalled in Section 5.2 is capable of constructing $STree'(S)$ in $O(\|S\|)$ time.

In this section, we give the first algorithm which builds $CDAWG'(S)$ in $O(\|S\|)$ time, and directly. Let $S = \{w_1, w_2, \dots, w_k\}$, where $k = |S|$. Then we consider set $S' = \{w_i\$i \mid w_i \in S \text{ and } \$i \notin Factor(S) \text{ for any } 1 \leq i \leq |S|\}$. Notice that S' has the prefix property, and thus, $CDAWG'(S') = CDAWG(S')$ for any S . $CDAWG'(S')$ can be constructed by a slight modification of the algorithm proposed in the previous section. We use a global variable e_i for each string in S' , where $1 \leq i \leq |S|$, which indicates the ending position of the open edges for each string. The set S' is input in the form of a single stream $s = w_1\$1w_2\$2 \cdots w_k\$k$. Whenever we encounter an end-marker $\$i$ in the input, we stop increasing the value of e_i . Then we create the new $(i + 1)$ -th sink node, and start increasing the value of e_{i+1} each time a new character is scanned until encountering $\$(i+1)$. Hereby we have the following:

Theorem 7 *Assume Σ is a fixed alphabet. For any set S of strings, the proposed algorithm directly constructs $CDAWG'(S')$ on-line and in $O(\|S'\|)$ time, using $O(\|S'\|)$ space.*

Fig. 19 shows construction of $CDAWG'(S')$, where $S' = \{\text{cocoa}\$1, \text{cola}\$2\}$.

8 Conclusion

The *compact directed acyclic word graph* ($CDAWG$) of a string w is a space-economical index structure that represents all factors of w . CDAWGs have so far been studied in literature such as [5,13].

In this paper we presented an on-line linear-time algorithm for constructing CDAWGs for a single string. The algorithm can easily be extended to construction of CDAWGs for a set of strings, running in linear time with respect to the total length of the strings in the set. A CDAWG can be obtained by compacting the corresponding DAWG or minimizing the corresponding suffix tree; however, our approach permits us to save time and space simultaneously, since the CDAWG can be built directly. Moreover, in our on-line manner the CDAWG for wa can be built simply by updating the CDAWG of w with the

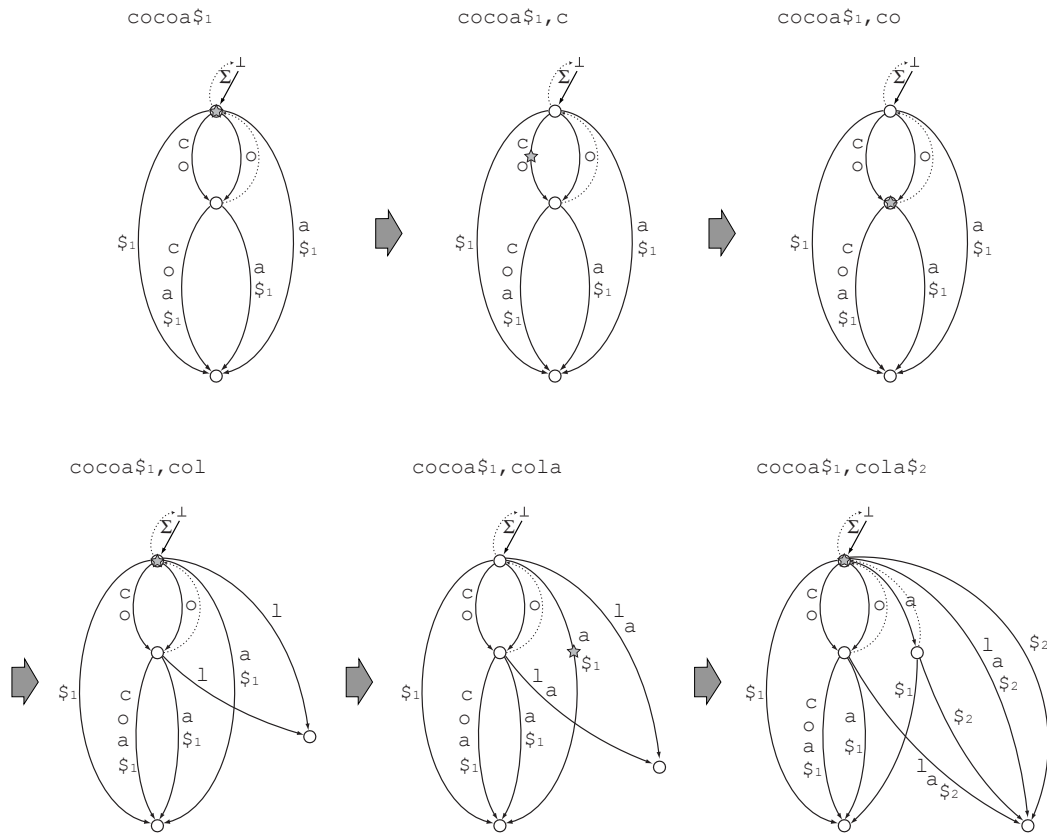


Fig. 19. Construction of $CDAWG'(S')$ for $S' = \{cocoa\$1, cola\$2\}$.

new character a , in contrast with Crochemore and V erin's off-line algorithm in [13] which needs to reconstruct it from scratch.

There have been made further work on CDAWGs. In [20,22] a *generic* algorithm to construct suffix tries, suffix trees, DAWGs, and CDAWGs was considered. The CDAWG of a *trie* was introduced in [19], where the trie is a tree structure representing a set of strings. An algorithm to construct the CDAWG for a given trie was also given in that paper, which runs in time proportional to the number of nodes in the trie. In [21], an on-line algorithm that builds *symmetric CDAWGs (SCDAWGs)* was introduced. The SCDAWG of a string w was first introduced in [5], which represents all factors of both w and its reversed string. Moreover, a liner-time algorithm to construct and maintain CDAWGs for a *sliding window* was given in [24]. An application of CDAWGs for a sliding window is the *prediction by partial matching (PPM)* style statistical data compression model [8,7].

Not only are CDAWGs attractive as an index structure, but also the underlying equivalence relation is useful in data mining and machine discovery from textual databases. In fact, the equivalence relation has played a central role in supporting expert researchers working on evaluating and interpreting literary knowledge mined from anthologies of classical Japanese poems [32].

References

- [1] A. Apostolico. The myriad virtues of subword trees. In *Combinatorial Algorithm on Words*, volume 12 of *NATO Advanced Science Institutes, Series F*, pages 85–96. Springer-Verlag, 1985.
- [2] A. Apostolico and S. Lonardi. A speed-up for the commute between subword trees and DAWGs. *Information Processing Letters*, 83(3):159–161, 2002.
- [3] M. Balík. Implementation of DAWG. In *Proc. The Prague Stringology Club Workshop '98 (PSCW'98)*, pages 26–35. Czech Technical University, 1998.
- [4] A. Blumer, J. Blumer, D. Haussler, A. Ehrenfeucht, M. T. Chen, and J. Seiferas. The smallest automaton recognizing the subwords of a text. *Theoretical Computer Science*, 40:31–55, 1985.
- [5] A. Blumer, J. Blumer, D. Haussler, R. McConnell, and A. Ehrenfeucht. Complete inverted files for efficient text retrieval and analysis. *J. ACM*, 34(3):578–595, 1987.
- [6] M. T. Chen and J. Seiferas. Efficient and elegant subword tree construction. In *Combinatorial Algorithm on Words*, volume 12 of *NATO Advanced Science Institutes, Series F*, pages 97–107. Springer-Verlag, 1985.
- [7] J. G. Cleary, W. J. Teahan, and I. H. Witten. Unbounded length contexts for PPM. In *Proc. Data Compression Conference '95 (DCC'95)*, pages 52–61. IEEE Computer Society, 1995.
- [8] J. G. Cleary and I. H. Witten. Data compression using adaptive coding and partial string matching. *IEEE Trans. Commun.*, 32(4):396–402, 1984.
- [9] M. Crochemore. Transducers and repetitions. *Theoretical Computer Science*, 45:63–86, 1986.
- [10] M. Crochemore. Reducing space for index implementation. *Theoretical Computer Science*, 292(1):185–197, 2003.
- [11] M. Crochemore and W. Rytter. *Text Algorithms*. Oxford University Press, New York, 1994.
- [12] M. Crochemore and W. Rytter. *Jewels of Stringology*. World Scientific, 2002.
- [13] M. Crochemore and R. Verin. On compact directed acyclic word graphs. In *Structures in Logic and Computer Science*, volume 1261 of *Lecture Notes in Computer Science*, pages 192–211. Springer-Verlag, 1997.
- [14] G. Gonnet, R. Baeza-Yates, and T. Snider. New indices for text: PAT trees and PAT arrays. *Information retrieval: data structures and algorithms*, pages 66–82, 1992.
- [15] R. Grossi and J. S. Vitter. Compressed suffix arrays and suffix trees with applications to text indexing and string matching. In *Proc. of 32nd ACM Symposium on Theory of Computing (STOC'00)*, pages 397–406, 2000.

- [16] D. Gusfield. *Algorithms on Strings, Trees, and Sequences*. Cambridge University Press, New York, 1997.
- [17] J. Holub and B. Melichar. Approximate string matching using factor automata. *Theoretical Computer Science*, 249:305–311, 2000.
- [18] S. Inenaga. Bidirectional construction of suffix trees. *Nordic Journal of Computing*, 10(1):52–67, 2003.
- [19] S. Inenaga, H. Hoshino, A. Shinohara, M. Takeda, and S. Arikawa. Construction of the CDAWG for a trie. In *Proc. The Prague Stringology Conference '01 (PSC'01)*, pages 37–48. Czech Technical University, 2001.
- [20] S. Inenaga, H. Hoshino, A. Shinohara, M. Takeda, and S. Arikawa. On-line construction of compact directed acyclic word graphs. Technical Report DOI-TR-CS-183, Department of Informatics, Kyushu University, 2001.
- [21] S. Inenaga, H. Hoshino, A. Shinohara, M. Takeda, and S. Arikawa. On-line construction of symmetric compact directed acyclic word graphs. In *Proc. of 8th International Symposium on String Processing and Information Retrieval (SPIRE'01)*, pages 96–110. IEEE Computer Society, 2001.
- [22] S. Inenaga, H. Hoshino, A. Shinohara, M. Takeda, and S. Arikawa. Unification of algorithms to construct index structures for texts. Technical Report DOI-TR-CS-196, Department of Informatics, Kyushu University, 2001.
- [23] S. Inenaga, H. Hoshino, A. Shinohara, M. Takeda, S. Arikawa, G. Mauri, and G. Pavesi. On-line construction of compact directed acyclic word graphs. In *Proc. 12th Annual Symposium on Combinatorial Pattern Matching (CPM'01)*, volume 2089 of *Lecture Notes in Computer Science*, pages 169–180. Springer-Verlag, 2001.
- [24] S. Inenaga, A. Shinohara, M. Takeda, and S. Arikawa. Compact directed acyclic word graphs for a sliding window. In *Proc. 9th International Symposium on String Processing and Information Retrieval (SPIRE'02)*, volume 2476 of *Lecture Notes in Computer Science*, pages 310–324. Springer-Verlag, 2002.
- [25] J. Kärkkäinen. Suffix cactus: A cross between suffix tree and suffix array. In *Proc. 6th Annual Symposium on Combinatorial Pattern Matching (CPM'95)*, volume 973 of *Lecture Notes in Computer Science*, pages 191–204. Springer-Verlag, 1995.
- [26] S. R. Kosaraju. Fast pattern matching in trees. In *Proc. 30th IEEE Symp. on Foundations of Computer Science*, pages 178–183, 1989.
- [27] S. Kurtz. Reducing the space requirement of suffix trees. *Software - Practice and Experience*, 29(13):1149–1171, 1999.
- [28] V. Mäkinen. Compact suffix array - a space-efficient full-text index. *Fundamenta Informaticae*, 56(1-2):191–210, 2003.
- [29] U. Manber and G. Myers. Suffix arrays: A new method for on-line string searches. *SIAM Journal on Computing*, 22(5):935–948, 1993.

- [30] E. M. McCreight. A space-economical suffix tree construction algorithm. *J. ACM*, 23(2):262–272, 1976.
- [31] K. Sadakane. Compressed text databases with efficient query algorithms based on the compressed suffix array. In *Proc. of 11th International Symposium on Algorithms and Computation (ISAAC'00)*, volume 1969 of *Lecture Notes in Computer Science*, pages 410–421. Springer-Verlag, 2000.
- [32] M. Takeda, T. Matsumoto, T. Fukuda, and I. Nanri. Discovering characteristic expressions from literary works. *Theoretical Computer Science*, 292(2):525–546, 2003.
- [33] E. Ukkonen. Approximate string matching over suffix trees. In *Proc. 4th Annual Symposium on Combinatorial Pattern Matching (CPM'93)*, volume 684 of *Lecture Notes in Computer Science*, pages 228–242. Springer-Verlag, 1993.
- [34] E. Ukkonen. On-line construction of suffix trees. *Algorithmica*, 14(3):249–260, 1995.
- [35] E. Ukkonen and D. Wood. Approximate string matching with suffix automata. *Algorithmica*, 10(5):353–364, 1993.
- [36] P. Weiner. Linear pattern matching algorithms. In *Proc. 14th Annual Symposium on Switching and Automata Theory*, pages 1–11, 1973.