

Compact Directed Acyclic Word Graphs for a Sliding Window

Shunsuke Inenaga¹, Ayumi Shinohara^{1,2},
Masayuki Takeda^{1,2}, and Setsuo Arikawa¹

¹ Department of Informatics, Kyushu University 33, Fukuoka 812-8581, Japan

² PRESTO, Japan Science and Technology Corporation (JST)

{s-ine, ayumi, takeda, arikawa}@i.kyushu-u.ac.jp

Abstract. The *suffix tree* is a well-known and widely-studied data structure that is highly useful for string matching. The suffix tree of a string w can be constructed in $O(n)$ time and space, where n denotes the length of w . Larsson achieved an efficient algorithm to maintain a suffix tree for a sliding window. It contributes to *prediction by partial matching (PPM)* style statistical data compression scheme. The *compact directed acyclic word graph (CDAWG)* is a more space-economical data structure for indexing a string. In this paper we propose a linear-time algorithm to maintain a CDAWG for a sliding window.

1 Introduction

String matching is one of the central tasks in data compression. A straightforward method for matching strings would be to construct an index structure for the full text to be compressed. It is, however, easy to imagine that indexing the whole string requires too much space for the storage. This fact implies that an index structure needs to be *dynamic* in order to be suitable for *processing part of the text*.

The *suffix tree* is a highly efficient data structure for string matching, therefore being extensively studied [18, 15, 1, 17, 5, 8]. $STree(w)$ represents all factors of w and can be constructed in linear time and space. The on-line algorithm by Ukkonen [17] processes a given string w from left to right, and at each j -th phase it maintains $STree(w[1 : j])$, where $w[1 : j]$ denotes the prefix of w of length j . Larsson [13] modified Ukkonen's algorithm so as to maintain $STree(w[i : j])$ with $0 \leq i \leq j \leq |w|$, for any factor of length $j - i + 1 = M$. The width M of indexed factors is called the *window size*. That is, Larsson presented an algorithm to maintain a suffix tree for a *sliding window* mechanism where the values of i and j are incremented.

An application of a suffix tree for a sliding window is the *prediction by partial matching (PPM)* style statistical data compression model [4, 3, 16]. PPM* [3] is an improvement that allows unbounded context length. PPM* employs a tree structure called the *context trie*, which supports indexes of the input string. The drawback of PPM* is, however, its too much computational resources in

both time and space, which weakens its practical usefulness. In particular, the context trie occupies major part of the space requirement. Larsson’s suffix tree for a sliding window offered a variant of PPM*, feasible in practice since its space requirement is bounded by the window size M and the running time is linear in the length of the input string w .

In this paper, we take another approach to reduce the space requirement in PPM*-style statistical compression. We propose an algorithm to maintain a *compact directed acyclic word graph* (CDAWG) for a sliding window, which performs in linear time and space. CDAWGs require less space than suffix trees in both theory and practice [2, 6]. In our previous work [11], we presented an on-line algorithm that constructs $CDAWG(w)$ in linear time and space. Moving the rightmost position of a sliding window can be accomplished by the algorithm. In case of a suffix tree, it is also rather straightforward to advance the leftmost position of a sliding window: basically we have only to remove the leaf node and its in-coming edge corresponding to the longest suffix. However, since a CDAWG is a graph, the matter is much more complex and technically difficult. Thus more detailed and precise discussions are necessary. In addition, we have to ensure that no edge labels refer to positions outside a sliding window. To guarantee it, Larsson utilized the technique of *credit issuing* first introduced in [7], which takes amortized constant time. We introduce an extended version of credit issuing that is modified to be suitable for treating CDAWGs.

2 Compact Directed Acyclic Word Graphs

2.1 Definitions

Let Σ be an *alphabet*. An element of Σ^* is called a *string*. The length of a string w is denoted by $|w|$. 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. The empty string is denoted by ε , that is, $|\varepsilon| = 0$. Let $\Sigma^+ = \Sigma^* - \{\varepsilon\}$. The factor of a string w that begins at position i and ends at position j is denoted by $w[i:j]$ for $1 \leq i \leq j \leq |w|$. For convenience, let $w[i:j] = \varepsilon$ for $j < i$. Let $S \subseteq \Sigma^*$. The cardinality of S is denoted by $|S|$. For any string $x \in \Sigma^*$, $Sx^{-1} = \{u \mid ux \in S\}$ and $x^{-1}S = \{u \mid xu \in S\}$. We define equivalence relations \equiv_w^L and \equiv_w^R on Σ^* by

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

respectively. Let $[x]_w^L$ and $[x]_w^R$ denote the equivalence classes of a string $x \in \Sigma^*$ under \equiv_w^L and \equiv_w^R , respectively. The longest elements in the equivalence classes $[x]_w^L$ and $[x]_w^R$ for $x \in Factor(w)$ are called their *representatives* and denoted by \overrightarrow{x}^w and \overleftarrow{x}^w , respectively. For any string $x \in Factor(w)$, there uniquely exist strings $\alpha, \beta \in \Sigma^*$ such that $\overrightarrow{x}^w = x\alpha$ and $\overleftarrow{x}^w = \beta x$.

We now introduce a relation X_w over Σ^* such that

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

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. It can be readily shown that \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 $x \in \text{Factor}(w)$, let \overrightarrow{x} denote the longest string in the equivalence class to which x belongs under the equivalence relation \equiv_w^L .

Note that \overrightarrow{x} and \overrightarrow{x} are not always equal. For example, consider the case that $w = \text{abab}$ and $x = \text{ab}$, where $\overrightarrow{x} = \text{ab}$ but $\overrightarrow{x} = \text{abab}$. More formally:

Proposition 1 ([10]). *Let $w \in \Sigma^*$. For any string $x \in \text{Factor}(w)$, \overrightarrow{x} is a prefix of \overrightarrow{x} . If $\overrightarrow{x} \neq \overrightarrow{x}$, then $\overrightarrow{x} \in \text{Suffix}(w)$.*

In the following, we define the suffix tree and the CDAWG of w , denoted by $\text{STree}(w)$ and $\text{CDAWG}(w)$, respectively. 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 give definitions of the *suffix links*, frequently used for time-efficient construction of the structures [18, 15, 17, 2, 6, 11, 9].

Definition 1. *$\text{STree}(w)$ is the tree (V, E) such that*

$$\begin{aligned} V &= \{\overrightarrow{x} \mid x \in \text{Factor}(w)\}, \\ E &= \{(\overrightarrow{x}, a\beta, \overrightarrow{x\beta}) \mid x, xa \in \text{Factor}(w), a \in \Sigma, \beta \in \Sigma^*, \overrightarrow{x\beta} = xa\beta, \overrightarrow{x} \neq \overrightarrow{x\beta}\}, \end{aligned}$$

and its suffix links are the set

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

Definition 2. *$\text{CDAWG}(w)$ is the dag (V, E) such that*

$$\begin{aligned} V &= \{[\overrightarrow{x}]_w^R \mid x \in \text{Factor}(w)\}, \\ E &= \{([\overrightarrow{x}]_w^R, a\beta, [\overrightarrow{x\beta}]_w^R) \mid x, xa \in \text{Factor}(w), a \in \Sigma, \beta \in \Sigma^*, \overrightarrow{x\beta} = xa\beta, \overrightarrow{x} \neq \overrightarrow{x\beta}\}, \end{aligned}$$

and its suffix links are the set

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

One can see that $\text{CDAWG}(w)$ is the “minimization” of $\text{STree}(w)$ due to “the $[(\cdot)]_w^R$ operation”.

The nodes $[\overrightarrow{\varepsilon}]_w^R = [\varepsilon]_w^R$ and $[\overrightarrow{w}]_w^R = [w]_w^R$ are called the *source* node and the *sink* node of $\text{CDAWG}(w)$, respectively. For any $x \in \text{Factor}(w)$ such that

$x = \overrightarrow{x}^w$, x is said to be represented by the *explicit* node $[\overrightarrow{x}^w]^R_w$. If $x \neq \overrightarrow{x}^w$, x is said to be on an *implicit* node. The implicit node is represented by a *reference* pair $([\overrightarrow{z}^w]^R_w, y)$ such that $z \in \text{Prefix}(x)$, $y \in \Sigma^*$ and $\overrightarrow{z}^w \cdot y = x$. When $|y|$ is minimum, the pair $([\overrightarrow{z}^w]^R_w, y)$ is called the *canonical* reference pair of x . Note that an explicit node can also be represented by a reference pair.

Proposition 1 implies that, for a string $x \in \text{Factor}(w)$, \overrightarrow{x}^w is not always represented on an explicit node in $CDAWG(w)$. Actually, in $CDAWG(\text{coco})$ displayed at Fig. 1, string $\overrightarrow{\text{co}}^w = \text{co}$ is on an implicit node, where $w = \text{coco}$.

To implement $CDAWG(w)$ using only $O(|w|)$ space, labels of edges are represented by two integers indicating their beginning and ending positions in w , respectively. Suppose that $([\overrightarrow{x}^w]^R_w, y, [\overrightarrow{z}^w]^R_w)$ is an edge of $CDAWG(w)$. Then the edge label y is actually represented by a pair (i, j) of integers such that $w[i : j] = y$. A reference pair can be represented in a similar way.

2.2 On-Line Algorithm to Construct CDAWGs

We here recall the on-line algorithm to construct $CDAWG(w)$, introduced in our previous work [11]. The algorithm is based on Ukkonen's on-line suffix tree construction algorithm [17]. It constructs $CDAWG(wa)$ from $CDAWG(w)$ by inserting suffixes of wa into $CDAWG(w)$ in decreasing order of their length. Let z be the longest string in $\text{Factor}(w) \cap \text{Suffix}(wa)$. Then z is called the *longest repeated suffix* of wa and denoted by $LRS(wa)$. Let $z' = LRS(w)$. Let $|wa| = l$ and $u_1, u_2, \dots, u_l, u_{l+1}$ be the suffixes of wa ordered in their length, that is, $u_1 = wa$ and $u_{l+1} = \varepsilon$. We categorize these suffixes into the following three groups.

- (Group 1)** u_1, \dots, u_{i-1}
- (Group 2)** u_i, \dots, u_{j-1} where $u_i = z'a$
- (Group 3)** u_j, \dots, u_{l+1} where $u_j = z$

Note all suffixes in Group 3 have already been represented in $CDAWG(w)$. Let v_1, \dots, v_{i-1} be the suffixes of w such that, for any $1 \leq k \leq i-1$, $v_k a = u_k$. Then we can insert all the suffixes of Group 1 into $CDAWG(w)$ by appending the character a at the end of edges leading to the sink node. Moreover, we can update those edges in constant time, by setting the ending position of the labels so to refer to a global variable e indicating the length of the scanned part of the input string. In this case, $e = l$. It therefore results in that we have only to care about those in Group 2. We start from the locus corresponding to $z' = v_i$ in $CDAWG(w)$, which is called the *active point*. If the active point is on an edge, a new explicit node is created and from there a new edge labeled by a is inserted leading to the sink node. The locus of v_{i+1} can be found by following the suffix link and possibly some downward edges. Assume v_{i+1} is also on an edge when it is found. Sometimes the edge can be *redirected* to the node created when u_i was inserted. For the detailed condition to examine whether the edge should be

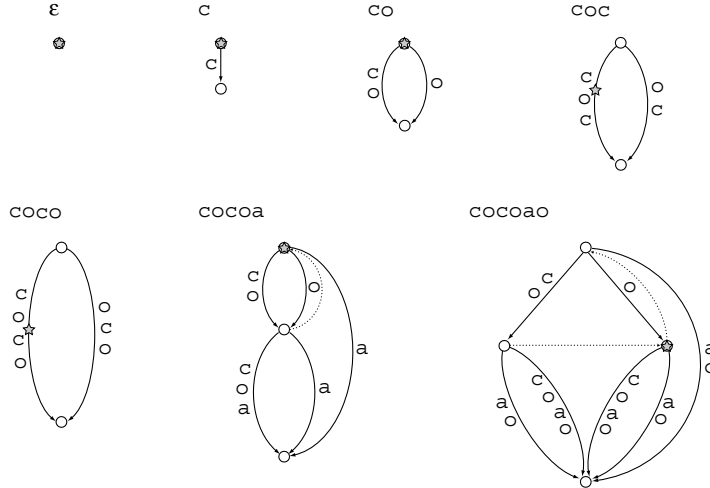


Fig. 1. The on-line construction of $CDAWG(w)$ with $w = cocoao$. The dotted lines represent the suffix links. The gray star indicates the active point of each step.

merged or not, see [11]. After the last suffix u_{j-1} is inserted to the CDAWG, all suffixes of wa are represented in the CDAWG. We now pay attention to $LRS(wa) = z = u_j$. Consider the case that $\xrightarrow{wa} u_j = u_j$, that is, the case that u_j is now represented in an explicit node. Then we have to check whether or not u_j is the representative of $[u_j]_{wa}^R$. If not, the explicit node is *separated* into two explicit nodes $[x]_{wa}^R$ and $[u_j]_{wa}^R$, where $x \in [u_j]_w^R$ and $x \neq u_j$. The on-line construction of $CDAWG(w)$ with $w = cocoao$ is shown in Fig. 1.

Theorem 1 ([11]). *For any string $w \in \Sigma^*$, $CDAWG(w)$ can be constructed on-line and in $O(|w|)$ time and space.*

3 Suffix Trees for a Sliding Window

In [13], Larsson introduced an algorithm for maintaining a suffix tree for a *sliding window*, whose width is $M > 0$. Let i (resp. j) be the leftmost (resp. rightmost) position of the window sliding in w , that is, $j - i + 1 = M$. To move the sliding window ahead, we need to increment i and j . Incrementing j can be accomplished by Ukkonen's on-line algorithm. On the other hand, incrementing i means to delete the leftmost character of the currently scanned string, that is, to convert $STree(bu)$ into $STree(u)$ with some $b \in \Sigma$ and $u, bu \in Factor(w)$. We focus on the path of $STree(bu)$ which spells out bu from the root node. This path is called the *backbone* of $STree(bu)$. Let x be the longest string in $Prefix(bu) - \{bu\}$ such that $\xrightarrow{bu} x = x$. The locus of x in $STree(bu)$ is called the *deletion point* and denoted by $DelPoint(bu)$. On the other hand, let z be the longest string in $Prefix(bu) - \{bu\}$

such that $\overset{bu}{\rightarrow}z = z$. The string z is called the *last node* in the backbone and denoted by $LastNode(bu)$.

When $DelPoint(bu) = LastNode(bu)$, there is an explicit node representing the string x in $STree(bu)$. Then there exists an edge $(\overset{bu}{\rightarrow}x, y, \overset{bu}{\rightarrow}bu)$ in $STree(bu)$ where $xy = bu$. Only by removing this edge, we can obtain $STree(u)$.

When $DelPoint(bu) \neq LastNode(bu)$, it follows from Proposition 1 that $x \in Suffix(bu)$. Moreover, $x = LRS(bu)$, as to be proven by Lemma 3 in Section 4.1.

Namely, the active point is on the locus for x in $STree(bu)$. Let $(\overset{bu}{\rightarrow}s, y, \overset{bu}{\rightarrow}bu)$ be the edge on which x is represented. Let $\overset{bu}{\rightarrow}s \cdot t = x$, where $t \in Prefix(y)$. We shorten the edge to $(\overset{u}{\rightarrow}s, t, \overset{u}{\rightarrow}x)$, and move the active point to the locus for the one-character shorter suffix of x .

Theorem 2 ([13]). *Let $w \in \Sigma^*$ and M be the window size. Larsson's algorithm runs in $O(|w|)$ time using $O(M)$ space.*

4 CDAWGs for a Sliding Window

In this section, we consider the maintenance of a CDAWG for a sliding window. Advancing the rightmost position of the window can be done by the on-line algorithm recalled in Section 2.2. Thus the matter is to move ahead the leftmost position of the window.

4.1 Edge Deletion

Given $CDAWG(w)$, we also focus on its backbone, the path spelling out w from the source node. Let $x = DelPoint(w)$. If $DelPoint(w) = LastNode(w)$, we remove the edge $([\overset{w}{\rightarrow}x]_w^R, y, [\overset{w}{\rightarrow}w]_w^R)$ such that $xy = w$. However, notice that this method might remove other suffixes of w from the CDAWG. More precise arguments follow.

Lemma 1. *Let $w \in \Sigma^+$, $x = DelPoint(w)$, and $z = LastNode(w)$. Assume $x = z$. Let s be any string in $[\overset{w}{\rightarrow}x]_w^R = [\overset{w}{\rightarrow}z]_w^R$. Then there uniquely exists a string $y \in \Sigma^+$ such that $sy \in Suffix(w)$.*

Proof. Since $x = DelPoint(w)$, there uniquely exists a character $a \in \Sigma$ such that $xa \in Factor(w)$ and $\overset{w}{\rightarrow}xa = w$. Let y be the string such that $xy = w$ with $y \in \Sigma^+$, where the first character of y is a . Let s be an arbitrary element in $[x]_w^R$. Since $x \in Prefix(w)$, $\overset{w}{\leftarrow}x = x$. Thus $s \in Suffix(x)$, which implies $sy \in Suffix(w)$. \square

For the case that $DelPoint(w) \neq LastNode(w)$, we have the following.

Lemma 2. Let $w \in \Sigma^+$, $x = DelPoint(w)$, and $z = LastNode(w)$. Assume $x \neq z$. Let s be any string in $[\overrightarrow{z}]_w^R$. Then there uniquely exist strings $t, u \in \Sigma^+$ such that $st \in Suffix(x)$ and $stu \in Suffix(w)$.

Proof. Since $\overrightarrow{z} = z$, $\overleftarrow{z} = z$. By the assumption that $z \neq x$, we have $z \in Prefix(x)$. Since $x = DelPoint(w)$, there uniquely exists a character $a \in \Sigma$ such that $\overleftarrow{z}a = x$. Thus there is a unique string $t \in \Sigma^+$ such that $zt = x$. Since $z \in Prefix(w)$, $z = \overleftarrow{z}$. Therefore, for any string $s \in [z]_w^R$ it holds that $st \in Suffix(x)$. Moreover, there uniquely exists a character $b \in \Sigma$ such that $\overrightarrow{x}b = w$. Let $u \in \Sigma^+$ be the string satisfying $\overrightarrow{x}b = xu$. Now we have $ztu = w$, and for any $s \in [z]_w^R$, it holds that $stu \in Suffix(w)$. \square

Lemma 3. Let $w \in \Sigma^+$, $x = DelPoint(w)$, and $z = LastNode(w)$. Assume $x \neq z$. Then $x = LRS(w)$.

Proof. Since $x \neq z$, $\overrightarrow{x} \neq \overleftarrow{x}$. Hence $\overleftarrow{x} = x \in Suffix(w)$ by Proposition 1. It is not difficult to show that x occurs in w just twice. Let $y = ax$ with $a \in \Sigma$, such that $y \in Suffix(w)$. Assume, for a contradiction, $y = LRS(w)$. On the assumption, y appears in w at least twice. If $y \notin Prefix(w)$, y must also occur in w as neither a prefix nor a suffix of w . It turns out that x appears three times in w : a contradiction. If $y \in Prefix(w)$, x is of the form a^ℓ . Then $y = DelPoint(w)$, which contradicts the assumption that $x = DelPoint(w)$. Consequently, $x = LRS(w)$. \square

According to the above three lemmas, we obtain the following theorem.

Theorem 3. Let $w \in \Sigma^+$, $x = DelPoint(w)$, and $z = LastNode(w)$. Let $k = |\overrightarrow{z}]_w^R|$. Suppose u_1, u_2, \dots, u_k be the suffixes of w arranged in decreasing order of their length, where $u_1 = w$.

1. When $x = z$: Let $xy = w$. Assume that the edge $([\overrightarrow{x}]_w^R, y, [\overrightarrow{w}]_w^R)$ is deleted from $CDAWG(w)$.
2. When $x \neq z$: Let $zt = x$ and $ztu = w$. Assume that the edge $([\overrightarrow{z}]_w^R, tu, [\overrightarrow{w}]_w^R)$ of $CDAWG(w)$ is shortened into the edge $([\overrightarrow{z}]_w^R, t, [\overrightarrow{x}]_w^R)$.

In both cases, the suffixes u_1, \dots, u_k are removed from the $CDAWG$.

What the above theorem implies is that after deleting or shortening the last edge in the backbone of $CDAWG(w)$, the leftmost position of a sliding window “skips” k characters at once. Let $DelSize(w) = k$. The next question is the exact upper bound of $DelSize(w)$. Fortunately, we achieve a reasonable result such that $DelSize(w)$ is at most about half of $|w|$. A more precise evaluation will be performed in Section 4.4.

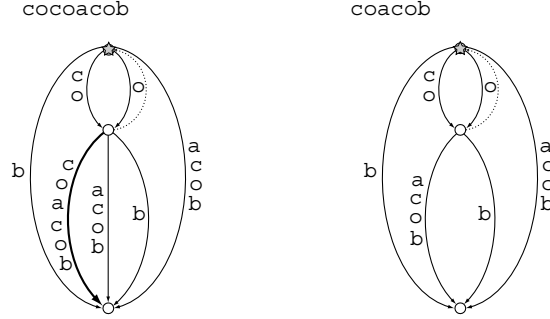


Fig. 2. On the left, $CDAWG(\text{cocoacob})$ is shown. The thick edge is to be deleted. The resulting structure is $CDAWG(\text{coacob})$, shown on the right. The gray star indicates the active point for each.

Can't we delete only the leftmost character of w in (amortized) constant time? We strongly believe the answer is “No”. The reason is as follows. Let $|w| = n$ where $w \in \Sigma^*$. Let u_1, u_2, \dots, u_{n+1} be all the suffixes of w arranged in decreasing order of their length. In [12], it has been proven that the total number of nodes necessary to keep $CDAWG(u_i)$ for every $1 \leq i \leq n + 1$ is $\Theta(n^2)$, even if we minimize the CDAWGs so to share as many nodes and edges as possible. Therefore, the amortized time complexity to delete the leftmost character of w would be proportional to n .

4.2 Maintaining the Structure of CDAWG

Suppose the last edge of the backbone of $CDAWG(w)$ is deleted or shortened right now. Let $k = DelSize(w)$. Let $u = w[k + 1 : n]$ where $n = |w|$. We sometimes need to modify the structure of the current graph, so that it exactly becomes $CDAWG(u)$. The *out-degree* of a node v is denoted by $OutDeg(v)$. Let $x = DelPoint(w)$ of $CDAWG(w)$.

Firstly, we consider when $OutDeg([\overrightarrow{x}]_w^R) \geq 3$ in the first cast of Theorem 3. In this case, $\overrightarrow{x} = x$ and $OutDeg([\overrightarrow{x}]_u^R) \geq 2$. It does not contradict Definition 2, and thus no more maintenance is required. An example of the case is shown in Fig. 2.

Secondly, we consider when $OutDeg([\overrightarrow{x}]_w^R) = 2$ in the first case of Theorem 3. Let $([\overrightarrow{r}]_w^R, s, [\overrightarrow{x}]_w^R)$ be an arbitrary in-coming edge of the node $[\overrightarrow{x}]_w^R$ in $CDAWG(w)$. Assume $\overrightarrow{r} = r$, that is, $rs \in Suffix(x)$. Let $([\overrightarrow{x}]_w^R, t, [\overrightarrow{w}]_w^R)$ be the edge which is to be the sole remaining out-going edge of the node $[\overrightarrow{x}]_w^R$ after the deletion. Notice that, however, $\overrightarrow{x} = u$. Thus the edge $([\overrightarrow{r}]_w^R, s, [\overrightarrow{x}]_w^R)$ is modified to $([\overrightarrow{r}]_u^R, st, [\overrightarrow{u}]_u^R)$. The total time required in the operations is pro-

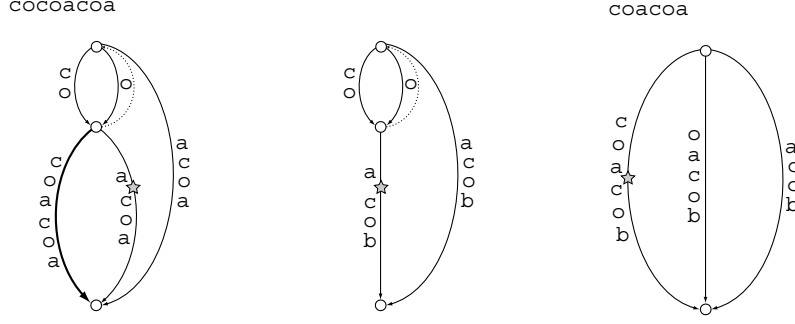


Fig. 3. On the left, $CDAWG(\text{cocoacoa})$ is shown, where the thick edge is to be deleted. The center is the intermediate structure in which the edge is deleted. After the modifications, we obtain $CDAWG(\text{coacoa})$, shown on the right. The gray star indicates the active point for each.

portional to the number of in-coming edges of the node $[\vec{x}]_w^R$ in $CDAWG(w)$. It is bounded by $DelSize(w)$.

Moreover, we might need a maintenance of the active point. Let $v = LRS(w)$. Supposing that $v \in Prefix(xt)$, v is represented on the edge $([\vec{x}]_w^R, t, [\vec{w}]_w^R)$ in $CDAWG(w)$. The active point is actually referred to as the pair $([\vec{x}]_w^R, p)$, where $p \in Prefix(t)$ and $xp = v$. The reference pair is modified to $([\vec{r}]_u^R, sp)$ in $CDAWG(u)$. Note that $\vec{r} \cdot sp = v$. An example of the case is shown in Fig. 3.

Thirdly, we consider the second case in Theorem 3. In this case, the last edge in the backbone, $([\vec{z}]_w^R, tu, [\vec{w}]_w^R)$, is shortened into $([\vec{z}]_u^R, t, [\vec{x}]_u^R) = ([\vec{z}]_u^R, t, [\vec{u}]_u^R)$ in $CDAWG(u)$. It implies that $x \neq LRS(u)$, although $x = LRS(w)$. The active point of $CDAWG(w)$ is represented by $([\vec{z}]_w^R, t)$, since $zt = x$ (by Lemma 3). Let $SufLink([\vec{z}]_w^R) = [\vec{s}]_w^R$. Assuming $\vec{s} = s$, s is the longest string such that $s \in Suffix(z)$ and $s \notin [\vec{z}]_w^R$. Notice that $LRS(u) = st$. Hereby, the reference pair of the active point is changed to $([\vec{s}]_u^R, t)$. If $[\vec{s}]_u^R$ is the explicit parent node nearest the locus of st , we are done. If not, the reference pair is canonized to $([\vec{r}]_u^R, p)$ such that $s \in Prefix(\vec{r})$, $st = \vec{r} \cdot p$, and $|p|$ is minimum. An example of the case is shown in Fig. 4.

4.3 Detecting $DelPoint(u)$

Suppose that after the edge deletion or shortening of $CDAWG(w)$, we got $CDAWG(u)$, where $u \in Suffix(w)$. The problem is how to locate $DelPoint(u)$ in $CDAWG(u)$. A naive solution is to traverse the backbone of $CDAWG(u)$ from

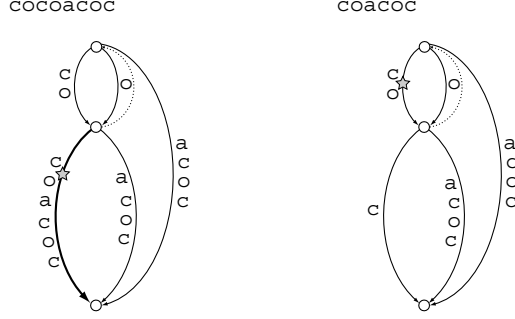


Fig. 4. On the left, $CDAWG(\text{cocoacoc})$ is shown. The thick edge is to be shortened. The resulting structure is $CDAWG(\text{coacoc})$, shown on the right. The gray star indicates the active point for each.

the source node. However, it takes $O(|u|)$ time, which leads to quadratic time complexity in total.

Our approach is to keep track of the *advanced point* that corresponds to the locus of $w[1 : n - 1]$, where $n = |w|$. Let $x = \text{LastNode}(w)$ and $xy = w$, that is, $([\overrightarrow{x}]_w^R, y, [\overrightarrow{w}]_w^R)$ is the edge for deletion or shortening. The canonical reference pair for the advanced point is $([\overrightarrow{x}]_w^R, t)$, where $t \in \text{Prefix}(y)$ and $\overrightarrow{x} \cdot t = w[1 : n - 1]$. We move to node $[\overrightarrow{s}]_w^R = \text{SufLink}([\overrightarrow{x}]_w^R)$. Suppose $CDAWG(w)$ has already been converted to $CDAWG(u)$. Assume $[\overrightarrow{s}]_w^R = [\overrightarrow{s}]_u^R$. Since $\overrightarrow{s} \cdot t = u[1 : m - 1]$ where $m = |u|$, $([\overrightarrow{s}]_u^R, t)$ is a reference pair of the next advanced point, and then it is canonized. Let $([\overrightarrow{s}']_u^R, t')$ be the canonical reference pair of the advanced point. Then $s' = \text{LastNode}(u)$. If $\text{LastNode}(u) \notin \text{Prefix}(LRS(u))$, that is, if the active point is not on the longest out-going edge from the node $[\overrightarrow{s}']_u^R$, $\text{DelPoint}(u) = \text{LastNode}(u)$. Otherwise, $\text{DelPoint}(u) = LRS(u)$. In the case that $[\overrightarrow{s}']_u^R \neq [\overrightarrow{s}]_w^R$, we perform the same procedure from its closest parent node (see Fig. 3). If Σ is fixed, the cost of canonizing the reference pair is only proportional to the number of nodes included in the path. The amortized number of such nodes is constant.

4.4 On Buffer Size

The main theorem of this section shows an exact estimation of the upperbound of $\text{DelSize}(w)$. For an alphabet Σ and an integer n , we define $\text{MaxDel}_\Sigma(n) = \max\{\text{DelSize}(w) \mid w \in \Sigma^*, |w| = n\}$.

Theorem 4. *If $|\Sigma| \geq 3$, $\text{MaxDel}_\Sigma(n) = \lceil \frac{n}{2} \rceil - 1$.*

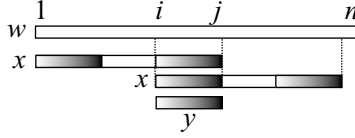


Fig. 5. The case $j > \frac{n}{2}$. x occurs at least twice in w , and the overlap y is in fact both a prefix and a suffix of x .

By this theorem, edge deletion or edge shortening can shrink the window size upto the half of the original size. Therefore, in order to keep the window size at least M , a buffer of size $2M + 1$ is necessary and sufficient.

We will prove the above theorem in the sequel.

Lemma 4. *Let $w \in \Sigma^*$. For any string $x \in \text{Factor}(w)$, let $\text{SufLink}(\overrightarrow{x}_w^R) = \overrightarrow{s}_w^R$. Then $|\overrightarrow{x}_w^R| = |\overrightarrow{x}| - |\overrightarrow{s}|$.*

Proof. $|\overrightarrow{x}_w^R| = |\text{Suffix}(\overrightarrow{x})| - |\text{Suffix}(\overrightarrow{s})| = (|\overrightarrow{x}| + 1) - (|\overrightarrow{s}| + 1) = |\overrightarrow{x}| - |\overrightarrow{s}|$. \square

Lemma 5. *Let $w \in \Sigma^*$ and $n = |w|$. For any $x \in \text{Prefix}(w) - \{w\}$ with $\overrightarrow{x} = x$, $|\overrightarrow{x}_w^R| \leq \lceil \frac{n}{2} \rceil - 1$.*

Proof. Let $j = |x| = |\overrightarrow{x}|$. Let $\text{SufLink}(\overrightarrow{x}_w^R) = \overrightarrow{s}_w^R$. We have the following three cases.

- (1) When $j < \frac{n}{2}$: Since j is an integer, $j \leq \lceil \frac{n}{2} \rceil - 1$, and $|\overrightarrow{x}_w^R| = |\overrightarrow{x}| - |\overrightarrow{s}| \leq \lceil \frac{n}{2} \rceil - 1$ by Lemma 4.
- (2) When $j > \frac{n}{2}$: (See Fig. 5.) The equivalences $x = w[1 : j]$ and $\overrightarrow{x} = x$ imply that $x = w[i : i + j - 1]$ for some $i \geq 2$ and $i + j - 1 \leq n$. Then $i - j \leq n - 2j + 1 < 1$, that is, $i \leq j$. Let $y = w[i : j]$. Its length is $|y| = j - i + 1 \geq 1$, and $y = x[i : j] \in \text{Suffix}(x)$. Since $\overrightarrow{x} = x$ and $y \in \text{Suffix}(x)$, $\overrightarrow{y} = y$. On the other hand, $y = w[i : j] = x[1 : j - i + 1] = w[1 : j - i + 1] \in \text{Prefix}(w)$, which implies $\overleftarrow{y} = y$. Thus y is the longest element of \overrightarrow{y}_w^R . Since $|x| > |y|$, $x \notin \overrightarrow{y}_w^R$. Therefore $|\overrightarrow{s}| \geq |y|$, which yields $|\overrightarrow{x}_w^R| = |\overrightarrow{x}| - |\overrightarrow{s}| \leq |x| - |y| = j - (j - i + 1) = i - 1 \leq n - j < n - \frac{n}{2} = \frac{n}{2}$. Thus $|\overrightarrow{x}_w^R| \leq \lceil \frac{n}{2} \rceil - 1$.
- (3) When $j = \frac{n}{2}$: Since $\overrightarrow{x} = x$, x occurs in w at least twice. If $x = w[i : i + j - 1]$ for some i with $2 \leq i \leq j$, we can show the inequality holds in the same way

as (2). Otherwise, $x = w[j + 1 : 2j] = w[j + 1 : n] = w[1 : j]$, that is $w = xx$. Then $\overrightarrow{x}^w = w \neq x$, which does not satisfy the precondition of the lemma.

In any cases, we have got the result. \square

We are ready to prove the upperbound of $MaxDel_\Sigma(n)$.

Lemma 6. $MaxDel_\Sigma(n) \leq \lceil \frac{n}{2} \rceil - 1$ for any Σ and any $n \geq 3$.

Proof. Let $x = DelPoint(w)$ and $z = LastNode(w)$. First we consider the case $x = z$. Since $\overrightarrow{x}^w = x$ and $x \in Prefix(w) - \{w\}$, $DelSize(w) = |\lceil \overrightarrow{z}^w \rceil_w^R| = |\lceil \overrightarrow{x}^w \rceil_w^R| \leq \lceil \frac{n}{2} \rceil - 1$ by Lemma 5.

We now assume $x \neq z$. Then $z \in Prefix(x)$, and $x = DelPoint(w)$ implies that $x \in Prefix(w) - \{w\}$, which yields $z \in Prefix(w) - \{w\}$. Thus by Lemma 5, $DelSize(w) = |\lceil \overrightarrow{z}^w \rceil_w^R| \leq \lceil \frac{n}{2} \rceil - 1$. \square

On the other hand, the lowerbound is given by the following lemma.

Lemma 7. If $|\Sigma| \geq 3$, $MaxDel_\Sigma(n) \geq \lceil \frac{n}{2} \rceil - 1$ for any $n \geq 1$.

Proof. For each $1 \leq n \leq 4$, the inequality trivially holds since $\lceil \frac{n}{2} \rceil - 1 \leq 1$. Let $\mathbf{a}, \mathbf{b}, \mathbf{c}$ be distinct symbols in Σ . For each odd $n \geq 5$, let $w_n = \mathbf{a}^k \mathbf{b} \mathbf{a}^k \mathbf{b} \mathbf{c}$, where $k = \frac{n-3}{2}$. Remark that $DelPoint(w_n) = \mathbf{a}^k \mathbf{b}$. Let $x = \mathbf{a}^k \mathbf{b}$. We can see that any suffix of x except ε belongs to $\lceil \overrightarrow{x}^w \rceil_w^R$, so that $SufLink(x) = \varepsilon$. Thus $DelSize(w_n) = |x| - |\varepsilon| = |\mathbf{a}^k \mathbf{b}| - 0 = k + 1 = \frac{n-1}{2} = \lceil \frac{n}{2} \rceil - 1$, since n is odd. For each even $n \geq 6$, let $w'_n = \mathbf{a}^{k-1} \mathbf{b} \mathbf{a}^k \mathbf{b} \mathbf{c}$, where $k = \frac{n}{2} - 1$, and we can verify that $DelSize(w'_n) = \lceil \frac{n}{2} \rceil - 1$ similarly. \square

Consequently, Theorem 4 is proved by Lemma 6 and Lemma 7. We note that for a binary alphabet $\Sigma = \{\mathbf{a}, \mathbf{b}\}$, two series of the strings $\mathbf{a}^k \mathbf{b} \mathbf{a}^k \mathbf{b} \mathbf{a} \mathbf{b}$ and $\mathbf{a}^{k-1} \mathbf{b} \mathbf{a}^k \mathbf{b} \mathbf{a} \mathbf{b}$ give the lowerbound $MaxDel_\Sigma(n) \geq \lceil \frac{n-3}{2} \rceil$.

On the on-line algorithm of [11], each node $\lceil \overrightarrow{x}^w \rceil_w^R$ of $CDAWG(w)$ stores the value of $|\lceil \overrightarrow{x}^w \rceil_w^R|$. By Lemma 4, it is guaranteed that we can calculate $DelSize(w)$ in constant time with no additional information.

4.5 Keeping Edge Labels Valid

As mentioned previously, an edge label is actually represented by a pair of integers indicating its beginning and ending positions in input string w , respectively. We must ensure that no edge label becomes “out of date” after the window slides, e.g., that no integer refers to a position outside the sliding window. In case of a suffix tree, when a new edge is created, we can guarantee the above regulation by traversing from the leaf node toward the root node while updating all edge labels encountered. However, this would yield quadratic time complexity in the aggregate. Larsson [13, 14] utilized *credit issuing*, an update-number-restriction

technique, originally proposed in [7], which takes in total $O(|w|)$ time and space. In the following, we introduce an extended credit issuing technique for CDAWGs. Our basic strategy is to show that we can handle the credit issuing as well as in case of suffix trees.

We assign each internal node s of $CDAWG(w)$ a binary counter called *credit*, denoted by $Cred(s)$. This credit counter is initially set to zero when s is created. When a node s receives a credit, we update the labels of in-coming edges of s . Then, if $Cred(s) = 0$, we set it to one, and stop. If $Cred(s) = 1$, after setting it to zero, we let the node s issue a credit to its parent nodes.

When s is newly created, $Cred(s) = 0$. The creation of the new node s implies that a new edge is to be inserted from s to the sink node. When the new edge is created leading to the sink node, the sink node *issues* a credit to the parent node s . Assume the new edge is labeled by pair (i, j) where i, j are some integers with $i < j$. Let ℓ be the length of the label of the in-coming edge of s . After s received a credit from the sink node, we reset its in-coming edge label to $(i - \ell, i - 1)$. Remember the edge redirection happening in the construction of a CDAWG (see Section 2.2 or [11]). If some edge is actually redirected to node s , its label is updated as well. Note that we need not change the value of $Cred(s)$ again.

Suppose a node r has right now received a credit from one of its child nodes. Assume $Cred(r)$ is currently one. We need to update all in-coming edge labels of r . We store a list in r to maintain its in-coming edges arranged in the order of the length of the path they correspond to. The maintenance of the list is an easy matter, since the on-line algorithm of [11] inserts edges to r in such order. Let t be an arbitrary parent node of r . Let k be the number of the in-coming edges of r connected from t . One might wonder that r must issue k credits to t , but there is the following time-efficient method. In case k is even, $Cred(t)$ need not to be changed because it is a *binary* counter. Contrarily, in case k is odd, we always change the value of $Cred(t)$. If $Cred(t)$ was one, we also have to update the in-coming edge of t . To do it, we focus on the *shortest* in-coming edge of r connected from t , which is in turn the shortest out-going edge of t leading to r . In updating the in-coming edges of t , we should utilize the label of the shortest edge, since the label corresponds to the possibly newest occurrence of the factors represented in node t . We continue updating edge labels by traversing the reversed graph rooted at r in width-first manner while issuing credits.

Recall the node separation in constructing a CDAWG (see Section 2.2 or [11]). Assume a node r has right now been created owing to the separation of a node s . The subgraph rooted at r is currently the same as the one rooted at s , since r was created as a clone of s . Thus we simply set $Cred(r) = Cred(s)$.

Now consider a node u to be deleted, corresponding to the second case of Section 4.2. It might have received a credit from its newest child node (that is not deleted), which has not been issued to its parent node yet. Therefore, when a node u is scheduled for deletion and $Cred(u) = 1$, node u issues credits to its parent nodes. However, this complicates the update of edge labels: several waiting credits may aggregate, causing nodes upper in the CDAWG to receive a credit *older* than the one it has already received from its another child node.

Therefore, before updating an edge label, we compare its previous value against the one associated with the received credit, and refer to the newer one. As well as the case of edge insertion mentioned in the above paragraph, we traverse the reversed graph rooted at u in width-first fashion to update edge labels. In the worst case, the updating cost is proportional to the number of paths from the source node to node u . Nevertheless, it is bounded by $DelSize(w)$.

By analogous arguments to [7, 13, 14], we can establish the following lemma.

Lemma 8. *All edge labels of a CDAWG can be kept valid in a sliding window, in linear time and space with respect to the length of an input string.*

As a conclusion of Section 4, we finally obtain the following.

Theorem 5. *Let $w \in \Sigma^*$ and M be the window size. The proposed algorithm runs in $O(|w|)$ time using $O(M)$ space.*

5 Conclusion

We introduced an algorithm to maintain CDAWGs for a sliding window, which runs in linear time and space. It can be an alternative of Larsson’s suffix tree algorithm in [13]. Moreover, CDAWGs are known to be more space-economical than suffix trees [2, 6], and thus our algorithm seems to contribute to reducing the space requirement in PPM*-style data compression scheme. We are currently implementing our algorithm to experimentally evaluate its efficiency in both time and space.

Acknowledgment

We wish to thank Ricardo Garcia and Veli Mäkinen who suggested us to apply our previous work [11] to the sliding window mechanism. We also appreciate Ricardo Garcia’s fruitful comments and suggestions in the early stage of this work.

References

1. A. Apostolico. The myriad virtues of subword trees. In A. Apostolico and Z. Galil, editors, *Combinatorial Algorithm on Words*, volume 12 of *NATO Advanced Science Institutes, Series F*, pages 85–96. Springer-Verlag, 1985.
2. 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.
3. 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.
4. J. G. Cleary and I. H. Witten. Data compression using adaptive coding and partial string matching. *IEEE Trans. Commun.*, 32(4):396–402, 1984.

5. M. Crochemore and W. Rytter. *Text Algorithms*. Oxford University Press, New York, 1994.
6. M. Crochemore and R. Verin. On compact directed acyclic word graphs. In J. Mycielski, G. Rozenberg, and A. Salomaa, editors, *Structures in Logic and Computer Science*, volume 1261 of *Lecture Notes in Computer Science*, pages 192–211. Springer-Verlag, 1997.
7. E. R. Fiala and D. H. Greene. Data compression with finite windows. *Commun. ACM*, 32(4):490–505, 1989.
8. D. Gusfield. *Algorithms on Strings, Trees, and Sequences*. Cambridge University Press, New York, 1997.
9. 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)*. Czech Technical University, 2001.
10. S. Inenaga, H. Hoshino, A. Shinohara, M. Takeda, and S. Arikawa. On-line construction of symmetric compact directed acyclic word graphs. In *Proc. 8th International Symposium on String Processing and Information Retrieval (SPIRE'01)*, pages 96–110. IEEE Computer Society, 2001.
11. 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 A. Amir and G. M. Landau, editors, *Proc. 12th Annual Symposium on Combinatorial Pattern Matching (CPM'01)*, volume 2089 of *Lecture Notes in Computer Science*, pages 169–180. Springer-Verlag, 2001.
12. S. Inenaga, A. Shinohara, M. Takeda, H. Bannai, and S. Arikawa. Space-economical construction of index structures for all suffixes of a string. In *Proc. 27th International Symposium on Mathematical Foundations of Computer Science (MFCS'02)*, Lecture Notes in Computer Science. Springer-Verlag, 2002. (to appear).
13. N. J. Larsson. Extended application of suffix trees to data compression. In *Proc. Data Compression Conference '96 (DCC'96)*, pages 190–199. IEEE Computer Society, 1996.
14. N. J. Larsson. *Structures of String Matching and Data Compression*. PhD thesis, Lund University, 1999.
15. E. M. McCreight. A space-economical suffix tree construction algorithm. *J. ACM*, 23(2):262–272, 1976.
16. A. Moffat. Implementing the PPM data compression scheme. *IEEE Trans. Commun.*, 38(11):1917–1921, 1990.
17. E. Ukkonen. On-line construction of suffix trees. *Algorithmica*, 14(3):249–260, 1995.
18. P. Weiner. Linear pattern matching algorithms. In *Proc. 14th Annual Symposium on Switching and Automata Theory*, pages 1–11, 1973.