

# Parameterized DAWGs: efficient constructions and bidirectional pattern searches

Katsuhito Nakashima<sup>1</sup>, Noriki Fujisato<sup>1</sup>, Diptarama Hendrian<sup>1</sup>, Yuto Nakashima<sup>2</sup>,  
Ryo Yoshinaka<sup>1</sup>, Shunsuke Inenaga<sup>2,3</sup>, Hideo Bannai<sup>4</sup>, Ayumi Shinohara<sup>1</sup>, and  
Masayuki Takeda<sup>2</sup>

<sup>1</sup>Graduate School of Information Sciences, Tohoku University, Japan

<sup>2</sup>Department of Informatics, Kyushu University, Japan

<sup>3</sup>PRESTO, Japan Science and Technology Agency, Japan

<sup>4</sup>M&D Data Science Center, Tokyo Medical and Dental University, Japan

## Abstract

Two strings  $x$  and  $y$  over  $\Sigma \cup \Pi$  of equal length are said to *parameterized match* ( $p$ -match) if there is a renaming bijection  $f : \Sigma \cup \Pi \rightarrow \Sigma \cup \Pi$  that is identity on  $\Sigma$  and transforms  $x$  to  $y$  (or vice versa). The  $p$ -matching problem is to look for substrings in a text that  $p$ -match a given pattern. In this paper, we propose *parameterized suffix automata* ( $p$ -suffix automata) and *parameterized directed acyclic word graphs* (PDAWG) which are the  $p$ -matching versions of suffix automata and DAWGs. While suffix automata and DAWGs are equivalent for standard strings, we show that  $p$ -suffix automata can have  $\Theta(n^2)$  nodes and edges but PDAWGs have only  $O(n)$  nodes and edges, where  $n$  is the length of an input string. We also give an  $O(n|\Pi| \log(|\Pi| + |\Sigma|))$ -time  $O(n)$ -space algorithm that builds the PDAWG in a left-to-right online manner. As a byproduct, it is shown that the *parameterized suffix tree* for the reversed string can also be built in the same time and space, in a right-to-left online manner. This duality also leads us to two further efficient algorithms for  $p$ -matching: Given the parameterized suffix tree for the reversal  $\bar{T}$  of the input string  $T$ , one can build the PDAWG of  $T$  in  $O(n)$  time in an offline manner; One can perform *bidirectional*  $p$ -matching in  $O(m \log(|\Pi| + |\Sigma|) + occ)$  time using  $O(n)$  space, where  $m$  denotes the pattern length and  $occ$  is the number of pattern occurrences in the text  $T$ .

## 1 Introduction

The *parameterized matching problem* ( $p$ -matching problem) [1] is a class of pattern matching where the task is to locate substrings of a text that have “the same structure” as a given pattern. More formally, we consider a parameterized string ( $p$ -string) over a union of two disjoint alphabets  $\Sigma$  and  $\Pi$  for static characters and for parameter characters, respectively. Two equal length  $p$ -strings  $x$  and  $y$  are said to *parameterized match* ( $p$ -match) if  $x$  can be transformed to  $y$  (and vice versa) by a bijection which renames the parameter characters. The  $p$ -matching problem is, given a text  $p$ -string  $T$  and pattern  $p$ -string  $P$ , to report the occurrences of substrings of  $T$  that  $p$ -match  $P$ .  $P$ -matching is well-motivated by plagiarism detection, software maintenance, RNA structural pattern matching, and so on [1, 2, 3, 4].

The *parameterized suffix tree* ( $p$ -suffix tree) [5] is the fundamental indexing structure for  $p$ -matching, which supports  $p$ -matching queries in  $O(m \log(|\Pi| + |\Sigma|) + occ)$  time, where  $m$  is the length of pattern  $P$ , and  $occ$  is the number of occurrences to report. It is known that the  $p$ -suffix tree of a text  $w$  of length  $n$  can be built in  $O(n \log(|\Pi| + |\Sigma|))$  time with  $O(n)$  space in an offline manner [6] and in a *left-to-right online* manner [2]. A *randomized*  $O(n)$ -time left-to-right online construction algorithm for  $p$ -suffix trees is also known [7]. Indexing  $p$ -strings has recently attracted much attention, and the

p-matching versions of other indexing structures, such as *parameterized suffix arrays* [8, 9, 10, 11], *parameterized BWTs* [12], and *parameterized position heaps* [13, 14, 15], have also been proposed.

This paper fills in the missing pieces of indexing structures for p-matching, by proposing the parameterized version of the *directed acyclic word graphs* (DAWGs) [16, 17], which we call the *parameterized directed acyclic word graphs* (PDAWGs).

For any standard string  $T$ , the following three data structures are known to be equivalent:

- (1) The *suffix automaton* of  $T$ , which is the minimum DFA that is obtained by merging isomorphic subtrees of the suffix trie of  $T$ .
- (2) The DAWG, which is the edge-labeled DAG of which each node corresponds to an equivalence class of substrings of  $T$  defined by the set of ending positions in  $T$ .
- (3) The *Weiner-link graph*, which is the DAG consisting of the nodes of the suffix tree of the reversal  $\bar{T}$  of  $T$  and the reversed suffix links (a.k.a. soft and hard Weiner links).

The equality of (2) and (3) in turn implies symmetry of suffix trees and DAWGs, namely:

- (a) The suffix links of the DAWG for  $T$  form the suffix tree for  $\bar{T}$ .
- (b) Left-to-right online construction of the DAWG for  $T$  is equivalent to right-to-left online construction of the suffix tree for  $\bar{T}$ .

Firstly, we present (somewhat surprising) combinatorial results on the p-matching versions of data structures (1) and (2). We show that the *parameterized suffix automaton* (*p-suffix automaton*), which is obtained by merging isomorphic subtrees of the *parameterized suffix trie* of a p-string  $T$  of length  $n$ , can have  $\Theta(n^2)$  nodes and edges in the worst case, while the PDAWG for any p-string has  $O(n)$  nodes and edges. On the other hand, the p-matching versions of data structures (2) and (3) are equivalent: After introducing the *parameterized Weiner links* on p-suffix trees, we show that the parameterized Weiner-link graph of the p-suffix tree for  $\bar{T}$  is equivalent to the PDAWG for  $T$ . As a corollary to this, symmetry (a) also holds: The suffix links of the PDAWG for  $T$  form the p-suffix tree for  $\bar{T}$ .

Secondly, we present algorithmic results on PDAWG construction. We propose left-to-right *online* construction of PDAWGs that works in  $O(n|\Pi| \log(|\Pi| + |\Sigma|))$  time with  $O(n)$  working space. In addition, as a byproduct of this algorithm, we obtain a right-to-left online construction of the p-suffix tree in  $O(n|\Pi| \log(|\Pi| + |\Sigma|))$  time with  $O(n)$  space. This can be seen as the p-matching version of symmetry (b). The complexities for our online algorithms are valid in the pointer machine, which is strictly weaker than the word RAM.

Thirdly, we propose an alternative *offline* algorithm which builds the PDAWG in  $O(n)$  time with  $O(n)$  working space, provided that the p-suffix tree of the reversal of the input string is given. While the proposed offline algorithm itself works in the pointer machine, there are two different complexities for p-suffix tree construction in the pointer machine and in the word RAM: The p-suffix tree can be built offline, in  $O(n \log(|\Pi| + |\Sigma|))$  time and  $O(n)$  space in the pointer machine [1] and  $O(n|\Pi|)$  time and  $O(n)$  space in the word RAM with word size  $\Omega(\log n)$  [11]. Putting these together, we obtain  $O(n \min\{\log(|\Pi| + |\Sigma|), |\Pi|\})$ -time  $O(n)$ -space offline construction of the PDAWG in the word RAM.

We also show that, using the PDAWG for a text string  $T$  and the p-suffix tree for its reversal  $\bar{T}$ , one can perform *bidirectional* p-matching, i.e., the pattern may grow in both forward and backward directions, in  $O(m \log(|\Pi| + |\Sigma|) + occ)$  time with  $O(n)$  space, where  $m$  denotes the pattern length and  $occ$  is the number of pattern occurrences in the text  $T$ . To our knowledge, this is the first index that allows bidirectional p-matching within linear space.

This paper is organized as follows. After defining basic mathematical notions used in this paper, we will briefly review existing indexing data structures for parameterized strings in Section 2. We propose three different data structures that can be thought to be the parameterized counterpart of DAWGs in Sections 3 to 5: parameterized suffix automata, pseudo-PDAWGs, and PDAWGs. The pseudo-PDAWGs are an intermediate data structure which makes the exposition of our pattern

matching algorithm with PDAWG easier to follow. Section 6 discusses the duality between PDAWGs and parameterized suffix trees, together with bidirectional pattern matching and offline construction algorithms. Section 7 presents an algorithm for constructing PDAWGs online. We conclude the paper in Section 8.

A preliminary version of this paper appeared in [18]. New materials given in this full version are complete proofs of the lemmas and theorems, introduction of pseudo-PDAWGs data structure, the bidirectional parameterized pattern matching algorithm, and more examples and figures.

## 2 Preliminaries

We first introduce definitions and notation used in this paper and then briefly review indexing data structures of parameterized strings.

### 2.1 Definitions and notation

We denote the set of strings over an alphabet  $A$  by  $A^*$ . For a string  $w = xyz \in A^*$ ,  $x$ ,  $y$ , and  $z$  are called *prefix*, *factor*, and *suffix* of  $w$ , respectively. The sets of the prefixes, factors, and suffixes of a string  $w$  are denoted by  $\text{Prefix}(w)$ ,  $\text{Factor}(w)$ , and  $\text{Suffix}(w)$ , respectively. The length of  $w$  is denoted by  $|w|$  and the  $i$ -th character of  $w$  is denoted by  $w[i]$  for  $1 \leq i \leq |w|$ . The factor of  $w$  that begins at position  $i$  and ends at position  $j$  is  $w[i : j]$  for  $1 \leq i \leq j \leq |w|$ . For convenience, we abbreviate  $w[1 : i]$  to  $w[:i]$  and  $w[i : |w|]$  to  $w[i:]$  for  $1 \leq i \leq |w|$ . The empty string is denoted by  $\varepsilon$ , that is  $|\varepsilon| = 0$ . Moreover, let  $w[i : j] = \varepsilon$  if  $i > j$ . The *reverse*  $\bar{w}$  of  $w \in A^*$  is inductively defined by  $\bar{\varepsilon} = \varepsilon$  and  $\overline{ax} = a\bar{x}$  for  $a \in A$  and  $x \in A^*$ .

Throughout this paper, we fix two disjoint ordered alphabets  $\Sigma$  and  $\Pi$ . We call elements of  $\Sigma$  *static* characters and those of  $\Pi$  *parameter* characters. Elements of  $\Sigma^*$  and  $(\Sigma \cup \Pi)^*$  are called *static strings* and *parameterized strings* (or *p-strings* for short), respectively.

Given two p-strings  $S_1$  and  $S_2$  of length  $n$ ,  $S_1$  and  $S_2$  are a *parameterized match* (*p-match*), denoted by  $S_1 \approx S_2$ , if there is a bijection  $f$  on  $\Sigma \cup \Pi$  such that  $f(a) = a$  for any  $a \in \Sigma$  and  $f(S_1[i]) = S_2[i]$  for all  $1 \leq i \leq n$  [1]. We use Kim and Cho's version of p-string encoding [19], which replaces 0 in Baker's prev-encoding [5] by  $\infty$ . Let  $\mathcal{N} = \mathbb{N} \cup \{\infty\} \setminus \{0\}$  where  $\mathbb{N}$  is the set of non-negative integers.<sup>1</sup> The *prev-encoding*  $\langle S \rangle$  of a p-string  $S$  is the string over  $\Sigma \cup \mathcal{N}$  of length  $|S|$  defined by

$$\langle S \rangle[i] = \begin{cases} S[i] & \text{if } S[i] \in \Sigma, \\ \infty & \text{if } S[i] \in \Pi \text{ and } S[i] \neq S[j] \text{ for } 1 \leq j < i, \\ i - j & \text{if } S[i] = S[j] \in \Pi, j < i \text{ and } S[i] \neq S[k] \text{ for any } j < k < i \end{cases}$$

for  $i \in \{1, \dots, |S|\}$ . We call a string  $x \in (\Sigma \cup \mathcal{N})^*$  a *pv-string* if  $x = \langle S \rangle$  for some p-string  $S$ . For any p-strings  $S_1$  and  $S_2$ ,  $S_1 \approx S_2$  if and only if  $\langle S_1 \rangle = \langle S_2 \rangle$  [1]. For example, given  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$  and  $\Pi = \{\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}\}$ ,  $S_1 = \mathbf{uvvauvb}$  and  $S_2 = \mathbf{xyyaxyb}$  are a p-match by  $f$  such that  $f(\mathbf{u}) = \mathbf{x}$  and  $f(\mathbf{v}) = \mathbf{y}$ , where  $\langle S_1 \rangle = \langle S_2 \rangle = \infty\infty\mathbf{1a43b}$ . For a p-string  $T$ , let  $\text{PFactor}(T) = \{\langle S \rangle \mid S \in \text{Factor}(T)\}$  and  $\text{PSuffix}(T) = \{\langle S \rangle \mid S \in \text{Suffix}(T)\}$  be the sets of prev-encoded factors (pv-factors) and suffixes (pv-suffixes) of  $T$ , respectively.

Let  $S$  be a string of length  $n$  over  $\Sigma \cup \Pi$ , and  $\pi$  denote the number of distinct parameter characters in  $S$ . The prev-encoding  $\langle S \rangle$  for string  $S$  can be computed in  $O(n \log \pi) \subseteq O(n \log |\Pi|)$  time using  $O(\pi) \subseteq O(n)$  space in an online fashion in the comparison model, or in  $O(n)$  time and space in the word RAM model when  $\Pi$  is an integer alphabet of polynomial size in  $n$ .

Not every factor of a pv-string is a pv-string. For example, for  $S_1 = \mathbf{uvvauvb}$  with  $\langle S_1 \rangle = \infty\infty\mathbf{1a43b}$ ,  $\langle S_1 \rangle[3 : 7] = \mathbf{1a43b}$  is not a pv-string, but we would like to obtain  $\langle S_1[3 : 7] \rangle = \langle \mathbf{vauvb} \rangle = \infty\mathbf{a}\infty\mathbf{3b}$

<sup>1</sup>Without loss of generality, we can assume  $\Sigma \cap \mathcal{N} = \Pi \cap \mathcal{N} = \emptyset$ . In the case where  $\Sigma$  and/or  $\Pi$  are integer alphabets, then for instance we can work on the modified alphabets  $\Sigma' = \{(0, c) \mid c \in \Sigma\}$ ,  $\Pi' = \{(1, x) \mid x \in \Pi\}$ , and  $\mathcal{N}' = \{(2, n) \mid n \in \mathcal{N}\}$ .

directly from  $\infty 1a43b$  without “decoding” the prev-encoding. For this end, we extend the notation  $\langle \cdot \rangle$  for factors  $x \in (\Sigma \cup \mathcal{N})^*$  of pv-strings. The *re-encoding*  $\langle x \rangle$  of  $x$  is the string of length  $|x|$  defined by

$$\langle x \rangle[i] = \begin{cases} \infty & \text{if } x[i] \in \mathcal{N} \text{ and } x[i] \geq i, \\ x[i] & \text{otherwise.} \end{cases}$$

We then have  $\langle \langle S \rangle[i : j] \rangle = \langle S[i : j] \rangle$  for any  $i, j$  and  $S \in (\Sigma \cup \Pi)^*$ . For example,  $\langle \langle S_1 \rangle[3 : 7] \rangle = \langle 1a43b \rangle = \infty a \infty 3b = \langle S_1[3 : 7] \rangle$  for  $S_1 = uvvauvb$ . We apply PFactor etc. to pv-strings  $w$  so that  $\text{PFactor}(w) = \{ \langle x \rangle \mid x \in \text{Factor}(w) \}$ . Here, we introduce an alternative definition of the re-encoding using the following notation:

$$\langle \langle a \rangle \rangle_i = \begin{cases} \infty & \text{if } a \in \mathcal{N} \text{ and } a > i, \\ a & \text{otherwise,} \end{cases}$$

for  $a \in \Sigma \cup \mathcal{N}$  and  $i \in \mathbb{N}$ . Then, the re-encoding can be defined inductively by  $\langle \varepsilon \rangle = \varepsilon$  and  $\langle xa \rangle = \langle x \rangle \langle \langle a \rangle \rangle_{|x|}$  for  $x \in (\Sigma \cup \mathcal{N})^*$  and  $a \in \Sigma \cup \mathcal{N}$ .

Let  $w, x, y \in (\Sigma \cup \mathcal{N})^*$ . The set of the *end positions* or *p-occurrences* of  $x$  in a pv-string  $w$  is defined by  $\text{RPos}_w(x) = \{ i \in \{0, \dots, |w|\} \mid x = \langle w[i - |x| + 1 : i] \rangle \}$ . We say  $x$  *occurs* in  $w$  at  $i$  if  $i \in \text{RPos}_w(x)$ . Note that  $0 \in \text{RPos}_w(x)$  iff  $x = \varepsilon$ . We write  $x \equiv_w^R y$  iff  $\text{RPos}_w(x) = \text{RPos}_w(y)$  and the equivalence class of  $x$  under  $\equiv_w^R$  as  $[x]_w^R$ . Note that for any  $x \notin \text{PFactor}(w)$ , which can be a non-pv-string,  $\text{RPos}_w(x) = \emptyset$ . For an equivalent class  $u = [x]_w^R$ , we may write  $\text{RPos}_w(u)$  to mean  $\text{RPos}_w(x)$  for  $x \in u$ . For a finite nonempty set  $X$  of strings which has no distinct elements of equal length, the shortest and longest elements of  $X$  are denoted by  $\lfloor X \rfloor$  and  $\lceil X \rceil$ , respectively.

**Example 1.** Let  $w = \langle xaxaya \rangle = \infty a 2a \infty a$  where  $\Sigma = \{a\}$  and  $\Pi = \{x, y\}$ . Then,  $\text{RPos}_w(a) = \text{RPos}_w(\infty a) = \{2, 4, 6\}$  and  $[a]_w^R = \{a, \infty a\}$ . On the other hand,  $\text{RPos}_w(2a) = \emptyset$  and  $[2a]_w^R = (\Sigma \cup \mathcal{N})^* \setminus \text{PFactor}(w)$ .

The *parameterized pattern matching problem* is to enumerate all the end positions of  $\langle P \rangle$  in  $\langle T \rangle$ , i.e., elements of  $\text{RPos}_{\langle T \rangle}(\langle P \rangle)$ , for given two p-strings  $T$  and  $P$ . The weaker version of the problem is to decide whether  $\text{RPos}_{\langle T \rangle}(\langle P \rangle) \neq \emptyset$ .

The basic properties on end positions of factors in static strings presented in [16] also hold for pv-strings.

**Lemma 1.** *Let  $x, y$  and  $w$  be pv-strings. If  $\text{RPos}_w(x) \cap \text{RPos}_w(y) \neq \emptyset$ , then either  $x \in \text{PSuffix}(y)$  or  $y \in \text{PSuffix}(x)$ . If  $x \in \text{PSuffix}(y)$ , then  $\text{RPos}_w(y) \subseteq \text{RPos}_w(x)$ . For any  $x \in \text{PFactor}(w)$ , there is  $k \in \mathbb{N}$  such that  $[x]_w^R = \{ \langle y[i : ] \rangle \mid 1 \leq i \leq k \}$  where  $y = \lceil [x]_w^R \rceil$ .*

We will use the above lemma implicitly in arguments in this paper.

We define notions symmetric to  $\text{RPos}$ ,  $\equiv^R$ , and  $[\cdot]_w^R$ . The start position set of  $x$  in a pv-string  $w$  is  $\text{LPos}_w(x) = \{ i \in \{0, \dots, |w|\} \mid \langle w[i : i + |x| - 1] \rangle = x \}$ . We write  $x \equiv_w^L y$  iff  $\text{LPos}_w(x) = \text{LPos}_w(y)$ . The equivalence class of  $x$  under  $\equiv_w^L$  is denoted by  $[x]_w^L$ .

## 2.2 Existing indexing structures for parameterized strings

A basic indexing structure of a p-string is a *parameterized suffix trie*. The parameterized suffix trie  $\text{PSTrie}(T)$  is the trie for  $\text{PSuffix}(T)$ . That is,  $\text{PSTrie}(T)$  is an edge-labeled tree  $(V, E)$  whose node set is  $V = \text{PFactor}(T)$  and edge set is  $E = \{ (x, a, xa) \in V \times (\Sigma \cup \mathcal{N}) \times V \}$ . We remark that each edge of  $\text{PSTrie}(T)$  is labeled by a single symbol from the prev-encoding  $\langle T \rangle$ , and that the out-going edge labels of each node are mutually distinct. An example of  $\text{PSTrie}(T)$  can be found in Figure 1(a). Like the standard suffix tries for static strings, the size of  $\text{PSTrie}(T)$  can be  $\Theta(|T|^2)$ . Obviously we can check whether  $T$  has a substring that p-matches  $P$  of length  $m$  in  $O(m \log(|\Pi| + |\Sigma|))$  time using  $\text{PSTrie}(T)$ , assuming that finding the edge to traverse for a given character takes  $O(\log(|\Pi| + |\Sigma|))$  time by, e.g.,

using balanced trees. We use the same assumption on other indexing structures considered in this paper.

A more compact representation of the suffix sets of p-strings is *parameterized suffix trees* [5]. The parameterized suffix tree  $\text{PSTree}(T)$  of a p-string  $T$  is the path-compacted (or Patricia) tree for  $\text{PSuffix}(T)$ . That is,  $\text{PSTree}(T)$  is an edge-labeled tree  $(V, E)$  of  $T$  where for  $w = \langle T \rangle$ ,

$$\begin{aligned} V &= \{ \llbracket [x]_w^L \rrbracket \mid x \in \text{PFactor}(T) \}, \\ E &= \{ (x, y, xy) \in V \times (\Sigma \cup \mathcal{N})^+ \times V \mid xy = \llbracket [xa]_w^L \rrbracket \in V \text{ for some } a \in \Sigma \cup \mathcal{N} \}. \end{aligned}$$

We remark that each edge of  $\text{PSTree}(T)$  is labeled by an element of  $\text{Factor}(\text{PSuffix}(T)) \setminus \{\varepsilon\}$ , and that the labels of the out-going edges of each node begin with mutually distinct symbols. An example of  $\text{PSTree}(T)$  can be found in Figure 1(b).

To store  $\text{PSTree}(T)$  in linear space, in an actual implementation, the label  $y$  of an edge  $(x, y, xy)$  is represented by two integers  $i$  and  $j$  such that  $y = \langle T[i - |x| : j] \rangle \llbracket [x]_w^L \rrbracket$ , where  $\langle T[i - |x| : j] \rangle = xy$ . In other words,  $y$  corresponds to  $T[i : j]$  but the prev-encoding is given relative to  $T[i - |x| : j]$ . The value  $|x|$  is stored in the node  $x$ , though it is possible to calculate  $|x|$  by reading edge labels from the root to  $x$ .

### 3 Parameterized suffix automata

Recall that the DAWG for a static string  $w \in \Sigma^*$  is isomorphic to a minimal deterministic finite automaton that accepts all the suffixes of  $w$ , which can be obtained by merging isomorphic subtrees of the suffix tries. A static string  $x$  occurs in  $w$  if and only if the automaton has a state that one can reach by reading  $x$ . One natural idea to define the parameterized counterpart of DAWGs, which we actually do not take, may be to have minimal deterministic finite automata for prev-encoded suffixes. This is equivalent to merging isomorphic subtrees of parameterized suffix tries. More formally, letting  $w = \langle T \rangle$  and  $[x]_w^N = \{ y \in \text{PFactor}(w) \mid xz \in \text{PSuffix}(w) \Leftrightarrow yz \in \text{PSuffix}(w) \text{ for all } z \in (\Sigma \cup \mathcal{N})^* \}$  (Nerode equivalence class),  $\text{PSAuto}(T)$  is defined as a directed acyclic graph  $(V, E)$  with

$$\begin{aligned} V &= \{ [x]_w^N \subseteq \text{PFactor}(w) \mid x \in \text{PFactor}(w) \}, \\ E &= \{ ([x]_w^N, a, [xa]_w^N) \in V \times (\Sigma \cup \mathcal{N}) \times V \mid xa \in \text{PFactor}(w) \}, \end{aligned}$$

where the initial state is  $[\varepsilon]_w^N$  and the final states are  $[x]_w^N$  for  $x \in \text{PSuffix}(w)$ . We then have  $x \in \text{PFactor}(w)$  if and only if one can reach some state in the automaton by reading  $x$ . Figure 1(c) shows an example of a parameterized suffix automaton, where we do not distinguish final states and other states. However, differently from the case of static strings, the size of  $\text{PSAuto}(T)$  can be  $\Theta(|T|^2)$  for  $T \in (\Sigma \cup \Pi)^*$ .

**Proposition 1.** *The size of  $\text{PSAuto}(T)$  is  $\Theta(|T|^2)$ .*

*Proof.* Let  $T_k = \mathbf{x}_1 \mathbf{a}_1 \dots \mathbf{x}_k \mathbf{a}_k \mathbf{x}_1 \mathbf{a}_1 \dots \mathbf{x}_k \mathbf{a}_k$  be a p-string over  $\Sigma_k = \{\mathbf{a}_1, \dots, \mathbf{a}_k\}$  and  $\Pi_k = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ , where  $|T_k| = 4k$ . For  $1 \leq i < j \leq k$ , we have  $y_{j,i} = \mathbf{a}_j \dots \mathbf{a}_k \mathbf{a}_1 \dots \mathbf{a}_i = \langle T_k[2j - 1 : 2k + 2i] \rangle \in \text{PFactor}(T_k)$ . We show that we reach different nodes by reading  $y_{j,i}$  and  $y_{j',i'}$  unless  $i = i'$  and  $j = j'$ . If  $i \neq i'$  or  $j \neq j'$ ,  $\mathbf{a}_{i+1} \dots \mathbf{a}_{j-1} (2k) \mathbf{a}_j$  can follow  $y_{j,i}$  but not  $y_{j',i'}$  to form an element of  $\text{PFactor}(T_k)$ . Therefore,  $\text{PSAuto}(T_k)$  must have at least  $k(k-1)/2 \in \Theta(k^2)$  nodes.  $\square$

For example, for  $k = 3$ ,  $i = i' = 1$ ,  $j = 2$ , and  $j' = 3$ , we have  $y_{2,1} = \langle \mathbf{x}_2 \mathbf{a}_2 \mathbf{x}_3 \mathbf{a}_3 \mathbf{x}_1 \mathbf{a}_1 \rangle = \mathbf{a}_2 \mathbf{a}_3 \mathbf{a}_1$  and  $y_{3,1} = \langle \mathbf{x}_3 \mathbf{a}_3 \mathbf{x}_1 \mathbf{a}_1 \rangle = \mathbf{a}_3 \mathbf{a}_1$ . Then,  $z = 6\mathbf{a}_2$  can follow  $y_{2,1}$ , i.e.,  $y_{2,1}z = \langle \mathbf{x}_2 \mathbf{a}_2 \mathbf{x}_3 \mathbf{a}_3 \mathbf{x}_1 \mathbf{a}_1 \mathbf{x}_2 \mathbf{a}_2 \rangle \in \text{PFactor}(T_3)$ . However,  $y_{3,1}z$  is not a pv-string, so not in  $\text{PFactor}(T_3)$ . We conclude  $[y_{2,1}]_{\langle T_3 \rangle}^N \neq [y_{3,1}]_{\langle T_3 \rangle}^N$ . We remark that Proposition 1 holds under binary alphabets, too, which can be shown by the standard binary encoding technique.

We will seek better ideas to define parameterized DAWGs in the following sections.

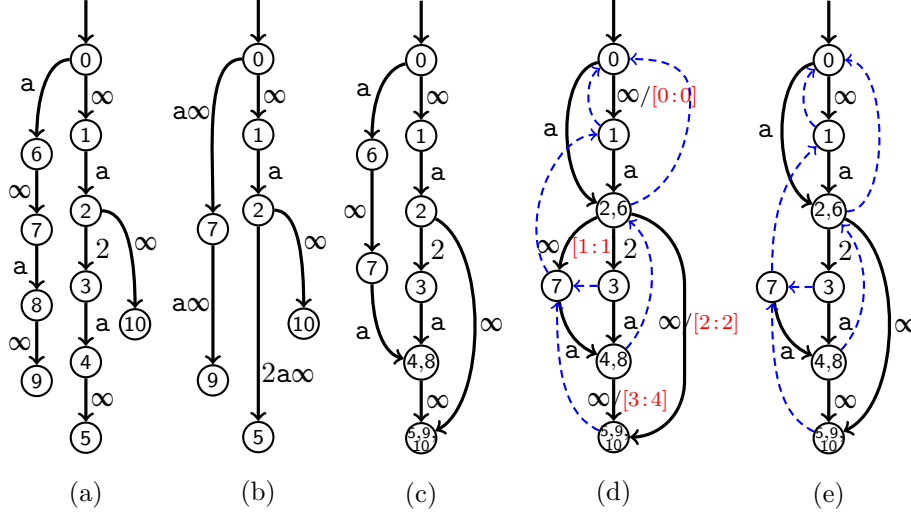


Figure 1: (a) The parameterized suffix trie  $\text{PSTrie}(T)$ , (b) the parameterized suffix tree  $\text{PSTree}(T)$ , (c) the parameterized suffix automaton  $\text{PSAuto}(T)$ , (d) the pseudo PDAWG  $\text{pPDAWG}(T)$ , and (e) the PDAWG  $\text{PDAWG}(T)$  for  $T = \text{xaxay}$  over  $\Sigma = \{a\}$  and  $\Pi = \{x, y\}$ , for which  $\langle T \rangle = w = \infty a 2 a \infty$ . Solid and broken arrows represent the edges and suffix links, respectively. Gate intervals of  $\infty$ -edges of  $\text{pPDAWG}(T)$  are shown with red letters. Some nodes of  $\text{PDAWG}(T)$  cannot be reached by following edges from the source node. The numbers in nodes illustrate how nodes in  $\text{PSTrie}(T)$  are skipped in  $\text{PSTree}(T)$  or merged in the other structures.

## 4 Pseudo parameterized directed acyclic word graphs

Another way to define a parameterized counterpart of DAWGs for static strings may be to merge nodes with the same end position sets in parameterized suffix tries. This approach and the one in the previous section result in the same structures for static strings, but it is not the case for p-strings. We call such a structure for a p-string  $T$  a *pseudo-PDAWG* and denote it by  $\text{pPDAWG}(T)$ , which can formally be defined as a direct acyclic graph  $(V, E)$  where for  $w = \langle T \rangle$

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

$$E = \{ ([x]_w^R, a, [xa]_w^R) \in V \times (\Sigma \cup \mathcal{N}) \times V \mid xa \in \text{PFactor}(w) \}.$$

The nodes  $[\varepsilon]_w^R$  and  $[w]_w^R$  are called the *source* and the *sink*, respectively. An edge  $(u, a, v) \in E$  is called the *a-edge* of  $u$ . In this section, we fix  $w = \langle T \rangle$ . Nodes of a parameterized suffix trie merged in the parameterized suffix automaton are also merged in the pseudo-PDAWG but not vice versa. Therefore, pseudo-PDAWGs can be smaller than parameterized suffix automata. In fact, as we will see later (Theorem 4), the number of nodes of a pseudo-PDAWG is linearly bounded in  $|T|$ . If  $T$  has no parameter characters, our pseudo-PDAWGs coincide with DAWGs for static strings [16].

However, this idea results in an apparent conflict for p-strings. Figure 1(d) shows the pseudo-PDAWG  $\text{pPDAWG}(T)$  obtained from  $\text{PSTrie}(T)$  of Figure 1(a) for  $T = \text{xaxay}$ . Concerning the two nodes  $a$  and  $\infty a$  in  $\text{PSTrie}(T)$ , we have  $\text{RPos}_w(a) = \text{RPos}_w(\infty a) = \{2, 4\}$ , so they shall be merged in  $\text{pPDAWG}(T)$ . However, the subtrees rooted at  $a$  and  $\infty a$  in  $\text{PSTrie}(T)$  have different shapes. The  $\infty$ -edges of those two nodes point at nodes  $a\infty$  and  $\infty a\infty$ , with  $\text{RPos}_w(a\infty) \neq \text{RPos}_w(\infty a\infty)$ , which shall not be merged. As a result of merging  $a$  and  $\infty a$ , the obtained node in  $\text{pPDAWG}(T)$  has two edges labeled with the same character  $\infty$  pointing at different nodes. Consequently,  $\text{pPDAWG}(T)$  has got a path labeled with  $\infty a \infty a \infty$ , which is not in  $\text{PFactor}(T)$ . This apparently obstructs pattern matching over  $\text{pPDAWG}(T)$ . Below we will explain how to resolve the problem and will show nice properties of pseudo-PDAWGs as indexing structures. Lemma 2 shows that if a node  $u$  in a pseudo-PDAWG has two or more edges labeled with the same character, the character must be  $\infty$ . We will determinize

$\infty$ -edges by giving them mutually exclusive “gate intervals” so that one can follow an  $\infty$ -edge only when the string read so far has a length in the interval, based on Lemma 3. In Figure 1(d), the gate intervals are shown in red beside those  $\infty$ -edges. One can follow an  $\infty$ -edge with a gate interval  $[i : j]$  only when  $\infty$  follows a prefix of length between  $i$  and  $j$ . This prevents one to follow the path labeled with  $\infty a \infty a \infty$  in Figure 1(d).

**Lemma 2.** *Given pv-strings  $xa$ ,  $ya$ , and  $w$  with  $a \in \Sigma \cup \mathcal{N}$ , if  $x \equiv_w^R y$  and  $xa \not\equiv_w^R ya$ , then  $a = \infty$ .*

*Proof.* Suppose that  $a \in \Sigma$  and  $x \equiv_w^R y$ . If  $i \in \text{RPos}_w(xa)$ , then  $i - 1 \in \text{RPos}_w(x) = \text{RPos}_w(y)$  and  $w[i] = a$ . This means  $i \in \text{RPos}_w(ya)$ . Therefore,  $\text{RPos}_w(xa) \subseteq \text{RPos}_w(ya)$  and symmetrically we can show  $\text{RPos}_w(ya) \subseteq \text{RPos}_w(xa)$ . That is,  $xa \equiv_w^R ya$ .

Suppose  $a \in \mathcal{N} \setminus \{\infty\}$  and  $x \equiv_w^R y$ . The facts that  $xa$  and  $ya$  are pv-strings and  $a \in \mathbb{N}$  imply  $|x|, |y| \geq a$ . If  $i \in \text{RPos}_w(xa)$ , then  $w[i] = a$ . Since  $i - 1 \in \text{RPos}_w(x) = \text{RPos}_w(y)$ , we have  $i \in \text{RPos}_w(ya)$ . Therefore,  $\text{RPos}_w(xa) = \text{RPos}_w(ya)$ .

Hence, only when  $a = \infty$ , it is possible that  $x \equiv_w^R y$  and  $xa \not\equiv_w^R ya$ .  $\square$

The following lemma means that  $\infty$ -edges of a node have respective gate intervals  $[i : j]$  according to which one can choose the valid one to follow.

**Lemma 3.** *Each  $\infty$ -edge  $(u, \infty, v)$  of  $\text{pPDAWG}(T)$  admits two integers  $i$  and  $j$  such that for any  $x \in u$ , we have  $x\infty \in v$  if and only if  $i \leq |x| \leq j$ .*

*Proof.* Let  $i = |y|$  and  $j = |z|$  for the shortest  $y$  and longest  $z$  such that  $y, z \in u$  and  $y\infty, z\infty \in v$ . Then, the “only if” direction is obvious. On the other hand, for every  $x$  such that  $x \in u$  and  $i \leq |x| \leq j$ , we have  $y \in \text{PSuffix}(x)$  and  $x \in \text{PSuffix}(z)$  by Lemma 1. This implies, again by Lemma 1,  $x\infty \in v$ .  $\square$

Now, we enhance  $\infty$ -edges of pseudo-PDAWGs with the gate intervals  $[i : j]$  given in Lemma 3. Suppose we have reached a node  $u$  by reading a pv-string  $x$ . If the next character  $a$  is not  $\infty$ , we simply follow the  $a$ -edge of  $u$ . If  $a = \infty$ , we follow the  $\infty$ -edge with interval  $[i : j]$  such that  $i \leq |x| \leq j$ . If  $u$  has no such edge, it means that  $xa \notin \text{PFactor}(T)$ . We remark that  $\text{pPDAWG}(T)$  in Figure 1(d) has a path  $a2$  which does not exist in  $\text{PSTrie}(T)$ , but it is harmless, since the non-pv-string  $a2$  can be obtained from no input pattern  $P$ .

**Proposition 2.** *Let  $w = \langle T \rangle$ . One can reach a node  $u$  from  $[\varepsilon]_w^R$  by reading  $x$  in  $\text{pPDAWG}(T)$  if and only if  $u = [x]_w^R$ .*

*Proof.* We show the proposition by induction on  $|x|$ . If  $x = \varepsilon$ , the conclusion is trivial. Suppose we have reached  $u = [x]_w^R$  by reading  $x$  and the next character is  $a \in \Sigma \cup \mathcal{N} \setminus \{\infty\}$ . If  $xa \in \text{PFactor}(T)$ , Lemma 2 implies that  $u$  has only one  $a$ -edge  $(u, a, v)$ , for which  $v = [xa]_w^R$  holds. If  $xa \notin \text{PFactor}(T)$ , then  $ya \notin \text{PFactor}(T)$  for all  $y \in u$ , which means that  $u$  has no  $a$ -edge by the definition of  $E$ . Suppose the next character is  $a = \infty$ . Lemma 3 implies that  $u$  has only one  $\infty$ -edge  $(u, \infty, v)$  with a length condition  $[i : j]$  satisfying  $i \leq |x| \leq j$ , for which  $v = [x\infty]_w^R$ , if and only if  $xa \in \text{PFactor}(T)$ .  $\square$

Therefore, one can decide whether  $\langle P \rangle \in \text{PFactor}(T)$  in time  $O(|P| \log(|\Pi| + |\Sigma|))$  using  $\text{pPDAWG}(T)$ .

In order to find all end positions of substrings of the text which p-matches with an input p-string, we further augment pseudo-PDAWGs, in the way analogous to the classical enhancement of DAWGs. One obvious idea might be to explicitly record  $\text{RPos}_w(u)$  in each node  $u \in V$  so that all the p-occurrences of  $x \in u$  can be found at the reached node, but it makes the data structure size non-linear. Instead, we assign the smallest number  $\ell_u$  in  $\text{RPos}_w(u)$  to each node  $u$ , i.e., the end position of the left most p-occurrence of  $x \in u$  in  $w$ . To find other p-occurrences, *suffix links* are useful, which are defined as

$$F = \{ (u, v) \in (V \setminus \{\varepsilon\}) \times V \mid v = [\langle x[2 : ] \rangle]_w^R \text{ for } x = [u] \}.$$

In other words, for two nodes  $u, v \in V$ , we have  $(u, v) \in F$  if and only if  $[v]$  is obtained by removing the first character of  $[u]$ . Thus, every pv-suffix of  $x \in u$  belongs to some  $v \in V$  which can be reached

---

**Algorithm 1:** Parameterized pattern matching algorithm based on pPDAWG( $T$ )

---

```
1 Let  $p \leftarrow \langle P \rangle$ ;  
2 Let  $u \leftarrow [\varepsilon]_w^R$ ;  
3 for  $i = 1$  to  $|p|$  do  
4   if  $p[i] \neq \infty$  and  $(u, p[i], v) \in E$  then Let  $u \leftarrow v$ ;  
5   else if  $u$  has an  $\infty$ -edge  $(u, \infty, v) \in E$  whose gate interval  $[j : k]$  satisfies  $j \leq i - 1 \leq k$   
6     then Let  $u \leftarrow v$ ;  
7   else return False;  
7 Traverse the reversed suffix link tree and  
   output  $\ell_v$  for all descendants  $v$  of  $u$ ;
```

---

from  $u$  by following the suffix links. Consequently, the reverse of suffix links  $\overline{F} = \{(v, u) \mid (u, v) \in F\}$  form a rooted tree, where  $[\varepsilon]_w^R$  is the root and  $v$  is a child of  $u$  if  $(u, v) \in \overline{F}$ . We call  $(V, \overline{F})$  the *reversed suffix link tree*. We will discuss in Section 6 that  $(V, \overline{F})$  is isomorphic to the parameterized suffix tree of  $\overline{T}$ , from which it is obvious that one can find all p-occurrences of  $x$  in  $w$  by visiting all the descendants of the reached node  $[x]_w^R$  in  $(V, \overline{F})$ . Here, independently of the duality arguments, we justify the p-matching procedure with reversed suffix links.

**Lemma 4.** *It holds that  $k \in \text{RPos}_w(v)$  if and only if  $k = \ell_u$  for some (not necessarily proper) descendant  $u$  of  $v$  in  $(V, \overline{F})$ .*

*Proof.* If  $u$  is a child of  $v$ , then  $\text{RPos}_w(u) \subsetneq \text{RPos}_w(v)$ . This implies the “if” direction.

To see the “only if” direction, let  $v_0 = [w[:k]]_w^R$ . Then, clearly  $\ell_{v_0} = k$  and  $x \in \text{PSuffix}(w[:k])$  for  $x \in v$ . There must be  $v_0, \dots, v_j$  such that  $(v_{i-1}, v_i) \in F$  for  $i = 1, \dots, j$  and  $v_j = v$ .  $\square$

Therefore, all the p-occurrence positions of the pattern can be found as  $\ell_u$  by traversing the descendants  $u$  in the reverse suffix link tree  $(V, \overline{F})$ . The matching procedure is summarized in Algorithm 1.

The following technical lemma can be used to show that the number of descendants of a node  $v$  is linearly bounded by  $\text{RPos}_w(v)$ .

**Lemma 5.** *Suppose that each node  $u$  of a rooted tree is assigned a nonempty finite set  $X_u$  so that*

- *if  $u$  is a child of  $v$ , then  $X_u \subsetneq X_v$ ,*
- *if  $u$  and  $v$  are siblings, then  $X_u \cap X_v = \emptyset$ .*

*Then,  $n \leq 2|X_r| - 1$  where  $n$  is the number of nodes of the tree and  $r$  is the root. Moreover,  $n = 2|X_r| - 1$  if and only if every leaf  $u$  has a singleton set  $X_u$  and every inner node  $v$  has exactly two children  $v_1$  and  $v_2$  such that  $X_{v_1}$  and  $X_{v_2}$  partition  $X_v$ .*

*Proof.* We show the lemma by induction on the number of nodes of the tree. Suppose  $r$  has children  $u_1, \dots, u_k$  and let  $n_i$  be the number of nodes of the subtree rooted by  $u_i$ . If  $k = 0$ , since  $X_r$  is nonempty, the lemma holds. If  $k = 1$ , by  $|X_{u_1}| < |X_r|$ ,  $n = n_1 + 1$  and the induction hypothesis  $n_1 \leq 2|X_{u_1}| - 1$ , we have  $n < 2|X_r| - 1$ . If  $k \geq 2$ ,  $n = 1 + \sum_{i=1}^k n_i \leq 1 + \sum_{i=1}^k (2|X_{u_i}| - 1) \leq 2|X_r| - k + 1 \leq 2|X_r| - 1$ . The equality signs hold only when  $k = 2$ ,  $|X_r| = |X_{u_1}| + |X_{u_2}|$ , and the subtrees rooted by  $u_i$  satisfy the stated condition.  $\square$

If  $(u_1, v), (u_2, v) \in F$ , then  $||u_1|| = ||u_2|| = ||v|| + 1$ . By Lemma 1,  $\text{RPos}_w(u_1) \cap \text{RPos}_w(u_2) = \emptyset$ , unless  $u_1 = u_2$ . This means that the reversed suffix link tree satisfies Lemma 5 for  $X_u = \text{RPos}_w(u)$ . Hence, the number of descendants of  $u$  in the reversed suffix link tree is at most  $2|\text{RPos}_w(u)| - 1$ . We obtain the following theorem.

**Theorem 1.** *Using pPDAWG( $T$ ) enhanced with the suffix links, we can find all substrings of  $T$  that p-match a given pattern  $P$  in  $O(|P| \log(|\Pi| + |\Sigma|) + \text{occ})$  time, where  $\text{occ}$  is the number of occurrences to report.*



The following theorem contrasts Proposition 1.

**Theorem 2.** *If  $n = |T| \geq 2$ ,  $\text{pPDAWG}(T)$  has at most  $2n - 1$  nodes. The bound is tight.*

*Proof.* Let  $aw = \langle T \rangle$  with  $a \in \Sigma \cup \{\infty\}$  and  $w \in (\Sigma \cup \mathcal{N})^+$ . In the reversed suffix link tree  $(V, \overline{F})$ ,  $[a]_{aw}^R = \{a\}$  is a child of the root  $[\varepsilon]_{aw}^R = \{\varepsilon\}$ . Suppose  $\{\varepsilon\}$  has children  $u_1, \dots, u_j$  in addition to  $\{a\}$  and  $\{a\}$  has children  $u_{j+1}, \dots, u_k$  in  $(V, \overline{F})$ , where  $\text{RPos}_w(u_1), \dots, \text{RPos}_w(u_k)$  are pairwise disjoint and  $k \geq 1$  by  $n \geq 2$ . Since only  $\varepsilon$  and  $a$  can end at the position 1,  $\text{RPos}_w(u_i) \subseteq \{2, \dots, n\}$ . Let us partition  $V$  into  $\{\{\varepsilon\}, \{a\}\}, V_1, \dots, V_k$  where each  $V_i$  consists of the nodes of the subtree rooted by  $u_i$ . By Lemma 5,  $|V_i| \leq 2|\text{RPos}_w(u_i)| - 1$  and therefore

$$|V| = 2 + \sum_{i=1}^k |V_i| \leq 2 + 2 \sum_{i=1}^k |\text{RPos}_w(u_i)| - k \leq 2 + 2|w| - 1 = 2n - 1.$$

The tightness is witnessed by a static string  $\mathbf{ab}^{n-1}$ , presented by Blumer et al. [16].  $\square$

Blumer et al. [16] have shown in addition that the DAWG for a static string of length  $n$  has at most  $3n - 4$  edges. We will later show in Corollary 1 a linear bound on the number of edges of pseudo-PDAWGs.

Our proposed pseudo-PDAWGs are compact enough and support efficient parameterized pattern matching. In the next section, we will present a modification of pseudo-PDAWGs as our main proposal indexing structure for parameterized strings.

## 5 Parameterized directed acyclic word graphs

In this section, we present a new indexing structure for parameterized strings, which we call *parameterized directed acyclic word graphs* (PDAWGs). A PDAWG is obtained from a pseudo-PDAWG by suppressing some  $\infty$ -edges and forgetting the assigned intervals of all  $\infty$ -edges. As compensation, we will make use of suffix links for matching. When two nodes  $x_1$  and  $x_2$  in  $\text{PSTrie}(T)$  are merged into  $u = [x_1]_w^R = [x_2]_w^R$  in  $\text{pPDAWG}(T)$ , the node  $u$  keeps all the outgoing edges of the original nodes  $x_1$  and  $x_2$ . In  $\text{PDAWG}(T)$ , we keep only the outgoing edges of  $[u]$ . Recall that if  $(x_1, a, y_1)$  is an edge of  $\text{PSTrie}(T)$  for some  $a \in \Sigma \cup \mathcal{N} \setminus \{\infty\}$ , then  $x_2$  also has an  $a$ -edge  $(x_2, a, y_2)$  such that  $y_1 \equiv_w^R y_2$ . Therefore, the difference of  $\text{pPDAWG}(T)$  and  $\text{PDAWG}(T)$  is only in  $\infty$ -edges. In this section, we fix a text  $T$  and its pv-encoding  $w = \langle T \rangle$ .

**Definition 1** (Parameterized directed acyclic word graphs). The *parameterized directed acyclic word graph* (PDAWG)  $\text{PDAWG}(T)$  of  $T$  is a triple  $(V, E, F)$  where

$$\begin{aligned} V &= \{ [x]_w^R \mid x \in \text{PFactor}(w) \}, \\ E &= \{ ([x]_w^R, a, [xa]_w^R) \in V \times (\Sigma \cup \mathcal{N}) \times V \mid x = [[x]_w^R] \text{ and } xa \in \text{PFactor}(w) \}, \\ F &= \{ (u, v) \in (V \setminus \{\{\varepsilon\}\}) \times V \mid v = [x[2 : ]]_w^R \text{ for } x = [u] \}. \end{aligned}$$

Elements of  $E$  are called edges and those of  $F$  are suffix links.

The sets  $V$  and  $F$  remain unchanged from pseudo-PDAWGs. Note that  $(V, E)$  is a directed acyclic graph and  $(V, \overline{F})$  is a tree. Since each  $u \in V \setminus \{\{\varepsilon\}\}$  has unique  $v$  such that  $(u, v) \in F$ , we often write  $F(u)$  for  $v$  regarding  $F$  as a function.

### 5.1 Parameterized matching based on PDAWGs

Figure 1(e) shows an example PDAWG. One may wonder how to find a p-occurrence of  $\langle \mathbf{axa} \rangle = \mathbf{a\infty a} \in \text{PFactor}(\mathbf{xaxay})$  using the PDAWG. Our transition function (Algorithm 2) is based on Lemma 7, which indicates the node we should visit when the succeeding character is  $a \in \Sigma \cup \mathcal{N}$  after reading  $x$ . The following lemma prepares for Lemma 7.

---

**Algorithm 2:** Function  $\text{trans}(u, i, a)$ 

---

```
1 if  $a \neq \infty$  then return  $\text{child}(u, a)$ ;  
2 else  
3   Let  $Z \leftarrow \{j \in \mathcal{N} \cap \text{Children}(u) \mid j > i\}$ ;  
4   if  $Z = \emptyset$  then return Null;  
5   else if  $Z = \{b\}$  then return  $\text{child}(u, b)$ ;  
6   else return  $F(\text{child}(u, b))$  for  $b = \min Z$ ;
```

---

**Lemma 6.** For  $x\infty \in \text{PFactor}(w)$ , let  $y = \lceil [x]_w^R \rceil$ . Then,

$$\text{RPos}_w(x\infty) = \bigcup \{ \text{RPos}_w(yj) \mid j \in \mathcal{N} \text{ and } j > |x| \}.$$

*Proof.* To show  $\text{RPos}_w(x\infty) \subseteq \bigcup_{j>|x|} \text{RPos}_w(yj)$ , suppose  $i \in \text{RPos}_w(x\infty)$ . This implies  $i-1 \in \text{RPos}_w(x) = \text{RPos}_w(y)$ ,  $w[i] \in \mathcal{N}$ , and  $w[i] > |x|$ . Thus,  $\langle w[i-|y| : i-1] \rangle = y$  and  $i \in \text{RPos}_w(yj)$  for  $j = \langle\langle w[i] \rangle\rangle_{|y|} \in \{w[i], \infty\}$ . Hence, we have  $j > |x|$ .

Conversely suppose  $i \in \text{RPos}_w(yj)$  for some  $j > |x|$ . Then,  $i-1 \in \text{RPos}_w(y) = \text{RPos}_w(x)$  and  $j = \langle\langle w[i] \rangle\rangle_{|y|} \in \{w[i], \infty\}$ . On the other hand,  $j > |x|$  implies  $\langle\langle w[i] \rangle\rangle_{|x|} = \infty$ . Thus,  $\langle w[i-|x| : i] \rangle = x\infty$ , i.e.,  $i \in \text{RPos}_w(x\infty)$ . This proves  $\text{RPos}_w(x\infty) \supseteq \bigcup_{j>|x|} \text{RPos}_w(yj)$ .  $\square$

**Lemma 7.** Suppose  $x \in \text{PFactor}(w)$  and  $a \in \Sigma \cup \mathcal{N}$ . Then, for  $y = \lceil [x]_w^R \rceil$ ,

$$[xa]_w^R = \begin{cases} [ya]_w^R & \text{if } a \neq \infty \text{ or } Z = \emptyset, \\ [yk]_w^R & \text{if } a = \infty \text{ and } |Z| = 1, \\ F([yk]_w^R) & \text{if } a = \infty \text{ and } |Z| \geq 2, \end{cases}$$

where  $Z = \{j \in \mathcal{N} \mid yj \in \text{PFactor}(w) \text{ and } j > |x| \}$  and  $k = \min Z$ .

*Proof.* If  $a \neq \infty$ , the lemma immediately follows from Lemma 2. We suppose  $a = \infty$ . If  $|Z| \leq 1$ , we obtain the lemma by Lemma 6.

Suppose  $|Z| \geq 2$ . In this case,  $|x| < k < \infty$ . By Lemma 6, we see that  $\text{RPos}_w(yk) \subsetneq \text{RPos}_w(x\infty)$ . Recall that in general  $\text{RPos}_w(u) \subseteq \text{RPos}_w(v)$  if and only if  $v$  is reachable from  $u$  by following a certain number (including zero) of suffix links in  $(V, F)$ . Hence, one can reach  $[x\infty]_w^R$  from  $[yk]_w^R$  by following at least one suffix link. To show that there is no other node between  $[x\infty]_w^R$  and  $[yk]_w^R$ , it suffices to show that for any  $z \in \text{PSuffix}(yk)$  such that  $|x\infty| < z < |yk|$  (and thus  $\text{RPos}_w(yk) \subseteq \text{RPos}_w(z) \subseteq \text{RPos}_w(x\infty)$ ), either  $\text{RPos}_w(z) = \text{RPos}_w(x\infty)$  or  $\text{RPos}_w(z) = \text{RPos}_w(yk)$ . Here,  $z$  must be of the form  $z = z' \langle\langle k \rangle\rangle_{|z'|}$ . The assumption implies  $\lceil [x]_w^R \rceil = \lceil [z']_w^R \rceil = \lceil [y]_w^R \rceil$ . Note that  $yk \in \text{PFactor}(w)$  and  $k < \infty$  implies  $|y| \geq k$ .

Suppose  $|x\infty| < |z| \leq k$ , i.e.,  $z = z'\infty$ . By Lemma 6,  $|z'| < k$  implies  $\text{RPos}_w(z'\infty) = \text{RPos}_w(x\infty) = \bigcup_{j \geq k} \text{RPos}_w(yj)$  by the choice of  $k$ .

Suppose otherwise,  $k < |z| < |yk|$ , i.e.,  $z = z'k$ . Lemma 2 implies  $\text{RPos}_w(z'k) = \text{RPos}_w(yk)$ .  $\square$

The function  $\text{trans}$  of Algorithm 2 is a straightforward realization of Lemma 7, where  $\text{Children}(u)$  denotes the set of labels of the outgoing edges of  $u$  and  $\text{child}(u, a)$  is the node that the  $a$ -edge of  $u$  points at. If  $u$  has no  $a$ -edge,  $\text{child}(u, a) = \text{Null}$ . In other words,  $\text{Children}(u) = \{a \in \mathcal{N} \cup \Sigma \mid (u, a, v) \in E \text{ for some } v \in V\}$  and  $\text{child}(u, a) = v \in V$  iff  $(u, a, v) \in E$ . The algorithm takes a node  $u \in V$ , a natural number  $i \in \mathbb{N}$ , and a character  $a \in \Sigma \cup \mathcal{N}$ , and returns the node where we should go by reading  $a$  from  $u$  assuming that we have read  $i$  characters so far. By Lemma 7,  $\text{trans}([x]_w^R, |x|, a) = [xa]_w^R$  for every  $xa \in \text{PFactor}(w)$ . On the other hand, suppose  $x \in \text{PFactor}(T)$  and  $xa \notin \text{PFactor}(T)$ . If  $a \neq \infty$ ,  $ya \notin \text{PFactor}(T)$  for  $y = \lceil [x]_w^R \rceil$ . If  $a = \infty$ , for any  $i \in \text{RPos}_w(x) = \text{RPos}_w(y)$ , either  $w[i+1] \in \Sigma$  or  $w[i+1] \leq |x|$ . Thus, the node  $[x]_w^R = [y]_w^R$  has no edge labeled with a character in  $\mathcal{N}$  greater than  $|x|$ . The algorithm returns **False**. Using  $\text{trans}$ , Algorithm 3 performs p-matching, where  $\ell_v = \min \text{RPos}(v)$ .

---

**Algorithm 3:** Parameterized pattern matching algorithm based on PDAWG( $T$ )

---

```
1  $p \leftarrow \langle P \rangle$ ;  
2 Let  $u$  be the source node of PDAWG( $T$ );  
3 for  $i = 1$  to  $|P|$  do  
4   | Let  $u \leftarrow \text{trans}(u, i - 1, p[i])$ ;  
5   | if  $u = \text{Null}$  then return False;  
6 Traverse the reversed suffix link tree and  
   output  $\ell_v$  for all descendants  $v$  of  $u$ ;
```

---

**Theorem 3.** Using PDAWG( $T$ ), we can find all  $p$ -occurrences of  $P$  in  $T$  in  $O(|P| \log(|\Pi| + |\Sigma|) + \text{occ})$  time, where  $\text{occ} = |\text{RPos}_{\langle T \rangle}(\langle P \rangle)|$  is the number of occurrences to report.

## 5.2 Size of PDAWGs

Blumer et al. [16] have shown that the DAWG for a static string of length  $n$  has at most  $2n - 1$  nodes and  $3n - 4$  edges. We show that PDAWGs have the same size bound.

**Theorem 4.** PDAWG( $T$ ) has at most  $2n - 1$  nodes and  $3n - 4$  edges when  $n = |T| \geq 3$ . Those bounds are tight.

*Proof.* Concerning the number of nodes, Theorem 2 holds for PDAWGs, since the node sets of PDAWGs and pseudo-PDAWGs are identical.

On the number of edges, we first give a weaker upper bound  $3n - 3$ , just like Blumer et al. [16] have done. Let  $\text{PDAWG}(T) = G = (V, E, F)$ ,  $\text{PSTrie}(T) = H = (U, D)$ ,  $V' \subsetneq V$  the set of non-sink nodes of  $G$ ,  $U' \subsetneq U$  the set of internal nodes of  $H$ , and  $d_G(v)$  denote the out-degree of node  $v$  in  $G$ . We have  $|E| = \sum_{v \in V'} d_G(v) = \sum_{v \in V'} d_H(\lceil v \rceil)$ , since  $d_G(v) = d_H(\lceil v \rceil)$  for all  $v \in V$ . Since  $H$  has at most  $n$  leaves,

$$n \geq |U \setminus U'| = 1 + \sum_{u \in U'} (d_H(u) - 1) \geq 1 + \sum_{v \in V'} (d_H(\lceil v \rceil) - 1) = 1 + |E| - |V'|,$$

which implies  $|E| \leq n + |V'| - 1 \leq 3n - 3$ .

This upper bound  $3n - 3$  could be achieved only when  $|V| = 2n - 1$ . We will show that if  $|V| = 2n - 1$ , then the skeleton (stripping off edge labels) of PDAWG( $T$ ) is isomorphic to that of PDAWG( $\mathbf{ab}^{n-1}$ ) for  $\mathbf{a}, \mathbf{b} \in \Sigma$  with  $\mathbf{a} \neq \mathbf{b}$ , where the source is the only branching node from which two paths of length  $n$  and  $n - 1$  reach the sink.

Let  $aw = \langle T \rangle$  with  $a \in \Sigma \cup \{\infty\}$  and  $\overline{F}(u) = \{v \mid (v, u) \in F\}$  for  $u \in V$ , i.e.,  $\overline{F}(u)$  is the set of children of  $u$  in the reversed suffix link tree. According to the proof of Theorem 2,  $|V| = 2n - 1$  can be achieved only when  $\overline{F}(\{\varepsilon\}) \cup \overline{F}(\{a\}) \setminus \{\{a\}\} = \{u\}$  for some  $u \in V$  such that  $\text{RPos}_{aw}(u) = \{2, \dots, n\}$ . Moreover, by Lemma 5, it must hold that  $\overline{F}(u) = \{u_1, u_2\}$  and that  $\text{RPos}_{aw}(u_1)$  and  $\text{RPos}_{aw}(u_2)$  partition  $\text{RPos}_{aw}(u) = \{2, \dots, n\}$ . Suppose  $\overline{F}(\{\varepsilon\}) = \{\{a\}\}$  and  $\overline{F}(\{a\}) = \{u\}$ . In this case, by  $\text{RPos}_{aw}(u) \subsetneq \text{RPos}_{aw}(a) \subsetneq \text{RPos}_{aw}(\varepsilon)$ , it must hold  $\text{RPos}_{aw}(a) = \{1, \dots, n\}$ . Since there are at most  $k + 1$  pv-strings  $x$  such that  $k \in \text{RPos}_{\langle S \rangle}(x)$  for any p-string  $S$  in general, there can be at most three nodes  $v$  in  $V$  such that  $2 \in \text{RPos}_{aw}(v)$ , which are actually  $\{\varepsilon\}$ ,  $\{a\}$  and  $u$ . Therefore,  $\text{RPos}_{aw}(u_1)$  and  $\text{RPos}_{aw}(u_2)$  cannot partition  $\text{RPos}_{aw}(u)$ . Consequently, it must hold  $\overline{F}(\{\varepsilon\}) = \{\{a\}, u\}$  and  $\overline{F}(\{a\}) = \emptyset$ . Since  $\text{RPos}_{aw}(\{a\}) \cap \text{RPos}_{aw}(u) = \emptyset$ , we must have  $\text{RPos}_{aw}(\{a\}) = \{1\}$ . By  $[w[k : k]]_{aw}^R \in \overline{F}(\{\varepsilon\})$  for all  $k$ , we must have either  $a \in \Sigma$  and  $w \in \mathcal{N}^{n-1}$  or  $w = b^{n-1}$  for some  $b \in \Sigma^{n-1}$ . In the latter case, it is easy to see that PDAWG( $T$ ) is isomorphic to PDAWG( $\mathbf{ab}^{n-1}$ ).

So, hereafter we assume  $a \in \Sigma$  and  $w \in \mathcal{N}^{n-1}$ . Then,  $\text{RPos}_w((aw)[k]) = \{k\}$  and  $[(aw)[k]]_w^R = \{(aw)[k]\}$  for all  $k \in \{1, \dots, n\}$ . We show by induction on  $k$  that  $\text{RPos}_w(\langle w[i : i + k - 1] \rangle) = \{k + 1, \dots, n\}$  for  $1 \leq i \leq n - k$  for  $1 \leq k \leq n - 1$ , i.e.,  $\text{PFactor}(w)$  has just one pv-string of length  $k$ . This is true for  $k = 1$ , since  $\langle w[i] \rangle = \infty$  for  $1 \leq i \leq n - 1$ . This is also trivially true for  $k = n - 1$ .

Suppose the claim holds up to  $k \leq n-3$ . According to the proof of Lemma 5, to achieve the tight upper bound  $|V| = 2n - 1$ , it must hold  $\overline{F}([w[:k]]_w^R) = \{u_1, u_2\}$  such that  $\text{RPos}(u_1)$  and  $\text{RPos}(u_2)$  partition  $\text{RPos}_w(w[:k]) = \{k+1, \dots, n\}$ . Since  $F([a(w[:k])]_w^R) = [w[:k]]_w^R$  and  $\text{RPos}_w(a(w[:k])) = \{k+1\}$ , it must hold  $\text{RPos}(u_2) = \{k+2, \dots, n\}$  for  $\overline{F}([w[:k]]_w^R) = \{[a(w[:k])]_w^R, u_2\}$ . Take an arbitrary element  $x \in u_2$ . The fact  $k+2, n \in \text{RPos}_w(x)$  means  $x = \langle w[i : k+2] \rangle$  for some  $i \geq 1$ . By induction hypothesis, any element of  $\text{PFactor}(w)$  of length  $j \leq k$  is  $\langle w[:j] \rangle$ . Therefore,  $x$  must have length  $k+1$ , i.e.,  $x = w[:k+1]$  and thus  $u_2 = \{w[:k+1]\}$ .  $\text{RPos}_{aw}(u_2) = \{k+1, \dots, n\}$  means  $w[:k+1] = \langle w[i : k+i-1] \rangle$  for  $1 \leq i < n-k$ . Summarizing above, all nodes of  $V$  are  $[w[:k]]_w^R$  and  $[a(w[:k])]_w^R$  for  $0 \leq k \leq n-2$  and the sink node  $\{w[:n-1], a(w[:n-1])\}$ . This PDAWG is isomorphic to  $\text{PDAWG}(\text{ab}^{n-1})$ , which has  $2n-1$  edges.

One can achieve  $|E| = 3n-4$  by the string  $\text{ab}^{n-2}\text{c}$ , given by Blumer et al. [16].  $\square$

**Corollary 1.**  $\text{pPDAWG}(T)$  has at most  $5n-7$  edges when  $n = |T| \geq 3$ .

*Proof.* By Theorem 4, it suffices to evaluate the number of  $\infty$ -edges of  $\text{pPDAWG}(T)$  which do not appear in  $\text{PDAWG}(T)$ . We show that each node of  $\text{pPDAWG}(T)$  has at most one such incoming  $\infty$ -edge. Suppose a (non-source) node  $v$  has two incoming  $\infty$ -edges  $(u_1, \infty, v)$  and  $(u_2, \infty, v)$  in  $\text{pPDAWG}(T)$ . This means that there are  $x_1 \in u_1$  and  $x_2 \in u_2$  such that  $x_1\infty, x_2\infty \in v$ . We assume without loss of generality that  $|x_1| < |x_2|$ , which implies  $|x_1| \leq \lceil |u_1| \rceil < \lfloor |u_2| \rfloor \leq |x_2|$ . Then,  $\lceil u_1 \rceil \infty \in v$  by Lemma 1, which means that  $\text{PDAWG}(T)$  has the edge  $(u_1, \infty, v)$ . That is, only  $(u_2, \infty, v)$  may be suppressed in  $\text{PDAWG}(T)$ . Therefore, the number of edges of  $\text{pPDAWG}(T)$  is at most  $|E| + |V| - 1 \leq 3n-4 + 2n-2-1 = 5n-7$  for  $\text{PDAWG}(T) = (V, E, F)$ .  $\square$

It is an open problem to give the tight upper bound on the number of edges of pseudo-PDAWGs.

## 6 Duality of PDAWGs and p-suffix trees

This section establishes the duality between parameterized suffix trees [5] and PDAWGs. This will give us the following two merits: the *first* bidirectional index for parametrized pattern matching (Section 6.1), and efficient offline construction of PDAWGs (Section 6.2).

### 6.1 Parameterized suffix trees and Weiner links

In this subsection, we first recall the basic properties of p-suffix trees (Section 6.1.1), and then the suffix links of p-suffix trees (Section 6.1.2). Then, we introduce our Weiner links of p-suffix trees (Section 6.1.3), and show that the Weiner links are equal to the PDAWG edges (Section 6.1.4). This immediately leads us to bidirectional indexing structure for p-matching.

#### 6.1.1 Basics of p-suffix trees

Let  $T$  be a p-string and consider its reversal  $\overline{T}$ . The *parameterized suffix tree* (*p-suffix tree*)  $\text{PSTree}(\overline{T})$  of  $\overline{T}$  is the path-compacted (or Patricia) tree for  $\text{PSuffix}(\overline{T})$ . Below let us recall the definition of  $\text{PSTree}(\overline{T})$ , which is the edge-labeled tree  $(V, E)$  of  $\overline{T}$  such that for  $w = \langle \overline{T} \rangle$

$$\begin{aligned} V &= \{ \lceil [x]_w^L \rceil \mid x \in \text{PFactor}(\overline{T}) \}, \\ E &= \{ (x, y, xy) \in V \times (\Sigma \cup \mathcal{N})^+ \times V \mid xy = \lceil [xa]_w^L \rceil \in V \text{ for some } a \in \Sigma \cup \mathcal{N} \}. \end{aligned}$$

Recall that each edge of  $\text{PSTree}(T)$  is labeled by an element of  $\text{Factor}(\text{PSuffix}(T)) \setminus \{\varepsilon\}$ . An example for  $\text{PSTree}(\overline{T})$  is given in Figure 2(a).

**Remark 1.** Since each node of  $\text{PSTree}(\overline{T})$  is defined as the longest member  $\lceil [x]_w^L \rceil$  of the equivalence class  $[x]_w^L$ ,  $\text{PSTree}(\overline{T})$  may contain an internal node that has only a single child. Such an internal node corresponds to a suffix of  $\overline{T}$  that has internal p-matching occurrences in  $\overline{T}$ . For instance,  $\text{a}\infty$  of the

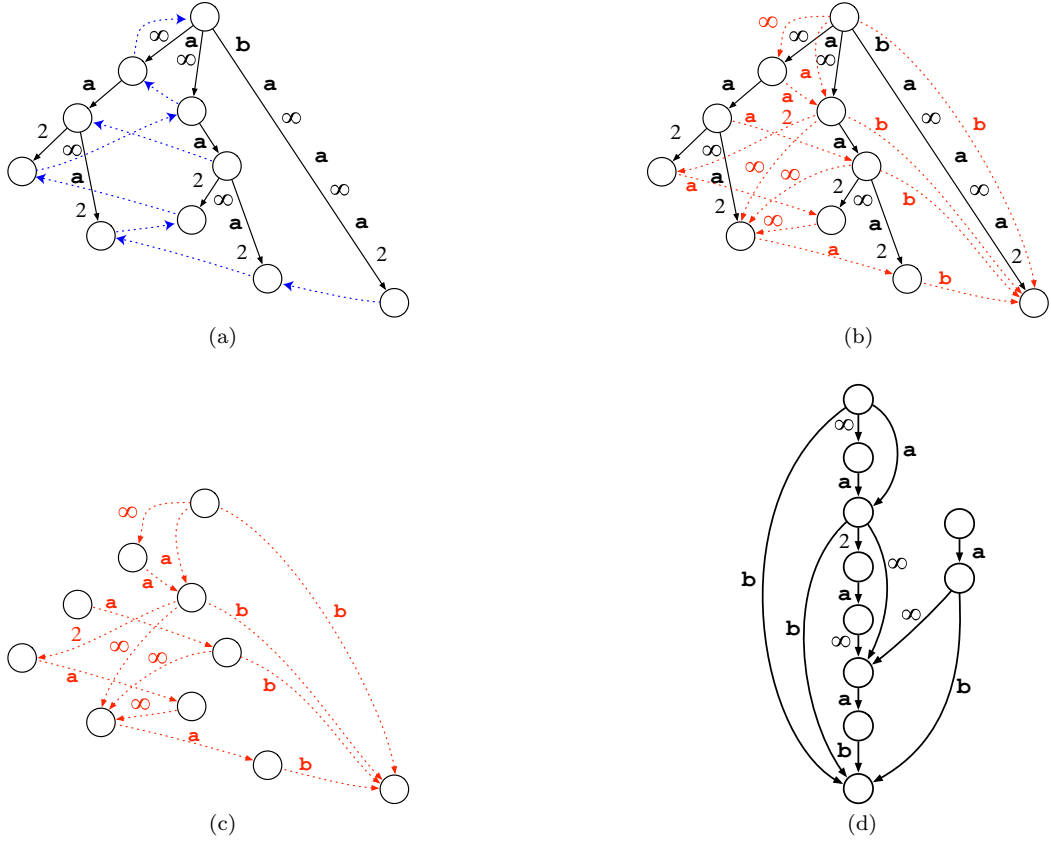


Figure 2: Consider string  $T = \text{yayaxab}$  over  $\Sigma = \{a, b\}$  and  $\Pi = \{x, y\}$ . (a) The parameterized suffix tree  $\text{PSTree}(\bar{T})$  with  $\bar{T} = \text{baxayay}$  augmented with the suffix links (dashed blue arcs). (b) The parameterized suffix tree  $\text{PSTree}(\bar{T})$  augmented with the Weiner links (dashed red arcs). (c) The DAG consisting of the p-suffix tree nodes and the Weiner-links. (d) The PDAWG  $\text{PDAWG}(T)$  with  $T = \text{yayaxab}$ . Observe that the edge-labeled graphs (c) and (d) are isomorphic.

parameterized suffix tree shown in Figure 2(a) has only a single child. This is because  $a\infty$  has three p-matching occurrences in  $\text{baxayay}$  at positions 2, 4, 6, where the last occurrence corresponds to the suffix  $\text{ay}$  and all of the other internal p-matching occurrences of  $a\infty$  are immediately followed by  $a$ . If we use a common convention that  $\bar{T}$  terminates with a unique character  $\$$ , all internal nodes of  $\text{PSTree}(\bar{T})$  become branching.

Recall that each node of  $\text{PSTree}(\bar{T})$  is a pv-string. The *string depth* of a node  $x$  of  $\text{PSTree}(\bar{T})$  is  $|x|$ . We store the string depth  $|x|$  in each node  $x$  of  $\text{PSTree}(\bar{T})$ . To represent  $\text{PSTree}(\bar{T})$  in linear space, in an actual implementation, the label  $y$  of an edge  $(x, y, xy)$  is represented by two integers  $i$  and  $j$  such that  $y = \langle \bar{T}[i - |x| : j] \rangle[|x| + 1 : j - i + |x| + 1]$ . In other words,  $y$  is the length- $(j - i + 1)$  suffix of the pv-encoding of  $\langle \bar{T}[i - |x| : j] \rangle = xy$  which labels the path from the root to the node  $xy$ .

### 6.1.2 Suffix links of p-suffix trees

Let us recall the *suffix links* of the p-suffix tree, which were first introduced by Baker [5]. Let  $u$  be a non-root node of  $\text{PSTree}(\bar{T})$  and let  $S$  be any substring of  $\bar{T}$  such that  $\langle S \rangle = u$ . The suffix link of  $u$ , denoted  $\text{sl}(u)$ , is a pointer from  $u$  to  $v = \langle S[2 : |S|] \rangle$ . Notice that  $v$  may not be a node of the p-suffix tree. For instance, see Figure 2(a). Consider node  $u = \infty a$ , in which we have  $S = \text{ya}$  or  $S = \text{xa}$ . In either case  $v = \langle S[2 : 2] \rangle = \langle a \rangle = a$ , however, there is no node representing  $a$ . In this paper, we define the suffix link of a node  $u = \langle S \rangle$  only if  $v = \langle S[2 : |S|] \rangle$  is a node.

We can characterize the target node  $v$  of the suffix link  $\text{sl}(u)$  depending on the first character  $u[1]$  of  $u$ , as follows:

- (i) If  $u[1] \in \Sigma$ , then  $v = u[2 : |u|]$ . This is the same as the suffix link of a standard suffix tree for exact matching.
- (ii) If  $u[1] = \infty$  and there exists a position  $a$  in  $u$  such that  $u[a] = a - 1$ , then  $u[a]$  “points” to the first position of  $u$  in the prev-encoding. Notice that such a position  $a$  is unique and that  $a$  is the smallest position in  $S$  that is larger than 1 with  $S[1] = S[a]$ . Now, the target node  $v$  is  $u[2 : a - 1] \cdot \infty \cdot u[a + 1 : |u|]$ , that is obtained by removing  $u[1]$  from  $u$  and replacing  $u[a]$  with  $\infty$ .
- (iii) If  $u[1] = \infty$  and there is no position  $a$  in  $u$  such that  $u[a] = a - 1$ , then  $v = u[2 : |u|]$ .

For examples of suffix links, see Figure 2(a). The suffix link of node  $\mathbf{a}\infty\mathbf{a}$  points to node  $\infty\mathbf{a}$ , which is Case (i). The suffix link of node  $\infty\mathbf{a}2$  points to node  $\mathbf{a}\infty$ , which is Case (ii). The suffix link of node  $\infty\mathbf{a}\infty\mathbf{a}2$  points to node  $\mathbf{a}\infty\mathbf{a}2$ , which is Case (iii).

### 6.1.3 Wiener links of p-suffix trees

We enhance p-suffix trees by introducing *Weiner links*, which are key tools in this whole section. The *reversals* of the suffix links of Section 6.1.2 are explicit Wiener links. In Case (i), the corresponding explicit Wiener link is labeled with the static character  $u[1] \in \Sigma$ . In Case (iii), the corresponding explicit Wiener link is labeled with  $u[1] = \infty$ . In Case (ii), we label the corresponding explicit Wiener link with the position  $a$ . In a unified manner for all these three cases, we can represent each explicit Wiener link by a triple  $(v, a, u)$ , where  $a \in \Sigma \cup \mathcal{N}$  and  $|u| = |v| + 1$ .

Now, we are to extend the notion of Wiener links to the case where there is no node  $u$  with  $|u| = |v| + 1$  for a given node  $v$  and  $a \in \Sigma \cup \mathcal{N}$ . In such a case, we use the shortest (i.e. shallowest) possible node as  $u$ , as follows. Let  $v$  be a node in  $\text{PSTree}(\bar{T})$  such that  $v = \langle S \rangle$  for some substring  $S$  of  $\bar{T}$ , and  $a \in \Sigma \cup \mathcal{N}$ . Let  $\alpha(a, v)$  be the pv-string such that

$$\alpha(a, v) = \begin{cases} av & \text{if } a \in \Sigma \cup \{\infty\} \text{ and } av \in \text{PFactor}(T), \\ \langle S[a] \cdot S \rangle & \text{if } a \in \mathcal{N} \setminus \{\infty\} \text{ and } \langle S[a] \cdot S \rangle \in \text{PFactor}(T), \\ \text{undefined} & \text{otherwise.} \end{cases} \quad (1)$$

This function  $\alpha$  corresponds to the reversals of the suffix links. When  $a \in \Sigma$ , it “prepends” label  $a$  to string (node)  $v$ . When  $a \in \mathcal{N}$ , it gives the pv-encoding of the p-string which is obtained by prepending the parameter character indicated by  $a$  to the p-string  $S$  whose pv-encoding is  $v$ . For  $a = \infty$ , the indicated parameter character is a fresh one that occurs nowhere in  $S$ . For  $a \in \mathcal{N} \setminus \{\infty\}$ , the parameter character is  $S[a]$ . Then, the Wiener link from node  $v$  to node  $u = \lceil [\alpha(a, v)]_{\langle \bar{T} \rangle}^L \rceil$  is labeled with  $a$ , which is represented by the triple  $(v, a, u)$ . Note that the operator  $\lceil [\cdot]_{\langle \bar{T} \rangle}^L \rceil$  brings us to the shallowest node  $u$  from the locus for  $\alpha(a, v)$  in  $\text{PSTree}(\bar{T})$ . Thus, each Wiener link increases the string depth by at least one, and hence  $|u| \geq |v| + 1$  always holds.

The Wiener link  $(v, a, u)$  is said to be *explicit* if  $|u| = |v| + 1$  (or equivalently  $u = \alpha(a, v)$ ), and *implicit* if  $|u| > |v| + 1$ . Namely,  $u = \alpha(a, v)$  if and only if  $v$  is obtained by simply removing the first character  $a$  from  $u = av$  (the first case in equation (1)), or by replacing  $u[a + 1] = a$  with  $\infty$  and removing the first character  $\infty$  from  $u$  (the second case in equation (1)). Basically the same arguments hold for implicit Wiener links, except in that we need to cut off the suffix of  $u$  to adjust the length to  $|v| + 1$ .

For examples of Wiener links of a p-suffix tree, see Figure 2(b). The node  $v = \mathbf{a}\infty$  has an explicit Wiener link which is labeled 2 and points to the node  $u = \infty\mathbf{a}2$ . This is because by replacing  $u[2+1] = 2$  with  $\infty$  and by removing the first character  $\infty$  from  $u$ , we obtain  $v = \mathbf{a}\infty$ . This corresponds to the second case in equation (1). For another example, consider the node  $u' = \infty\mathbf{a}\infty\mathbf{a}2$  which has three in-coming Wiener links all labeled  $\infty$ . The Wiener link from the node  $v' = \mathbf{a}\infty\mathbf{a}2$  is explicit, while the other two from the node  $v'' = \mathbf{a}\infty\mathbf{a}$  and  $v''' = \mathbf{a}\infty$  are implicit. Note that all these Wiener links correspond to the first case in equation (1), namely, one can obtain  $v'$  by simply removing the first  $\infty$  from  $u$ , and can obtain  $v''$  and  $v'''$  by removing the first  $\infty$  from  $u$  and removing the suffixes of  $u$  accordingly.

### 6.1.4 Duality between PDAWG and p-suffix trees and bidirectional searches

Our Weiner links for  $\text{PSTree}(\overline{T})$  permit us to design a Weiner-style [20] right-to-left online construction of  $\text{PSTree}(\overline{T})$ , which turns out to be equivalent to our left-to-right online construction of  $\text{PDAWG}(T)$  to be presented in Section 7. This observation is based on the following duality between  $\text{PSTree}(\overline{T})$  and  $\text{PDAWG}(T)$ .

To establish the correspondence between  $\text{PDAWG}(T)$  and  $\text{PSTree}(\overline{T})$ , we define the *p-reverse*  $\tilde{x}$  of a pv-string  $x$  so that  $\tilde{x} = \langle \overline{S} \rangle$  iff  $x = \langle S \rangle$  for any p-string  $S \in (\Sigma \cup \Pi)^*$ . For example, for  $T = \mathbf{xaxy}$  with  $\mathbf{a} \in \Sigma$  and  $\mathbf{x}, \mathbf{y} \in \Pi$ , we have  $\langle \overline{T} \rangle = \langle \overline{\mathbf{a}2\mathbf{x}\mathbf{a}2\mathbf{y}} \rangle = \langle \mathbf{y}\mathbf{x}\mathbf{a}\mathbf{x} \rangle = \langle \overline{T} \rangle$ .

For technical convenience, we rename the nodes  $[x]_{\langle T \rangle}^R$  of  $\text{PDAWG}(T)$  to be  $\lceil [x]_{\langle T \rangle}^R \rceil$  in this section. Moreover, we call an edge  $(x, a, y)$  of  $\text{PDAWG}(T)$  *primary* if  $y = xa$ , and *secondary* otherwise.

**Theorem 5.** *The following correspondence between  $\text{PDAWG}(T) = (V_D, E_D)$  and  $\text{PSTree}(\overline{T}) = (V_T, E_T)$  holds.*

- (1)  $\text{PDAWG}(T)$  has a node  $x \in V_D$  iff  $\text{PSTree}(\overline{T})$  has a node  $\tilde{x} \in V_T$ .
- (2)  $\text{PDAWG}(T)$  has a primary edge  $(x, a, y) \in E_D$  iff  $\text{PSTree}(\overline{T})$  has an explicit Weiner link  $(\tilde{x}, a, \tilde{y})$ .
- (3)  $\text{PDAWG}(T)$  has a secondary edge  $(x, a, y) \in E_D$  iff  $\text{PSTree}(\overline{T})$  has an implicit Weiner link  $(\tilde{x}, a, \tilde{y})$ .
- (4)  $\text{PDAWG}(T)$  has a suffix link from  $xy$  to  $y$  iff  $\text{PSTree}(\overline{T})$  has an edge  $(\tilde{y}, \tilde{x}, \tilde{xy}) \in E_T$ .

*Proof.* To make the arguments simpler, we assume for now that  $T$  begins with a unique character  $\$$  that does not occur elsewhere in  $T$ . The case without  $\$$  can be shown similarly.

- (1) By the symmetry  $\text{RPos}_{\langle T \rangle}(x) = \{n + 1 - k \mid k \in \text{LPos}_{\langle \overline{T} \rangle}(\tilde{x})\}$ , we have  $x = \lceil [x]_{\langle T \rangle}^R \rceil$  if and only if  $\tilde{x} = \lceil [\tilde{x}]_{\langle \overline{T} \rangle}^L \rceil$ . To see why this holds more intuitively, let  $\text{parent}(\tilde{x})$  be the parent of  $\tilde{x}$  in  $\text{PSTree}(\overline{T})$ , and let  $\ell$  be the edge label from  $\text{parent}(\tilde{x})$  to  $\tilde{x}$ . Then, for any locus on this edge representing  $\tilde{z}_i = \text{parent}(\tilde{x}) \cdot \ell[1 : i]$ , with  $1 \leq i \leq |\ell|$ , there are the same leaves below it. Since each leaf of  $\text{PSTree}(\overline{T})$  corresponds to a distinct position in  $\langle \overline{T} \rangle$ , every  $\tilde{z}_i$  has the same set of beginning positions in  $\langle \overline{T} \rangle$  (note that  $\tilde{z}_\ell = \tilde{x}$ ). By symmetry, this in turn means that  $z_i$  has the same set of ending positions in  $\langle \overline{T} \rangle = \langle T \rangle$ , i.e.  $\{z_i \mid 1 \leq i \leq |\ell|\} = [x]_{\langle T \rangle}^R$ . The other way (from PDAWG nodes to p-suffix tree nodes) can be shown analogously.
- (2) Because  $(\tilde{x}, a, \tilde{y})$  is an explicit Weiner link,  $\alpha(a, \tilde{x}) = \tilde{y}$ . By the definition of operator  $\sim$ , we obtain  $xa = y$ . Hence there is a primary edge from node  $x$  to  $y$  labeled  $a$  in  $\text{PDAWG}(T)$ . The other way (from PDAWG primary edges to p-suffix tree explicit Weiner links) can be shown analogously.
- (3) Similar to (2).
- (4) Immediately follows from the proof for (1).

□

Properties (2) and (3) imply that we can use the Weiner links of  $\text{PSTree}(\overline{T})$  for parameterized matching for the forward string  $T$  as in Algorithm 3. Further, we obtain the following corollary that allows for bidirectional parameterized pattern searches:

**Corollary 2.** *Using a pair of  $\text{PDAWG}(T)$  and  $\text{PSTree}(\overline{T})$  with pv-encodings  $\langle T \rangle$  and  $\langle \overline{T} \rangle$ , one can perform forward and backward search to find all substrings of  $T$  that p-match a given pattern  $P$  in  $O(m \log(|\Pi| + |\Sigma|) + \text{occ})$  time, where  $m$  is the length of pattern  $P$  and  $\text{occ}$  is the number of occurrences to report.*

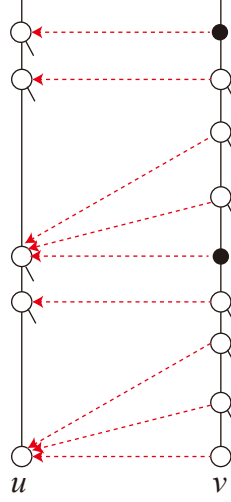


Figure 3: Illustration for our algorithm of Theorem 6 that propagates the Weiner links between the ancestors of  $v$  and  $u$ , in a bottom up manner. The white circles represent nodes of  $\text{PSTree}(\bar{T})$  and the black circles represent imaginary nodes. The red arcs represent Weiner links.

*Proof.* For a query pattern  $P$  that grows in both directions, one can easily maintain  $\langle P \rangle$  in  $O(\log |\Pi|)$  time per added character using  $O(|\Pi|)$  working space.

It is known that  $\text{PSTree}(\bar{T})$  allows for the navigation of a (reversed) pattern  $\bar{P}$  of length  $m$  in  $O(m \log(|\Pi| + |\Sigma|))$  time [5]. This can be translated to an amortized  $O(\log(|\Pi| + |\Sigma|))$ -time navigation per input character that is added to the *left end* of  $P$  (which is the right-end of  $\bar{P}$ .) A symmetric argument holds for our  $\text{PDAWG}(T)$ , which leads to an amortized  $O(\log(|\Pi| + |\Sigma|))$ -time navigation per input character that is added to the *right end* of  $P$ .

After locating the locus for the whole pattern  $P$ , we can report all the *occ* occurrences in  $O(\text{occ})$  time using Theorem 3.  $\square$

## 6.2 Offline construction of PDAWGs via p-suffix trees

In this section, we present a fast offline construction algorithm for  $\text{PDAWG}(T)$ , provided that  $\text{PSTree}(\bar{T})$  (without suffix links) has already been built.

By the definition of our Weiner links on parameterized suffix trees, the following monotonicity holds.

**Lemma 8.** *Suppose that a node  $v$  in  $\text{PSTree}(\bar{T})$  has an (implicit or explicit) Weiner link  $(v, k, u)$  with label  $k \in \Sigma \cup \mathcal{N}$ . Then, any ancestor  $v'$  of  $v$  has an (implicit or explicit) Weiner link  $(v', \langle\langle k \rangle\rangle_{|v'|}, u')$  where  $u'$  is the shallowest ancestor of  $u$  with  $|u'| \geq |v'| + 1$ .*

It follows from Theorem 5 and Lemma 8 that there is a simple *offline* algorithm that builds the PDAWG by computing Weiner links in a bottom-up manner over the PST for the reversed text string.

**Theorem 6.** *Let  $T$  be a p-string of length  $n$ . Given  $\text{PSTree}(\bar{T})$  (without suffix links), we can build  $\text{PDAWG}(T)$  in  $O(n)$  time and space.*

*Proof.* We first compute the (reversed) suffix links of the nodes of  $\text{PSTree}(\bar{T})$  that corresponds to  $\langle S \rangle$  for all suffixes of  $T$ , together with their labels which are the characters of  $\langle T \rangle$ .

For each Weiner link  $L = (v, k, u)$  between two nodes corresponding to two consecutive suffixes of  $T$ , perform the following:

- (1) Perform the following  $\text{update}(v, u)$  function:



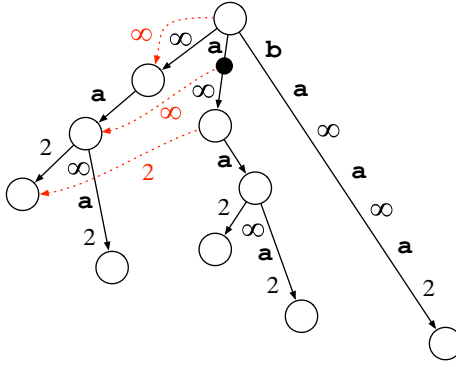


Figure 4: An example for how the Weiner links are propagated in a bottom-up manner, with the same string as in Figure 2. We pick the Weiner link  $(a\infty, 2, \infty a2)$  between two nodes  $a\infty = \langle ay \rangle$  and  $\infty a2 = \langle yay \rangle$ . An imaginary node (black circle) is created at string depth  $|\infty a| - 1 = 1$ . By removing the imaginary nodes and the Weiner links from them, we obtain  $\text{PDAWG}(T)$  (see Figure 2).

- (a) If  $|\text{parent}(u)| = \text{parent}(v) + 1$ , set  $v \leftarrow \text{parent}(v)$  and  $u \leftarrow \text{parent}(u)$ .
  - (b) If  $|\text{parent}(u)| < \text{parent}(v) + 1$ , set  $v \leftarrow \text{parent}(v)$ .
  - (c) If  $|\text{parent}(u)| > \text{parent}(v) + 1$ , then create an imaginary node  $w$  at string depth  $|\text{parent}(u)| - 1$  between  $\text{parent}(v)$  and  $v$ . Set  $v \leftarrow w$  and  $u \leftarrow \text{parent}(u)$ .
- (2) (a) If there is no Weiner link between  $v$  and  $u$ , create a new Weiner link  $(v, \langle\langle k \rangle\rangle_{|v|}, u)$ . Go to (1).
  - (b) Otherwise, stop the propagation for  $L$ .

See also Figure 3 that illustrates our algorithm. The correctness of this algorithm is immediate from Lemma 8.

It is clear from Lemma 8 that the complexity of this algorithm is linear in the number of Weiner links created. It follows from Theorem 4 and Theorem 5 that the number of Weiner links  $(v, k, u)$  such that  $v$  is *not* an imaginary node is  $O(n)$ . It also follows from our duality discussion in Section 6.2 and the definition of  $\text{pPDAWG}(T)$ , that each Weiner link  $(w, k, u)$  such that  $w$  is an imaginary node corresponds to an edge in  $E' \setminus E$ , where  $E'$  is the set of edges of  $\text{pPDAWG}(T)$  and  $E$  is the set of edges of  $\text{PDAWG}(T)$ . Since  $|E'| = O(n)$  by Lemma 1, the total time complexity of this algorithm is  $O(n)$ .  $\square$

See also Figure 4 for a concrete example of our offline algorithm that computes  $\text{PDAWG}(T)$  and  $\text{pPDAWG}(T)$  from  $\text{PSTree}(\overline{T})$ .

## 7 Online algorithm for constructing PDAWGs

This section proposes an algorithm constructing PDAWGs online. Our algorithm is based on the one by Blumer et al. [16] for constructing DAWGs of static strings. In fact, if the input string is static, the behavior of our algorithm coincides with theirs. In this section, we refer to  $\text{PDAWG}(T)$  as  $\text{PDAWG}(\langle T \rangle)$  for  $T \in (\Sigma \cup \Pi)^*$  and consider updating  $\text{PDAWG}(w)$  to  $\text{PDAWG}(wa)$  for a pv-string  $wa$  where  $a \in \Sigma \cup \mathcal{N}$ . To distinguish sets of nodes, edges, and suffix links of  $\text{PDAWG}(w)$  and  $\text{PDAWG}(wa)$ , we add subscripts  $w$  and  $wa$  to respective sets.

We first consider the difference of the node sets  $V_w$  of  $\text{PDAWG}(w)$  and  $V_{wa}$  of  $\text{PDAWG}(wa)$ . Recall that  $V_w$  is a partition of  $\text{PFactor}(w)$ . Concerning strings in  $\text{PFactor}(wa) \setminus \text{PFactor}(w) \subseteq \text{PSuffix}(wa)$ , we create a new sink node  $[wa]_{wa}^R = \text{PFactor}(wa) \setminus \text{PFactor}(w)$  in  $\text{PDAWG}(wa)$  for which  $\text{RPos}_{wa}([wa]_{wa}^R) = \{|wa|\}$ . Other strings  $z$ , which are already in  $\text{PFactor}(w)$ , may or may not get a new end position  $|wa|$ , depending on whether they are in  $\text{PSuffix}(wa)$ . If all or no members of  $[z]_w^R$  get the new position,

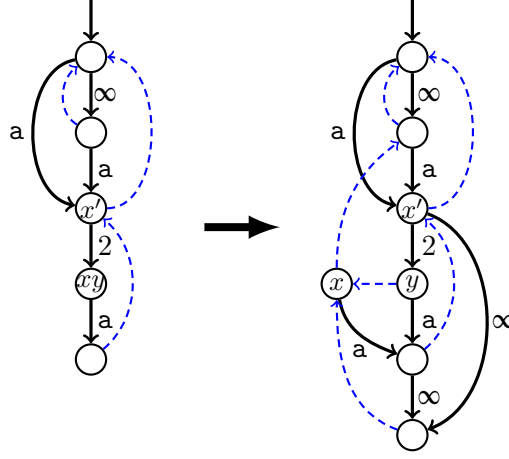


Figure 5: PDAWG( $w$ ) and PDAWG( $wa$ ) for  $w = \infty a 2 a$  and  $a = \infty$ . The strings  $\text{LRS}(wa) = x = a\infty$ ,  $\text{preLRS}(wa) = x' = a$ , and  $y = \lceil [x]_w^R \rceil = \infty a 2$  are shown on the respective nodes where they belong.

we have  $[z]_w^R = [z]_{wa}^R$ . We keep those nodes. However, it is possible that some but not all members of  $[z]_w^R$  get the new end position. In this case, the node  $[z]_w^R$  will be split into two. Let us define the *longest repeated suffix (LRS)* of  $wa \in (\Sigma \cup \mathcal{N})^+$  to be  $\text{LRS}(wa) = \lceil \text{PSuffix}(wa) \cap \text{PFactor}(w) \rceil$ . Then, a pv-string  $z \in \text{PFactor}(w)$  will have  $|wa| \in \text{RPos}_{wa}(z)$  if and only if  $z$  is a pv-suffix of  $\text{LRS}(wa)$ . That is, a node  $[z]_w^R$  will be split if and only if some of  $[z]_w^R$  are pv-suffix of  $\text{LRS}(wa)$  and some are not. In other words, the split node in  $V_w$  includes  $\text{LRS}(wa)$  and longer pv-strings. On the other hand, if  $\lceil [\text{LRS}(wa)]_w^R \rceil = \text{LRS}(wa)$ , then no node will be split in the update. We call  $[\text{LRS}(wa)]_w^R \in V_w$  the *LRS node* (w.r.t.  $wa$ ). The following lemma for node splits on PDAWGs is an analog to that for DAWGs.

**Lemma 9** (Node update). *For  $x = \text{LRS}(wa)$  and  $y = \lceil [x]_w^R \rceil$ ,*

$$V_{wa} = V_w \setminus \{[x]_w^R\} \cup \{[x]_{wa}^R, [y]_{wa}^R, [wa]_{wa}^R\}.$$

*If  $x = y$ , then  $[x]_w^R = [x]_{wa}^R = [y]_{wa}^R$ , i.e.,  $V_{wa} = V_w \cup \{[wa]_{wa}^R\}$ . Otherwise,  $[x]_w^R = [x]_{wa}^R \cup [y]_{wa}^R$  and  $[x]_{wa}^R \neq [y]_{wa}^R$ .*

*Proof.* First remark that  $\text{RPos}_{wa}(z) = \text{RPos}_w(z) \cup \{|wa|\}$  for all  $z \in \text{PSuffix}(wa)$  and  $\text{RPos}_{wa}(z) = \text{RPos}_w(z)$  for all  $z \notin \text{PSuffix}(wa)$ . For those  $z \in \text{PSuffix}(wa) \setminus \text{PFactor}(w)$ , we have  $\text{RPos}_{wa}(z) = \{|wa|\}$  and  $[wa]_{wa}^R = \text{PSuffix}(wa) \setminus \text{PFactor}(w) \in V_{wa} \setminus V_w$ . For  $z \in \text{PFactor}(w)$ , if  $[z]_w^R \neq [z]_{wa}^R$ , some elements of  $[z]_w^R$  are in  $\text{PSuffix}(wa)$  and some are not. That is,  $[z]_w^R$  is partitioned into two non-empty equivalence classes  $\{z' \in [z]_w^R \mid z' \in \text{PSuffix}(wa)\}$  and  $\{z' \in [z]_w^R \mid z' \notin \text{PSuffix}(wa)\}$ . By definition, the longest of the former is  $x = \text{LRS}(wa)$  and the longest of the latter is  $y = \lceil [x]_w^R \rceil$ . Otherwise,  $[z]_w^R = [z]_{wa}^R \in V_w \cap V_{wa}$ .  $\square$

**Example 2** (Figure 5). Let  $w = \infty a 2 a$  and  $a = \infty$ . Then,  $\text{LRS}(wa) = \langle w[2 : 3] \rangle = \langle wa[4 : 5] \rangle = a\infty$ . We have  $\text{LRS}(wa) \neq \lceil [\text{LRS}(wa)]_w^R \rceil = \infty a 2$ , where  $\text{RPos}_w(a\infty) = \text{RPos}_w(\infty a 2) = \{3\}$ . On the other hand,  $\text{RPos}_{wa}(a\infty) = \{3, 5\} \neq \text{RPos}_{wa}(\infty a 2) = \{3\}$ . Therefore, PDAWG( $wa$ ) has two more nodes than PDAWG( $w$ ).

When updating PDAWG( $w$ ) to PDAWG( $wa$ ), all edges that do not touch the LRS node  $[\text{LRS}(wa)]_w^R$  are kept by definition. What we have to do is to manipulate incoming edges for the new sink node  $[wa]_{wa}^R$ , and, if necessary, to split the LRS node into two and to manipulate incoming and outgoing edges of them. Therefore, it is very important to identify the LRS node  $[\text{LRS}(wa)]_w^R$  and to decide whether  $\text{LRS}(wa) = \lceil [\text{LRS}(wa)]_w^R \rceil$ . We remark that  $\text{LRS}(wa) = \varepsilon$  if and only if  $wa \in \Sigma^* \{\infty\} \cup (\mathcal{N} \cup \Sigma \setminus \{a\})^* \Sigma$ . This special case where  $\text{LRS}(wa) = \varepsilon$  is easy to handle, since the LRS node will never be split by  $[\varepsilon]_w^R = \{\varepsilon\}$ . Hereafter, we assume that  $\text{LRS}(wa) \neq \varepsilon$  and define  $\text{preLRS}(wa)$  to be the prefix

of  $\text{LRS}(wa)$  obtained by removing the last character:  $\text{preLRS}(wa) = \text{LRS}(wa)[ : |\text{LRS}(wa) - 1| ]$ . Here,  $\text{LRS}(wa)$  and  $\text{preLRS}(wa)$  are pv-suffixes of  $wa$  and  $w$ , respectively. The LRS node can be reached from the pre-LRS node  $[\text{preLRS}(wa)]_w^R$  by reading one more character, and the pre-LRS node can be found by following suffix links from the sink node  $[w]_w^R$  of  $\text{PDAWG}(w)$ . This appears quite similar to online construction of DAWGs for static strings, but there are nontrivial differences. Main differences from the DAWG construction are in the following points:

- Our PDAWG construction uses  $\text{trans}_w(u, i, \langle\langle a \rangle\rangle_i)$  with an appropriate  $i$ , when the original DAWG construction refers to  $\text{child}_w(u, a)$ ,
- While  $\text{preLRS}(wa)$  is the longest of its equivalence class for static strings in  $\text{DAWG}(w)$ , it is not necessarily the case for p-strings (like the one in Figure 5). This affects the procedure to find the node of  $\text{LRS}(wa)$ . As a consequence, if  $\text{preLRS}(wa) \neq \lceil [\text{preLRS}(wa)]_w^R \rceil$  and the LRS node is split,  $\text{PDAWG}(wa)$  has no edge from  $[\text{preLRS}(wa)]_{wa}^R$  to  $[\text{LRS}(wa)]_{wa}^R$ . Even it is possible that  $[\text{LRS}(wa)]_{wa}^R$  has no incoming edges in  $\text{PDAWG}(wa)$ .
- When the LRS node of  $\text{PDAWG}(w)$  is split into two in  $\text{PDAWG}(wa)$ , the outgoing edges of the two obtained nodes are identical in the DAWG construction, while it is not necessarily the case anymore in the PDAWG construction.

Let us call a sequence  $\langle v_0, v_1, \dots, v_k \rangle$  of nodes of  $V_w$  the *suffix link chain* of  $v_0$  if  $v_{i+1} = F_w(v_i)$  for  $0 \leq i < k$  and  $v_k = \{\varepsilon\}$ . In DAWGs, the pre-LRS node is the first node with an  $a$ -edge on the suffix link chain of the old sink  $[w]_w^R$ . However, it is not necessarily the case for PDAWGs. Lemma 10 below suggests how to find the pre-LRS and LRS nodes, how to obtain the length of the LRS, and how to decide whether the LRS node should be split. The first item of the lemma describes how to find the pre-LRS node. The pre-LRS node is on the suffix link chain  $\langle [w]_w^R = u_0, u_1, u_2, \dots, [\varepsilon]_w^R \rangle$  of the old sink node  $[w]_w^R$ . It is the firstly found node  $u_j$  from which one can read  $\langle\langle a \rangle\rangle_{\lceil u_j \rceil}$ . The second item will be used to identify the length of the pre-LRS  $x'$  and the LRS  $x = x'a'$ , where  $a' = \langle\langle a \rangle\rangle_{|x'|}$ . If the pre-LRS node  $u_j$  has an edge labeled with  $\langle\langle a \rangle\rangle_{\lceil u_j \rceil}$ , then the longest element  $\lceil u_j \rceil$  is the pre-LRS  $x'$  and thus  $\langle \lceil u_j \rceil \cdot a \rangle = \lceil u_j \rceil \cdot \langle\langle a \rangle\rangle_{\lceil u_j \rceil}$  is the LRS. Otherwise, the lengths of the pre-LRS and the LRS can be determined by  $a$  and the largest number (or  $\infty$ ) labeling an outgoing edge of  $u_j$ . The third and fourth items are immediate consequences of Lemmas 7 and 9, respectively. We can reach the LRS node  $[x'a']_w^R$  from the pre-LRS node by the transition function  $\text{trans}_w(u_j, |x'|, a')$ . The LRS node shall be split if and only if the LRS is not the longest of the node. Hereafter, throughout this section, we fix the following variables:  $x = \text{LRS}(wa)$ ,  $x' = \text{preLRS}(wa)$ ,  $a' = \langle\langle a \rangle\rangle_{|x'|}$ , (i.e.,  $x = x'a'$ ),  $y = \lceil [x]_w^R \rceil$ ,  $u_0 = [w]_w^R$ ,  $u_i = F_w(u_{i-1})$  for  $i \geq 1$  as long as  $F_w(u_{i-1})$  is defined, and  $a_i = \langle\langle a \rangle\rangle_{\lceil u_i \rceil}$ .

**Lemma 10.** *We have*

1.  $x' \in u_j$  for the least  $j$  such that  $\text{trans}_w(u_j, \lceil u_j \rceil, \langle\langle a \rangle\rangle_{\lceil u_j \rceil}) \neq \text{Null}$ ,
2.  $|x'a'| = \begin{cases} \lceil u_j \rceil + 1 & \text{if } \text{child}_w(u_j, a_j) \neq \text{Null}, \\ \min\{a, \max(\text{Children}_w(u_j) \cap \mathcal{N})\} & \text{otherwise,} \end{cases}$   
for  $j$  such that  $x' \in u_j$ ,
3.  $[x'a']_w^R = \text{trans}_w(u_j, |x'|, a')$  for  $j$  such that  $x' \in u_j$ ,
4.  $[x'a']_w^R \neq [x'a']_{wa}^R$  if and only if  $|x'a'| \neq \lceil [x'a']_w^R \rceil$ .

*Proof.* Suppose  $x' \in u_j$ . Note that  $\text{PSuffix}(w) = \bigcup_{i \geq 0} u_i$ .

(1) Since every string  $z \in u_i \subseteq \text{PSuffix}(w)$  with  $i < j$  is properly longer than  $x'$ , particularly for  $z = \lceil u_i \rceil$ ,  $\langle za \rangle \in \text{PSuffix}(wa)$  is properly longer than the LRS  $x = x'a'$ . Thus,  $\langle za \rangle = z\langle\langle a \rangle\rangle_{|z|} \notin \text{PFactor}(w)$ .

<sup>2</sup>This corrects an error in the conference version [18].

Therefore,  $\text{trans}_w(u_i, \lfloor u_i \rfloor, \langle a \rangle_{\lfloor u_i \rfloor}) = \text{Null}$ . On the other hand, the fact  $\lfloor u_j \rfloor \in \text{PSuffix}(x')$  implies  $\langle \lfloor u_j \rfloor a \rangle = \lfloor u_j \rfloor \langle a \rangle_{\lfloor u_j \rfloor} \in \text{PSuffix}(x'a') \subseteq \text{PFactor}(w)$ . Hence,  $\text{trans}_w(u_j, \lfloor u_j \rfloor, \langle a \rangle_{\lfloor u_j \rfloor}) \neq \text{Null}$ .

(2) Recall that for any element  $z$  of the pre-LRS node  $u_j$ ,  $\langle za \rangle \in \text{PFactor}(w)$  if and only if  $|\langle za \rangle| \leq |x'a'|$ . Hence,  $\text{child}_w(u_j, a_j) \neq \text{Null}$  if and only if the longest element  $\lfloor u_j \rfloor$  of the pre-LRS node  $u_j$  is the pre-LRS, i.e.,  $|x'a'| = \lfloor u_j \rfloor + 1$ .

Now suppose  $\text{child}_w(u_j, a_j) = \text{Null}$ , but  $\text{trans}_w(u_j, \lfloor u_j \rfloor, \langle a \rangle_{\lfloor u_j \rfloor}) \neq \text{Null}$ . Then,  $\lfloor u_j \rfloor$  is properly longer than  $x'$ . Let  $x'' \in u_j$  be the element of  $u_j$  of length  $|x'| + 1$  and  $a'' = \langle a \rangle_{|x''|}$ . Since  $x'a'$  is the LRS, the longer p-suffix  $x''a'' \in \text{PSuffix}(wa)$  does not occur in  $w$ , so  $x'a' \not\equiv_w^R x''a''$ , whereas  $x' \equiv_w^R x''$ . Here, one can see that  $a' = \infty$  as follows. If  $a' = a''$ , then  $a' = \infty$  by Lemma 2. If  $a' \neq a''$ , then it can happen only when  $|x'| < a'' = a < \infty$  and  $a' = \infty$ . Let  $Z = \text{Children}_w(u_j) \cap \mathcal{N}$  and  $Z_i = \{k \in Z \mid k > i\}$ . Lemma 7 and the fact  $x'a' \in \text{PFactor}(w)$  imply  $Z_{|x'|} \neq \emptyset$ . On the other hand,  $x''a'' \notin \text{PFactor}(w)$  implies either  $a'' = \infty$  and  $Z_{|x''|} = \emptyset$  or  $0 < a'' = a < \infty$ . In the former case,  $Z_{|x'|} \neq Z_{|x''|} = \emptyset$  implies that  $Z_{|x'|}$  is the singleton set with the element  $\max Z = |x''|$ . By  $a' = \langle a \rangle_{|x'|} = \infty$ , we have  $|x'| < a$ . That is,  $\max Z = |x''| = |x'a'| \leq a$ . In the latter case,  $a' = \langle a \rangle_{|x'|} = \infty \neq a'' = \langle a \rangle_{|x''|} \in \mathbb{N}$  implies  $a = |x''|$ . By  $Z_{|x'|} \neq \emptyset$ , we have  $|x'| < \max Z$ . That is,  $a = |x''| = |x'a'| \leq \max Z$ . Summarizing these cases,  $|x'a'| = \min\{a, \max Z\}$ .

(3) By Lemma 7. (4) By Lemma 9.  $\square$

To update the PDAWG based on Lemma 10, we need to know the lengths of the longest and shortest elements of each node. Accordingly, our algorithm maintains the value  $\text{len}(u) = \lfloor u \rfloor$  for each node  $u$ . The length of the shortest element can be calculated by  $\lfloor u \rfloor = \lfloor F(u) \rfloor + 1 = \text{len}(F(u)) + 1$ .

**Example 3.** See Figure 5, where  $wa = \infty a 2 a \infty$ , and the pre-LRS and LRS are  $x' = a$  and  $x = a\infty$ , respectively. The pre-LRS node  $[x']_w^R$  is the first node  $u$  on the suffix link chain of the old sink  $[w]_w^R$  such that  $\text{trans}_w(u, \lfloor u \rfloor, \infty) \neq \text{Null}$ . This is how we find the pre-LRS node by Lemma 10.1. The lengths of the pre-LRS and LRS can be known by Lemma 10.2. In this case,  $\text{child}_w([x']_w^R, \infty) = \text{Null}$ . Hence,  $|x| = \min\{\infty, \max \text{Children}_w([x']_w^R) \cap \mathcal{N}\} = 2$  and thus  $|x'| = 1$ . From the identified pre-LRS node, one can reach the LRS node by  $\text{trans}_w([x']_w^R, 1, \infty) = [x]_w^R$ . Here,  $\lfloor [x']_w^R \rfloor = 3 \neq |x| = 2$ . So, the LRS node should be split.

Figure 6 shows another example, with a step-by-step illustration of our algorithm, which will be explained in more detail later. Compare the initial PDAWG (a) and the goal PDAWG (d). For  $wa = \infty a 2 a \infty a a$ , the pre-LRS and LRS are  $x' = \varepsilon$  and  $x = a$ , respectively. Indeed, the source node  $[\varepsilon]_w^R$  is the only node on the suffix link chain of the old sink  $[w]_w^R$  that has an  $a$ -edge. Since  $\text{child}_w([\varepsilon]_w^R, a) \neq \text{Null}$ , we know  $|x| = \lfloor [\varepsilon]_w^R \rfloor + 1 = 1$  and  $|x'| = 0$ . The LRS node can be found by  $\text{trans}_w([\varepsilon]_w^R, 0, a) = [x]_w^R$ . Since  $\lfloor [x]_w^R \rfloor = 2 > |x| = 1$ , the LRS node should be split.

Edges are created or replaced in accordance with the definition of a PDAWG. The incoming edges for the new sink node  $[wa]_{wa}^R$  of  $\text{PDAWG}(wa)$  are given as Lemma 11, except for when we have an edge from  $[\text{LRS}(wa)]_{wa}^R$  to the new sink in the case of node split. This case will be discussed later in Lemma 12.

**Lemma 11** (Incoming edges of the new sink). *Let  $u'_i = \lfloor u_i \rfloor_{wa}^R$ , which coincides with  $u_i$  unless  $u_i$  is the LRS node and split. The incoming edges for the new sink  $[wa]_{wa}^R$  are exactly those  $(u'_i, a_i, [wa]_{wa}^R)$  such that  $\text{child}_w(u_i, a_i) = \text{Null}$ , except for an edge from  $[\text{LRS}(wa)]_{wa}^R$  in the case of node split.*

*Proof.* Recall that  $\text{RPos}_{wa}(wa) = \{|wa|\}$ . For  $[wa]_{wa}^R$  to have an incoming edge from a node  $u \in V_{wa} \cap V_w$ ,  $\lfloor u \rfloor$  must occur as a pv-suffix of  $w$ . Then, the node  $u$  must be in the suffix link chain of the old sink node, i.e.,  $u = u_i$  for some  $i$ . Here,  $\lfloor u_i \rfloor a_i \in \text{PSuffix}(wa)$ . On the other hand,  $\lfloor u_i \rfloor a_i$  should not occur anywhere in  $w$ , since it occurs only at  $|wa|$  in  $wa$ . That is,  $u_i$  should not have an edge  $a_i$  in  $\text{PDAWG}(w)$ . Therefore, we must have  $(u_i, a_i, [wa]_{wa}^R) \in E_{wa}$  if and only if  $\text{child}_w(u_i, a_i) = \text{Null}$ .

It remains to consider whether we should have  $([y]_{wa}^R, \langle a \rangle_{|y|}, [wa]_{wa}^R) \in E_{wa}$  in the case of node split. This can be discussed in essentially the same way as above. We have  $([y]_{wa}^R, \langle a \rangle_{|y|}, [wa]_{wa}^R) \in E_{wa}$  if and only if  $[y]_w^R = u_i$ , i.e.,  $[y]_{wa}^R = u'_i$ , for some  $i$  and  $\text{child}_w(u_i, a_i) = \text{Null}$ .  $\square$

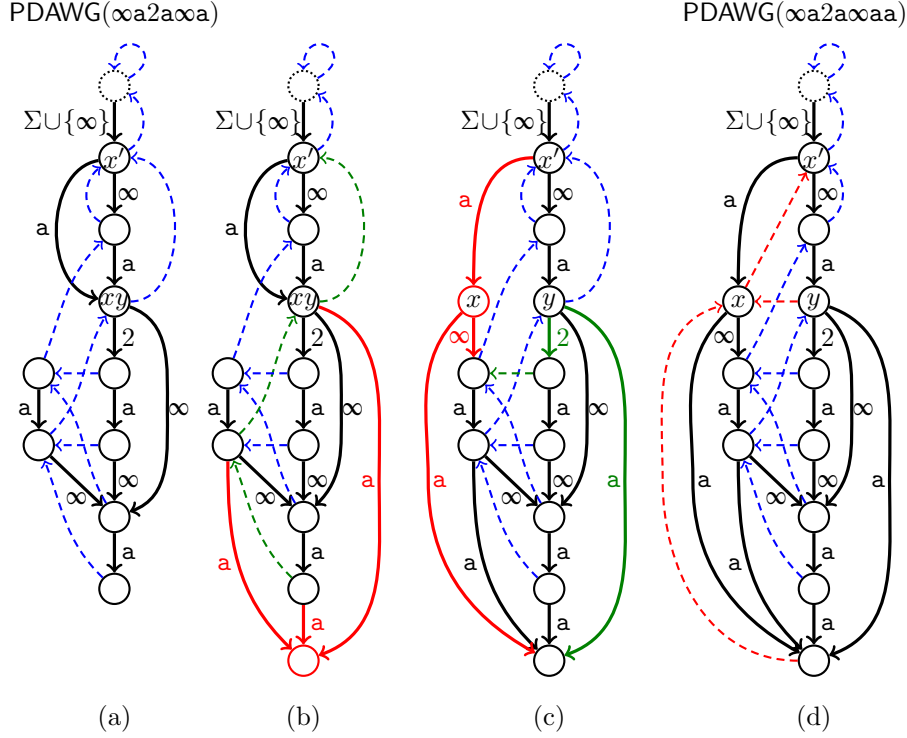


Figure 6: Updating  $\text{PDAWG}(w)$  (a) to  $\text{PDAWG}(wa)$  (d) for  $w = \langle \mathbf{xaxaya} \rangle$  and  $wa = \langle \mathbf{xaxayaa} \rangle$ . The subfigures (b) and (c) illustrate the structures right after Lines 10 and 19 of Algorithm 4, respectively. The circle with a broken line represents the dummy node  $\top$ . We present  $x = \text{LRS}(wa) = \mathbf{a}$ ,  $y = \lceil [\text{LRS}(wa)]_w^R \rceil = \infty \mathbf{a}$ , and  $x' = \text{preLRS}(wa) = \varepsilon$  in the respective nodes where they belong. Algorithm 4 modifies (a) into (b) by creating the new sink, adding incoming  $\mathbf{a}$ -edges to the new sink from the nodes reachable by suffix links from the old sink, until arriving at a node that already has an  $\mathbf{a}$ -edge. This node is the pre-LRS node  $[x']_w^R$ . In (b), the newly added node and edges are colored red and the suffix links we followed are green. We now know the lengths of the pre-LRS and the LRS are 0 and 1, respectively. The LRS node  $v$  is  $\text{trans}([x']_w^R, 0, \mathbf{a}) = [x]_w^R$ . Since  $\text{len}(v) > 1$ , the LRS node shall be split into  $[y]_{wa}^R$  and  $[x]_{wa}^R$ , as shown in (c). Our algorithm recycles the LRS node  $v$  for  $[y]_{wa}^R$  and create  $v'$  for  $[x]_{wa}^R$ . Then, the outgoing  $\mathbf{a}$ -edge of the pre-LRS node is redirected from  $v$  to  $v'$ . The  $\mathbf{a}$ -edge of  $v$  is copied for  $v'$ , which has just been added in the previous step in this particular example. In addition,  $v'$  gets an  $\infty$ -edge toward  $\text{trans}_w(v, 1, \infty)$ . The new node  $v'$  and its incoming and outgoing edges are colored red, and the edges and the suffix link referenced for making those outgoing edges are colored green in (c). Finally, we obtain (d) by creating suffix links from the new sink to the new node, and from the new node  $v'$  to the node that the old suffix link of  $v$  points at, and by redirecting that of  $v$  toward  $v'$ . Those suffix links are colored red.

This is not much different from DAWG update, except that the pre-LRS node has an edge toward the new sink when the pre-LRS is not the longest in the pre-LRS node.

**Example 4.** In Figure 5, the new sink gets an incoming edge from the pre-LRS node  $[x']_{wa}^R$  in addition to the one from the old sink. In Figure 6, the new sink gets incoming edges from the nodes on the suffix link chain of the old sink located *before* the pre-LRS node, including the one from  $[y]_{wa}^R$ .

If the LRS node  $[\text{LRS}(wa)]_{wa}^R$  is not split, we have nothing more to do on edges.

Hereafter, we suppose that the LRS node shall be split. That is,  $x \neq y$  for  $x = \text{LRS}(wa)$  and  $y = \lceil [\text{LRS}(wa)]_w^R \rceil$ . By definition, all edges of  $\text{PDAWG}(w)$  that do not involve the LRS node  $[\text{LRS}(wa)]_w^R$

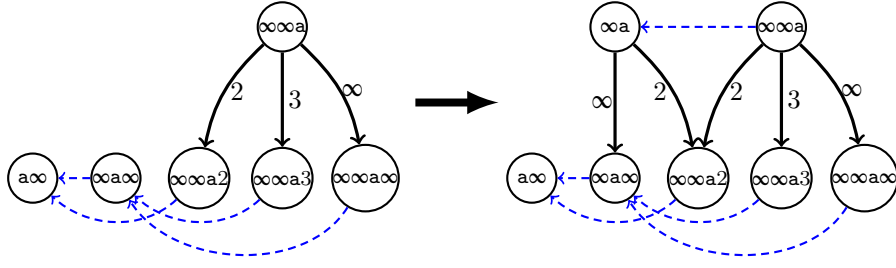


Figure 7: Subgraphs of  $\text{PDAWG}(w)$  and  $\text{PDAWG}(wa)$  for  $w = \infty\infty a3\infty a2\infty a\infty$  and  $a = a$ , where  $\text{LRS}(wa) = \infty a$  and  $\lceil [\infty a]_w^R \rceil = \infty\infty a$ . For explanation, we show  $\lceil u \rceil$  in each node  $u$ .  $[\infty a]_{wa}^R$  does not inherit the outgoing edges of  $[\infty a]_w^R = [\infty\infty a]_w^R$  labeled with 3 and  $\infty$ . Instead, the 3-edge and  $\infty$ -edge are bundled into a single  $\infty$ -edge which points at  $\text{trans}_w([\infty a]_w^R, 2, \infty) = F([\infty\infty a3]_w^R) = [\infty a\infty]_w^R$ . On the other hand, the 2-edge of  $[\infty\infty a]_w^R$  is simply copied for  $[\infty a]_{wa}^R$ .

will be inherited to  $\text{PDAWG}(wa)$ . First, we make it clear that we will have no edge between  $[x]_{wa}^R$  and  $[y]_{wa}^R$ . To have  $([x]_{wa}^R, b, [y]_{wa}^R) \in E_{wa}$  for some  $b \in \Sigma \cup \mathcal{N}$ , it requires  $\text{RPos}_{wa}(xb) = \text{RPos}_{wa}(y)$ . Particularly for  $i = \min \text{RPos}_{wa}(xb)$ , we have  $i-1 \in \text{RPos}_{wa}(x)$ , which contradicts  $i = \min \text{RPos}_{wa}(y) = \min \text{RPos}_w(y) = \min \text{RPos}_w(x) = \min \text{RPos}_{wa}(x)$ . Similarly we conclude  $([y]_{wa}^R, b, [x]_{wa}^R) \notin E_{wa}$ .

Lemma 12 below is concerned with the outgoing edges of  $[y]_{wa}^R$  and  $[x]_{wa}^R$ . In the DAWG construction, those two nodes  $[x]_{wa}^R$  and  $[y]_{wa}^R$  simply inherit the outgoing edges of the LRS node  $[x]_w^R = [y]_w^R$ . However, in the PDAWG construction, due to the prev-encoding rule on parameter characters, the node  $[x]_{wa}^R$  will lose edges whose labels are integers greater than  $|x|$ , as demonstrated in Figure 7. Those edges are “bundled” into a single  $\infty$ -edge which points at  $\text{trans}_w([y]_w^R, |x|, \infty)$ . This is described as the second item of Lemma 12.

**Lemma 12** (Outgoing edges of the LRS node). *For  $u \in V_w \cap V_{wa}$ ,*

- $([y]_{wa}^R, b, u) \in E_{wa}$  if and only if  $([y]_w^R, b, u) \in E_w$ ,
- $([x]_{wa}^R, b, u) \in E_{wa}$  if and only if  $\text{trans}_w([y]_w^R, |x|, b) = u$  if and only if either  $([y]_w^R, b, u) \in E_w$  and  $\langle\langle b \rangle\rangle_{|x|} \neq \infty$  or  $\text{trans}_w([y]_w^R, |x|, \infty) = u$  and  $\langle\langle b \rangle\rangle_{|x|} = \infty$ .

Moreover,

- $([x]_{wa}^R, \langle\langle a \rangle\rangle_{|x|}, [wa]_{wa}^R) \in E_{wa}$  if and only if  $([y]_{wa}^R, \langle\langle a \rangle\rangle_{|y|}, [wa]_{wa}^R) \in E_{wa}$ .
- If  $([x]_{wa}^R, \infty, [wa]_{wa}^R) \in E_{wa}$ , then  $Z = \{i \in \mathcal{N} \cap \text{Children}_w([y]_w^R) \mid i > |x|\}$  is empty.

*Proof.* The claims on edges toward nodes  $u \in V_w$  is an immediate consequence of the definition of PDAWG edges and the soundness of the function  $\text{trans}$  (Algorithm 2 and Lemma 7). We prove the third and fourth claims.

Suppose  $([x]_{wa}^R, \langle\langle a \rangle\rangle_{|x|}, [wa]_{wa}^R) \in E_{wa}$ , i.e.,  $\text{RPos}_{wa}(\langle xa \rangle) = \{|wa|\}$ . Then,  $x \in \text{PSuffix}(w)$  and  $x \equiv_w^R y$  imply  $y \in \text{PSuffix}(w)$ . Hence,  $\langle ya \rangle \in \text{PSuffix}(wa)$ . Since  $\text{RPos}_{wa}(\langle ya \rangle) \subseteq \text{RPos}_{wa}(\langle xa \rangle)$ , we have  $\text{RPos}_{wa}(\langle ya \rangle) = \{|wa|\}$  and thus  $([y]_{wa}^R, \langle\langle a \rangle\rangle_{|y|}, [wa]_{wa}^R) \in E_{wa}$ . Conversely, suppose  $([y]_{wa}^R, \langle\langle a \rangle\rangle_{|y|}, [wa]_{wa}^R) \in E_{wa}$ , i.e.,  $\text{RPos}_{wa}(\langle ya \rangle) = \{|wa|\}$ . We have  $\langle xa \rangle \in \text{PSuffix}(\langle ya \rangle) \subseteq \text{PSuffix}(wa)$ . Since  $\langle xa \rangle$  is longer than the LRS  $x$ , it cannot occur in  $w$ . Thus,  $|wa|$  is the only end position of  $\langle xa \rangle$  in  $wa$ . Therefore,  $([x]_{wa}^R, \langle\langle a \rangle\rangle_{|x|}, [wa]_{wa}^R) \in E_{wa}$ .

Suppose  $([x]_{wa}^R, \infty, [wa]_{wa}^R) \in E_{wa}$  and  $Z \neq \emptyset$ . The fact  $([x]_{wa}^R, \infty, [wa]_{wa}^R) \in E_{wa}$  implies  $x\infty \in \text{PSuffix}(wa)$ . The fact  $Z \neq \emptyset$  implies, by Lemma 7,  $x\infty \in \text{PFactor}(w)$ , which contradicts that  $x$  is the LRS.  $\square$

Let us turn our attention to the incoming edges of the LRS node. Lemma 13 is no more than a direct implication of the definition of edges of PDAWGs. If a node  $u$  has an outgoing edge pointing

at the LRS node, the edge will point at  $[y]_{wa}^R$  in  $\text{PDAWG}(wa)$  if and only if  $\lceil u \rceil$  is longer than the pre-LRS  $x'$ . Otherwise, it will point at  $[x]_{wa}^R$  possibly changing the label to  $\infty$  if necessary.

**Lemma 13** (Incoming edges of the LRS node). *We have*

- $(u, b, [y]_{wa}^R) \in E_{wa}$  if and only if  $(u, b, [y]_w^R) \in E_w$  and  $|\lceil u \rceil| > |x'|$ ,
- $(u, b, [x]_{wa}^R) \in E_{wa}$  if and only if  $(u, b, [y]_w^R) \in E_w$  and  $|\lceil u \rceil| \leq |x'|$ , where  $b = \langle\langle a \rangle\rangle_{|\lceil u \rceil|}$  and  $\lceil u \rceil \in \text{PSuffix}(x')$ .

*Proof.* Recall that  $(u, b, [y]_w^R) \in E_w$  if and only if  $\lceil u \rceil b \in [y]_w^R = [x]_{wa}^R \cup [y]_{wa}^R$ . Here, for each  $(u, b, [y]_w^R) \in E_w$ , clearly  $\lceil u \rceil b \in [x]_{wa}^R$  if and only if  $|\lceil u \rceil| \leq |x'|$ . In this case,  $\lceil u \rceil b \in [x]_{wa}^R$  implies  $b = \langle\langle a \rangle\rangle_{|\lceil u \rceil|}$  and  $\lceil u \rceil b \in \text{PSuffix}(x')$ , so  $\lceil u \rceil \in \text{PSuffix}(x')$ .  $\square$

Therefore, the incoming edges of  $[x]_{wa}^R$  can be constructed by visiting nodes  $u$  on the suffix link chain of the pre-LRS node and checking their  $\langle\langle a \rangle\rangle_{|\lceil u \rceil|}$ -edges. Note that in the online construction of a DAWG, the edge from the pre-LRS node  $[x']_w^R$  to the LRS node  $[y]_w^R$  in the old DAWG will be inherited to the new node  $[x]_{wa}^R$  in the new DAWG, since always  $\lceil [x']_w^R \rceil = x'$  holds. However, it is not necessarily the case in the PDAWG construction.

**Example 5.** Let us come back to the example in Figure 5, where  $x = a\infty$ ,  $x' = a$ , and  $y = \infty a2$  for  $wa = \infty 2 a \infty a \infty$ . The outgoing edge of the LRS node  $[x]_w^R \in V_w$  is inherited to both  $[x]_{wa}^R$  and  $[y]_w^R$  in  $\text{PDAWG}(wa)$ .

The LRS node has one incoming edge in  $\text{PDAWG}(w)$ . Since  $|\lceil [x']_w^R \rceil| = |\infty a| > |x'| = |a|$ , the 2-edge from  $[x']_w^R$  to  $[y]_w^R$  in  $\text{PDAWG}(w)$  is kept as the 2-edge from  $[x']_{wa}^R$  to  $[y]_{wa}^R$  in  $\text{PDAWG}(wa)$  and, as a result, the new node  $[x]_{wa}^R$  has no incoming edges.

Recall that in Figure 6, we have  $x = a$ ,  $x' = \varepsilon$ , and  $y = \infty a$  for  $wa = \infty a 2 a \infty a a$ . The LRS node has two outgoing edges:  $\infty$ -edge and 2-edge. The new node  $[y]_{wa}^R \in V_{wa}$  inherits those two. In addition,  $[y]_{wa}^R$  gets an  $a$ -edge pointing at the new sink  $[wa]_{wa}^R$ . On the other hand,  $[x]_{wa}^R \in V_{wa}$  has the  $\infty$ -edge pointing at the node  $\text{trans}_w([x]_w^R, |x|, \infty)$ , but it has no 2-edge. In addition,  $[x]_{wa}^R$  gets an  $a$ -edge toward  $[wa]_{wa}^R$ , too.

The LRS node has two incoming edges labeled with  $a$ . The one from  $[\varepsilon]_w^R$  is inherited to  $[x]_{wa}^R$  by  $|\varepsilon| \leq |x'| = 0$ , and the one from  $[\infty]_w^R$  is inherited to  $[y]_{wa}^R$  by  $|\infty| = 1 > |x'|$ .

Updates of suffix links simply follow the definition.

**Lemma 14** (Suffix link update). *Suppose  $V_{wa} = V_w \cup \{[wa]_{wa}^R\}$ , i.e., no node split happens. Then, for each  $u \in V_{wa}$ ,*

$$F_{wa}(u) = \begin{cases} [\text{LRS}(wa)]_{wa}^R & \text{if } u = [wa]_{wa}^R, \\ F_w(u) & \text{otherwise.} \end{cases}$$

*Suppose  $[x]_w^R \neq [x]_{wa}^R$  for  $x = \text{LRS}(wa)$ , where  $V_{wa} = V_w \setminus \{[y]_w^R\} \cup \{[wa]_{wa}^R, [x]_{wa}^R, [y]_{wa}^R\}$  for  $y = \lceil [x]_w^R \rceil$ . Then, for each  $u \in V_{wa}$ ,*

$$F_{wa}(u) = \begin{cases} [x]_{wa}^R & \text{if } u \in \{[wa]_{wa}^R, [y]_{wa}^R\}, \\ F_w([y]_w^R) & \text{if } u = [x]_{wa}^R, \\ F_w(u) & \text{otherwise.} \end{cases}$$

*Proof.* Recall that  $\text{LRS}(wa)$  is the longest p-suffix of  $wa$  that occurs in  $w$ . That is, all the pv-suffixes of  $wa$  longer than  $\text{LRS}(wa)$  belongs to the sink  $[wa]_{wa}^R$ . Therefore, in both cases,  $F([wa]_{wa}^R) = [x]_{wa}^R$ . If  $[x]_w^R = [x]_{wa}^R$ , all the suffix links in  $\text{PDAWG}(w)$  are kept. If  $[x]_w^R \neq [x]_{wa}^R$ ,  $F([y]_{wa}^R) = [x]_{wa}^R$ , since the longest p-suffix of  $y$  not in  $[y]_{wa}^R$  is  $x$ . On the other hand, the longest p-suffix of  $x$  not in  $[x]_{wa}^R$  is the longest p-suffix of  $y$  not in  $[y]_w^R$ , i.e.,  $F_w([y]_w^R)$ .  $\square$



---

**Algorithm 4:** Constructing PDAWG( $T$ )

---

```
1 Let  $V \leftarrow \{\top, \rho\}$ ,  $E \leftarrow \{(\top, a, \rho) \mid a \in \Sigma \cup \{\infty\}\}$ ,  $F(\top) \leftarrow \top$ ,  $F(\rho) \leftarrow \top$ ,  $\text{len}(\top) = -1$ ,  
    $\text{len}(\rho) \leftarrow 0$ ,  $\text{sink} \leftarrow \rho$ , and  $t \leftarrow \langle T \rangle$ ;  
2 for  $i \leftarrow 1$  to  $|t|$  do  
3   Let  $a \leftarrow t[i]$  and  $u \leftarrow \text{sink}$ ;  
4   Create a new node and let  $\text{sink}$  be that node with  $\text{len}(\text{sink}) = i$ ;  
5   while  $\text{trans}(u, \text{len}(F(u)) + 1, \langle\langle a \rangle\rangle_{\text{len}(F(u))+1}) = \text{Null}$  do  
6     Let  $\text{child}(u, \langle\langle a \rangle\rangle_{\text{len}(u)}) \leftarrow \text{sink}$  and  $u \leftarrow F(u)$ ;  
     //  $u$  corresponds to  $[\text{preLRS}(t[:i])]_{t[:i-1]}^R$   
7   if  $\langle\langle a \rangle\rangle_{\text{len}(u)} \in \text{Children}(u)$  then //  $\text{preLRS}(t[:i]) = [[\text{preLRS}(t[:i])]_{t[:i-1]}^R]$   
8     Let  $k \leftarrow \text{len}(u) + 1$  and  $v \leftarrow \text{child}(u, \langle\langle a \rangle\rangle_{\text{len}(u)})$   
9   else //  $\text{preLRS}(t[:i]) \neq [[\text{preLRS}(t[:i])]_{t[:i-1]}^R]$   
10    Let  $k \leftarrow \min\{a, \max(\text{Children}(u) \cap \mathcal{N})\}$ ,  $v \leftarrow \text{trans}(u, k - 1, \infty)$ ,  
     $\text{child}(u, \langle\langle a \rangle\rangle_{\text{len}(u)}) \leftarrow \text{sink}$ , and  $u \leftarrow F(u)$ ;  
    //  $u$  corresponds to  $F_{[t[:i-1]]}([\text{preLRS}(t[:i])]_{t[:i-1]}^R)$   
    //  $v$  corresponds to  $[\text{LRS}(t[:i])]_{t[:i-1]}^R$  and  $k = |\text{LRS}(t[:i])|$   
11   if  $\text{len}(v) = k$  then Let  $F(\text{sink}) \leftarrow v$ ; // No node split  
12   else // Node split  
13     Create a new node  $v'$ ; //  $v'$  corresponds to  $[\text{LRS}(t[:i])]_{t[:i]}^R$   
14     Let  $\text{len}(v') \leftarrow k$ ;  
     // Outgoing edges of the new node  
15     for each  $b \in \text{Children}(v)$  such that  $\langle\langle b \rangle\rangle_k \neq \infty$  do  
16       Let  $\text{child}(v', b) \leftarrow \text{child}(v, b)$ ;  
17     if  $\text{trans}(v, k, \infty) \neq \text{Null}$  then Let  $\text{child}(v', \infty) \leftarrow \text{trans}(v, k, \infty)$ ;  
     // Incoming edges of the new node  
18     while  $\text{child}(u, \langle\langle a \rangle\rangle_{\text{len}(u)}) = v$  do  
19       Let  $\text{child}(u, \langle\langle a \rangle\rangle_{\text{len}(u)}) \leftarrow v'$  and  $u \leftarrow F(u)$ ;  
       // Suffix links  
20     Let  $F(v') \leftarrow F(v)$ ,  $F(v) \leftarrow v'$  and  $F(\text{sink}) \leftarrow v'$ ;  
21 return  $(V, E, F)$ ;
```

---

Algorithm 4 constructs PDAWGs based on the above lemmas, where  $\rho$  and  $t[:i]$  in the pseudo code corresponds to the source node  $\{\varepsilon\}$  and  $wa$  in the main body of the text, respectively. Figure 6 illustrates an example run. For technical convenience, like the standard DAWG construction algorithm, we add a dummy node  $\top$  to the PDAWG that has edges to the source node labeled with all elements of  $\Sigma \cup \{\infty\}$  and let  $F(\rho) = F(\top) = \top$ . This trick allows us to uniformly treat the special case where the LRS node is  $\rho$ , in which case  $\top$  is regarded as the pre-LRS node. Each node  $u$  does not remember the elements of  $u$  but it remembers  $\text{len}(u) = |u|$ . For the dummy node  $\top$ , we let  $\text{len}(\top) = -1$ . Note that  $|u| = |\text{len}(F(u))| + 1$  (used in Line 5). Hereafter, we use functions  $F$ ,  $\text{child}$ ,  $\text{trans}$ , etc. without a subscript specifying a text, to refer to the data structure that the algorithm is manipulating, rather than the mathematical notion relative to the text. Of course, we design our algorithm so that those functions coincide with the corresponding mathematical notions.

Suppose we have constructed  $\text{PDAWG}(w)$  and want to obtain  $\text{PDAWG}(wa)$ . The sink node of  $\text{PDAWG}(w)$ , which we denote as  $\text{oldsink}$ , corresponds to  $[w]_w^R$ . We first make a new sink node  $\text{newsink} = [wa]_{wa}^R$  and let  $\text{len}(\text{newsink}) = |wa|$ . Then, in the **while** loop of Line 5, we visit  $u_0, u_1, \dots, u_j$  on the suffix link chain of  $\text{oldsink}$ , until we find the pre-LRS node  $u_j = [\text{preLRS}(wa)]_w^R$ . By Lemma 10.1, we identify  $u_j$  when exiting the **while** loop. For each node  $u_i$  with  $i < j$ , we make an edge labeled with  $\langle\langle a \rangle\rangle_{\text{len}(u_i)}$  pointing at  $\text{newsink}$  by Lemma 11. The algorithm identifies the length  $k$  of the LRS



in Lines 7–10 based on Lemma 10.2. If  $k - 1 < \text{len}(u_j)$ , the pre-LRS node  $u_j$  also has an edge pointing at *newsink* by Lemma 11. This is done at Line 10. At this moment, we have created all the incoming edges of *newsink*, possibly except the one from  $[\text{preLRS}(wa)]_{wa}^R$ , if necessary, because that node  $[\text{preLRS}(wa)]_{wa}^R$  has not yet been created in the case of node split. We then reach the LRS node  $v = [\text{LRS}(wa)]_w^R = \text{trans}(u_j, k - 1, \langle\langle a \rangle\rangle_{k-1})$  by Lemma 10.3. At the moment entering Line 11, the variable  $u$  represents the first node on the suffix link chain of the pre-LRS node such that  $\lceil u \rceil$  is not longer than the pre-LRS. It is just the pre-LRS node  $u_j$  if  $k - 1 = \text{len}(u_j)$ , and it is  $F(u_j)$  otherwise. Then, in both cases,  $u$  is the first node whose edge toward the LRS node will be redirected to  $[x]_{wa}^R$  in the **while** loop of Line 18 in accordance with Lemma 13, if it has one.

We compare  $k$  and  $\text{len}(v)$  to decide whether the LRS node should be split based on Lemma 10.4. If  $|\text{LRS}(wa)| = \text{len}(v)$ , the node  $v$  will not be split, in which case we obtain  $\text{PDAWG}(wa)$  by making  $F(\text{newsink}) = v$  at Line 11 by Lemma 14. Suppose  $k < \text{len}(v)$ . In this case, the LRS node  $v$  must be split. We reuse the old node  $v$ , which used to correspond to  $[\text{LRS}(wa)]_w^R$ , as a new node corresponding to  $[[[\text{LRS}(wa)]_w^R]]_{wa}^R$ , and create another new node  $v'$  for  $[\text{LRS}(wa)]_{wa}^R$  with  $\text{len}(v') = k$ . Edges are determined in accordance with Lemmas 12 and 13. The outgoing edges from  $v$  shall be kept. We create outgoing edges of  $v'$  referring to the corresponding transitions from  $v$ . If  $(v, b, u) \in E$  with  $\langle\langle b \rangle\rangle_k \neq \infty$ , then we add  $(v', b, u)$  to  $E$  in the **for** loop of Line 15. In addition, we add  $(v', \infty, \text{trans}(v, k, \infty))$  to  $E$  if  $\text{trans}(v, k, \infty) \neq \text{Null}$  at Line 17. We note that thanks to the third and fourth claims of Lemma 12, this correctly creates an edge  $(v', \infty, \text{newsink})$  if necessary. We should have  $(v', \infty, \text{newsink}) \in E$  if and only if  $(v, \langle\langle a \rangle\rangle_{\text{len}(v)}, \text{newsink}) \in E \wedge \langle\langle a \rangle\rangle_k = \infty$  if and only if  $\text{trans}(v, k, \langle\langle a \rangle\rangle_k) = \text{newsink}$ . All incoming edges of  $v$  from nodes on the suffix link chain of  $u$  in  $\text{PDAWG}(w)$  are redirected to  $v'$  in the **while** loop of Line 18, where  $u$  is the pre-LRS node  $u_j$  if  $\lceil u_j \rceil$  is the pre-LRS, and it is  $F(u_j)$  otherwise.

At last, the suffix link from *newsink* to  $v$  and the suffix link from  $v$  to  $v'$  are determined in accordance with Lemma 14.

**Remark 2.** One can compute  $\ell_u = \min \text{RPos}(u)$  online simply by letting  $\ell_{[wa]_{wa}^R} = |wa|$  when creating the new sink node  $[wa]_{wa}^R$ . When the LRS node  $[y]_w^R$  is split into  $[y]_{wa}^R$  and  $[x]_{wa}^R$ , we have  $\min \text{RPos}([y]_w^R) = \min \text{RPos}([y]_{wa}^R) = \min \text{RPos}([x]_{wa}^R)$ . So, it is enough to copy the  $\ell$  value.

## 7.1 Time complexity analysis

Let us call an edge  $(u, a, v)$  *primary* if  $\lceil v \rceil = \lceil u \rceil \cdot a$ , and *secondary* otherwise. The following lemma is an adaptation of the corresponding one for DAWGs by Blumer et al. [16].

**Lemma 15.** *Let  $\text{SC}_w(u)$  be the set of nodes on the suffix link chain of a node  $u$ . If  $\text{PDAWG}(w)$  has a primary edge from  $u$  to  $v$ , then the total number of secondary edges from nodes in  $\text{SC}_w(u)$  to nodes in  $\text{SC}_w(v)$  is bounded by  $|\text{SC}_w(u)| - |\text{SC}_w(v)| + |\Pi| + 1$ .*

*Proof.* Let us count the number of edges from nodes in  $\text{SC}_w(u)$  to  $\text{SC}_w(v)$ . Baker [1, Lemma 1] showed that in a parameterized suffix tree, each path from the root to a leaf has at most  $|\Pi|$  nodes with *bad suffix links*. Through the duality of PDAWGs and parameterized suffix trees stated in Lemma 5, this means that  $\text{SC}_w(v)$  contains at most  $|\Pi| + 1$  nodes which have no incoming primary edges, where the additional one node is the root of the PDAWG. Since each node has at most one incoming primary edge, the number of primary edges in concern is at least  $|\text{SC}_w(v)| - |\Pi| - 1$  in total. Since each node in  $\text{SC}_w(u)$  has just one outgoing edge to  $\text{SC}_w(v)$ , we obtain the lemma.  $\square$

**Theorem 7.** *Given a string  $T$  of length  $n$ , Algorithm 4 constructs  $\text{PDAWG}(T)$  in  $O(n|\Pi| \log(|\Pi| + |\Sigma|))$  time and  $O(n)$  space online, by reading  $T$  from left to right.*

*Proof.* Since the size of a PDAWG is bounded by  $O(n)$  (Theorem 4) and nodes are monotonically added, it is enough to bound the number of edges and suffix links that are deleted. In each iteration of the **for** loop, at most one suffix link is deleted. So at most  $n$  suffix links are deleted in total. We count the number of edges whose target is altered from  $v = [\text{LSR}(wa)]_w^R$  to  $v' = [\text{LSR}(wa)]_{wa}^R$  on Line 18 when updating  $\text{PDAWG}(w)$  to  $\text{PDAWG}(wa)$ . Let  $k_i$  be the number of such edges at the  $i$ -th iteration

of the **for** loop. Note that those are all secondary edges from a node in  $\text{SC}_w(u_0)$  for the pre-LRS node  $u_0$ . By Lemma 15,

$$\begin{aligned}
\sum_{i=1}^n k_i &\leq \sum_{i=1}^n (|\text{SC}_{wa}(w)| - |\text{SC}_{wa}(wa)| + |\Pi| + 1) \\
&\leq \sum_{i=1}^n (|\text{SC}_w(w)| - |\text{SC}_{wa}(wa)| + |\Pi| + 1) \\
&= |\text{SC}_\varepsilon(\varepsilon)| - |\text{SC}_t(t)| + (|\Pi| + 1)n \in O(|\Pi|n). \quad \square
\end{aligned}$$

Since the suffix links of  $\text{PDAWG}(T)$  form the p-suffix tree of  $\bar{T}$  (see Subsection 6.2), the following corollary is immediate from Theorem 7.

**Corollary 3.** *The p-suffix tree of a string  $S$  of length  $n$  can be constructed in  $O(n|\Pi| \log(|\Pi| + |\Sigma|))$  time and  $O(n)$  space online, by reading  $S$  from right to left.*

Differently from the online DAWG construction algorithm [16], we have the factor  $|\Pi|$  in our algorithm complexity analysis. Actually, our algorithm takes time proportional to the difference of the old and new PDAWGs modulo logarithmic factors, as long as the difference is defined so that the split node  $[\text{LRS}(wa)]_w^R$  automatically becomes  $[[[\text{LRS}(wa)]_w^R]]_{wa}^R$  rather than  $[\text{LRS}(wa)]_{wa}^R$ . In this sense, our algorithm is optimal. It is open whether we could improve the analysis.

## 8 Concluding remarks

In this paper, we proposed a new indexing structure for parameterized pattern matching—the PDAWGs. We showed that  $\text{PDAWG}(T)$  for an input text string  $T$  of length  $n$  over a static alphabet  $\Sigma$  and a parameterized alphabet  $\Pi$  can be built in  $O(n|\Pi| \log(|\Pi| + |\Sigma|))$  time with  $O(n)$  space, in a right-to-left online manner. The duality of our PDAWGs and parameterized suffix trees [5] permits us  $O(n)$ -time offline construction of  $\text{PDAWG}(T)$  provided that the p-suffix tree of  $\bar{T}$  has already been built. It also gives us a linear-space bidirectional index for parameterized pattern matching in  $O(m \log(|\Pi| + |\Sigma|) + occ)$  time.

The major open question is whether our upper bound  $O(n|\Pi| \log(|\Pi| + |\Sigma|))$  for the time complexity of online construction of the PDAWGs is tight. Our construction algorithm is “optimal” in the sense that, given a new character  $a$  to append, the time for updating  $\text{PDAWG}(T)$  to  $\text{PDAWG}(Ta)$  is proportional to the difference of the two DAWGs, ignoring the logarithmic factors which are the costs for searching for edges and/or suffix links. We have not found an instance that requires  $\Omega(n|\Pi|)$  changes into the DAWG during the whole online construction.

## Acknowledgment

The authors are deeply grateful to the anonymous reviewer for many valuable comments that have improved the paper’s readability and preciseness. This work was supported by JSPS KAKENHI Grant Numbers JP19K20208, JP18K18002, JP18H04091, JP18K11150, JP17H01697, JP16H02783, JP20H04141, JP15H05706, JP21K11745, JP18H04098, and JST PRESTO Grant Number JPMJPR1922.

## References

- [1] B. S. Baker, Parameterized pattern matching: Algorithms and applications, *Journal of Computer and System Sciences* 52 (1) (1996) 28–42.
- [2] T. Shibuya, Generalization of a suffix tree for RNA structural pattern matching, *Algorithmica* 39 (1) (2004) 1–19.

- [3] J. Mendivelso, Y. Pinzón, Parameterized matching: Solutions and extensions, in: Proceedings of the Prague Stringology Conference 2015, 2015, pp. 118–131.
- [4] J. Mendivelso, S. V. Thankachan, Y. Pinzón, A brief history of parameterized matching problems, *Discrete Applied Mathematics* 274 (2020) 103–115.
- [5] B. S. Baker, A theory of parameterized pattern matching: algorithms and applications, in: Proceedings of the twenty-fifth annual ACM symposium on Theory of Computing, 1993, pp. 71–80.
- [6] S. R. Kosaraju, Faster algorithms for the construction of parameterized suffix trees, in: Proceedings of IEEE 36th Annual Foundations of Computer Science, 1995, pp. 631–638.
- [7] T. Lee, J. C. Na, K. Park, On-line construction of parameterized suffix trees for large alphabets, *Information Processing Letters* 111 (5) (2011) 201–207.
- [8] S. Deguchi, F. Higashijima, H. Bannai, S. Inenaga, M. Takeda, Parameterized suffix arrays for binary strings, in: Proceedings of the Prague Stringology Conference 2008, 2008, pp. 84–94.
- [9] T. I, S. Deguchi, H. Bannai, S. Inenaga, M. Takeda, Lightweight parameterized suffix array construction, in: Proceedings of the 20th International Workshop on Combinatorial Algorithms, 2009, pp. 312–323.
- [10] R. Beal, D. A. Adjeroh, p-suffix sorting as arithmetic coding, *Journal of Discrete Algorithms* 16 (2012) 151–169.
- [11] N. Fujisato, Y. Nakashima, S. Inenaga, H. Bannai, M. Takeda, Direct linear time construction of parameterized suffix and LCP arrays for constant alphabets, in: Proceedings of the 26th International Symposium on String Processing and Information Retrieval, 2019, pp. 382–391.
- [12] A. Ganguly, R. Shah, S. V. Thankachan, pBWT: Achieving succinct data structures for parameterized pattern matching and related problems, in: Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms, 2017, pp. 397–407.
- [13] Diptarama, T. Katsura, Y. Otomo, K. Narisawa, A. Shinohara, Position heaps for parameterized strings, in: Proceedings of the 28th Annual Symposium on Combinatorial Pattern Matching, 2017, pp. 8:1–8:13.
- [14] N. Fujisato, Y. Nakashima, S. Inenaga, H. Bannai, M. Takeda, Right-to-left online construction of parameterized position heaps, in: Proceedings of the Prague Stringology Conference 2018, 2018, pp. 91–102.
- [15] N. Fujisato, Y. Nakashima, S. Inenaga, H. Bannai, M. Takeda, The parameterized position heap of a trie, in: Proceedings of the International Conference on Algorithms and Complexity, 2019, pp. 237–248.
- [16] A. Blumer, J. Blumer, D. Haussler, A. Ehrenfeucht, M.-T. Chen, J. Seiferas, The smallest automation recognizing the subwords of a text, *Theoretical Computer Science* 40 (1985) 31–55.
- [17] M. Crochemore, Transducers and repetitions, *Theoretical Computer Science* 45 (1) (1986) 63–86.
- [18] K. Nakashima, N. Fujisato, D. Hendrian, Y. Nakashima, R. Yoshinaka, S. Inenaga, H. Bannai, A. Shinohara, M. Takeda, DAWGs for parameterized matching: Online construction and related indexing structures, in: Proceedings of the 31st Annual Symposium on Combinatorial Pattern Matching, 2020, pp. 26:1–26:14.
- [19] S. Kim, H. Cho, Simpler FM-index for parameterized string matching, *Information Processing Letters* 165 (2021) 106026.

- [20] P. Weiner, Linear pattern matching algorithm, in: Proceedings of the 14th IEEE Symposium on Switching and Automata Theory, 1973, pp. 1–11.