

Enhanced Covers of Regular & Indeterminate Strings using Prefix Tables

Ali Alatabbi¹, A. S. Sohidull Islam², M. Sohel Rahman^{*1,3},
Jamie Simpson^{4,5}, and W. F. Smyth^{**2,5}

¹ Department of Informatics, King's College London
ali.alatabbi@kcl.ac.uk

² Algorithms Research Group, Department of Computing & Software
McMaster University
sohansayed@gmail.com, smyth@mcmaster.ca

³ Department of Computer Science & Engineering
Bangladesh University of Engineering & Technology
msrahman@cse.buet.ac.bd

⁴ Department of Mathematics and Statistics, Curtin University of Technology
Jamie.Simpson@curtin.edu.au

⁵ School of Engineering & Information Technology
Murdoch University

Abstract. A *cover* of a string $x = x[1..n]$ is a proper substring u of x such that x can be constructed from possibly overlapping instances of u . A recent paper [12] relaxes this definition — an *enhanced cover* u of x is a border of x (that is, a proper prefix that is also a suffix) that covers a *maximum* number of positions in x (not necessarily all) — and proposes efficient algorithms for the computation of enhanced covers. These algorithms depend on the prior computation of the *border array* $\beta[1..n]$, where $\beta[i]$ is the length of the longest border of $x[1..i]$, $1 \leq i \leq n$. In this paper, we first show how to compute enhanced covers using instead the *prefix table*: an array $\pi[1..n]$ such that $\pi[i]$ is the length of the longest substring of x beginning at position i that matches a prefix of x . Unlike the border array, the prefix table is robust: its properties hold also for *indeterminate strings* — that is, strings defined on *subsets* of the alphabet Σ rather than individual elements of Σ . Thus, our algorithms, in addition to being faster in practice and more space-efficient than those of [12], allow us to easily extend the computation of enhanced covers to indeterminate strings. Both for regular and indeterminate strings, our algorithms execute in expected linear time. Along the way we establish an important theoretical result: that the expected maximum length of any border of any prefix of a regular string x is approximately 1.64 for binary alphabets, less for larger ones.

* Supported in part by a Commonwealth Academic Fellowship and a ACU Titular Fellowship. Currently on a sabbatical leave from BUET.

** Supported in part by a grant from the Natural Sciences & Engineering Research Council (NSERC) of Canada.

1 Introduction

The concept of *periodicity* is fundamental to combinatorics on words and related algorithms: it is difficult to imagine a research contribution that does not somehow involve periods of strings. But periodicity alone may not be the best descriptor of a string; for example, $\mathbf{x} = \text{abaababab}$, a string of length $n = 9$, has period 7 and corresponding *generator*⁶ abaabab , but it could well be more interesting that every position but one in \mathbf{x} lies within an occurrence of ab . In 1990 Apostolico & Ehrenfeucht [3] introduced the idea of quasiperiodicity: a *quasiperiod* or *cover* of a string \mathbf{x} is a proper substring \mathbf{u} of \mathbf{x} such that any position in \mathbf{x} is contained in an occurrence of \mathbf{u} ; \mathbf{u} is then said to *cover* \mathbf{x} , which is said to be *quasiperiodic*. Thus, for example, $\mathbf{u} = \text{aba}$ is a cover of $\mathbf{x} = \text{ababaaba}$. Several linear-time algorithms were proposed for the computation of covers [4,8,17,18], culminating in an algorithm [16] to compute the *cover array* γ , where $\gamma[i]$ gives the length j of the longest cover of $\mathbf{x}[1..i]$. Since the longest cover of $\mathbf{x}[1..j]$ is also a cover of $\mathbf{x}[1..i]$, γ implicitly specifies all the covers of every prefix of \mathbf{x} . A recent paper [2] extends the computation of γ to “indeterminate strings” (see below for definition).

Even though the cover of a string can provide useful information, quasiperiodic strings are on the other hand infrequent among strings in general. Another approach to string covering was therefore proposed in [15]: a set $U_k = U_k(\mathbf{x})$ of strings, each of length k , is said to be a *minimum k -cover* of \mathbf{x} if every position in \mathbf{x} lies within some occurrence of an element of U_k , and no smaller set of k -strings has this property. Thus $U_2(\text{abaababab}) = U_2(\text{ababaaba}) = \{ab, ba\}$. In [10] the computation of U_k was shown to be NP-complete, though an approximate polynomial-time algorithm was presented in [14].

Recall that a *border* of \mathbf{x} is a possibly empty proper prefix of \mathbf{x} that is also a suffix: every nonempty string has a border of length zero. Recently the promising idea of an *enhanced cover* was introduced [12]; that is, a border \mathbf{u} of $\mathbf{x} = \mathbf{x}[1..n]$ that covers a maximum number $m \leq n$ of positions in \mathbf{x} . Then the *minimum enhanced cover* $\text{mec}(\mathbf{x})$ is the *shortest* border \mathbf{u} that covers m positions, and [12] presented an algorithm to compute $\text{mec}(\mathbf{x})$ in $\Theta(n)$ time. Thus for $\mathbf{x} = \text{abaababab}$, $\text{mec}(\mathbf{x}) = \text{ab}$. Further, on the analogy of the cover array defined above, the authors proposed the *minimum enhanced cover array* $\text{MEC}_{\mathbf{x}}$ — for every $i \in 1..n$, $\text{MEC}_{\mathbf{x}}[i] = |\text{mec}(\mathbf{x}[1..i])|$, the length of the minimum enhanced cover of $\mathbf{x}[1..i]$ — and showed how to compute it in $\mathcal{O}(n \log n)$ time. In this

⁶ Notation and terminology generally follow [19].

paper we introduce in addition the **CMEC** array, where $\text{CMEC}[i]$ specifies the number of positions in \mathbf{x} covered by the border of length $\text{MEC}[i]$. Thus, for example, $\text{MEC}_{\text{abaababab}} = 001123232$ and $\text{CMEC}_{\text{abaababab}} = 002346688$.

In order to compute $\text{MEC}\mathbf{x}$, the authors of [12] made use of a variant of the **border array** — that is, an integer array $\beta[1..n]$ in which for every $i \in 1..n$, $\beta[i]$ is the length of the longest border of $\mathbf{x}[1..i]$. In this paper we adopt a different approach to the computation of $\text{MEC}\mathbf{x}$, using, instead of a border array, the **prefix table** $\pi = \pi[1..n]$, where for every $i \in 1..n$, $\pi[i]$ is the length of the longest substring at position i of \mathbf{x} that equals a prefix of \mathbf{x} . It has long been folklore that β and π are “equivalent”, but it has only recently been made explicit [6] that each can be computed from the other in linear time. However, this equivalence holds only for **regular** strings \mathbf{x} in which each entry $\mathbf{x}[i]$ is constrained to be a single element of the underlying alphabet Σ .

We say that a letter λ is **indeterminate** if it is any nonempty subset of Σ , and thus a string \mathbf{x} is said to be **indeterminate** if some constituent letter $\mathbf{x}[i]$ is indeterminate. The idea of an indeterminate string was first introduced in [11] — with letters constrained to be either regular (single elements of Σ) or Σ itself — and the properties of these strings have been much studied by Blanchet-Sadri [7] and her collaborators as “partial words” or “strings with holes”. Indeterminate strings can model DNA sequences on $\Sigma = \{A, C, G, T\}$ when ambiguities arise in determining individual nucleotides (letters).

Two indeterminate letters λ and μ are said to **match** (written $\lambda \approx \mu$) whenever $\lambda \cap \mu \neq \emptyset$, a relation that is in general nontransitive [13,22]: $a \approx \{a, b\}$ and $\{a, b\} \approx b$, but $a \not\approx b$. An important consequence of this nontransitivity is that the border array no longer correctly describes *all* of the borders of \mathbf{x} : it is no longer necessarily true, as for regular strings, that if \mathbf{u} is the longest border of \mathbf{v} , in turn the longest border of \mathbf{x} , then \mathbf{u} is a border of \mathbf{x} . On the other hand, the prefix array retains all its properties for indeterminate strings \mathbf{x} and, in particular, correctly identifies all the borders of every prefix of \mathbf{x} [6]. [20] describes algorithms to compute the prefix table of an indeterminate string; conversely, [9] proves that there exists an indeterminate string corresponding to every feasible prefix table, while [1] describes an algorithm to compute the lexicographically least indeterminate string determined by any given feasible prefix table. There is thus a many-many relationship between the set of all indeterminate strings and the set of all prefix tables. Consequently, computing $\text{MEC}\mathbf{x}$ (or simply **MEC** when there is no ambiguity) from the prefix table $\pi = \pi\mathbf{x}$

rather than from a variant of the border array allows us to extend the computation to indeterminate strings.

In Section 2 we outline the basic methodology and data structures used to compute the minimum enhanced cover array from the prefix table, while illustrating the ideas with an example. Then Section 3 provides a proof of the algorithm’s correctness, as well as an analysis of its complexity, both worst and average case. In Section 4 we discuss the practical application of our algorithms, in terms of time and space requirements, and compare our prefix-based implementation with the border-based implementation of [12]. Section 5 extends the enhanced cover array algorithm to indeterminate strings (for rooted covers) and outlines various other extensions, particularly to generalizations of MECs.

2 Methodology

In this section we describe the computation of $\text{MEC}_{\mathbf{x}}$, the enhanced cover array of \mathbf{x} , based on the prefix array $\boldsymbol{\pi}$. Since every minimum enhanced cover of \mathbf{x} is also a border of \mathbf{x} , we are initially interested in the covers of prefixes of \mathbf{x} . For this purpose we need arrays whose size is B , the maximum length of any border of any prefix of \mathbf{x} . Noting that B must be the maximum entry in the prefix array $\boldsymbol{\pi}$, we write $B = \max_{2 \leq i \leq n} \boldsymbol{\pi}[i]$.

Definition 1 *In the **maximum no cover** array $\text{MNC} = \text{MNC}[1..B]$, for every $q \in 1..B$, $\text{MNC}[q] = q'$, where q' is the maximum integer in $1..q$ such that the prefix $\mathbf{x}[1..q']$ has no cover — that is, such that $\gamma[q'] = 0$.*

As shown in Figure 1, once B is computed in $\Theta(n)$ time from the prefix array $\boldsymbol{\pi}$, MNC can be easily computed in $\Theta(B)$ time using the cover array $\gamma[1..B]$ of $\mathbf{x}[1..B]$. Note that the entries in MNC are monotone nondecreasing with $1 \leq \text{MNC}[q] \leq q$ for every $q \in 1..B$. The following is fundamental to the execution of our main algorithm:

Observation 2 *If a prefix $\mathbf{v} = \mathbf{x}[1..q]$ of \mathbf{x} has a cover \mathbf{u} , then $\mathbf{v} \neq \text{mec}(\mathbf{x})$ (since $|\mathbf{u}| < q$ and \mathbf{u} covers every position covered by \mathbf{v}).*

Thus $\text{MNC}[q]$ specifies an upper bound $q' \in 1..q$ on the length of a minimum enhanced cover of \mathbf{x} . Two other arrays are required for the computation, both of length B :

Definition 3 *For every $q \in 1..B$:*

- $\text{PR}[q]$ *is the rightmost position in \mathbf{x} at which the prefix $\mathbf{x}[1..q]$ occurs;*

```

procedure Compute_MNC( $n, \pi; B, \gamma, \text{MNC}$ )
 $B \leftarrow \pi[2]$ 
for  $i \leftarrow 3$  to  $n$  do
     $B \leftarrow \max(B, \pi[i])$ 
     $\triangleright$  Compute  $\gamma[1..B]$  of  $\mathbf{x}[1..B]$  using
     $\triangleright$  the algorithm Compute_PCR of [2].
    Compute_PCR( $B, \pi; \gamma$ )
     $\triangleright$  Note that MNC can overwrite  $\gamma$ .
    for  $q \leftarrow 1$  to  $B$  do
        if  $\gamma[q] = 0$  then  $\text{MNC}[q] \leftarrow q$ 
        else  $\text{MNC}[q] \leftarrow \text{MNC}[q-1]$ 

```

Fig. 1. Computing MNC from the prefix array $\pi[1..n]$ and the cover array $\gamma[1..B]$.

- $\text{CPR}[q]$ is the number of positions in \mathbf{x} covered by occurrences of $\mathbf{x}[1..q]$.

Here is an example of the arrays introduced thus far:

```

      1 2 3 4 5 6 7 8 9 10
 $\mathbf{x} = a b a b a a b a b a$ 
 $\pi = 10 0 3 0 1 5 0 3 0 1$ 
 $\gamma = 0 0 0 2 3$ 
 $\text{MNC} = 1 2 3 3 3$ 
 $\text{PR} = 10 8 8 6 6$ 
 $\text{CPR} = 6 8 10 8 10$ 
 $\text{MEC} = 0 0 1 2 3 1 2 3 2 3$ 
 $\text{CMEC} = 0 0 2 4 5 4 6 8 8 10$ 

```

Note that for $\mathbf{x}[1..9]$ and $\mathbf{x}[1..10]$, there are actually *two* borders that cover a maximum number of positions; in each case the border of minimum length is identified in MEC.

The algorithm Compute_MEC is shown in Figure 2. In the first stage, B and MNC are computed and the arrays CMEC , PR and CPR are initialized. Then every position $i > 1$ such that $q = \gamma[i] > 0$ is considered. Using MNC , the longest prefix $\mathbf{Q}' = \mathbf{x}[1..q']$ of $\mathbf{x}[1..q]$ that does *not* have a cover is identified; for prefixes of $\mathbf{x}[1..q]$ that do have a cover, the appropriate PR and CPR values have already been updated. There are two main steps in the processing of \mathbf{Q}' :

- Since i has now become the rightmost occurrence of \mathbf{Q}' in $\mathbf{x}[1..i]$, we must set $\text{PR}[q'] \leftarrow i$ and increment the corresponding number $\text{CPR}[q']$ of positions covered.

```

procedure Compute_MEC( $\pi$ ; MEC, CMEC)
 $n \leftarrow |\pi|$ 
Compute_MNC( $n, \pi$ ; B,  $\gamma$ , MNC)
MEC  $\leftarrow 0^n$ ; CMEC  $\leftarrow 0^n$ ; PR  $\leftarrow 1^B$ 
for  $q \leftarrow 1$  to B do CPR[ $q$ ]  $\leftarrow q$ 
for  $i \leftarrow 2$  to  $n$  do
     $q \leftarrow \pi[i]$ 
     $\triangleright \mathbf{x}[i..i+q-1] = \mathbf{x}[1..q]$ .
    while  $q > 0$  do
     $\triangleright \mathbf{x}[1..q']$  is the longest prefix of  $\mathbf{x}[1..q]$  without a cover.
         $q' \leftarrow \text{MNC}[q]$ 
     $\triangleright \mathbf{x}[1..q']$  also occurs at  $i$ : update CPR[ $q'$ ] & PR[ $q'$ ].
        if  $i - \text{PR}[q'] < q'$  then
            CPR[ $q'$ ]  $\leftarrow \text{CPR}[q'] + i - \text{PR}[q']$ 
        else
            CPR[ $q'$ ]  $\leftarrow \text{CPR}[q'] + q'$ 
            PR[ $q'$ ]  $\leftarrow i$ 
     $\triangleright$  Update CMEC & MEC accordingly for interval  $i..i+q'-1$ .
        if CPR[ $q'$ ]  $\geq$  CMEC[ $i+q'-1$ ] then
            MEC[ $i+q'-1$ ]  $\leftarrow q'$ 
            if CPR[ $q'$ ]  $>$  CMEC[ $i+q'-1$ ] then
                CMEC[ $i+q'-1$ ]  $\leftarrow \text{CPR}[q']$ 
         $q \leftarrow q' - 1$ 

```

Fig. 2. Computing MEC and CMEC from the prefix array π .

- If the number CPR[q'] of positions covered by occurrences of Q' exceeds CMEC[$i+q-1$], then CMEC and MEC must be updated accordingly.

These steps are repeated recursively for the longest proper prefix of Q' that does not have a cover.

3 Correctness & Complexity of Compute_MEC

We begin by proving the correctness of Compute_MEC, which depends on the prior computation of $\pi = \pi_{\mathbf{x}}$ [6]. Consider first procedure Compute_MNC, where B is computed, followed by the cover array $\gamma[1..B]$. Then for every $q \in 1..B$, MNC[q] $\leftarrow q$ whenever there is no cover of $\mathbf{x}[1..q]$, with MNC[q] $\leftarrow \text{MNC}[q-1]$ otherwise, an easy and straightforward calculation.

Compute_MEC then independently considers positions $i = 2, 3, \dots, n$ for which $\pi[i] > 0$; that is, such that a border of \mathbf{x} of length $q = \pi[i]$

begins at i . The internal **while** loop then processes in decreasing order of length the prefixes $Q' = x[1..q']$ of $x[1..q]$ that have no cover — and that therefore, by Observation 2, can possibly be minimum enhanced covers of $x[1..i+q'-1]$. Thus, for every $i \in 2..n$, all such borders $x[1..q] = x[i..i+q-1]$ are considered and, for each one, all such prefixes Q' . For each q' :

- the number $\text{CPR}[q']$ of positions covered by Q' is updated, as well as the position $\text{PR}[q'] = i$ of rightmost occurrence of Q' ;
- $\text{MEC}[i+q'-1]$ and $\text{CMEC}[i+q'-1]$ are updated accordingly for sufficiently large $\text{CPR}[q']$.

We claim therefore that

Theorem 4 *For a given string x , Compute_MEC correctly computes the minimum enhanced cover array MEC_x and the number CMEC_x of positions covered by it, based solely on the prefix array π_x .*

We have seen that in aggregate Compute_MEC processes a subset of the nonempty borders of every prefix $x[1..i]$, devoting $\mathcal{O}(1)$ time to each one. As we have seen, each border Q' in each such subset is constrained to have no cover. We say that a string v is **strongly periodic** if it has a border u such that $|u| \geq |v|/2$; otherwise v is said to be **weakly periodic**. Observe that the borders Q' must all be weakly periodic; if not, then they would have a cover u with $|u| \geq |v|/2$. In [12] the following result is proved:

Lemma 5 *There are at most $\log_2 n$ weakly periodic borders of a string of length n .*

It follows then that for each $i \in 2..n$, there are at most $\log_2 i$ borders considered, thus overall $\mathcal{O}(n \log n)$ time.

The space requirement of Compute_MEC , apart from the π , MEC and CMEC arrays, each of length n , consists of three integer arrays (MNC (overwriting γ), PR , CPR), each of length $B < n$. Thus

Theorem 6 *In the worst case, Compute_MEC computes MEC and CMEC from π using*

- $\mathcal{O}(n \log n)$ time;
- three additional arrays $1..B$ of integers $1..n$, thus $\Theta(B \log n)$ bits of space.

Now consider the expected (average) case behaviour of Compute_MEC . This depends critically on the expected length of the maximum border of

$\mathbf{x}[1..n]$; that is, the expected value of \mathbf{B} . We show in the Appendix that for a given alphabet size, \mathbf{B} approaches a limit as n goes to infinity. The limit is approximately 1.64 for binary alphabets, 0.69 for ternary alphabets, and monotone decreasing in alphabet size. Thus

Theorem 7 *In the average case, Compute_MEC requires $\mathcal{O}(n)$ time and $\Theta(\log n)$ additional bits of space.*

4 Comparing Border-Based and Prefix-Based Algorithms

As has been mentioned above, in order to compute $\text{MEC}_{\mathbf{x}}$, the authors of [12] made use of the *border array*. On the other hand Compute_MEC is based on the *prefix table*. We have already highlighted the advantage Compute_MEC has because of the use of a *prefix table* in lieu of a *border array* especially in the context of indeterminate strings. Additionally, the simplicity and low space usage of Compute_MEC encourage us to compare its practical performance with the algorithm of [12]. To this end this can be seen as a comparison between a border-based algorithm (i.e., the algorithm of [12]) for computing $\text{MEC}_{\mathbf{x}}$ and a prefix-based algorithm (i.e, Compute_MEC of the current paper) for doing the same. In what follows we will refer to the former algorithm as ECB and the latter as ECP.

We have implemented ECP (i.e., Compute_MEC) in C# using Visual Studio 2010. We got the implementation of ECB from the authors of [12]. However, ECB was implemented in C. To ensure a level playing ground, we re-implemented ECB in C# following their implementation. Then we have run both the algorithms on all binary strings of lengths 2 to 30. The experiments have been carried out on a Windows Server 2008 R2 64-bit Operating System, with Intel(R) Core(TM) i7 2600 processor @ 3.40GHz having an installed memory (RAM) of 8.00 GB. The results are illustrated in Figure 3 and 4, where the maximum number of operations carried out by each algorithm is reported in Figure 3. Figure 4 shows the ratio of the total number of operations performed by the Border-Based (ECB) [12] and Prefix-Based (ECP) algorithm to the length n of string, for all strings on the binary alphabet. As is evident from Figure 3 and 4, ECP outperforms ECB and in fact it does show a linear behaviour verifying the claim in Theorem 7 above.

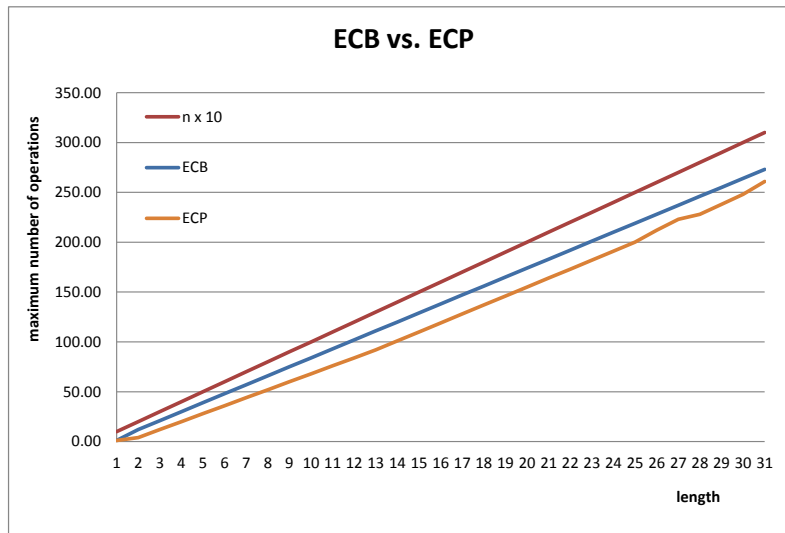


Fig. 3. The maximum number of operations performed by the Border-Based (ECB) [12] and Prefix-Based (ECP) algorithm (i.e., Compute_MEC) to compute the Minimum Enhanced Cover array, for all strings on the binary alphabet.

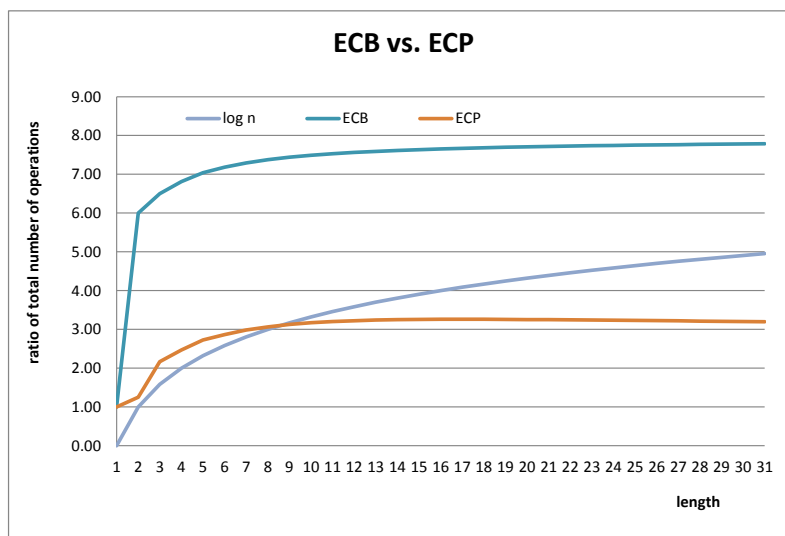


Fig. 4. Ratio of the total number of operations performed by the Border-Based (ECB) [12] and Prefix-Based (ECP) algorithm to the length n of string, for all strings on the binary alphabet.

5 Extensions

In Sections 2 and 3 we describe an algorithm to compute the minimum enhanced cover array $\text{MEC}_{\mathbf{x}}$ of a given string \mathbf{x} , based only on the prefix array $\boldsymbol{\pi}_{\mathbf{x}}$. As noted in the Introduction, since the prefix array can be computed also for indeterminate strings [20], this immediately raises the possibility of extending the MEC calculation to indeterminate strings.

In [2] two definitions of “cover” for an indeterminate string are proposed: a *sliding cover* where adjacent or overlapping covering substrings of \mathbf{x} must match, and a *rooted cover* where each covering substring is constrained only to match a prefix of \mathbf{x} . The nontransitivity of matching (see Section 1) inhibits implementation of a sliding cover, but [2] shows how to compute all the rooted covers of indeterminate \mathbf{x} from its prefix array in $\mathcal{O}(n^2)$ worst case time, $\Theta(n)$ in the average case. Thus it becomes possible to execute `Compute_MNC` for rooted covers, simply by replacing the function call to `Compute_PCR` by a function call to `PCInd` of [2]; that is, to compute the rooted cover array $\boldsymbol{\gamma}_{\mathbf{R}}[1..B]$, hence `MNC[1..B]` and thus $\text{MEC}_{\mathbf{x}}$, all for indeterminate strings. Let us call this new algorithm `Compute_MEC_Ind`. We recall now a lemma from [5] stating that the expected number of borders in an indeterminate string is bounded above by a constant, approximately 29. Therefore, also for indeterminate strings, B can be treated as a constant, and we have the following remarkable result:

Theorem 8 *In the average case, `Compute_MEC_Ind` requires $\mathcal{O}(n)$ time and $\Theta(\log n)$ additional bits of space.*

We note further that the prefix array can be efficiently computed in a compressed form [20], taking advantage of the fact that for $i \in 1..n$, $\boldsymbol{\pi}[i] \neq 0$ if and only if $\mathbf{x}[i] = \mathbf{x}[1]$. Thus we can use two arrays `POS` and `LEN` to store nonzero positions in $\boldsymbol{\pi}$ and the values at those positions, respectively, thus saving much space in cases that arise in practice. We have designed a `POS/LEN` version of `Compute_MEC` that space restrictions do not allow us to describe here.

Finally, [12] describes extensions of the minimum enhanced cover array calculation, as follows:

- computation of the enhanced left-cover array of \mathbf{x} ;
- computation of the enhanced left-seed array of \mathbf{x} .

Our prefix array approach yields efficient algorithms for these problems also, that may similarly be extended to rooted covers of indeterminate strings.

References

1. Ali Alatabbi, M. Sohel Rahman & W. F. Smyth, **Inferring an indeterminate string from a prefix graph**, *J. Discrete Algorithms* (2014) to appear.
2. Ali Alatabbi, M. Sohel Rahman & W. F. Smyth, **Computing covers using prefix tables**, <http://arxiv.org/abs/1412.3016>.
3. Alberto Apostolico & Andrzej Ehrenfeucht, *Efficient Detection of Quasi-periodicities in Strings*, Tech. Report No. 90.5, The Leonardo Fibonacci Institute, Trento, Italy (1990).
4. Alberto Apostolico, Martin Farach & C. S. Iliopoulos, Optimal superprimitivity testing for strings, *Inform. Process. Lett.* 39-1 (1991) 17-20.
5. Md. Faizul Bari, Mohammad Sohel Rahman, Rifat Shahriyar, **Finding All Covers of an Indeterminate String in $O(n)$ Time on Average**. *Proc. Prague Stringology Conference* (2009) 263–271
6. Widmer Bland, Gregory Kucherov & W. F. Smyth, **Prefix table construction & conversion**, *Proc. 24th IWOCA*, Springer Lecture Notes in Computer Science LNCS 8288 (2013) 41–53.
7. Francine Blanchet-Sadri, *Algorithmic Combinatorics on Partial Words*, Chapman & Hall/CRC (2008) 385 pp.
8. D. Breslauer, An on-line string superprimitivity test, *Inform. Process. Lett.* 44-6 (1992) 345-347.
9. Manolis Christodoulakis, P, J. Ryan, W. F. Smyth & Shu Wang, **Indeterminate strings, prefix arrays & undirected graphs**, submitted for publication (2013).
10. Richard Cole, C. S. Iliopoulos, Manal Mohamed, W. F. Smyth & Lu Yang, **The complexity of the minimum k-cover problem**, *J. Automata, Languages & Combinatorics* 10-5/6 (2005) 641–653.
11. Michael J. Fischer & Michael S. Paterson, **String-matching and other products**, *Complexity of Computation, Proc. SIAM-AMS* 7 (1974) 113-125.
12. Tomáš Flouri, C. S. Iliopoulos, Tomasz Kociumaka, Solon P. Pissis, Simon J. Puglisi, W. F. Smyth & Wojciech Tyczyński, **Enhanced string covering**, *Theoret. Comput. Sci.* 506 (2013) 102–114.
13. Jan Holub & W. F. Smyth, **Algorithms on indeterminate strings**, *Proc. 14th Australasian Workshop on Combinatorial Algs.* (2003) 36–45.
14. C. S. Iliopoulos, Manal Mohamed & W. F. Smyth, **New complexity results for the k-covers problem**, *Inform. Sciences* 181 (2011) 2571–2575.
15. C. S. Iliopoulos & W. F. Smyth, **On-line algorithms for k-covering**, *Proc. Ninth Australasian Workshop on Combinatorial Algs.* (1998) 64–73.
16. Yin Li & W. F. Smyth, **An optimal on-line algorithm to compute all the covers of a string**, *Algorithmica* 32-1 (2002) 95–106.
17. Dennis Moore & W. F. Smyth, An optimal algorithm to compute all the covers of a string, *Inform. Process. Lett.* 50 (1994) 239-246.
18. Dennis Moore & W. F. Smyth, Correction to: An optimal algorithm to compute all the covers of a string, *Inform. Process. Lett.* 54 (1995) 101-103.
19. Bill Smyth, *Computing Patterns in Strings*, Pearson Addison-Wesley (2003) 423 pp.
20. W. F. Smyth & Shu Wang, **New perspectives on the prefix array**, *Proc. 15th String Processing & Inform. Retrieval Symp.*, Springer Lecture Notes in Computer Science LNCS 5280 (2008) 133–143.
21. W. F. Smyth & Shu Wang, **A new approach to the periodicity lemma on strings with holes**, *Theoret. Comput. Sci.* 410-43 (2009) 4295–4302.

22. W. F. Smyth & Shu Wang, **An adaptive hybrid pattern-matching algorithm on indeterminate strings**, *Internat. J. Foundations of Computer Science* 20–6 (2009) 985–1004.

APPENDIX

We write $|\mathbf{x}|$ for the length of string \mathbf{x} . Here we show that the expected length of the longest border of a string \mathbf{x} approaches a limit as $|\mathbf{x}|$ tends to infinity, the limit depending on the alphabet size. For a binary alphabet it is approximately 1.64. We use the following notation. $\sigma = |\Sigma|$ is the alphabet size, $B(w)$ is length of longest border of string w and $B_k(w)$ is length of longest border of string w which has length at most k (ie, ignoring any borders longer than k). Thus if $\mathbf{x} = babaabababbabaabab$ then $B(\mathbf{x}) = 8$ since \mathbf{x} has longest border $babaabab$ and $B_4(\mathbf{x}) = 3$ since the longest border of \mathbf{x} which has length at most 4 is aba . W_n is the set of all strings of length n on an alphabet of size σ . Since W_0 contains only the empty string we have $|W_0| = 1$.

Lemma 9 *The number of strings of length n on an alphabet of size σ which have a border of length k (not necessarily the longest border) is σ^{n-k} .*

Proof. A string with border of length k is periodic with period $n - k$ and so is determined by its length $n - k$ prefix. This prefix can be chosen in σ^{n-k} ways. \square

We also need the following formula (which can be obtained using a computer algebra system).

Lemma 10 $\sum_{i=a}^b m\sigma^m = \frac{\sigma^{b+1}(\sigma(b+1) - \sigma - b - 1)}{(\sigma - 1)^2} - \frac{\sigma^a(a\sigma - a - \sigma)}{(\sigma - 1)^2}$.

Clearly $|W_n| = \sigma^n$. The expected size of the longest border of a string of length n on an alphabet of size σ is therefore

$$\bar{B}(n) = \frac{1}{\sigma^n} \sum_{w \in W_n} B(w). \quad (1)$$

Similarly, the expected size of the longest border not exceeding k is

$$\bar{B}_k(n) = \frac{1}{\sigma^n} \sum_{w \in W_n} B_k(w). \quad (2)$$

Clearly $B(w) \geq B_k(w)$ so

$$\overline{B}(n) \geq \overline{B}_k(n). \quad (3)$$

Note that if $n \geq 2k$ then $W_n = \{uvw : u \in W_k, x \in W_{n-2k}, v \in W_k\}$ and so

$$\overline{B}_k(n) = \frac{1}{\sigma^n} \sum_{u \in W_k} \sum_{x \in W_{n-2k}} \sum_{v \in W_k} B_k(uxv). \quad (4)$$

Now $B_k(uxv) = B_k(uv)$ so if $n \geq 2k$,

$$\begin{aligned} \overline{B}_k(n) &= \frac{1}{\sigma^n} \sum_{u \in W_k} \sum_{v \in W_k} B_k(uv) \sum_{x \in W_{n-2k}} 1 \\ &= \frac{\sigma^{n-2k}}{\sigma^n} \sum_{u \in W_k} \sum_{v \in W_k} B_k(uv) \\ &= \frac{1}{\sigma^{2k}} \sum_{w \in W_{2k}} B_k(w) \\ &= \overline{B}_k(2k). \end{aligned} \quad (5)$$

With (3) we then have, for $n \geq 2k$,

$$\overline{B}(n) \geq \overline{B}_k(2k). \quad (6)$$

Now any border that is counted in the right hand side of (1) but not counted on the right hand side of (2) has length at least $k+1$. The sum of the lengths of such borders is, by Lemma 9,

$$\sum_{m=k+1}^n m\sigma^{n-m}.$$

So, by Lemma 10 and (5),

$$\begin{aligned} \overline{B}(n) &\leq \frac{1}{\sigma^n} \left(\sum_{w \in W_n} B_k(w) + \sum_{m=k+1}^n m\sigma^{n-m} \right) \\ &= \overline{B}_k(n) + \frac{1}{\sigma^n} \left(\frac{\sigma^{n-k+1}k + \sigma^{n-k+1} - \sigma^{n-k}k - \sigma n - \sigma + n}{(\sigma - 1)^2} \right) \\ &< \overline{B}_k(n) + \frac{\sigma^{-k+1}k + \sigma^{-k+1} - \sigma^{-k}k - \sigma}{(\sigma - 1)^2} \\ &= \overline{B}_k(2k) + O(k\sigma^{-k}). \end{aligned} \quad (7)$$

Thus for all $n \geq 2k$

$$\overline{B}_k(2k) \leq \overline{B}(n) \leq \overline{B}_k(2k) + O(\sigma^{-k})$$

so they're contained in an arbitrarily small interval. Call this interval I_k and define $J_1 = I_1$ and for $i \geq 2$ $J_i = I_i \cap J_{i-1}$. Then J_1, J_2, \dots is a sequence of nested intervals whose lengths have limit 0. By the Nested Intervals Theorem this means the limit of \overline{B}_n exists.

Using (3) and (7) with $k = 11$ we find that $\lim_{n \rightarrow \infty} \overline{B}(n)$ lies in the interval $(1.6356, 1.6420)$ for binary alphabets. For ternary alphabets using $k = 6$ the limit lies in $(0.6811, 0.6864)$.