

# Simple KMP Pattern-Matching on Indeterminate Strings<sup>\*</sup>

Neerja Mhaskar<sup>1</sup> and W. F. Smyth<sup>1,2</sup>

<sup>1</sup> Algorithms Research Group, Department of Computing & Software  
McMaster University, Canada

pophlin@mcmaster.ca, smyth@mcmaster.ca

<sup>2</sup> School of Engineering & Information Technology  
Murdoch University, Western Australia

**Abstract.** In this paper we describe a simple, fast, space-efficient approach to finding all matches of an indeterminate pattern  $\mathbf{p} = \mathbf{p}[1..m]$  in an indeterminate string  $\mathbf{x} = \mathbf{x}[1..n]$ , where both  $\mathbf{p}$  and  $\mathbf{x}$  are defined on a “small” ordered alphabet  $\Sigma$  — say,  $\sigma = |\Sigma| \leq 9$ . A preprocessing phase replaces  $\Sigma$  by an integer alphabet  $\Sigma_I$  of size  $\sigma_I = \sigma$  that (reversibly, in time linear in string length) maps both  $\mathbf{x}$  and  $\mathbf{p}$  into equivalent regular strings  $\mathbf{y}$  and  $\mathbf{q}$ , respectively, on  $\Sigma_I$ , whose maximum (indeterminate) letter can be expressed in a 32-bit word (for  $\sigma \leq 4$ , thus for DNA sequences, an 8-bit representation suffices). We then describe an efficient version KMP\_INDET of the venerable Knuth-Morris-Pratt algorithm to find all occurrences of  $\mathbf{q}$  in  $\mathbf{y}$  (that is, of  $\mathbf{p}$  in  $\mathbf{x}$ ), but, whenever necessary, using the prefix array, rather than the border array, to control shifts of the transformed pattern  $\mathbf{q}$  along the transformed string  $\mathbf{y}$ . Although requiring  $\mathcal{O}(m^2n)$  time in the theoretical worst case, in cases of practical interest KMP\_INDET executes in  $\mathcal{O}(n)$  time. A noteworthy feature is the very small additional space requirement:  $\Theta(m)$  words in all cases. We conjecture that a similar approach may yield practical and efficient indeterminate equivalents to other well-known pattern-matching algorithms, especially Boyer-Moore and its variants.

**Keywords:** indeterminate, degenerate, conservative degenerate, pattern-matching, KMP, indeterminate encoding

## 1 Introduction

Given a fixed finite alphabet  $\Sigma = \{\lambda_1, \lambda_2, \dots, \lambda_\sigma\}$ , a *regular letter*, also called a *character*, is any single element of  $\Sigma$ , while an *indeterminate letter* is any subset of  $\Sigma$  of cardinality greater than one. A *regular string*  $\mathbf{x} = \mathbf{x}[1..n]$  on  $\Sigma$  is an array of regular letters drawn from  $\Sigma$ . An *indeterminate string*  $\mathbf{x}[1..n]$  on  $\Sigma$  is an array of letters drawn from  $\Sigma$ , of which at least one is indeterminate. Whenever entries  $\mathbf{x}[i]$  and  $\mathbf{x}[j]$ ,  $1 \leq i, j \leq n$ , both contain the same character (possibly other characters as well), we say that  $\mathbf{x}[i]$  *matches*  $\mathbf{x}[j]$  and write  $\mathbf{x}[i] \approx \mathbf{x}[j]$ .

In this paper we describe a simple transformation of  $\Sigma$  that permits all subsets of  $\Sigma$  to be replaced by single integer values, while maintaining matches and non-matches between all transformed entries  $\mathbf{x}[i_1]$  and  $\mathbf{x}[i_2]$ ,  $1 \leq i_1, i_2 \leq n$ , in  $\mathbf{x}$ . The method is effective on small alphabets (say  $|\Sigma| = \sigma \leq 9$ ), including in particular the important case of DNA sequences ( $\Sigma_{DNA} = \{a, c, g, t\}$ ). Thus, in many cases, cumbersome and time-consuming matches of indeterminate letters can be efficiently handled. For background on pattern-matching in indeterminate strings, see [6,1,8,11,9,2,16,17,4,5].

<sup>\*</sup> Supported by Grant No. 105–36797 from the Natural Sciences & Engineering Research Council of Canada (NSERC). Also the authors thank anonymous reviewers for several valuable suggestions.

We will make use of a mapping  $\mu_j \leftarrow \lambda^{(j)}$ ,  $j = 1, 2, \dots, \sigma$ , of the letters  $\lambda^{(j)}$  of  $\Sigma$  chosen in some order, where  $\mu_j$  is the  $j^{\text{th}}$  prime number ( $\mu_1 = 2, \mu_2 = 3$ , and so on). Then, given  $\mathbf{x} = \mathbf{x}[1..n]$  on  $\Sigma$  (the **source string**), we can apply the mapping to compute  $\mathbf{y} = \mathbf{y}[1..n]$  (the **mapped string**) according to the following rule:

(R) For every  $\mathbf{x}[i] = \{\lambda_1, \lambda_2, \dots, \lambda_k\}$ ,  $1 \leq k \leq \sigma$ ,  $1 \leq i \leq n$ , where  $\lambda_h \in \Sigma$ ,  $1 \leq h \leq k$ , set

$$\mathbf{y}[i] \leftarrow \prod_{h=1}^k \mu_{\lambda_h}.$$

When  $k = \sigma$ ,  $\mathbf{y}$  achieves the maximum value, which we denote by  $P_\sigma = \prod_{j=1}^{\sigma} \mu_j$  (often called a **hole** [2]). More generally, since the mapping  $\pi$  yields all possible products of the first  $\sigma$  prime numbers, it imposes an order on indeterminate letters drawn from  $\Sigma$ :  $\mathbf{x}[i_1] < \mathbf{x}[i_2] \Leftrightarrow \mathbf{y}[i_1] < \mathbf{y}[i_2]$ .

For example, consider a DNA source string  $\mathbf{x} = a\{a, c\}g\{a, t\}t\{c, g\}$ , over  $\Sigma_{DNA}$ . Then  $\sigma = 4$ , and applying (R) for  $1 \leq k \leq 4$  (based on the mapping  $\mu : 2 \leftarrow a, 3 \leftarrow c, g \leftarrow 5, 7 \leftarrow t$ ), we compute a mapped string  $\mathbf{y} = 2/6/5/14/7/15$ , so that

$$a < g < \{a, c\} < t < \{a, t\} < \{c, g\}.$$

On the other hand, a different mapping (say,  $\mu : 2 \leftarrow t, 3 \leftarrow c, 5 \leftarrow a, 7 \leftarrow g$ ) would yield  $\mathbf{y} = 5/15/7/10/2/21$  and a quite different ordering

$$t < a < g < \{a, t\} < \{a, c\} < \{c, g\}.$$

**Lemma 1** *Let  $k_i$  denote the number of letters in  $\mathbf{x}[i]$ . Then Rule (R) computes  $\mathbf{y}$  in time  $\Theta(K\mathbf{x})$ , where  $K\mathbf{x} = \sum_{i=1}^n k_i$ .*

Note that when the letters in  $\mathbf{x}$  are strongly indeterminate — that is,  $K\mathbf{x} \in \Theta(\sigma n)$  —, then the approach proposed here (replacing  $\mathbf{x}$  by  $\mathbf{y}$ ) has the advantage that subsequent processing of  $\mathbf{y}$  requires access only to a single integer at each position.

**Lemma 2** *If  $\mathbf{y}$  is computed from  $\mathbf{x}$  by Rule (R), then for every  $i_1, i_2 \in 1..n$ ,  $\mathbf{x}[i_1] \approx \mathbf{x}[i_2]$  if and only if  $\gcd(\mathbf{y}[i_1], \mathbf{y}[i_2]) > 1$ .*

*Proof.*

( $\Rightarrow$ ) By contradiction. Suppose  $\mathbf{x}[i_1] \approx \mathbf{x}[i_2]$ ,  $1 \leq i_1, i_2 \leq n$ , but  $\gcd(\mathbf{y}[i_1], \mathbf{y}[i_2]) = 1$ ; that is,  $\mathbf{y}[i_1]$  and  $\mathbf{y}[i_2]$  have no common divisor. Since for every  $i$ , the letter  $\mathbf{y}[i]$  is a product of the prime numbers assigned to the characters in  $\mathbf{x}[i]$ , we see that therefore  $\mathbf{x}[i_1]$  and  $\mathbf{x}[i_2]$  can have no character in common; that is,  $\mathbf{x}[i_1] \not\approx \mathbf{x}[i_2]$ , a contradiction.

( $\Leftarrow$ ) By the reverse argument. □

Two strings  $\mathbf{x}_1$  and  $\mathbf{x}_2$  of equal length  $n$  are said to be **isomorphic** if and only if for every  $i, j \in \{1, \dots, n\}$ ,

$$\mathbf{x}_1[i] \approx \mathbf{x}_1[j] \iff \mathbf{x}_2[i] \approx \mathbf{x}_2[j]. \quad (1)$$

We thus have:

**Observation 3** *If  $\mathbf{x}$  is an indeterminate string on  $\Sigma$ , and  $\mathbf{y}$  is the numerical string constructed by applying Rule (R) to  $\mathbf{x}$ , then  $\mathbf{x}$  and  $\mathbf{y}$  are isomorphic.*

**Observation 4** By virtue of Lemma 2 and (1),  $\mathbf{y}$  can overwrite the space required for  $\mathbf{x}$  (and vice versa) with no loss of information.

**Observation 5** Suppose  $\ell_1$  and  $\ell_2$  are integers representable in at most  $B$  bits. Then  $\gcd(\ell_1, \ell_2)$  can be computed in time bounded by  $\mathcal{O}(M_B \log B)^1$ , where  $M_B$  denotes the maximum time required to compute  $\ell_1 \ell_2$  over all such integers.

**Observation 6** For  $\sigma = 9$  corresponding to the first nine prime numbers

$$2, 3, 5, 7, 11, 13, 17, 19, 23$$

$P_\sigma = 223, 092, 870$ , a number representable in less than  $B = 32$  bits, a single computer word. Thus by Observation 5, the time required to match any two indeterminate letters is bounded by  $\mathcal{O}(5M_{32})$ . When  $\sigma = 4$ , corresponding to  $\Sigma_{DNA}$ ,  $2 \times 3 \times 5 \times 7 = 210 < 256$ , and so  $B = 8$  and the matching time reduces to  $\mathcal{O}(3M_8)$ .

**Observation 7** We assume therefore that, for  $\sigma \leq 9$ , computing a match between  $\mathbf{x}[i_1]$  and  $\mathbf{x}[i_2]$  on  $\Sigma$  (that is, between  $\mathbf{y}[i_1]$  and  $\mathbf{y}[i_2]$  computed using Rule (R)) requires time bounded above by a (small) constant.

Other models to represent indeterminate strings have been proposed [14,10]. For example, the model proposed in [14] maps all the non-empty letters (both regular and indeterminate) over the DNA alphabet  $\Sigma_{DNA} = \{A, C, G, T\}$  to the IUPAC symbols  $\Sigma_{IUPAC} = \{A, C, G, T, R, Y, S, W, K, M, B, D, H, V, N\}$ . Then given an indeterminate string over  $\Sigma_{DNA}$ , it constructs an isomorphic regular string over  $\Sigma_{IUPAC}$ . In [10], the model proposed maps each symbol in the DNA alphabet to a 4-bit integer power of 2; that is,  $\{A, C, G, T\}$  is mapped to  $\{2^0, 2^1, 2^2, 2^3\}$ . Then a non-empty indeterminate letter over  $\Sigma_{DNA}$  is represented as  $\Sigma_{\{s \in \mathcal{P}(\Sigma)\}}^s$  of maximum size  $15 = 1111_2$ . Also, instead of using the natural order on integers, [10] uses a Gray code [7] to order indeterminate letters over  $\Sigma_{DNA}$ . Note that with the Gray code two successive values differ by only one bit, such as 1100 and 1101, which enables minimizing the number of separate intervals associated with each of the four symbols of  $\Sigma_{DNA}$ .

## 2 Pattern Matching Algorithm for Indeterminate Strings

In this section we describe a simple, fast, space-efficient algorithm KMP\_INDDET that, in order to compute all occurrences of a source pattern  $\mathbf{p} = \mathbf{p}[1..m]$  in a source string  $\mathbf{x} = \mathbf{x}[1..n]$ , computes all the positions at which the corresponding mapped pattern  $\mathbf{q} = \mathbf{q}[1..m]$  occurs in the mapped string  $\mathbf{y} = \mathbf{y}[1..n]$ . We begin with the following result:

**Lemma 8** For alphabet  $\Sigma$  of size  $\sigma \leq 9$ , the positions of occurrence of  $\mathbf{p}$  in  $\mathbf{x}$  can be computed in  $O(mn)$  time.

*Proof.* By Lemma 2 and Observation 5, the positions of occurrence of  $\mathbf{q}$  in  $\mathbf{y}$  can be trivially computed in  $O(mnM_B \log B) = O(mn)$  time, with constant of proportionality  $M_B \log B$ .  $\square$

<sup>1</sup> [https://en.wikipedia.org/wiki/Greatest\\_common\\_divisor#Complexity](https://en.wikipedia.org/wiki/Greatest_common_divisor#Complexity)

As noted earlier, many pattern matching algorithms have been proposed for indeterminate strings. In [12] Iliopoulos and Radoszewski propose an  $\mathcal{O}(n \log m)$  algorithm for a constant alphabet. This is the best theoretical bound known so far for pattern matching algorithms on indeterminate strings over a constant alphabet. Recently, pattern matching algorithms for *conservative* indeterminate strings, where the number of indeterminate letters in text  $\mathbf{x}$  and pattern  $\mathbf{p}$  is bounded above by a constant  $k$ , have been proposed [4,5]. In [4], Crochemore et. al present an  $\mathcal{O}(nk)$  algorithm which uses suffix trees and other auxiliary data structures to search for  $\mathbf{p}$  in  $\mathbf{x}$ . In [5], Daykin et. al propose a pattern matching algorithm by first constructing the Burrows Wheeler Transform (BWT) of  $\mathbf{x}$  in  $\mathcal{O}(mn)$  time, and use it to find all occurrences of  $\mathbf{p}$  in  $\mathbf{x}$  in  $\mathcal{O}(km^2 + q)$  time, where  $q$  is the number of occurrences of the pattern in  $\mathbf{x}$ , and  $\mathcal{O}(km^2)$  is the time required to compute it.

## 2.1 Definitions

We give here a few essential definitions.

Given  $\mathbf{x}[1..n]$ , then for  $1 \leq i \leq n$  and  $1 \leq j \leq n$ ,  $\mathbf{u} = \mathbf{x}[i..j]$  is called a *substring* of  $\mathbf{x}$ , an *empty* substring  $\varepsilon$  if  $j < i$ . If  $i = 1$ ,  $\mathbf{u}$  is a *prefix* of  $\mathbf{x}$ , a *suffix* if  $j = n$ . A string  $\mathbf{x}$  has a *border*  $\mathbf{u}$  if  $|\mathbf{u}| < |\mathbf{x}|$  and  $\mathbf{x}$  has both prefix and suffix equal to  $\mathbf{u}$ . Note that a border of  $\mathbf{x}$  may be empty.

A *border array*  $\beta_{\mathbf{x}} = \beta_{\mathbf{x}}[1..n]$  of  $\mathbf{x}$  is an integer array where for every  $i \in [1..n]$ ,  $\beta_{\mathbf{x}}[i]$  is the length of the longest border of  $\mathbf{x}[1..i]$ . A *prefix array*  $\pi_{\mathbf{x}} = \pi_{\mathbf{x}}[1..n]$  of  $\mathbf{x}$  is an integer array where for every  $i \in [1..n]$ ,  $\pi_{\mathbf{x}}[i]$  is the length of the longest substring starting at position  $i$  that matches a prefix of  $\mathbf{x}$ . See Figure 1 for an example.

	1	2	3	4	5	6	7	8	9	10	11	12	13
$\mathbf{x}$	a	a	b	a	a	b	a	a	{a, b}	b	a	a	{a, c}
$\beta_{\mathbf{x}}$	0	1	0	1	2	3	4	5	6	3	4	5	2
$\pi_{\mathbf{x}}$	13	1	0	6	1	0	3	5	1	0	2	2	1

**Figure 1.** Border array  $\beta_{\mathbf{x}}$ , and Prefix array  $\pi_{\mathbf{x}}$  computed for the string  $\mathbf{x} = aabaabaa\{a, b\}baa\{a, c\}$ .

In Lemmas 9 and 10, we rephrase earlier results on running times for computing the border array and prefix array of a string of length  $n$ .

**Lemma 9 ([15,16])** *The border array and prefix array of a regular string of length  $m$  can be computed in  $\mathcal{O}(m)$  time.*

**Lemma 10 ([15,16])** *The border array and prefix array of an indeterminate string of length  $m$  can be computed in  $\mathcal{O}(m^2)$  time in the worst-case,  $\mathcal{O}(m)$  in the average case.*

For completeness we give in Figure 2 the KMP algorithm for regular strings  $\mathbf{x} = \mathbf{x}[1..n]$ . In case of a mismatch or after a full match, KMP computes the shift of the pattern  $\mathbf{p} = \mathbf{p}[1..m]$  along  $\mathbf{x}$  by using the border array of  $\mathbf{p}$ , which as we have seen is computable in  $\mathcal{O}(m)$  time. Thus KMP runs in  $\mathcal{O}(n)$  time.

```

function KMP( $\mathbf{x}, n, \mathbf{p}, m$ ) : Integer List
 $i \leftarrow 0$ ;  $j \leftarrow 0$ 
 $indexlist \leftarrow \emptyset$   $\triangleright$  List of indices where  $\mathbf{p}$  occurs in  $\mathbf{x}$ 
 $\beta_{\mathbf{p}} \leftarrow$  Border array of pattern  $\mathbf{p}$ 
while  $i < n$  do
  if  $\mathbf{p}[j + 1] = \mathbf{x}[i + 1]$  then
     $j \leftarrow j + 1$ ;  $i \leftarrow i + 1$ 
    if  $j = m$  then
       $indexlist \leftarrow indexlist \cup \{i - j + 1\}$ 
       $j \leftarrow \beta_{\mathbf{p}}[j]$ 
    else
      if  $j = 0$  then  $i \leftarrow i + 1$ 
      else
         $j \leftarrow \beta_{\mathbf{p}}[j]$ 
return  $indexlist$ 

```

**Figure 2.** KMP checks whether the regular pattern  $\mathbf{p}$  occurs in the regular text  $\mathbf{x}$ . If it does, then it outputs the set of indices at which  $\mathbf{p}$  occurs in  $\mathbf{x}$ ; otherwise returns an empty set.

## 2.2 KMP Algorithm for Indeterminate Strings

We now describe KMP\_INDET (see Figure 3), which searches for pattern  $\mathbf{q} = \mathbf{q}[1..m]$  in text  $\mathbf{y} = \mathbf{y}[1..n]$ , outputting the indices at which  $\mathbf{q}$  occurs in  $\mathbf{y}$  (thus, at which  $\mathbf{p}$  occurs in  $\mathbf{x}$ ). Thus our algorithm implements the KMP algorithm [13] on indeterminate strings that have been transformed using Rule (R). However, note that this transformation is not necessary for the algorithm to work: we use it to improve space and time efficiency. The algorithm also works with other indeterminate string encoding/transformations mentioned in the previous section. While scanning  $\mathbf{y}$  from left to right and performing letter comparisons, KMP\_INDET checks whether the prefix of  $\mathbf{q}$  and the substring of  $\mathbf{y}$  currently being matched are both regular. If so, then it uses the border array  $\beta_{\mathbf{q}_\ell}$  of the longest regular prefix  $\mathbf{q}_\ell$  of  $\mathbf{q}$  of length  $\ell$ , to compute the shift; if not, it constructs a new string  $\mathbf{q}'$ , which is a concatenation of the longest proper prefix of the matched prefix of  $\mathbf{q}$  and the longest proper suffix of the matched substring of  $\mathbf{y}$ , using the prefix array  $\pi_{\mathbf{q}'}$  of  $\mathbf{q}'$  to compute the shift. The COMPUTE\_SHIFT function given in Figure 4 implements this computation.

In order to determine whether or not indeterminate letters are included in any segment  $\mathbf{q}' = \mathbf{q}[1..j-1]\mathbf{y}[i-j+2..i]$ , two variables are employed:  $indet_y$  and the length  $\ell$  of the longest regular prefix  $\mathbf{q}_\ell$  of  $\mathbf{q}$ .  $indet_y$  is a Boolean variable that is true if and only if the current segment  $\mathbf{y}[i-j+2..i]$  contains an indeterminate letter;  $\ell$  is pre-computed in  $\mathcal{O}(m)$  time as a byproduct of the one-time calculation of  $\mathbf{q}_\ell$ .

If  $\mathbf{y}$  and  $\mathbf{q}$  are both regular, then KMP\_INDET reduces to the KMP algorithm [13]. Otherwise, it checks whether indeterminate letters exist in the matched prefix of  $\mathbf{q} = \mathbf{q}[1..j-1]$ , or the matched substring of  $\mathbf{y} = \mathbf{y}[i-j+2..i]$ . If they do, then the shift in  $\mathbf{q}$  is equal to the maximum length of the prefix of  $\mathbf{q}[1..j-1]$  that matches with a suffix of  $\mathbf{y}[i-j+2..i]$ . To compute this length, the algorithm first builds a new string  $\mathbf{q}' = \mathbf{q}[1..j-1]\mathbf{y}[i-j+2..i]$  and, based on an insight given in [16], computes its *prefix* array  $\pi_{\mathbf{q}'}$  rather than its border array. To compute the shift only the last  $j$  entries of  $\pi_{\mathbf{q}'}$  are examined; that is, entries  $k = j+1$  to  $2(j-1)$ . Note that we need to consider only those entries  $k$  in  $\pi_{\mathbf{q}'}[j+1..2(j-1)]$ , where a prefix of  $\mathbf{q}'$  matches the suffix at

```

function KMP_INDET( $\mathbf{y}, n, \mathbf{q}, m$ ) : Integer List
 $i \leftarrow 0$ ;  $j \leftarrow 0$ ;  $indet_y \leftarrow \mathbf{false}$ 
 $indexlist \leftarrow \emptyset$   $\triangleright$  List of indices where  $\mathbf{q}$  occurs in  $\mathbf{y}$ 
 $\mathbf{q}_\ell \leftarrow$  longest regular prefix of  $\mathbf{q}$  of length  $\ell$ 
 $\beta_{\mathbf{q}} \leftarrow$  Compute_ $\beta(\mathbf{q}_\ell)$   $\triangleright$  Border Array of  $\mathbf{q}_\ell$ 
while  $i < n$  do
  if  $\mathbf{q}[j+1] \approx \mathbf{y}[i+1]$  then
    if INDET( $\mathbf{y}[i+1]$ ) then  $indet_y \leftarrow \mathbf{true}$ 
     $j \leftarrow j+1$ ;  $i \leftarrow i+1$ 
    if  $j = m$  then
       $indexlist \leftarrow indexlist \cup \{i-j+1\}$ 
       $j \leftarrow$  Compute_Shift( $indet_y, \mathbf{y}, \mathbf{q}, i, j, \beta_{\mathbf{q}}, \ell$ )
       $indet_y \leftarrow \mathbf{false}$ 
  else
    if  $j = 0$  then  $i \leftarrow i+1$ 
    else
       $j \leftarrow$  Compute_Shift( $indet_y, \mathbf{y}, \mathbf{q}, i, j, \beta_{\mathbf{q}}, \ell$ )
       $indet_y \leftarrow \mathbf{false}$ 
return  $indexlist$ 

```

**Figure 3.** KMP\_INDET checks whether the pattern  $\mathbf{q}$  occurs in the text  $\mathbf{y}$  (both possibly indeterminate). If it does, then it outputs the set of indices at which  $\mathbf{q}$  occurs in  $\mathbf{y}$ ; otherwise returns an empty set.

$k$  ( $\mathbf{q}'[k..2(j-1)]$ ); that is, the entries where  $\pi_{\mathbf{q}'}[k] = 2j - k - 1$ . The shift is simply the maximum over such entries in  $\pi_{\mathbf{q}'}$ . (Recall that computing the border array for an indeterminate string is not useful as the matching relation  $\approx$  is not transitive [8].)

```

function COMPUTE_SHIFT( $indet_y, \mathbf{y}, \mathbf{q}, i, j, \beta_{\mathbf{q}}, \ell$ ) : Integer
 $\triangleright$   $\ell$  is length of longest regular prefix of  $\mathbf{q}$ .
if  $indet_y$  or  $j > \ell$  then
   $\mathbf{q}' = \mathbf{q}[1..j-1]\mathbf{y}[i-j+2..i]$ 
   $\pi_{\mathbf{q}'} \leftarrow$  Compute_ $\pi(\mathbf{q}')$   $\triangleright$  Prefix Array of pattern  $\mathbf{q}'$ 
   $max \leftarrow 0$ 
  for  $k = j$  to  $2(j-1)$ 
    if  $max < \pi_{\mathbf{q}'}[k]$  and  $\pi_{\mathbf{q}'}[k] = 2j - k - 1$  then
       $max \leftarrow \pi_{\mathbf{q}'}[k]$ 
   $j \leftarrow max$ 
else  $\triangleright$  prefix of  $\mathbf{q}$  & substring of  $\mathbf{y}$  are regular
   $j \leftarrow \beta_{\mathbf{q}}[j]$ 
return  $j$ 

```

**Figure 4.** COMPUTE\_SHIFT computes the shift in the pattern when a mismatch occurs or the end of pattern is reached.

Figure 5 represents the processing of the text  $\mathbf{x} = aabaabaa\{a, b\}baa\{a, c\}$  and pattern  $\mathbf{p} = aabaa$  corresponding to the processing of  $\mathbf{y}$  and  $\mathbf{q}$  by KMP\_INDET. KMP\_INDET first computes  $\beta_{\mathbf{p}} = (0, 1, 0, 1, 2)$  and  $\ell = 5$ . Initially the pattern is aligned with  $\mathbf{x}$  at position 1. Since it matches with the text ( $j = 5$ ), and  $indet_x = \mathbf{false}$  and  $5 \leq (\ell = 5)$ , we compute the shift from  $\beta_{\mathbf{p}}[5] = 2$ . Therefore, the pattern

is then aligned with  $\mathbf{x}$  at position  $i = 4$ . Analogously, the pattern is next aligned with  $\mathbf{x}$  at position  $i = 7$ . Since a mismatch occurs at  $i + 1 = 10, j + 1 = 4$ , and because  $\text{indet}_x = \text{true}$ , we construct  $\mathbf{p}' = \mathbf{p}[1..2]\mathbf{x}[8..9] = \text{aba}\{a, b\}$  and compute  $\pi_{\mathbf{p}'} = (4, 0, 2, 1)$ . Then shift is equal to 2. Therefore the pattern is aligned with  $\mathbf{x}$  at 8. Since it matches (and because it is the last match),  $\text{KMP\_INDET}$  returns the list  $\{1, 4, 8\}$ .

$\text{KMP\_INDET}$  contains a function  $\text{INDET}$  that determines whether or not the current position  $\mathbf{y}[i+1]$  is indeterminate. To enable this query to be answered efficiently, we suppose that an array  $P = P[1..9]$  has been created with  $P[t]$  equal to the  $t^{\text{th}}$  prime number in the range 2..23 (for  $\sigma = 9$ ). Then  $\mathbf{y}[i+1]$  is indeterminate if and only if it exceeds 23 or else does not occur in  $P$ . The worst case time requirement for  $\text{INDET}$  is therefore  $\log_2 9$  times a few microseconds, the time for a binary search.

	1	2	3	4	5	6	7	8	9	10	11	12	13
$\mathbf{x}$	a	a	b	a	a	b	a	a	{a, b}	b	a	a	{a, c}
	a	a	b	a	a								
				a	a	b	a	a					
							a	a	b	x			
								a	a	b	a	a	

**Figure 5.** The figure simulates the execution of  $\text{KMP\_INDET}$  on the text  $\mathbf{x} = \text{aabaabaa}\{a, b\}\text{baa}\{a, c\}$  and pattern  $\mathbf{p} = \text{aaba}$ . After execution,  $\text{KMP\_INDET}$  returns the list of positions  $\{1, 4, 8\}$  at which  $\mathbf{p}$  occurs in  $\mathbf{x}$ . ‘x’ in the third alignment identifies a mismatch.

Now we discuss the running time of algorithm  $\text{KMP\_INDET}$ . It is clear that the running time of  $\text{KMP\_INDET}$  for a regular pattern and regular text is linear. Otherwise, when a matched prefix of  $\mathbf{q}$  or a matched substring of  $\mathbf{y}$  contains an indeterminate letter, then the algorithm constructs the prefix array of a new string  $\mathbf{q}'$  which is a concatenation of the matched strings. In the worst case we might need to construct the prefix array of  $\mathbf{q}'$  for each iteration of the **while** loop. By Lemma 10 and because  $\mathbf{q}'$  can be of length at most  $2(m-1)$ , in the worst case the total time required for the execution of  $\text{KMP\_INDET}$  is  $\mathcal{O}(m^2n)$ . Theorem 11 states these conclusions:

**Theorem 11** *Given text  $\mathbf{y} = \mathbf{y}[1..n]$  and pattern  $\mathbf{q} = \mathbf{q}[1..m]$  on an alphabet of constant size  $\sigma$ ,  $\text{KMP\_INDET}$  executes in  $\mathcal{O}(n)$  time when  $\mathbf{y}$  and  $\mathbf{q}$  are both regular; otherwise, when both are indeterminate, the worst-case upper bound is  $\mathcal{O}(m^2n)$ . The algorithm’s additional space requirement is  $\mathcal{O}(m)$ , for the pattern  $\mathbf{q}'$  and corresponding arrays  $\beta_{\mathbf{q}'}$  and  $\pi_{\mathbf{q}'}$ .*

An improved theoretical bound to compute the prefix array for a string over a constant alphabet is given in [12], and is summarized in Lemma 12. Using Lemma 12 we restate Theorem 11 resulting in an improved run time complexity for  $\text{KMP\_INDET}$ .

**Lemma 12** ([12]) *The prefix array of an indeterminate string of length  $n$  over a constant-sized alphabet can be computed in  $\mathcal{O}(n\sqrt{n})$  time and  $\mathcal{O}(n)$  space.*

**Theorem 13** *Given text  $\mathbf{y} = \mathbf{y}[1..n]$  and pattern  $\mathbf{q} = \mathbf{q}[1..m]$  on an alphabet of constant size  $\sigma$ ,  $\text{KMP\_INDET}$  executes in  $\mathcal{O}(n)$  time when  $\mathbf{y}$  and  $\mathbf{q}$  are both regular; otherwise, when both are indeterminate, the worst-case upper bound is  $\mathcal{O}(nm\sqrt{m})$ .*

The algorithm's additional space requirement is  $\mathcal{O}(m)$ , for the pattern  $\mathbf{q}'$  and corresponding arrays  $\beta_{\mathbf{q}'}$  and  $\pi_{\mathbf{q}'}$ .

We provide context for the result given in Theorems 11 and 13 by the following:

**Remark 14** *One of the features that makes this algorithm truly practical is that, apart from the  $\mathcal{O}(n)$  time in-place mapping of  $\mathbf{x}$  into  $\mathbf{y}$  and  $\mathbf{p}$  into  $\mathbf{q}$ , there is no preprocessing and no auxiliary data structure requirement. As a result, processing is direct and immediate, requiring negligible additional storage.*

**Remark 15** *The worst case time requirement is predicated on a requirement for  $\mathcal{O}(n)$  (short) shifts of  $\mathbf{q}$  along  $\mathbf{y}$ , each requiring a worst-case  $\mathcal{O}(m^2)$  prefix array calculation. For example, this circumstance could occur with  $\mathbf{p} = \{a, b\}c^{m-1}$  and  $\mathbf{x} = a^n$  or with  $\mathbf{p} = ab$  and  $\mathbf{x} = \{a, c\}^n$ .*

**Remark 16** *Indeed, given a regular pattern and a string  $\mathbf{x}$  containing  $Q$  indeterminate letters (a case considered in both [4] and [5]), KMP\_INDET may make as many as  $mQ$  shifts, each requiring  $\mathcal{O}(m^2)$  processing, thus  $\mathcal{O}(m^3Q)$  overall. Therefore, if  $m^3Q$  is small with respect to  $m^2n$  —  $Q$  small with respect to  $n/m$  — then KMP\_INDET will execute in  $\mathcal{O}(n)$  time.*

### 3 Conclusion

We have described a simple procedure, based on the KMP algorithm, to do pattern-matching on indeterminate strings that is very time-efficient in cases that arise in practice and moreover uses negligible  $\Theta(m)$  space in all cases. We conjecture that a similar approach is feasible for the Boyer-Moore algorithm [3], together with its numerous variants (BM-Horspool, BM-Sunday, BM-Galil, Turbo-BM): see [15, Ch. 8] and

<https://www-igm.univ-mlv.fr/~lecroq/string/>

It would also be of interest to optimize KMP\_INDET for the conservative indeterminate strings mentioned in Section 2. And we look forward to experimental comparison of the running times of existing indeterminate pattern-matching algorithms with those of KMP\_INDET, assuming various frequencies of indeterminate letters.

### References

1. K. ABRAHAMSON: *Generalized string matching*. SIAM Journal of Computing, 16(6) 1987, pp. 1039–1051.
2. F. BLANCHET-SADRI: *Algorithmic Combinatorics on Partial Words*, Chapman & Hall CRC, 2008.
3. R. S. BOYER AND J. S. MOORE: *A fast string searching algorithm*. Communications of the ACM, 20(10) 1977, pp. 762–772.
4. M. CROCHEMORE, C. S. ILIOPOULIS, R. KUNDU, M. MOHAMED, AND F. VAYANI: *Linear algorithm for conservative degenerate pattern matching*. Eng. Appls. of Artificial Intelligence, 51 2016, pp. 109–114.
5. J. W. DAYKIN, R. GROULT, Y. GUESNET, T. LECROQ, A. LEFEBVRE, M. LÉONARD, L. MOUCHARD, E. PRIEUR-GASTON, AND B. WATSON: *Efficient pattern matching in degenerate strings with the Burrows-Wheeler transform*. Information Processing Letters, 147 2019, pp. 82–87.



6. M. FISCHER AND M. PATERSON: *String matching and other products*, in Complexity of Computation,, R. Karp, ed., American Mathematical Society, 1974, pp. 113–125.
7. F. GRAY: *Pulse code communication*. Hughes Aircraft Company, U.S. Patent no. 2632058, 1953.
8. J. HOLUB AND W. F. SMYTH: *Algorithms on indeterminate strings*. Proc. 14th Australasian Workshop on Combinatorial Algs. (AWOCA), 2003, pp. 36–45.
9. J. HOLUB, W. F. SMYTH, AND S. WANG: *Hybrid pattern-matching algorithms on indeterminate strings*, in London Algorithmics and Stringology, J. W. Daykin, M. Mohamed, and K. Steinhofel, eds., King’s College Texts in Algorithmics, 2006, pp. 115–133.
10. L. HUANG, V. POPIC, AND S. BATZOGLOU: *Short read alignment with populations of genomes*. Bioinformatics, 29(13) 06 2013, pp. i361–i370.
11. C. S. ILIOPOULOS, M. MOHAMED, L. MOUCHARD, W. F. SMYTH, K. G. PERDIKURI, AND A. K. TSAKALIDIS: *String regularities with don’t cares*. Nordic J. Computing, 10(1) 2003, pp. 40–51.
12. C. S. ILIOPOULOS AND J. RADOSZEWSKI: *Truly subquadratic-time extension queries and periodicity detection in strings with uncertainties*, in CPM, 2016.
13. D. E. KNUTH, J. H. MORRIS, AND V. R. PRATT: *Fast pattern matching in strings*. SIAM Journal of Computing, 6(2) 1977, pp. 323–350.
14. P. PROCHÁZKA AND J. HOLUB: *On-line searching in IUPAC nucleotide sequences*, in Proceedings of the 12th International Joint Conference on Biomedical Engineering Systems and Technologies (BIOSTEC 2019) - Volume 3: BIOINFORMATICS, Prague, Czech Republic, February 22-24, 2019, E. D. Maria, A. L. N. Fred, and H. Gamboa, eds., SciTePress, 2019, pp. 66–77.
15. B. SMYTH: *Computing Patterns in Strings*, Pearson/Addison–Wesley, 2003.
16. W. F. SMYTH AND S. WANG: *New perspectives on the prefix array*. Proc. 15th String Processing & Inform. Retrieval Symp. (SPIRE), 5280 2008, pp. 133–143.
17. W. F. SMYTH AND S. WANG: *An adaptive hybrid pattern-matching algorithm on indeterminate strings*. Internat. J. Foundations of Computer Science, 20(6) 2009, pp. 985–1004.