# Indeterminate strings, prefix arrays & undirected graphs

Manolis Christodoulakis [a], P.J. Ryan [b], W.F. Smyth [b,c,*,1], Shu Wang [d]

[a] *Department of Electrical & Computer Engineering, University of Cyprus, PO Box 20537, 1678 Nicosia, Cyprus*
[b] *Algorithms Research Group, Department of Computing & Software, McMaster University, Hamilton, Ontario, L8S 4K1, Canada*
[c] *School of Engineering & Information Technology, Murdoch University, 90 South Street, Murdoch, WA 6150, Australia*
[d] *IBM Toronto Software Lab, 8200 Warden Avenue, Markham, Ontario, L6G 1C7, Canada*

## A B S T R A C T

An integer array $y = y[1..n]$ is said to be ***feasible*** if and only if $y[1] = n$ and, for every $i \in 2..n$, $i \le i + y[i] \le n + 1$. A string is said to be ***indeterminate*** if and only if at least one of its elements is a subset of cardinality greater than one of a given alphabet $\Sigma$; otherwise it is said to be ***regular***. A feasible array $y$ is said to be ***regular*** if and only if it is the prefix array of some regular string. We show using a graph model that every feasible array of integers is a prefix array of some (indeterminate or regular) string, and for regular strings corresponding to $y$, we use the model to provide a lower bound on the alphabet size. We show further that there is a 1–1 correspondence between labelled simple graphs and indeterminate strings, and we show how to determine the minimum alphabet size $\sigma$ of an indeterminate string $x$ based on its ***associated graph*** $\mathcal{G}_x$. Thus, in this sense, indeterminate strings are a more natural object of combinatorial interest than the strings on elements of $\Sigma$ that have traditionally been studied.

© 2015 Elsevier B.V. All rights reserved.

## 1. Introduction

Pattern matching in strings — that is, locating all the occurrences of a given pattern in a given text — has been studied for at least half a century. A major breakthrough was the realization that preprocessing the pattern would allow the problem to be solved significantly faster. Perhaps the first form of preprocessing was proposed in the seminal paper by Morris & Pratt [24], which computed the ***border array*** of the pattern; that is, an array $\beta$, of the same length as the pattern $p$, such that $\beta[i]$ is the length of the longest proper prefix of $p[1..i]$ that is also a suffix.

In recent years, a generalization of the classical string pattern matching problem has been introduced, where either the pattern or the text, or both, contain *sets* of symbols at each position, as opposed to a single symbol per position in regular strings. These types of sequences are known as ***indeterminate strings*** and were first introduced in a famous paper by Fischer & Paterson [12], then later studied by Abrahamson [1]. In the last ten years or so, much work has been done by Blanchet-Sadri and her associates (for example, [4]) on "strings with holes" — that is, strings on an alphabet $\Sigma$ augmented by a single letter, a "hole" or "wildcard", that matches all other symbols in $\Sigma$. The monograph [3] summarizes much of the pioneering work in this area. For indeterminate strings in their full generality, the third and fourth authors of this paper

---

* Corresponding author at: Algorithms Research Group, Department of Computing & Software, McMaster University, Hamilton, Ontario, L8S 4K1, Canada.
*E-mail addresses:* christodoulakis.manolis@ucy.ac.cy (M. Christodoulakis), ryanpj@mcmaster.ca (P.J. Ryan), smyth@mcmaster.ca (W.F. Smyth), wangs@ca.ibm.com (S. Wang).
*URL:* http://www.cas.mcmaster.ca/cas/research/algorithms.htm (W.F. Smyth).

have collaborated in several papers, especially in the contexts of pattern-matching [16–18,29] and extensions to periodicity [27,28].

In the search for a preprocessing approach to speed up the pattern matching problem on indeterminate strings, it soon became clear that the border array is of limited use. For regular strings $x$, the border array has the desirable property that any border of a border of $x$ is also a border of $x$ − thus $\beta$ implicitly specifies every border of every prefix of $x$. For indeterminate $x$, however, due to the nontransitivity of the match operation, this is not true [29,28]. Hence border arrays cannot be used to speed up pattern matching on indeterminate strings. However, it turns out to be possible to make use of another data structure, the **prefix array** $\pi$, in which $\pi[i]$ is the length of the longest substring beginning at position $i$ of $x$ that matches a prefix of $x$.

Apparently the first algorithm for computing the prefix array occurred as a routine in the repetitions algorithm of Main & Lorentz [22]; see also [26, pp. 340–347]. A slightly improved algorithm is given in [21, Section 8.4], and two algorithms for computing a "compressed" prefix array are described in [27]. A comprehensive treatment of prefix array construction algorithms can be found in [5]. As noted above, for regular strings the border array and the prefix array are equivalent: it is claimed in [9,10], and demonstrated in detail in [5], that there are $\Theta(n)$-time algorithms to compute one from the other. On the other hand, as shown in [27], for indeterminate strings the prefix array actually allows all borders of every prefix to be specified, while the border array does not [16,19]. Thus the prefix array provides a more compact and more general mechanism for identifying borders, hence for describing periodicity, in indeterminate strings.

Ref. [27] describes an algorithm that computes the prefix array of any indeterminate string. In this paper we consider the "reverse engineering" problem of computing a string corresponding to a given "feasible" array $y$ − that is, any array that could conceivably be a prefix array. The first reverse engineering problem was introduced in [13,14], where a linear-time algorithm was described to compute a lexicographically least string whose border array was a given integer array − or to return the result that no such string exists. There have been many such results published since, corresponding to other data structures and other conditions; for example, [2,11,15]. In [8] a linear-time algorithm is described to compute a lexicographically least regular string $x$ corresponding to a given feasible array $y$, or to return an error if $y$ corresponds to no regular string.

In this paper we solve the more general reverse engineering problem for any feasible array $y$, regardless of whether it corresponds to a regular string or not. Moreover, we establish a remarkable connection between labelled graphs and indeterminate strings. The remainder of the paper is organized as follows. Section 2 provides preliminary information and all the necessary definitions that are used throughout the paper. In Section 3 we prove the surprising result that every feasible array is in fact a prefix array of some string (regular or indeterminate); further, we characterize the minimum alphabet size of a regular string corresponding to a given prefix array in terms of the largest clique in the negative "prefix" graph $\mathcal{P}^-$. We go on to give necessary and sufficient conditions that a given prefix array is regular. Section 4 establishes the duality between strings (whether regular or indeterminate) and labelled undirected graphs; also it provides a characterization of the minimum alphabet size of an indeterminate string $x$ in terms of the number of "independent" maximal cliques in the "associated graph" $\mathcal{G}_x$. Section 5 outlines future work.

## 2. Preliminaries

Traditionally, a string is a sequence of letters taken from some alphabet $\Sigma$. Since we discuss "indeterminate strings" in this paper, we begin by generalizing the definition as follows:

**Definition 1.** A **string** with base alphabet $\Sigma$ is either empty or else a sequence of nonempty subsets of $\Sigma$. A 1-element subset of $\Sigma$ is called a **regular** letter; otherwise it is **indeterminate**. Similarly, a nonempty string consisting only of regular letters is **regular**, otherwise **indeterminate**. The empty string $\varepsilon$ is regular.

All alphabets and all strings discussed in this paper are finite. We denote by $\Sigma'$ the set of all nonempty subsets of $\Sigma$, with $\sigma = |\Sigma|$ and $\sigma' = |\Sigma'| = 2^{\sigma} - 1$. On a given alphabet $\Sigma$, there are altogether $(\sigma')^n$ distinct nonempty strings of length $n$, of which $\sigma^n$ are regular.

**Definition 2.** Two elements $\lambda, \mu$ of $\Sigma'$ are said to **match** (written $\lambda \approx \mu$) if they have nonempty intersection. Two strings $x$, $y$ **match** ($x \approx y$) if they have the same length and all corresponding letters match.

Thus two regular letters match if and only if they are equal. But note that for indeterminate letters $\lambda, \mu, \nu$, it may be that $\lambda \approx \mu$ and $\lambda \approx \nu$, while $\mu \not\approx \nu$: for example, $\lambda = \{1, 2\}, \mu = 1, \nu = 2$.

**Definition 3.** If a string $x$ can be written $x = u_1 v$ and $x = w u_2$ for nonempty strings $v$, $w$, where $u_1 \approx u_2$, then $x$ is said to have a **border** of length $|u_1| = |u_2|$.

Note that choosing $v = w = x$ yields the empty border $\varepsilon$ of length 0.