



Murdoch
UNIVERSITY

MURDOCH RESEARCH REPOSITORY

This is the author's final version of the work, as accepted for publication following peer review but without the publisher's layout or pagination.

The definitive version is available at :

<http://dx.doi.org/10.1016/j.tcs.2009.07.010>

Smyth, W.F. and Wang, Shu (2009) A new approach to the periodicity lemma on strings with holes. *Theoretical Computer Science*, 410 (43). pp. 4295-4302.

<http://researchrepository.murdoch.edu.au/id/eprint/28087/>

Copyright: © 2009 Elsevier B.V.

It is posted here for your personal use. No further distribution is permitted.

A New Approach to the Periodicity Lemma on Strings with Holes ^{☆,☆☆}

W. F. Smyth*

*Algorithms Research Group, Department of Computing & Software, McMaster University,
Hamilton, Ontario, Canada, L8S 4K1*

*Digital Ecosystems & Business Intelligence Institute, Curtin University of Technology,
Perth WA 6845, Australia*

Shu Wang

*Algorithms Research Group, Department of Computing & Software, McMaster University,
Hamilton, Ontario, Canada, L8S 4K1*

Abstract

We first give an elementary proof of the periodicity lemma for strings containing one hole (variously called a “wild card” or a “don’t-care” or an “indeterminate letter” in the literature). The proof is modelled on Euclid’s algorithm for the greatest common divisor and is simpler than the original proof given in [BB99]. We then study the two hole case, where our result agrees with the one given in [BSH02] but is more easily proved and enables us to identify a maximum-length prefix or suffix of the string to which the periodicity lemma does apply. Finally we extend our result to three or more holes using elementary methods and state a version of the periodicity lemma that applies to all strings with or without holes. We describe an algorithm that, given the locations of the holes in a string, computes maximum length substrings to which the periodicity lemma applies, in time proportional to the number of holes. Our approach is quite different from the one in [BSH02, BS04] and also simpler.

Key words: periodicity, periodicity lemma, indeterminate string, hole

[☆]Supported in part by grants from the Natural Sciences & Engineering Research Council of Canada.

^{☆☆}The authors express their gratitude to three anonymous referees, whose comments have materially improved the quality of this paper.

*Corresponding author

Email addresses: smyth@mcmaster.ca (W. F. Smyth), shuw@mcmaster.ca (Shu Wang)

1. Introduction

Over the last few years researchers have shown interest [BB99, IMM⁺03, BSH02] in strings that may contain *don't-care* letters; that is, letters $*$ that match every letter in a given alphabet Σ . More generally, several papers [HS03, HSW06, HSW08] have studied “indeterminate” strings that may contain “indeterminate” letters — those that match various subsets of Σ . In this article we study the more general model.

Let Σ be an alphabet and let λ_i , $|\lambda_i| \geq 2$, $1 \leq i \leq m$, be pairwise distinct subsets of Σ . We form a new alphabet $\Sigma' = \Sigma \cup \{\lambda_1, \lambda_2, \dots, \lambda_m\}$ and define a new relation *match* (\approx) on Σ' as follows:

- for every $\mu_1, \mu_2 \in \Sigma$, $\mu_1 \approx \mu_2$ if and only if $\mu_1 = \mu_2$;
- for every $\mu \in \Sigma$ and every $\lambda \in \Sigma' - \Sigma$, $\mu \approx \lambda$ and $\lambda \approx \mu$ if and only if $\mu \in \lambda$;
- for every $\lambda_i, \lambda_j \in \Sigma' - \Sigma$, $\lambda_i \approx \lambda_j$ if and only if $\lambda_i \cap \lambda_j \neq \emptyset$.

This idea seems to have first been mentioned in [FP74].

We observe that *match* is reflexive and symmetric but not necessarily transitive; for example, if $\lambda = \{a, b\}$, then $a \approx \lambda$ and $b \approx \lambda$ does not imply $a \approx b$. In this paper $\mathbf{x} = \mathbf{x}[1..n]$ is always a nonempty string on Σ' that may therefore contain some $\lambda \in \Sigma' - \Sigma$ at some position $h \in 1..n$; that is, $\mathbf{x}[h] = \lambda$. We refer to an occurrence of λ in \mathbf{x} as a *hole*, generalizing the usage in [BB99, BSH02, BS04], where always $\Sigma' = \Sigma \cup \{\Sigma\}$. Here a hole is equivalent to an *indeterminate letter* as defined in [HS03]. We also sometimes refer to the position h itself as a hole.

A string \mathbf{x} has *period* (*strong period*) p if and only if for every $i, j \in 1..n$ such that $i \equiv j \pmod{p}$, $\mathbf{x}[i] \approx \mathbf{x}[j]$; \mathbf{x} has *weak period* p if and only if for every $i, j \in 1..n$ such that $j = i+p$, $\mathbf{x}[i] \approx \mathbf{x}[j]$. For example, in the following table \mathbf{x} has a weak period but not a strong period of length 2.

$$\begin{array}{cccccc} & 1 & 2 & 3 & 4 & 5 & 6 \\ \mathbf{x} = & a & b & a & * & a & c \end{array}$$

On strings without holes, periodicity and weak periodicity are equivalent.

2. Strings With One Hole

We first consider strings with exactly one hole. In [BB99] a variant of the periodicity lemma [FW65] for such strings was stated, proved, and shown to be sharp:

Lemma 1. *If \mathbf{x} with one hole has weak periods p and $q > p$, and $n \geq p+q$, then \mathbf{x} has strong period $d = \gcd(p, q)$.*

We prove this lemma here based on the Euclidean algorithm, extending the proof given in [Smy03] for the original periodicity lemma. As observed in [BB99], it suffices to establish the case $n = p+q$, since therefore for larger n , the lemma holds for every factor of length $p+q$, hence for \mathbf{x} itself. We first prove a preliminary result:

Lemma 2. *Suppose $\mathbf{x} = \mathbf{x}[1..p+q]$ has weak periods p and $q > p$ with a single hole $\mathbf{x}[h] = \lambda$.*

(a) $h \in 1..q \Rightarrow \mathbf{x}[1..q]$ has weak periods p and $q-p$;

(b) $h \in p+1..p+q \Rightarrow \mathbf{x}[p+1..p+q]$ has weak periods p and $q-p$.

Proof. We prove (a); the proof of (b) is analogous. Since \mathbf{x} has weak periods p and $q > p$, therefore $\mathbf{x}[1..q]$ has weak period p . Since for $i > p$, $i+(q-p) > q$, we need consider only $i \in 1..p$. For these values of i , it follows from weak q periodicity that $\mathbf{x}[i] \approx \mathbf{x}[i+q]$ and from weak p periodicity that $\mathbf{x}[i+q] \approx \mathbf{x}[i+q-p]$. Since $h \leq q$, we know that $\mathbf{x}[i+q] \neq \lambda$, hence that $\mathbf{x}[i] \approx \mathbf{x}[i+q-p]$. Therefore $\mathbf{x}[1..q]$ also has weak period $q-p$, as required. \square

Since h satisfies the hypothesis of either Lemma 2(a) or Lemma 2(b) (or both), we can always reduce \mathbf{x} with a single hole, whose length $p+q$ is the sum of its distinct weak periods p and q , to a substring \mathbf{y} with a single hole whose length q is the sum of its (not necessarily distinct) weak periods p and $q-p$: \mathbf{y} is either a prefix $\mathbf{x}[1..q]$ or a suffix $\mathbf{x}[p+1..p+q]$ of \mathbf{x} . If $p = q-p$, we have computed $p = \gcd(p, q) = q/2$; if not, we can perform another reduction. Let us write $\mathbf{x}^{(0)} = \mathbf{x}$ and for $r \geq 0$, let $\mathbf{x}^{(r+1)}$ be the reduction (hence a substring) of $\mathbf{x}^{(r)}$. By the correctness of the Euclidean algorithm, a finite number $k \geq 1$ of reductions yields a string $\mathbf{x}^{(k)} = \mathbf{x}^{(k)}[1..2d]$ that contains one hole and has weak period $d = \gcd(p, q)$. But then, since $\mathbf{x}^{(k)}$ takes the form $\mathbf{u}\mathbf{u}$, where $\mathbf{u} = \mathbf{x}[1..d]$, it actually has strong period d . We illustrate this reduction process with an example in Tables 1–4. Starting with a string $\mathbf{x}^{(0)}$ that has weak periods $q^{(0)} = 8$ and $p^{(0)} = 6$, we recursively reduce it to $\mathbf{x}^{(3)}$ that has a strong period 2.

$\mathbf{x}^{(0)} =$	1	2	3	4	5	6	7	8	9	10	11	12	13	14
	a	b	a	b	a	b	*	b	(a	b	a	b	a	b)

Table 1: $|\mathbf{x}^{(0)}| = 14, q^{(0)} = 8, p^{(0)} = 6, q^{(0)} - p^{(0)} = 2$

Lemma 3. *If for some $r \in 1..k$, $\mathbf{x}^{(r)}$ has strong period d , then $\mathbf{x}^{(r-1)}$ also has strong period d .*

Proof. According to the nature of a reduction, $\mathbf{x}^{(r-1)}$ has weak periods p and $q > p$ that are divisible by $d = q-p$, and $|\mathbf{x}^{(r-1)}| = p+q$. We want to prove that for every $i, j \in 1..p+q$ such that $i \equiv j \pmod{d}$, $\mathbf{x}^{(r-1)}[i] \approx \mathbf{x}^{(r-1)}[j]$. We consider three cases:

$$\mathbf{x}^{(1)} = \begin{array}{cccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 \\ (a & b) & a & b & a & b & * & b & & & & & & & & \end{array}$$

Table 2: $|\mathbf{x}^{(1)}| = 8, q^{(1)} = 6, p^{(1)} = 2, q^{(1)} - p^{(1)} = 4$

$$\mathbf{x}^{(2)} = \begin{array}{cccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 \\ & & & (a & b) & a & b & * & b & & & & & & & \end{array}$$

Table 3: $|\mathbf{x}^{(2)}| = 6, q^{(2)} = 4, p^{(2)} = 2, q^{(2)} - p^{(2)} = 2$

1. both i and j lie in $\mathbf{x}^{(r)}$;
2. one position (say i) lies in $\mathbf{x}^{(r)}$, but not j ;
3. neither i nor j lies in $\mathbf{x}^{(r)}$.

Case (1) is straightforward since $\mathbf{x}^{(r)}$ is strongly d periodic.

In case (2), assume without loss of generality that $\mathbf{x}^{(r)} = \mathbf{x}^{(r-1)}[1..q]$ — the proof for suffix $\mathbf{x}^{(r)} = \mathbf{x}^{(r-1)}[p+1..p+q]$ is analogous. By the weak periodicity of $\mathbf{x}^{(r-1)}$, $\mathbf{x}^{(r-1)}[j-q] \approx \mathbf{x}^{(r-1)}[j]$ and $\mathbf{x}^{(r-1)}[j-p] \approx \mathbf{x}^{(r-1)}[j]$, where $j-q < j-p \leq q$, so that both $j-q$ and $j-p$ are positions in $\mathbf{x}^{(r)}$. Since there is exactly one hole in $\mathbf{x}^{(r)}$, we may denote by j^* any one of $j-q, j-p$ that is *not* a hole. Since $i \equiv j \pmod{d}$ and d divides both p and q , $i \equiv j^* \pmod{d}$. Then by the strong d periodicity of $\mathbf{x}^{(r)}$,

$$\mathbf{x}^{(r-1)}[i] \approx \mathbf{x}^{(r-1)}[j^*] \approx \mathbf{x}^{(r-1)}[j].$$

Since j^* is not a hole, $\mathbf{x}^{(r-1)}[i] \approx \mathbf{x}^{(r-1)}[j]$, as required.

In case (3) we again need only consider prefix $\mathbf{x}^{(r)} = \mathbf{x}^{(r-1)}[1..q]$. Using the same argument as in case (2), we can find $j^* < q$, not a hole, such that $\mathbf{x}^{(r-1)}[j^*] \approx \mathbf{x}^{(r-1)}[j]$. But now the same construction applies also to $i > q$, allowing us to find $i^* < q$, not a hole, such that $\mathbf{x}^{(r-1)}[i^*] \approx \mathbf{x}^{(r-1)}[i]$. Since $i \equiv j \pmod{d}$, it follows that $i^* \equiv j^* \pmod{d}$, so that by the strong d periodicity of $\mathbf{x}^{(r)}$, $\mathbf{x}^{(r-1)}[i^*] \approx \mathbf{x}^{(r-1)}[j^*]$. Thus $\mathbf{x}^{(r-1)}[i] \approx \mathbf{x}^{(r-1)}[j]$. (In fact, in this case, $\mathbf{x}^{(r-1)}[i] = \mathbf{x}^{(r-1)}[j]$.) \square

Lemma 3 allows us to reconstruct \mathbf{x} by reversing the reduction, and shows that every intermediate substring $\mathbf{x}^{(r)}$ has the same strong period. Using again the example in Tables 4–1, we see that starting with $\mathbf{x}^{(3)}$ of strong period 2, every intermediate substring $\mathbf{x}^{(2)}$, $\mathbf{x}^{(1)}$, and eventually $\mathbf{x}^{(0)}$ will have the same strong period 2.

Therefore, Lemmas 2–3 imply Lemma 1, the periodicity lemma for strings with one hole.

$\mathbf{x}^{(3)} =$	1	2	3	4	5	6	7	8	9	10	11	12	13	14
					a	b	$*$	b						

Table 4: $|x^{(3)}| = 4, q^{(3)} = 2, p^{(3)} = 2, q^{(3)} - p^{(3)} = 0$

3. Strings With Two Holes

Let $\mathbf{x} = \mathbf{x}[1..n]$ be a string with two holes that is weakly p, q periodic with $q > p$, where $n \geq 2(p+q) - d$, $d = \gcd(p, q)$. Let $L_0 = p+q-d, L_1 = p+q$, and observe that $L_1 > L_0 \geq q$. Consider the prefix $\mathbf{x}_1 = \mathbf{x}[1..L_0]$ of length L_0 and the suffix $\mathbf{x}_2 = \mathbf{x}[n-L_1+1..n]$ of length L_1 . Since there are only two holes, no matter where they lie at least one of \mathbf{x}_1 and \mathbf{x}_2 must, by the periodicity lemmas for no-hole and one-hole strings, be d periodic. Of course the same statement holds for $\mathbf{x}_1 = \mathbf{x}[1..L_1]$ and $\mathbf{x}_2 = \mathbf{x}[n-L_0+1..n]$.

Since part of \mathbf{x} is strongly d periodic, we are encouraged to investigate whether there is a way to extend the d periodic portion(s), perhaps to all of \mathbf{x} . The following definition provides one basis for such an extension:

Definition 4. Suppose that $\mathbf{x} = \mathbf{x}[1..n]$ is a string with at most two holes that is weakly p, q periodic, $q > p$. For $i \in L_0+1..n$, we say that $\mathbf{x}[1..i-1]$ is *right-extendible* (RE) if at least one of the following conditions holds:

1. $\mathbf{x}[i-p] \in \Sigma$;
2. $\mathbf{x}[i-q] \in \Sigma$;
3. $i+p \leq n$ and $\mathbf{x}[i+p-q] \in \Sigma$;

For example, in Table 5, x has weak periods $q = 6$ and $p = 4$. Since $d = \gcd(6, 4) = 2$, $L_0 = 6+4-2 = 8$ and $L_1 = 6+4 = 10$. There is no hole in $x[1..L_0]$, therefore according to the original periodicity lemma, $x[1..L_0]$ is (strongly) d periodic. Furthermore, according to Definition 4, for all $i \in 9..13$, $x[1..i]$ is right-extendible, while $x[1..14]$ is not right-extendible.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
$x = a$	b	a	b	a	b	a	b	$*$	b	$*$	b	a	b	c	b	a	b

Table 5: Example: Right extendibility of a string with two holes

We remark that if neither condition (1) nor (2) in Definition 4 is satisfied, then both $i-p$ and $i-q$ are holes; since \mathbf{x} contains at most two holes, therefore for $i+p \leq n$, $\mathbf{x}[i+p] \in \Sigma$, and so condition (3) can fail to hold only in the case that $q = 2p$ — thus $i+p-q = i-p$. This is the “special” case described in [BSH02].

We shall see in the next section that for strings with an arbitrary number of holes, a weaker (and more general) definition of RE suffices. Based on the

RE property, the following lemma allows us to extend a d periodic prefix to the right:

Lemma 5. *Suppose that a string \mathbf{x} on Σ' with at most two holes is weakly p, q periodic, $q > p$, and let $d = \gcd(p, q)$. If $\mathbf{x}[1..i-1]$ is d periodic and RE, then $\mathbf{x}[1..i]$ is d periodic.*

Proof. We need only prove that for every $j \in 1..i$ such that $j \equiv i \pmod{d}$, $\mathbf{x}[j] \approx \mathbf{x}[i]$.

Suppose condition (1) of Definition 4 holds. By d periodicity, for every $j \in 1..i-1$ such that $j \equiv (i-p) \pmod{d}$, $\mathbf{x}[j] \approx \mathbf{x}[i-p]$. By weak p periodicity we know that $\mathbf{x}[i] \approx \mathbf{x}[i-p]$. Because $\mathbf{x}[i-p]$ is not a hole, it follows that for every $j \in 1..i$ such that $j \equiv i \pmod{d}$, $\mathbf{x}[j] \approx \mathbf{x}[i]$, so that $\mathbf{x}[1..i]$ is d periodic.

The proof for condition (2) is analogous.

Suppose that neither condition (1) or condition (2) holds, but that (3) is true. By d periodicity, for every $j \in 1..i-1$ such that $j \equiv (i+p-q) \pmod{d}$, $\mathbf{x}[j] \approx \mathbf{x}[i+p-q]$. Since there are at most two holes, $\mathbf{x}[i+p] \in \Sigma$ and so $\mathbf{x}[i] = \mathbf{x}[i+p]$; by weak q periodicity, $\mathbf{x}[i+p] \approx \mathbf{x}[i+p-q]$; since moreover $\mathbf{x}[i+p-q] \in \Sigma$, in fact $\mathbf{x}[i] = \mathbf{x}[i+p-q]$. It follows that for every $j \in 1..i$ such that $j \equiv i \pmod{d}$, $\mathbf{x}[j] \approx \mathbf{x}[i]$, so that again $\mathbf{x}[1..i]$ is d periodic. \square

A symmetrical definition and lemma enable us to extend a d periodic suffix to the left:

Definition 6. Suppose that $\mathbf{x} = \mathbf{x}[1..n]$ is a string with zero or more holes that is weakly p, q periodic, $q > p$. For $i \in 1..n-L_0$, we say that $\mathbf{x}[i+1..n]$ is **left-extendible** (LE) if at least one of the following conditions holds:

1. $\mathbf{x}[i+p] \in \Sigma$;
2. $\mathbf{x}[i+q] \in \Sigma$;
3. $i > p$ and $\mathbf{x}[i-p+q] \in \Sigma$;

Lemma 7. *Suppose that a string \mathbf{x} on Σ' with at most two holes is weakly p, q periodic, $q > p$, and let $d = \gcd(p, q)$. If $\mathbf{x}[i+1..n]$ is d periodic and LE, then $\mathbf{x}[i..n]$ is d periodic.* \square

We see that under specified conditions, we can extend a strongly d periodic prefix/suffix of \mathbf{x} by one to the right/left, respectively. If this process can be iterated to cover all of \mathbf{x} , then \mathbf{x} is d periodic. We summarize our results as

Lemma 8. *Suppose that $\mathbf{x} = \mathbf{x}[1..n]$ is a string with two holes and weak periods p and $q > p$, where $n \geq L_0+L_1$, $d = \gcd(p, q)$. Then:*

- (a) *At least one of $\mathbf{x}[1..L_0]$ and $\mathbf{x}[n-L_1+1..n]$ is d periodic.*
- (b) *If $\mathbf{x}[1..L_0]$ is d periodic and for every $i \in L_0+1..n$, $\mathbf{x}[1..i-1]$ is RE, then \mathbf{x} is d periodic.*

- (c) If $\mathbf{x}[n-L_1+1..n]$ is d periodic and for every $i \in 1..n-L_1$, $\mathbf{x}[i+1..n]$ is LE, then \mathbf{x} is d periodic. \square

As suggested earlier, this result can also be stated in terms of $\mathbf{x}[1..L_1]$ and $\mathbf{x}[n-L_0+1..n]$; note also that it applies to strings with any form of hole, not only don't-cares. Lemma 8 basically agrees with the result given in [BSH02], where d periodicity of \mathbf{x} is shown to depend on \mathbf{x} being “not $(2, p, q)$ -special”. However, the iterative approach given here is simpler and leads directly to a straightforward $\Theta(n)$ -time algorithm to compute the maximum-length d periodic suffix/prefix of $\mathbf{x}[1..n]$ with two holes.

To understand this better, again we consider the weakly 4, 6 periodic two-hole string of Table 5. By Lemma 5 the 2 periodic prefix $\mathbf{x}[1..8]$ can be iteratively extended to the right, yielding the conclusion that $\mathbf{x}[1..14]$ is 2 periodic. Since none of the conditions (1)-(4) of Definition 4 is satisfied in position 15, no further extension is possible. This makes sense since $\mathbf{x}[15] = c$, so that $\mathbf{x}[1..15]$ is not 2 periodic. Observe however that even if we transform \mathbf{x} into \mathbf{x}' by changing position 15 from c to a , $\mathbf{x}'[1..14]$ can still not be right-extended, because of the definition. Nevertheless \mathbf{x}' is in fact 2 periodic.

In order to resolve such situations, we state a more precise version of Lemma 8, as follows:

Corollary 9. *Suppose that $\mathbf{x} = \mathbf{x}[1..n]$ is a string with two holes h_1 and $h_2 > h_1$ and weak periods p and $q > p$, where $n \geq L_0 + L_1$, $d = \gcd(p, q)$.*

- (a) *If $h_2 - h_1 \neq q - p$, then \mathbf{x} is d -periodic.*
- (b) *If $h_2 - h_1 = q - p$, then*
- (i) *$h_2 + p > n$ or $h_1 \leq p \Rightarrow \mathbf{x}$ is d periodic;*
 - (ii) *otherwise, $\mathbf{x}[h_2 + p] = \mathbf{x}[h_1 - p] \Leftrightarrow \mathbf{x}$ is d periodic.*

Proof.

- (a) If the gap between the holes is never $q - p$, then either condition (1) or condition (2) of both Definitions 4 and 6 will hold for every i . Thus one of Lemmas 5 and 7 can be used to extend the d periodic segment of \mathbf{x} to the full range $1..n$.
- (b) Suppose then that the gap between holes is exactly $q - p$. Even so, if $h_2 + p > n$ (respectively, $h_1 \leq p$), there can exist no i such that conditions (1)-(3) of Definition 4 (respectively, 6) all fail to hold. Again, the d periodic segment can be extended, either right or left, to the full range.

Suppose then that $h_2 + p \leq n$ and $h_1 > p$. Since $n \geq L_0 + L_1$, either $\mathbf{x}[1..h_2 + p - 1]$ or $\mathbf{x}[h_1 - p + 1..n]$ is d periodic. In both cases, to establish whether the d periodic range can be extended (to $\mathbf{x}[1..h_2 + p]$ or to $\mathbf{x}[h_1 - p..n]$), it suffices to perform the single comparison

$$\mathbf{x}[h_2 + p] : \mathbf{x}[h_1 - p],$$

where, since two holes are accounted for, both must be regular letters in Σ . If unequal, then the d periodic range cannot be extended; if equal, then since the remainder of the string contains no holes, the entire string is d periodic.

□

This result yields the following simple constant-time algorithm:

```

function  $d$ -range( $\mathbf{x}, n, p, q, h_1, h_2$ )
if  $h_2 - h_1 \neq q - p$  or  $h_2 + p > n$  or  $h_1 \leq p$  then
    return  $1, n$ 
elseif  $\mathbf{x}[h_2 + p] = \mathbf{x}[h_1 - p]$  then
    return  $1, n$ 
elseif  $h_1 + h_2 > n$  then
    return  $1, h_2 + p - 1$ 
else
    return  $h_1 - p + 1, n$ 

```

Figure 1: For weakly p, q periodic $\mathbf{x}[1..n]$, $q > p$, $n \geq L_0 + L_1$, identify the maximum d periodic range that contains holes h_1 and $h_2 > h_1$.

Our methodology extends easily and naturally to three or more holes, as discussed in the next section.

4. Strings With Zero or More Holes

For a string \mathbf{x} with three holes and length $n \geq 2L_1$, again we consider a prefix $\mathbf{x}_1 = \mathbf{x}[1..L_1]$ and a suffix $\mathbf{x}_2 = \mathbf{x}[n - L_1 + 1..n]$: now both of them have length L_1 . Note that since there are only three holes, at least one of these substrings has no more than one hole. If at least two holes lie in \mathbf{x}_1 , so that at most one hole lies in \mathbf{x}_2 , then by Lemma 1 we know that \mathbf{x}_2 is d periodic; otherwise \mathbf{x}_1 is d periodic. In either case, at least a substring (prefix or suffix) of \mathbf{x} is d periodic. Figure 2 shows possible positions of these three holes, where in this case \mathbf{x}_1 is d periodic.

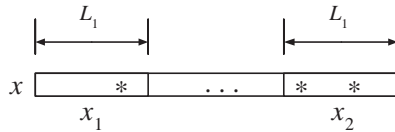


Figure 2: Possible positions of three holes

We can extend this result to any number of holes. For $d = \gcd(p, q)$, in addition to $L_0 = p + q - d$, $L_1 = p + q$, for $k \geq 2$ define $L_k = L_{k-2} + L_1$. Thus for odd k , $L_k = \lceil (k+1)/2 \rceil (p+q)$, while for even k , $L_k = L_{k+1} - d$. We claim that the following lemma holds:

Lemma 10. For a string \mathbf{x} with $k \geq 0$ holes, if \mathbf{x} is weakly p, q periodic and $|\mathbf{x}| \geq L_k$, then a substring of \mathbf{x} of length at least L_0 is d periodic, where $d = \gcd(p, q)$.

Proof. We prove this result by induction. For $k = 0$ and $k = 1$, the lemma holds by the periodicity lemmas for zero hole and one hole. If it holds for $k - 2$, then for a string \mathbf{x} with $|\mathbf{x}| \geq L_k$, we consider its prefix $\mathbf{x}_1 = x[1..L_{i-2}]$ and its suffix $\mathbf{x}_2 = x[n - L_1 + 1..n]$ of length L_1 . If the number of holes in \mathbf{x}_1 is less than or equal to $k - 2$, then by the inductive assumption \mathbf{x}_1 has a d periodic substring of length L_0 . Otherwise the number of holes in \mathbf{x}_1 is greater than $k - 2$, so that the number of holes in \mathbf{x}_2 is at most 1, implying by Lemma 1 that \mathbf{x}_2 is d periodic. \square

Note that unlike the 2-hole and 3-hole cases, in a string \mathbf{x} with more than three holes the substring of \mathbf{x} (let's call it \mathbf{x}_d) that may initially be d periodic is not necessarily a prefix or a suffix of \mathbf{x} . Therefore if \mathbf{x}_d can be extended both to the left and to the right until all of \mathbf{x} is covered, we may still claim that all of \mathbf{x} is d periodic. Observe that \mathbf{x}_d must itself contain a substring of length d without holes:

- * in the case that $|\mathbf{x}_d| = L_0$, \mathbf{x}_d contains no holes and $L_0 \geq 2d$;
- * if $|\mathbf{x}_d| = L_1$, \mathbf{x}_d contains at most one hole and $L_1 \geq 3d$.

Figure 3 demonstrates a possible position of \mathbf{x}_d and a substring of \mathbf{x}_d without holes.

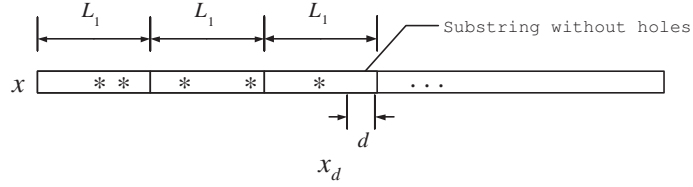


Figure 3: Possible position of \mathbf{x}_d

To accommodate three or more holes, we give a more general definition of RE and LE as follows:

Definition 11. Suppose a string \mathbf{x} with zero or more holes is weakly p, q periodic, $q > p$, with a substring $\mathbf{x}_d = \mathbf{x}[i..j]$, $j - i \geq p - 1$, that is d periodic, $d = \gcd(p, q)$.

- (a) \mathbf{x}_d is said to be RE iff $\mathbf{x}[j + 1] = \{\Sigma\}$ (hole) or there exists an integer sequence s_1, s_2, \dots, s_t , $t \geq 2$, such that

- * $s_1 = j + 1 \leq n$ and $s_t \in i..j$;
- * for every $\ell \in 2..t$, $\mathbf{x}[s_\ell] \in \Sigma$ and $|s_\ell - s_{\ell-1}| = p$ or q .

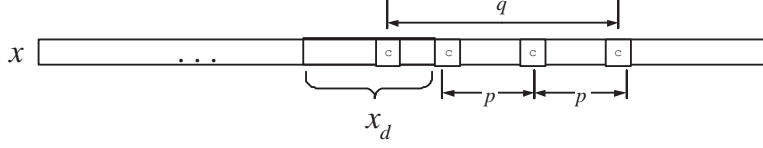


Figure 4: Example of RE and a path

(b) Symmetrically, \mathbf{x}_d is LE iff $\mathbf{x}[i-1] = \{\Sigma\}$ or there exists an integer sequence s_1, s_2, \dots, s_t , $t \geq 2$, such that

- * $s_1 = i - 1 \geq 1$ and $s_t \in i..j$;
- * for every $\ell \in 2..t$, $\mathbf{x}[s_\ell] \in \Sigma$ and $|s_\ell - s_{\ell-1}| = p$ or q .

Intuitively, this definition means that if we can find a path starting from $\mathbf{x}[j+1]$ that at each step identifies a next position p or q positions away and not a hole, terminating at a position that lies between i and j — then $\mathbf{x}[i..j]$ is RE (similarly for LE). Figure 4 illustrates an example of RE and such a path.

Note that Definitions 4 and 6 given in the previous section are special cases of this general definition.

Lemma 12. *Suppose that a string \mathbf{x} with zero or more holes is weakly p, q periodic, $q > p$, with $d = \gcd(p, q)$. If there exist i and $j \geq i + p - 1$ such that $\mathbf{x}[i..j]$ is d periodic and RE (respectively, LE), then $\mathbf{x}[i..j+1]$ (respectively, $\mathbf{x}[i-1..j]$) is d periodic.*

Proof. We prove the RE case only. If $\mathbf{x}[j+1] = \{\Sigma\}$ then certainly for every $\ell \in i..j$ such that $\ell \equiv (j+1) \pmod{d}$, $\mathbf{x}[\ell] \approx \mathbf{x}[j+1]$. Otherwise there exists a sequence s_1, s_2, \dots, s_t as described in Definition 11(a). We see that

$$\mathbf{x}[j+1] \approx \mathbf{x}[s_2] \approx \mathbf{x}[s_3] \approx \dots \approx \mathbf{x}[s_t],$$

and since every $\mathbf{x}[s_\ell] \in \Sigma$, $2 \leq \ell \leq t$, it follows that $\mathbf{x}[j+1] \approx \mathbf{x}[s_t]$. Since moreover $j+1 \equiv s_\ell \pmod{d}$ for every $\ell \in 2..t$, we conclude in particular that $j+1 \equiv s_t \pmod{d}$. Since $s_t \in i..j$ and $\mathbf{x}[i..j]$ is d periodic, therefore $\mathbf{x}[j+1] \approx \mathbf{x}[r]$ for every $r \in i..j$ such that $r \equiv (j+1) \pmod{d}$. Thus $\mathbf{x}[i..j+1]$ is d periodic, as required. \square

We now define functions **Right-Extend** and **Left-Extend** as follows:

Definition 13. Suppose that \mathbf{x} is weakly p, q periodic, $q > p$, with a d periodic substring $\mathbf{x}[i..j]$, where $d = \gcd(p, q)$ and $j - i \geq p - 1$. The function **Right-Extend** maps the pair (i, j) to $(i, j+1)$ if $\mathbf{x}[i..j]$ is RE and to (i, j) otherwise. The function **Left-Extend** maps the pair (i, j) to $(i-1, j)$ if $\mathbf{x}[i..j]$ is LE and to (i, j) otherwise.

Using these functions, we can state a general characterization of the left and right extensions that guarantee that \mathbf{x} is d periodic.

Lemma 14. *If \mathbf{x} with $k \geq 0$ holes has weak periods p and $q > p$, and $|\mathbf{x}| \geq L_k$, then at least a substring $\mathbf{x}[i..j]$ of length L_0 is d periodic, where $d = \gcd(p, q)$. If there exists a concatenation of functions $E = E_1 \circ E_2 \circ \dots \circ E_t$ where for every $\ell \in 1..t$, $E_\ell \in \{\mathbf{Right-Extend}, \mathbf{Left-Extend}\}$, and such that $E(i, j) = (1..n)$, then \mathbf{x} is d periodic. \square*

This is a statement of the periodicity lemma that applies to all strings with or without holes. However, as in the two-hole case (Corollary 9), we can be more precise: we now describe a straightforward algorithm that identifies a maximum-length d periodic substring of \mathbf{x} that contains a substring initially known to be d periodic. The algorithm uses a list of the k holes in \mathbf{x} and executes in $O(k)$ time.

Consider $\mathbf{x} = \mathbf{x}[1..n]$, $n \geq L_k$, with $k \geq 0$ holes. Suppose an array $H[1..k]$ gives the locations of all the holes in \mathbf{x} in ascending order. We add $H[0] = 0$ and $H[k+1] = n+1$. By Lemma 10 we may suppose that a $\Theta(k)$ scan of H has yielded a range $i..j$ in \mathbf{x} such that $\mathbf{x}[i..j]$ is d periodic, as well as a position s in H such that $H[s] < j$, $H[s+1] > j$, where in addition one of the following holds:

- * $j-i > L_0$ and $H[s] < i$;
- * $j-i > L_1$ and $H[s-1] < i$, $H[s] \in i..j$.

In either of these cases $\mathbf{x}[i..j]$ contains a substring $\mathbf{x}[\ell..\ell+d-1]$ such that for every $i' \in \ell..\ell+d-1$, $\mathbf{x}[i'] \in \Sigma$ (i' not a hole).

In addition to H , it is convenient also to compute a Boolean array $N[1..k]$ defined as follows: for every $s \in 1..k$, $N[s] = \mathbf{TRUE}$ if $\mathbf{x}[H[s]+q-p]$ is a hole, $N[s] = \mathbf{FALSE}$ otherwise. Figure 5 describes the preprocessing that computes N in $\Theta(k)$ time.

```

N[k] ← FALSE; r ← 2; δ ← q-p
for s ← 1 to k-1 do
  START ← H[s]
  while r ≤ k and H[r]-START < δ do
    r ← r+1
  if r > k or H[r]-START ≠ δ or H[r]+p > n then
    N[s] ← FALSE
  else
    N[s] ← TRUE; r ← r+1

```

Figure 5: Preprocessing: compute $N = N[1..k]$ in $\Theta(k)$ time from the array H of holes.

We are now in a position to describe an algorithm that extends a d periodic range $i..j$ in \mathbf{x} to the right by processing H and N from left to right, with minimal access to \mathbf{x} itself. The function **right-extend** shown in Figure 6 uses a current hole s to extend the current range: it returns $s+1$ and an extended right boundary j if further extension to the right may be possible; otherwise, it returns $s = k+1$ and the absolute rightmost boundary j of the

d periodic substring. It executes in constant time for each position s in H . (Note that here we assume the mathematical `mod` operation can be performed in constant time, since $(a \bmod b = a - \lfloor a/b \rfloor \cdot b)$; thus the complexity of `mod` is equivalent to that of division and multiplication.) A corresponding algorithm `left-extend` deals with left extension of range $i..j$. Overall, repeated execution of `right-extend` and `left-extend` will yield a maximum-length d periodic substring that contains the original d periodic range $i..j$, thus generalizing the algorithm described in Figure 1 for the two-hole case.

```

function right-extend( $H, N, s, k, \mathbf{x}, i, j, \ell, n, p, q, d$ )
if  $j - H[s] \geq q$  or not  $N[s]$  then
     $s \leftarrow s + 1$ ;  $j \leftarrow \max\{j, \min\{H[s] + q - 1, n\}\}$ 
else
     $j \leftarrow H[s] + q$ 
    if  $\mathbf{x}[j] \approx \mathbf{x}[\ell + (j - i) \bmod d]$  then
         $s \leftarrow s + 1$ 
    else
         $j \leftarrow j - 1$ ;  $s \leftarrow k + 1$ 
return  $j, s$ 

```

Figure 6: This function uses a single hole $H[s]$ to extend the d periodic range $i..j$ to the right.

We remark that a little further preprocessing may be done to form an array $\mathbf{z}[1..d] = \mathbf{x}[\ell..\ell + d - 1]$. Apart from H , N and \mathbf{z} , at most one reference to $\mathbf{x}[j]$ is then required in order to right-extend range $i..j$.

For a string \mathbf{x} with multiple holes and with weak periods $p = 4$ and $q = 6$, we illustrate the right extend process in Tables 6–8. Starting in Table 6, we first identify a periodic substring $\mathbf{x}[1..10]$ of length $p + q = 10$ with strong period $d = \gcd(4, 6) = 2$. As we already know, the existence of such a substring is guaranteed by Lemma 10. Let $[\ell..\ell + d - 1]$ be $\mathbf{x}[1..2]$. Since the position of the first hole $H[s] = 9$, we immediately know that $\mathbf{x}[1..9 + q - 1]$ is d periodic. Because every position in $\mathbf{x}[9 + 1..9 + q - 1]$ is RE according to Definition 11, Lemma 12 tells us that $\mathbf{x}[1..9 + q - 1]$ is d periodic. Since $N[s] = \text{TRUE}$ indicates that both $\mathbf{x}[15 - p]$ and $\mathbf{x}[15 - q]$ are holes, we have to compare $\mathbf{x}[15]$ with $\mathbf{x}[2]$. Since they match, we right-extend j from position 10 to 15.

Next we consider $H[s] = 11$ in Table 7. Since $N[s] = \text{FALSE}$, without any comparison we know that $\mathbf{x}[1..11 + q]$ is d periodic and therefore right-extend j to position 17.

Finally we consider $H[s] = 12$ in Table 8. Because $N[s] = \text{TRUE}$ and $\mathbf{x}[12 + q]$ does not match $\mathbf{x}[2]$, the algorithm correctly returns the maximum d periodic range $1..17$.

5. Summary and Future Work

The periodicity lemma is perhaps the fundamental result of stringology. In this paper we extend this result to strings with holes, an increasingly important

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
$\mathbf{x} =$	a	b	a	b	a	b	a	b	*	b	*	*	a	*	a	b	a	c	a	b
	i									j										

Table 6: $H[s] = 9, N[s] = TRUE, x[15] \approx x[1], j \leftarrow 15$

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
$\mathbf{x} =$	a	b	a	b	a	b	a	b	*	b	*	*	a	*	a	b	a	c	a	b
	i														j					

Table 7: $H[s] = 11, N[s] = FALSE, j \leftarrow 17$

algorithmic topic. Throughout this paper we have used elementary and simple methods independent of number theory. In the case that the number of holes is arbitrary, we have taken a quite different approach than the graph-theoretical one of [BS04]. Our Lemma 14 is very general, covering indeterminate strings whose holes are not necessarily don't-cares; it leads to the algorithm that identifies maximum-length d periodic substrings of \mathbf{x} . We would like to extend other important results in stringology to strings with holes (indeterminate strings).

References

- [BB99] J. Berstel and L. Boasson. Partial words and a theorem of Fine and Wilf. *Theoret. Comput. Sci.*, 218:135–141, 1999.
- [BS04] F. Blanchet-Sadri. Periodicity on partial words. *Computers and Mathematics with Applications*, 47:71–82, 2004.
- [BSH02] F. Blanchet-Sadri and Robert A. Hegstrom. Partial words and a theorem of fine and wilf revisited. *Theor. Comput. Sci.*, 270(1-2):401–419, 2002.
- [FP74] M. J. Fischer and M. S. Paterson. String matching and other products. In R.M.Karp, editor, *SIAM-AMS Proceedings*, number 7 in Complexity of Computation, pages 113–125, 1974.
- [FW65] N. J. Fine and H. S. Wilf. Uniqueness theorems for periodic functions. In *Proc. Amer. Math. Soc.*, volume 16, pages 109–114, 1965.
- [HS03] Jan Holub and W. F. Smyth. Algorithms on indeterminate strings. In Mirka Miller and Kunsoo Park, editors, *Proc. 14th Australasian Workshop on Combinatorial Algorithms*, pages 36–45, 2003.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
$\mathbf{x} =$	a	b	a	b	a	b	a	b	$*$	b	$*$	$*$	a	$*$	a	b	a	c	a	b
	i																			j

Table 8: $H[s] = 12$ $N[s] = TRUE$, $not(x[18] \approx x[2])$, **return**

- [HSW06] Jan Holub, W. F. Smyth, and Shu Wang. Hybrid pattern-matching algorithms on indeterminate strings. *London Algorithmics and Stringology 2006*, J. Daykin, M. Mohamed and K. Steinhofel (eds.), King's College London Series Texts in Algorithmics, pages 115–133, 2006.
- [HSW08] Jan Holub, W. F. Smyth, and Shu Wang. Fast pattern-matching on indeterminate strings. *J. Discrete Algorithms*, 6(1):37–50, 2008.
- [IMM⁺03] C.Š. Iliopoulos, Manal Mohamed, Laurent Mouchard, Katerina G. Perdikuri, W.F. Smyth, and Athanasios K. Tsakalidis. String regularities with don't cares. *Nordic J. Computing*, 10(1):40–51, 2003.
- [Smy03] Bill Smyth. *Computing Patterns in Strings*. Pearson Addison Wesley, 2003.