



HAL
open science

On the Cyclic Regularities of Strings

Oluwole Ajala, Miznah Alshammary, Mai Alzamel, Jia Gao, Costas Iliopoulos,
Jakub Radoszewski, Wojciech Rytter, Bruce Watson

► **To cite this version:**

Oluwole Ajala, Miznah Alshammary, Mai Alzamel, Jia Gao, Costas Iliopoulos, et al.. On the Cyclic Regularities of Strings. 15th IFIP International Conference on Artificial Intelligence Applications and Innovations (AIAI), May 2019, Hersonissos, Greece. pp.219-224, 10.1007/978-3-030-19909-8_19 . hal-02363851

HAL Id: hal-02363851

<https://inria.hal.science/hal-02363851>

Submitted on 14 Nov 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

On the cyclic regularities of strings

Oluwole Ajala¹, Miznah Alshammary¹, Mai Alzamel¹, Jia Gao¹, Costas Iliopoulos¹, Jakub Radoszewski², Wojciech Rytter² and Bruce Watson³

¹ Faculty of Natural and Mathematical Sciences, King's College London, United Kingdom

{oluwole.ajala, miznah.alshammary, mai.alzamel, jia.gao, c.iliopoulos}@kcl.ac.uk

² Faculty of Mathematics, Informatics and Mechanics, University of Warsaw, Poland

{jrad, rytter}@mimuw.edu.pl

³ Faculty of Informatics Science, Stellenbosch University, South Africa

{bwwatson}@sun.ac.za

Abstract. Regularities in strings are often related to periods and covers, which have extensively been studied, and algorithms for their efficient computation have broad application. In this paper we concentrate on computing cyclic regularities of strings, in particular, we propose several efficient algorithms for computing: (i) cyclic periodicity; (ii) all cyclic periodicity; (iii) maximal local cyclic periodicity; (iv) cyclic covers.

Keywords: Cyclic regularities, Periods, Covers

1 Introduction and Related Work

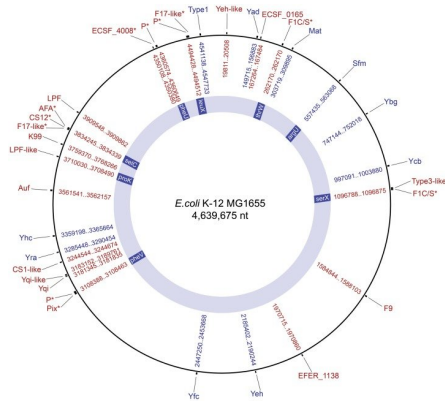
A fundamental concept of repeating patterns or *regularities* is that of periods (also known as powers). A period of order k is defined by a concatenation of k identical blocks of symbols. The study of periods can be traced to as far back as the early 1900s with the work of [9], who researched a set of strings that

do not contain any substrings that are periods. Periods in diverse forms gained prominence, when they became key structures in computational biology, where they are associated with various regulatory mechanisms and play an important role in genomic fingerprinting [6].

In regularities in strings, one of the most general notions is related to period or power, for instance, given a string x of length n , a period k of a string x is a sub string of x , if it can be decomposed into equal-length blocks of symbols, such that $x = u^k u'$, where u' is a prefix of u . However, for simplicity we will discard u' and only consider u^k .

So far, regularities in strings related with periods and powers, which have been extensively studied, [3] [4] [5] [8] and algorithms for their efficient computation have broad applications. In this paper, we study cyclic factors of strings. The motivation of cyclic factors comes from viruses. The viruses are circular strings, for example Escherichia coli (E.coli) has 154 bases and it is circular [12](Fig. 1). Formally, the viruses break up at any point of the circle, for example, that can appear in the DNA sequence as $x_\delta \dots x_n x_1 \dots x_{\delta-1}$ breaking up at position δ (Fig. 2). Now, we propose several efficient algorithms for computing: (i) cyclic periodicity; (ii) all cyclic periodicity; (iii) maximal local cyclic periodicity; (iv) cyclic covers.

Fig. 1. The E. coli K-12 MG1655 chromosome (outer black ring) was used as a reference map to visualise the locus position of 30 chromosome-borne CU fimbrial types. Types highlighted in blue are present in E. coli K-12 MG1655, types in red are absent in this strain. Fimbrial types associated with PAIs are indicated by an asterisk. A number of PAI associated fimbrial gene clusters occupy different locus positions relative to the MG1655 genome. tRNA sites that flank CU-containing PAIs are indicated on the inner blue ring. [12]



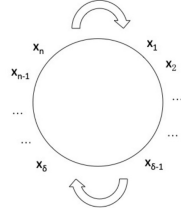


Fig. 2. Circular Pattern.

2 Preliminaries

A *string* x of length $|x| = n$ over an alphabet Σ of size σ can be denoted as $x[1..n] = x[1]x[2]\dots x[i]\dots x[n]$, where $1 \leq i \leq n$ and the i -th letter of x is denoted by $x[i] \in \Sigma$. The empty string ϵ is the string of length 0. The string x^R is the *reverse* of string x . And $x[i..j]$, $1 \leq i \leq j \leq n$, denotes the contiguous substring (or factor) of letters, such as $x[i]x[i+1]x[i+2]\dots x[j]$. A substring $x[i..j]$ is a suffix of x , if $j = n$ and a substring $x[i..j]$ is a prefix of x , if $i = 1$. Given a cyclic factor u of length k , $1 \leq k \leq n$, we denote by $c(u)$, for instance, $u = ababc$, $c(u)$ is one of the following rotations: $ababc$, $babca$, $abcab$, $bcaba$, $cabab$. Moreover, we say $u = u_1u_2\dots u_n$, $c_\delta(u) = u_\delta\dots u_nu_1\dots u_{\delta-1}$.

In this paper, suffix trees are used extensively as computation tools. For a general introduction to suffix trees, see [2], [7], [10], [11].

K-CYCLIC PERIOD

Input: Given a string x of length n , and an integer $k < n$, compute k-cyclic-period ℓ of x , where $x = u_1u_2\dots u_\ell$, $u_i = c(u_j)$, $|u_i| = |u_j| = k$, $\forall i, j$, $1 \leq i \leq \ell$, $1 \leq j \leq \ell$, and $k \times \ell = n$.

Output: k-cyclic-period ℓ of x

Example 1. Consider a string $x = aaabaabaabaabaa =: u_1u_2u_3u_4$, where $u_1 = aaab$, $u_2 = aaba$, $u_3 = abaa$, $u_4 = baaa$ and $k = 4$, $\ell = 4$. Therefore x has a period of length ℓ .

Definition 1. A cyclic periodic array A of a string x of length n is defined to be as follows: $A[i] := \ell$, $1 \leq i \leq n$, if and only if $x[1..i]$ has cyclic periodicity ℓ by a string u and there no u' , with $|u'| \leq |u|$ that is a cyclic period of $x[1..i]$.

Example 2. Consider a string $x = aababa$ of length 6, a cyclic periodic array A as follows:

$$\begin{array}{ll} x[1] = a \implies A[1] := 1 & x[1..4] = aaba \implies A[4] := 1 \\ x[1..2] = aa \implies A[2] := 2 & x[1..5] = aabab \implies A[5] := 1 \\ x[1..3] = aab \implies A[3] := 1 & x[1..6] = aababa \implies A[6] := 2 \end{array}$$

Definition 2. We define maximal local k -cyclic periodicity of a string x , if a substring y is cyclic periodic and y is not a substring of another cyclic periodic strings.

Example 3. Consider a string $x = aaaababaaa$, $\Sigma = \{a, b\}$, and a substring $y = aabababaa$ is 3-cyclic periodic and substring $y\alpha = aabababaaa$, $\alpha \in \Sigma$, is not cyclic periodic, and substring $\beta y = aaabababaa$, $\beta \in \Sigma$ is not cyclic periodic. Therefore, the substring $y = aabababaa$ is maximal local 3-cyclic periodic in string $x = aaaabababaaa$.

Definition 3. We say that a string x of length n is cyclic-coverable by a string u of length k' , if and only if, for every position i of x , the following condition holds $x[\beta..\gamma] = c(u)$, $1 \leq \beta \leq i \leq \gamma \leq n$.

Example 4. Consider a string $x = aababaa = u_1u_2$, $u_1 = aaba$, $u_2 = abaa$, $k' = 4$, $\gamma = 2$, is cyclic coverable by a string u , for every position i of x , $x[1..4] = x[4..7] = c(u)$.

Definition 4. Compute all cyclic covers of a given string x , that is for all possible length cyclic covers.

Example 5. Consider a string $x = ababbaba$, then ab , $abab$, $ababb$, $ababbab$ are cyclic covers of x .

3 Computing k -cyclic periodicity

Theorem 1. *Given a string x of length n and an integer k , $1 \leq k \leq n$, test whether it is k -cyclic periodic; this can be determined in $\mathcal{O}(n/k)$ time and $\mathcal{O}(n)$ space.*

Proof. We construct the suffix tree of x (see [7], [10], [11]). We let $u = x[1 \dots k]$, then let ℓ_m denote the depth of the lowest common ancestor of $x[1 \dots n]$ (see [1]), and $x[i_m \dots n]$. We compute the LCA ℓ_m of $x[1 \dots n]$ and $x[i_m \dots n]$ for $i_m = 2k, 3k, \dots$, and $\ell_k = n - k$, if $\ell_m = 1$ for some m , then x is not k -cyclic periodic string. Now consider $C_m^{right} = (u_{\ell_m+1} \dots u_k)^R$, compute the ℓ'_m the LCA of u^R and C_m^{right} . If $\ell'_m \geq \ell_m$ for all m , then x is k -cyclic periodic. \square

4 Computing all cyclic periodicities

Theorem 2. *Given a string x of length n , test whether it is k -cyclic periodicity for all $1 \leq k \leq n$, this can be determined in $\mathcal{O}(n \log n)$ time and $\mathcal{O}(n)$ space.*

Proof. We apply the algorithm of Theorem 1 for $k = 1, 2, \dots, n$ and we test all cyclic periods of length k . The construction of the suffix tree of string x and x^R is done once costing $\mathcal{O}(n)$. The total cost is

$$\mathcal{O}(n) + \mathcal{O}\left(\sum_{k=1}^n n/k\right) = \mathcal{O}(n \log n)$$

.

\square

Lemma 1. *Compute the cyclic period of x .*

.

Proof. The smallest cyclic-period of x is the cyclic-period of x . \square

5 Computing maximal local k -cyclic periodicity

Theorem 3. *We can compute all k -cyclic periodicity of x in $\mathcal{O}(n \log n)$ time.*

Proof. We apply the algorithm for $k = 1, 2, \dots, n$ and in this case, extend it to cyclic periods of length $k + 1$, where $|y| = m$ is cyclic periodic and $y\alpha = m + 1$ is not cyclic periodic. Next, we perform this algorithm on string x^R as $\mathcal{T}(x^R)$, where $|y| = m$, again is cyclic periodic and $\beta y = m + 1$ is not cyclic periodic.

The construction of the suffix tree of string x is done once. The total cost is

$$\mathcal{O}\left(\sum_{k=1}^n n/k\right) = \mathcal{O}(n \log n)$$

.

□

Lemma 2. *Compute maximal local k -cyclic periodicity of x .*

.

Proof. We compute and merge the arrays for $y\alpha$ and βy of x . That is the maximal local k -cyclic periodicity of x . □

6 Computing k' -cyclic coverability

Theorem 4. *Given a string x of length n and an integer k' , $1 \leq k' \leq n$, test whether it is k' -cyclic coverable, this can be determined in $\mathcal{O}(n)$ time and $\mathcal{O}(n)$ space.*

Proof. We compute the suffix tree of string x as $\mathcal{T}(x)$, and also we compute the suffix tree of string x^R as $\mathcal{T}(x^R)$.

Then we check $x[1, k']$ with each one of $x[n - k' + 1, n]$, $x[n - k', n - 1]$, $x[n - k' - 1, n - 2] \dots x[2, k' + 1]$, together with the reverse pairs in $\mathcal{T}(x^R)$. This way we build a collection of cyclic covers if there is one.

The construction of the suffix tree costs $\mathcal{O}(n)$; checking of equality costs $\mathcal{O}(1)$ and there are n factors. The total time is $\mathcal{O}(n)$. \square

7 Computing all cyclic coverability

Theorem 5. *Given a string x of length n , test whether it is k' -cyclic coverable for $1 \leq k' \leq n$, this can be determined in $\mathcal{O}(n^2)$ time and $\mathcal{O}(n)$ space.*

Proof. We apply the algorithm for $k' = 1, 2, \dots, n$ and we compare all cyclic coverable of length k' . The construction of the suffix tree of string x is done once. The total cost is

$$\mathcal{O}\left(\sum_{k'=1}^n n\right) = \mathcal{O}(n^2)$$

.

\square

Lemma 3. *Compute the cyclic coverability of x .*

.

Proof. The smallest cyclic coverable of x is the all the cyclic coveralbe of x . \square

8 Conclusions and open problems

In this paper, we defined k -cyclic periodicity, we presented several efficient algorithms for computing: (i) cyclic periodicity; (ii) all cyclic periodicity; (iii) maximal local cyclic periodicity; (iv) cyclic covers.

Future work will be focused on computing the cyclic-periodic array, that is the cyclic periodicity of every prefix of string x and computing the cyclic-coverability array, that is testing each prefix of x , for cyclic-coverability. Finally, we will extend the cyclic periodicity to cover the case $u_1u_2u_2 \dots u_ku^1$, where $u_i=c(u_j) \forall i, j$ and u^1 is a substring of some u_i .

References

1. Bender, M.A., Farach-Colton, M.: The LCA problem revisited. In: Latin American Symposium on Theoretical Informatics. LNCS, vol. 1776, pp. 88–94. Springer-Verlag (2000)
2. Crochemore, M., Hancart, C., Lecroq, T.: Algorithms on Strings. Cambridge University Press, New York, NY, USA (2007)
3. Defant, C.: Anti-power prefixes of the thue-morse word. arXiv preprint arXiv:1607.05825 (2016)
4. Erdős, P., et al.: Anti-ramsey theorems (1973)
5. Fujita, S., Magnan, C., Ozeki, K.: Rainbow generalizations of ramsey theory: a survey. *Graphs and Combinatorics* 26(1), 1–30 (2010)
6. Kolpakov, R., Bana, G., Kucherov, G.: mreps: efficient and flexible detection of tandem repeats in dna. *Nucleic acids research* 31(13), 3672–3678 (2003)
7. McCreight, E.M.: A space-economical suffix tree construction algorithm. *Journal of the ACM (JACM)* 23(2), 262–272 (1976)
8. Narayanan, S.: Functions on antipower prefix lengths of the thue-morse word. arXiv preprint arXiv:1705.06310 (2017)
9. Thue, A.: Uber unendliche zeichenreihen. *Norske Vid Selsk. Skr. I Mat-Nat Kl.(Christiana)* 7, 1–22 (1906)
10. Ukkonen, E.: Constructing suffix trees on-line in linear time. In: Proceedings of the IFIP 12th World Computer Congress on Algorithms, Software, Architecture-Information Processing'92, Volume 1-Volume I. pp. 484–492. North-Holland Publishing Co. (1992)
11. Weiner, P.: Linear pattern matching algorithms. In: 14th Annual Symposium on Switching and Automata Theory (swat 1973). pp. 1–11. IEEE (1973)
12. Wurple, D.J., Beatson, S.A., Totsika, M., Petty, N.K., Schembri, M.A.: Chaperone-usher fimbriae of escherichia coli. *PloS one* 8(1), e52835 (2013)