# Finding all covers of an indeterminate string in $O(n)$ time on average

Md. Faizul Bari    M. Sohel Rahman    Rifat Shahriyar

Department of Computer Science and Engineering
Bangladesh University of Engineering and Technology

The Prague Stringology Conference 2009

# Outline

# Outline

## Basic Definitions

For a strings $x = uwv$:

- $|x|$ is the **length** of $x$
- $\epsilon$ is the **empty** string
- $x[i]$ is the $i$-**th symbol** of $x$
- $w$ is a **substring** of $x$ and $x$ is a **superstring** of $w$
- $u(v)$ is a **prefix (suffix)** of $x$
- $x[i \ldots j]$ denotes the substring of $x$ starting at position $i$ and ending at $j$

Introduction
Our Contributions
Summary

Definitions
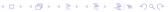The Problems That We Studied
Previous Works

## Basic Definitions

For strings $x = x[1 \ldots n]$ and $y = y[1 \ldots m]$:

- $xy$ denotes the **concatenation** of strings $x$ and $y$.
- $x^k$ denotes the concatenation of $k$ copies of $x$.
- If $x[n - i + 1 \ldots n] = y[1 \ldots i]$ for some $i \geq 1$, the string $x[1 \ldots n]y[i + 1 \ldots m]$ is a **superposition** of $x$ and $y$. We also say that $x$ overlaps $y$.

Introduction
Our Contributions
Summary

Definitions
The Problems That We Studied
Previous Works

# Border and Border Array

Border and Border Array:

- A **border** $u$ of $x$ is a prefix of $x$ that is also a suffix of $x$.

- That is $u = x[1 \dots b] = x[n - b + 1 \dots n]$ for some $b \in \{0 \dots n - 1\}$.

- The border array of $x$ is an array $\beta$ such that for all $i \in \{1 \dots n\}$, $\beta[i]$ = length of the **longest border** of $x[1 \dots i]$.

## Cover and Cover Array

Cover and Cover Array:

- A substring *w* of *x* is called a *cover* of *x*, if *x* can be constructed by **concatenating** or **overlapping** copies of *w*. We also say that *w* covers *x*.
- For example, if *x* = *ababaaba*, then *aba* and *x* are covers of *x*.
- The *cover array* $\gamma$, is a data structure used to store the length of the **longest proper cover** of every prefix of *x*;
- That is for all $i \epsilon 1 \ldots n$, $\gamma[i]$ = length of the longest proper cover of $x[1 \ldots i]$ or 0.

Introduction
Our Contributions
Summary

Definitions
The Problems That We Studied
Previous Works

## Indeterminate Strings

Indeterminate Strings:

- An indeterminate string is a sequence $T = T[1]T[2]\ldots T[n]$, where $T[i] \subseteq \Sigma$ for each $i$, and $\Sigma$ is a given alphabet of fixed size.
- If at any position in an indeterminate string, $|T[i]| = 1$, we call this a **solid symbol**. However, when $|T[i]| \geq 1$, we call this a **non-solid symbol**.
- In an indeterminate string a non-solid position can contain up to $|\Sigma|$ symbols.

Introduction
Our Contributions
Summary

Definitions
The Problems That We Studied
Previous Works

# Outline

Introduction
Our Contributions
Summary

Definitions
The Problems That We Studied
Previous Works

## The Problems That We Studied

The problems that we studied here are

- Finding all the covers of an indeterminate string
- Finding the cover array of an indeterminate string

Introduction
Our Contributions
Summary

Definitions
The Problems That We Studied
Previous Works

# The Problems That We Studied

The problems that we studied here are

- Finding all the covers of an indeterminate string
- Finding the cover array of an indeterminate string

Introduction
Our Contributions
Summary

Definitions
The Problems That We Studied
Previous Works

# Outline

Introduction
Our Contributions
Summary

Definitions
The Problems That We Studied
Previous Works

## Previous Works

Regularities of *conservative* indeterminate strings

- In [1], the authors investigated the regularities of *conservative* indeterminate strings.
- In a conservative indeterminate string the number indeterminate positions is bounded by a constant.
- The authors presented algorithms for finding
  - The smallest *conservative cover* (number of indeterminate position in the cover is bounded by a given constant)
  - $\lambda$-conservative covers (conservative covers having a fixed length $\lambda$)
  - $\lambda$-conservative seeds.

Introduction
Our Contributions
Summary

Definitions
The Problems That We Studied
Previous Works

## Previous Works

Regularities of *conservative* indeterminate strings

- In [1], the authors investigated the regularities of *conservative* indeterminate strings.

- In a conservative indeterminate string the number indeterminate positions is bounded by a constant.

- The authors presented algorithms for finding
  - The smallest *conservative cover* (number of indeterminate position in the cover is bounded by a given constant)
  - $\lambda$-conservative covers (conservative covers having a fixed length $\lambda$)
  - $\lambda$-conservative seeds.

## Previous Works

Regularities of *conservative* indeterminate strings

- In [1], the authors investigated the regularities of *conservative* indeterminate strings.
- In a conservative indeterminate string the number indeterminate positions is bounded by a constant.
- The authors presented algorithms for finding
  - The smallest *conservative cover* (number of indeterminate position in the cover is bounded by a given constant)
  - $\lambda$-conservative covers (conservative covers having a fixed length $\lambda$)
  - $\lambda$-conservative seeds.

Introduction
Our Contributions
Summary

Definitions
The Problems That We Studied
Previous Works

## Previous Works

Regularities on indeterminate strings (without any restriction)

- Antoniou et al. presented an $O(n \log n)$ algorithm to find the smallest cover of an indeterminate string in [2].
- They showed that their algorithm can be easily extended to compute all the covers of $x$. The later algorithm runs in $O(n^2 \log n)$ time.

Introduction
Our Contributions
Summary

Definitions
The Problems That We Studied
Previous Works

## Previous Works

Regularities on indeterminate strings (without any restriction)

- Antoniou et al. presented an $O(n \log n)$ algorithm to find the smallest cover of an indeterminate string in [2].
- They showed that their algorithm can be easily extended to compute all the covers of $x$. The later algorithm runs in $O(n^2 \log n)$ time.

# Outline

# Main Results

- We devise an algorithm for computing all the covers of an indeterminate string $x$ of length $n$ in $O(n^2)$ time in the worst case.

- We also show that our algorithm works in $O(n)$ time on the average.

- We extend our algorithm to compute the cover array of $x$ in $O(n^2)$ time and $O(n)$ space complexity in the worst case.

- Notably, our algorithm, unlike the algorithm of [1], does not enforce the restriction that the cover or the input string $x$ must be conservative.

# Main Results

- We devise an algorithm for computing all the covers of an indeterminate string $x$ of length $n$ in $O(n^2)$ time in the worst case.

- We also show that our algorithm works in $O(n)$ time on the average.

- We extend our algorithm to compute the cover array of $x$ in $O(n^2)$ time and $O(n)$ space complexity in the worst case.

- Notably, our algorithm, unlike the algorithm of [1], does not enforce the restriction that the cover or the input string $x$ must be conservative.

## Main Results

- We devise an algorithm for computing all the covers of an indeterminate string $x$ of length $n$ in $O(n^2)$ time in the worst case.
- We also show that our algorithm works in $O(n)$ time on the average.
- We extend our algorithm to compute the cover array of $x$ in $O(n^2)$ time and $O(n)$ space complexity in the worst case.
- Notably, our algorithm, unlike the algorithm of [1], does not enforce the restriction that the cover or the input string $x$ must be conservative.

# Main Results

- We devise an algorithm for computing all the covers of an indeterminate string $x$ of length $n$ in $O(n^2)$ time in the worst case.
- We also show that our algorithm works in $O(n)$ time on the average.
- We extend our algorithm to compute the cover array of $x$ in $O(n^2)$ time and $O(n)$ space complexity in the worst case.
- Notably, our algorithm, unlike the algorithm of [1], does not enforce the restriction that the cover or the input string $x$ must be conservative.

# Outline

# Definition of the The First Problem

We start with a formal definition of the first problem we handle in this paper.

## Problem

*Computing All Covers of an Indeterminate String over a fixed alphabet.*
**Input:** *We are given an indeterminate string x, of length n on a fixed alphabet Σ.*
**Output:** *We need to compute all the covers of x.*

# Our Algorithm Depends on the Following Facts

### Fact

*Every cover of string x is also a border of x.*

### Fact

*If u and c are covers of x and $|u| < |c|$ then u must be a cover of c.*

# Our Algorithm Depends on the Following Lemma

## Lemma

*The expected number of borders of an indeterminate string is bounded by a constant.*

# The Algorithm

- In the first step, the deterministic border array of $x$ is computed.
- In the second step, we check each border whether it is a cover of $x$ or not.

# The Algorithm

- In the first step, the deterministic border array of $x$ is computed.
- In the second step, we check each border whether it is a cover of $x$ or not.

# The First Step

- Here, we utilize the algorithm provided by Holub and Smyth [3] for computing the deterministic border of an indeterminate string.

- The output of the algorithm is a two dimensional list $\beta$.

- Each entry $\beta_i$ of $\beta$ contains a list of pair $(b, \nu_a)$, where $b$ is the length of the border and $\nu_a$ represents the required assignment.

- This list is kept sorted in decreasing order of border lengths of $x[1 \ldots i]$.

# The First Step

- Here, we utilize the algorithm provided by Holub and Smyth [3] for computing the deterministic border of an indeterminate string.

- The output of the algorithm is a two dimensional list $\beta$.

- Each entry $\beta_i$ of $\beta$ contains a list of pair $(b, \nu_a)$, where $b$ is the length of the border and $\nu_a$ represents the required assignment.

- This list is kept sorted in decreasing order of border lengths of $x[1 \ldots i]$.

# The First Step

- Here, we utilize the algorithm provided by Holub and Smyth [3] for computing the deterministic border of an indeterminate string.
- The output of the algorithm is a two dimensional list $\beta$.
- Each entry $\beta_i$ of $\beta$ contains a list of pair $(b, \nu_a)$, where $b$ is the length of the border and $\nu_a$ represents the required assignment.
- This list is kept sorted in decreasing order of border lengths of $x[1 \ldots i]$.

# The First Step

- Here, we utilize the algorithm provided by Holub and Smyth [3] for computing the deterministic border of an indeterminate string.
- The output of the algorithm is a two dimensional list $\beta$.
- Each entry $\beta_i$ of $\beta$ contains a list of pair $(b, \nu_a)$, where $b$ is the length of the border and $\nu_a$ represents the required assignment.
- This list is kept sorted in decreasing order of border lengths of $x[1 \ldots i]$.

# The First Step: Running Time Analysis

- If we assume that the maximum number of borders of any prefix of $x$ is $m$, then the worst case running time of the algorithm is $O(nm)$.

- But from Lemma 3 we know that the expected number of borders of an indeterminate string is bounded by a constant.

- As a result the expected running time of the above algorithm is $O(n)$.

# The First Step: Running Time Analysis

- If we assume that the maximum number of borders of any prefix of $x$ is $m$, then the worst case running time of the algorithm is $O(nm)$.

- But from Lemma 3 we know that the expected number of borders of an indeterminate string is bounded by a constant.

- As a result the expected running time of the above algorithm is $O(n)$.

# The First Step: Running Time Analysis

- If we assume that the maximum number of borders of any prefix of *x* is *m*, then the worst case running time of the algorithm is $O(nm)$.

- But from Lemma 3 we know that the expected number of borders of an indeterminate string is bounded by a constant.

- As a result the expected running time of the above algorithm is $O(n)$.

## The Second Step

- Now, we find out the covers of string $x$. Here we need only the last entry of the border array, $\beta_n$, where $n = |x|$.
- To identify a border as a cover of $x$ we use the pattern matching technique of an Aho-Corasick automaton.
- We build an Aho-Corasick automaton with the dictionary containing the border of $x$ and parse $x$ through the automaton to find out whether $x$ can be covered by the it or not.

# The Second Step

- Now, we find out the covers of string $x$. Here we need only the last entry of the border array, $\beta_n$, where $n = |x|$.

- To identify a border as a cover of $x$ we use the pattern matching technique of an Aho-Corasick automaton.

- We build an Aho-Corasick automaton with the dictionary containing the border of $x$ and parse $x$ through the automaton to find out whether $x$ can be covered by the it or not.

# The Second Step

- Now, we find out the covers of string $x$. Here we need only the last entry of the border array, $\beta_n$, where $n = |x|$.
- To identify a border as a cover of $x$ we use the pattern matching technique of an Aho-Corasick automaton.
- We build an Aho-Corasick automaton with the dictionary containing the border of $x$ and parse $x$ through the automaton to find out whether $x$ can be covered by the it or not.

# The Second Step: Building the Aho-Corasick automaton

- Suppose in iteration $i$, we have the length of the $i$th border of $\beta_n$ equal to $b$.
- In this iteration, we build an Aho-Corasick automaton for the following dictionary:
- 

$$D = \{x[1]x[2]\dots x[b]\}, \qquad \text{where } \forall j \in \ 1 \text{ to b }, \ x[j] \in \Sigma$$

# The Second Step: Building the Aho-Corasick automaton

- Suppose in iteration $i$, we have the length of the $i$th border of $\beta_n$ equal to $b$.
- In this iteration, we build an Aho-Corasick automaton for the following dictionary:

  ○

  $$D = \{x[1]x[2]\dots x[b]\}, \qquad \text{where } \forall j \in \text{ 1 to b }, \ x[j] \in \Sigma$$

# The Second Step: Building the Aho-Corasick automaton

- Suppose in iteration $i$, we have the length of the $i$th border of $\beta_n$ equal to $b$.
- In this iteration, we build an Aho-Corasick automaton for the following dictionary:
- 

  $$D = \{x[1]x[2]\dots x[b]\}, \qquad \text{where } \forall j \in \text{ 1 to b }, \ x[j] \in \Sigma$$

## The Second Step: Function isCover

Algorithm 1 formally presents the steps of a function *isCover*(),
which is the heart of the second step.

1: Construct the Aho-Corasick automaton for *c*
2: parse *x* and compute the positions where *c* occurs in *x* and
   put the positions in the array *Pos*
3: **for** $i = 2$ to $|Pos|$ **do**
4:    **if** $Pos[i] - Pos[i-1] > |c|$ **then**
5:       Return FALSE
6:    **end if**
7: **end for**
8: Return TRUE

**Algorithm 1**: Function isCover(x, c)

# The Second Step: Running Time Analysis of Algorithm 1

- Clearly, Steps 3 and 2 run in $O(n)$.
- Now, the complexity of Step 1 is linear in the size of the dictionary on which the automaton is build.
- Here the length of the string in the dictionary can be $n - 1$ in the worst case. So, the time and space complexity of this algorithm is $O(n)$.

# The Second Step: Running Time Analysis of Algorithm 1

- Clearly, Steps 3 and 2 run in $O(n)$.
- Now, the complexity of Step 1 is linear in the size of the dictionary on which the automaton is build.
- Here the length of the string in the dictionary can be $n - 1$ in the worst case. So, the time and space complexity of this algorithm is $O(n)$.

# The Second Step: Running Time Analysis of Algorithm 1

- Clearly, Steps 3 and 2 run in $O(n)$.
- Now, the complexity of Step 1 is linear in the size of the dictionary on which the automaton is build.
- Here the length of the string in the dictionary can be $n - 1$ in the worst case. So, the time and space complexity of this algorithm is $O(n)$.

# The Second Step: A Further Improvement

- According to Fact 2, if $u$ and $c$ are covers of $x$ and $|u| < |c|$ then $u$ must be a cover of $c$.

- Now if $\beta_n = \{b_1, b_2, \ldots, b_m\}$ then from the definition of border array $b_1 > b_2 > \ldots > b_n$.

- Now if in any iteration we find a $b_i$ that is a cover of $x$ then from Fact 2, we can say that for all $j \in i + 1 \ldots m$, if $b_j$ is a cover of $x$ if and only if $b_j$ is a cover of $b_i$.

- So instead of parsing $x$ we can parse $b_i$ for the subsequent automatons and as $|b_i| \leq |x|$.

# The Second Step: A Further Improvement

- According to Fact 2, if $u$ and $c$ are covers of $x$ and $|u| < |c|$ then $u$ must be a cover of $c$.
- Now if $\beta_n = \{b_1, b_2, \ldots, b_m\}$ then from the definition of border array $b_1 > b_2 > \ldots > b_n$.
- Now if in any iteration we find a $b_i$ that is a cover of $x$ then from Fact 2, we can say that for all $j \in i + 1 \ldots m$, if $b_j$ is a cover of $x$ if and only if $b_j$ is a cover of $b_i$.
- So instead of parsing $x$ we can parse $b_i$ for the subsequent automatons and as $|b_i| \leq |x|$.

## The Second Step: A Further Improvement

- According to Fact 2, if $u$ and $c$ are covers of $x$ and $|u| < |c|$ then $u$ must be a cover of $c$.
- Now if $\beta_n = \{b_1, b_2, \ldots, b_m\}$ then from the definition of border array $b_1 > b_2 > \ldots > b_n$.
- Now if in any iteration we find a $b_i$ that is a cover of $x$ then from Fact 2, we can say that for all $j \in i + 1 \ldots m$, if $b_j$ is a cover of $x$ if and only if $b_j$ is a cover of $b_i$.
- So instead of parsing $x$ we can parse $b_i$ for the subsequent automatons and as $|b_i| \leq |x|$.

## The Second Step: A Further Improvement

- According to Fact 2, if $u$ and $c$ are covers of $x$ and $|u| < |c|$ then $u$ must be a cover of $c$.
- Now if $\beta_n = \{b_1, b_2, \ldots, b_m\}$ then from the definition of border array $b_1 > b_2 > \ldots > b_n$.
- Now if in any iteration we find a $b_i$ that is a cover of $x$ then from Fact 2, we can say that for all $j \in i + 1 \ldots m$, if $b_j$ is a cover of $x$ if and only if $b_j$ is a cover of $b_i$.
- So instead of parsing $x$ we can parse $b_i$ for the subsequent automatons and as $|b_i| \leq |x|$.

## The Second Step: Overall Algorithms

1: $k \leftarrow n$
2: $AC \leftarrow \phi$ {$AC$ is a list used to store the covers of $x$}
3: **for all** $b \, \epsilon \, \beta_n$ **do**
4:    **if** $isCover(x[1..k], x[1..b]) = true$ **then**
5:       $m \leftarrow b$
6:       $AC.Add(k)$
7:    **end if**
8: **end for**

**Algorithm 2**: Computing all covers of $x$
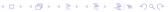
## The Second Step: Running Time Analysis

- The running time of Algorithm 2 is $O(nm)$, where $m$ is number of borders of $x$ or alternatively number of entries in $\beta_n$.

- Again, from Lemma 3 we can say that the number of borders of an indeterminate string is bounded by a constant on average.

- Hence, the expected running time of Algorithm 2 is $O(n)$.

# The Second Step: Running Time Analysis

- The running time of Algorithm 2 is $O(nm)$, where $m$ is number of borders of $x$ or alternatively number of entries in $\beta_n$.
- Again, from Lemma 3 we can say that the number of borders of an indeterminate string is bounded by a constant on average.
- Hence, the expected running time of Algorithm 2 is $O(n)$.

# The Second Step: Running Time Analysis

- The running time of Algorithm 2 is $O(nm)$, where $m$ is number of borders of $x$ or alternatively number of entries in $\beta_n$.
- Again, from Lemma 3 we can say that the number of borders of an indeterminate string is bounded by a constant on average.
- Hence, the expected running time of Algorithm 2 is $O(n)$.

## The Second Step: Running Time Analysis

- It follows from above that our algorithm for finding all the covers of an indeterminate string of length $n$ runs in $O(n)$ time on the average.

- The worst case complexity of our algorithm is $O(nm)$, i.e., $O(n^2)$.

- Which is also an improvement since the current best known algorithm [2] for finding all covers requires $O(n^2 \log n)$ in the worst case.

# The Second Step: Running Time Analysis

- It follows from above that our algorithm for finding all the covers of an indeterminate string of length $n$ runs in $O(n)$ time on the average.

- The worst case complexity of our algorithm is $O(nm)$, i.e., $O(n^2)$.

- Which is also an improvement since the current best known algorithm [2] for finding all covers requires $O(n^2 \log n)$ in the worst case.

## The Second Step: Running Time Analysis

- It follows from above that our algorithm for finding all the covers of an indeterminate string of length $n$ runs in $O(n)$ time on the average.
- The worst case complexity of our algorithm is $O(nm)$, i.e., $O(n^2)$.
- Which is also an improvement since the current best known algorithm [2] for finding all covers requires $O(n^2 \log n)$ in the worst case.

## Definition of the The Second Problem

We start with a formal definition of the second problem we handle in this paper.

### Problem

*Computing the Cover array of an Indeterminate String over a fixed alphabet.*
**Input:** *We are given an indeterminate string x, of length n on a fixed alphabet* Σ.
**Output:** *We need to compute the cover array of x.*

## Algorithm for Computing the Cover Array

- Here we only need the length of the largest border of each prefix of $x$. This information is stored in the first entry of each $\beta_i$ of the border array.

- Let us assume that $\beta_i[1]$ denotes the first entry of the list $\beta_i$ that is $\beta_i[1]$ is the length of the largest border of $x[1 \ldots i]$.

## Algorithm for Computing the Cover Array

- Here we only need the length of the largest border of each prefix of $x$. This information is stored in the first entry of each $\beta_i$ of the border array.
- Let us assume that $\beta_i[1]$ denotes the first entry of the list $\beta_i$ that is $\beta_i[1]$ is the length of the largest border of $x[1\ldots i]$.

# Algorithm for Computing the Cover Array

1: $\gamma[i] \leftarrow 0 \qquad \forall\, i\, \epsilon\, \{1\ldots n\}$
2: **for** $i \leftarrow 1$ to $n$ **do**
3:     **if** *isCover*$(x, x[1\ldots\beta_{i[1]}])$ = *true* **then**
4:         $\gamma[i] \leftarrow \beta_{i[1]}$
5:     **end if**
6: **end for**

        **Algorithm 3**: Computing cover array $\gamma$ of *x*

# Running Time Analysis

- As the worst case running time of the *isCover*$(x, c)$ function is $O(n)$ and the algorithm iterates over the $n$ lists of the border array $\beta$, the running time of Algorithm 3 is $O(n^2)$.

# Summary

- In this paper we have presented an average case $O(n)$ time and space complex algorithm for computing all the covers of a given indeterminate string $x$ of length $n$.

- We have also presented an algorithm for computing the cover array $\gamma$ of an indeterminate string. This algorithm requires $O(n^2)$ time and $O(n)$ space in the worst case.

- Both of these algorithms are improvement over existing algorithms.

# **Thank You**

## References I

📄 P. ANTONIOU, M. CROCHEMORE, C. S. ILIOPOULOS,
I. JAYASEKERA, AND G. M. LANDAU:
*Conservative string covering of indeterminate strings*.
Proceedings of the Prague Stringology Conference 2008,
2008, pp. 108–115.

📄 P. ANTONIOU, C. S. ILIOPOULOS, I. JAYASEKERA, AND
W. RYTTER:
*Computing repetitive structures in indeterminate strings*.
Proceedings of the 3rd IAPR International Conference on
Pattern Recognition in Bioinformatics (PRIB 2008), 2008.

# References II

📄 J. HOLUB AND W. F. SMYTH:
*Algorithms on indeterminate strings.*
Miller, M., Park, K. (eds.): Proceedings of the 14th
Australasian Workshop on Combinatorial Algorithms
AWOCA'03, 2003, pp. 36–45.