

Importance of Aho-Corasick String Matching Algorithms in Real World Applications

Franks Jhon Colque Zegarra¹ Luis Mamani Chirinos²

National University of San Agustin, School of Computer Science

czfranks@gmail.com

January 3, 2018

Overview

1 Introduction

- Motivation
- Set Matching Problem
- Trie

2 Aho-Corasick

- Aho-Corasick Automaton
- Functions Of Aho-Corasick
- Aho-Corasick algorithm
 - Phase I
 - Phase II

3 Applications Of Aho-Corasick

- Digital Forensics and Text Mining
- Intrusion Detecting
- Detecting Plagiarism
- Bioinformatics

Motivation

- String Matching Algorithms we try to find the position where patterns are found within a larger string or text.
- String Matching can be performed in the text String through Single Pattern and Multipattern Pattern occurrences.
- Aho-Corasick solve the multipattern matching efficiently.
- Aho-Corasick was invented by Alfred V. Aho and Margaret J. Corasick

Set Matching Problem.

In the Set Matching Problem we locate occurrences of any pattern of set $\mathcal{P} = \{P_0, \dots, P_k\}$, in target $\mathcal{T}[1 \dots m]$

Set Matching can be solved in time

$$O(|P_1| + m + \dots + |P_k| + m) = O(n + km)$$

where

$$n = \sum_{i=1}^k |P_i|$$

Aho-Corasick algorithm(AC) is a good solution to set Matching
Now set Matching can be solved in time

$$O(n + m + z)$$

where

$z = \text{number of pattern occurrences in } \mathcal{T}$

Trie

A **Trie** (or a **keyword Tree**) is an ordered tree data structure that is used to store a dynamic set or associative array where the keys are usually strings.

A Trie for a set of patterns \mathcal{P} is a rooted tree \mathcal{K} such that:

1. Each edge of \mathcal{K} is labeled by a character.
2. Any two nodes out of a node have different labels.

$L(v)$: Label of node v as concatenation on the path from the root to v .

3. For each $P \in \mathcal{P}$ there's a node v with $L(v) = P$
4. The label $L(v)$ of any leaf v equals some $P \in \mathcal{P}$.

Example Trie

A Trie for $\mathcal{P} = \{he, she, his, hers\}$:

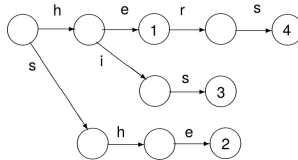


Figure: A Trie is an efficient implementation of a **dictionary** of strings.

Trie: Construction

Construction for $\mathcal{P} = \{he, she, his, hers\}$:

- 1) Begin with a root node only
- 2) Insert each pattern P_i
- 3) Starting at the root, follow the path labeled by chars of P_i
- 4) If the path ends before P_i , continue it by adding new edges and nodes for the remaining characters of P_i
- 5) Store identifier i of P_i at the terminal node of the path

This takes $O(|P_1| + \dots + |P_k|) = O(n)$ time.

Lookup of a String P : ...

Aho Corasick(AC) automaton

Aho Corasick(AC) is the multipattern matching which locates all the occurrences of set of patterns in a text of string. It first creates Deterministic automata for all predefined patterns.

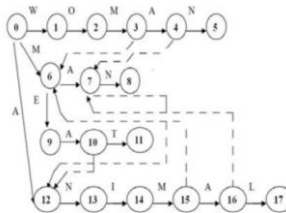


Figure: AC Automaton of patterns $\mathcal{P} = \{woman, man, meat, animal\}$

Goto Function

The **goto function** $g(q, a)$ gives the state entered from current state q by matching target char a

if edge(q, v) is labeled by a , then $g(q, a) = v$

$g(0, a) = 0$ for each a that does not label an edge out of the root

Otherwise $g(q, v) = \phi$

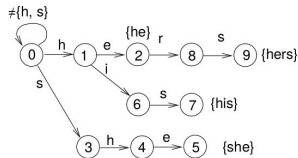


Figure: Goto function

Failure Function

The **failure function** $f(q)$ for $q \neq 0$:

$f(q)$ is the node labeled by the longest proper suffix w of $\mathcal{L}(q)$, w is a prefix of some pattern

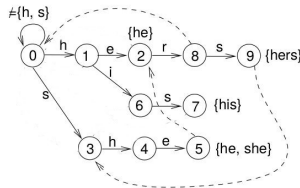


Figure: Fail transition of nodes 5, 8, 9

Output Function

The **output function** $out(q)$ gives the set patterns recognized when entering state q $f(q)$ is the node labeled by the longest proper suffix w of $\mathcal{L}(q)$, w is a prefix of some pattern

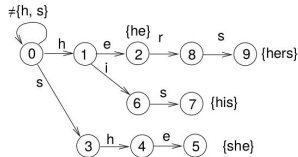


Figure: **Out** function

Aho-Corasick Automaton search of target $T[1 \cdots m]$

Example (code AC Search)

```
q := 0;  
for i := 1 to m do  
  while  $g(q, T[i]) = \phi$  do  
    q := f(q);  
  q := g(q, T[i]);  
  if  $out(q) = \phi$  then  
    print i, out(q);  
endfor;
```

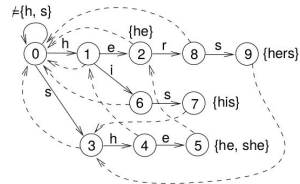


Figure: Search text "Ushers"

Complexity of Search in Aho-Corasick Automaton

Searching target $\mathcal{T}[1 \cdots m]$ with an AC automaton takes time $O(m + z)$

- Each **goto** either stay at the root or increased the depth of q by 1
- Each **fail** moves q closer to the root
- The occurrences can be reported in $z \times O(1) = O(z)$

Constructing an Aho-Corasick automaton

The Aho-corasick automaton can be constructed in two phases.

- Phase 1
- Phase 2

Constructing Phase 1

Phase 1

- 1 Construct the Trie Tree for \mathcal{P}
- 2 Complete the goto function for the root

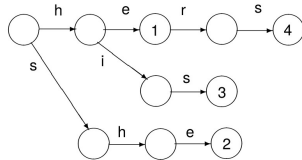


Figure: Trie of $\mathcal{P} = \{he, she, his, hers\}$

This Phase 1 of construction takes time $O(n)$

Constructing Phase 2

```

Q := emptyQueue();
for a ∈ Σ do
  if q(0, a) = q ≠ 0 then
    f(q) := 0; enqueue(q, Q);
while not isEmpty(Q) do
  r := dequeue(Q);
  for a ∈ Σ do
    if g(r, a) = u ≠ ∅ then
      enqueue(u, Q); v := f(r);
      while g(v, a) = ∅ do v := f(v); // (*)
      f(u) := g(v, a);
      out(u) := out(u) ∪ out(f(u));
  
```

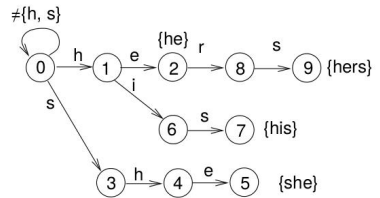


Figure: Trie of $\mathcal{P} = \{he, she, his, hers\}$

Figure: Algorithm phase 2

Constructing Phase 2

```

Q := emptyQueue();
for a ∈ Σ do
  if q(0, a) = q ≠ 0 then
    f(q) := 0; enqueue(q, Q);
while not isEmpty(Q) do
  r := dequeue(Q);
  for a ∈ Σ do
    if g(r, a) = u ≠ ∅ then
      enqueue(u, Q); v := f(r);
      while g(v, a) = ∅ do v := f(v); // (*)
      f(u) := g(v, a);
      out(u) := out(u) ∪ out(f(u));
  
```

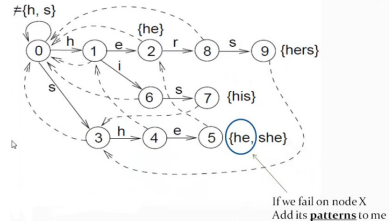


Figure: Trie of $\mathcal{P} = \{he, she, his, hers\}$

Figure: Algorithm phase 2

Constructing Phase 2

```

Q := emptyQueue();
for a ∈ Σ do
  if q(0, a) = q ≠ 0 then
    f(q) := 0; enqueue(q, Q);
while not isEmpty(Q) do
  r := dequeue(Q);
  for a ∈ Σ do
    if g(r, a) = u ≠ ∅ then
      enqueue(u, Q); v := f(r);
      while g(v, a) = ∅ do v := f(v); // (*)
      f(u) := g(v, a);
      out(u) := out(u) ∪ out(f(u));
  
```

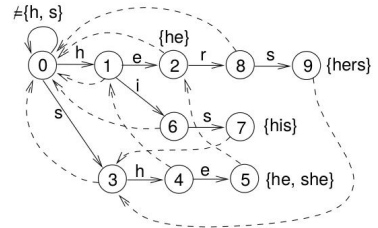


Figure: Trie of $\mathcal{P} = \{he, she, his, hers\}$

Figure: Algorithm phase 2

Constructing Phase 2

- Functions fail and output are computed for the nodes of the trie in a breadth-first order
- Then Phase 2 can be run in time $O(n)$
- Is also an $O(n)$ bound for the number that the f transitions on line that contains (*)?

Digital Forensics and Text Mining

- **Digital Forensics:** is a mathematical scheme for demonstrating the authenticity of a digital message or document.
- **Text Mining:** refers to the process of deriving high-quality information from text. Deriving patterns within the structured data, and finally evaluation and interpretation of the output.



Figure: Term map of the Journal of the American Society for Information Science and Technology

Intrusion Detecting

- **Digital Signature:** is to compare one brand with a file virus database to identify matches
- **Heuristic detection:** is the scan code looking for patterns that resemble those used in the virus.

```
franks@franks-HP:~$ clamscan ~/Descargas/test2.zip
/home/franks/Descargas/test2.zip: ClamAV-Test-Signature FOUND
----- SCAN SUMMARY -----
Known viruses: 4205080
Engine version: 0.98.7
Scanned directories: 0
Scanned files: 1
Infected files: 1
Data scanned: 0.00 MB
Data read: 0.00 MB (ratio 0.00:1)
Time: 9.091 sec (0 m 9 s)
```

Figure: Digital Signature

Zem_Bug	Zhurekov (Clam)	Zhengxi.7313 (Clam)	Zhekov-1882
Zhekov-1915	Zhekov-1915	Zhekov-2435	Zhekov-2968
Zhekov-2970	Zhekov-A	Zhekov-B	ZigZag
ZigZag-127 (B)	Zodac #1	Zohos.4516 (Clam)	Zombie.2553 (Clam)
Zombie.667 (Clam)	Zombie.PM.4506 (Clam)	Zombie.PM.4502 (Clam)	.0017.0001.000
.0017.0001.001	.0017.0001.002	.0017.0001.003	.0017.0001.004
.0021.0004.003	.0023.0004.004	.0023.0005.000	.0023.0005.001
.0023.0005.002	.0025.0006.000	.0025.0006.001	.0026.0006.000

Figure: Index of virus detected by clamAV

Detecting Plagiarism

- Plagiarism Detecting is process of finding within a work or document.
- Plagiarism is the act of copying the idea of someone else.
- There are so many algorithms have been proposed for plagiarism detection.

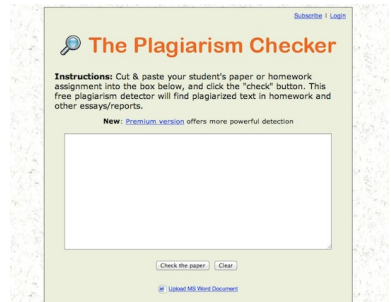


Figure: The Plagiarism Checker

Bioinformatics

Bioinformatics is the study of biological science which deals with the methods of storing, **retrieving** and analyzing biological data, such nucleic acid(DNA,RNA) and protein sequence.

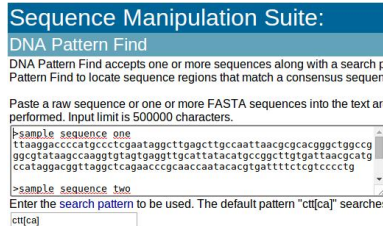


Figure: DNA Pattern Find

References



Saima Hasib, Mahak Motwani, Amit Saxema (2013)

Importance of Aho-Corasick String Matching Algorithm in Real World Applications

(IJCSIT)International journal of Computer Science and Information Technologies, Vol. 4(3), 2013, 467-469



Pekka Kilpelainen (2005)

Set Matching and Aho-Corasick Algorithm

Biosequence Algorithms, Spring 2005.

The End