COMP20007 Design of Algorithms

Input Enhancement Part 2: String Searching

Daniel Beck

Lecture 18

Semester 1, 2020

String Search - Recap

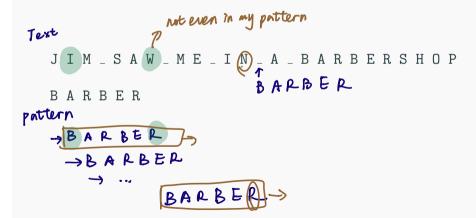
String Search - Recap

• Goal: given text of size *n*, find a string (pattern) of size *m*.

String Search - Recap

- Goal: given text of size *n*, find a string (pattern) of size *m*.
- Brute force algorithm: $O(m \times n)$.

String Search - Brute Force



• Longer shifts can be made if we know statistics about the pattern.

- Longer shifts can be made if we know statistics about the pattern.
- Same way we could sort arrays more efficiently by knowing statistics about the array (Counting Sort).

- Longer shifts can be made <u>if we know statistics about the</u> pattern.
- Same way we could sort arrays more efficiently by knowing statistics about the array (Counting Sort).
- The Horspool's algorithm use this idea to make string search faster.

- Longer shifts can be made <u>if we know statistics about the</u> pattern.
- Same way we could sort arrays more efficiently by knowing statistics about the array (Counting Sort).
- The Horspool's algorithm use this idea to make string search faster.
- Key idea: scan the text from left to right but scan the pattern from right to left.

The last character is not in the pattern.

Shift the whole pattern.

BARBER

Shift the pattern until the last occurence of the character.

The last character matches but one of the m-1 characters does not match and the last character is unique.

JIM_SAW_ME_IN_A_BARBERSHOP

SEESAW

SEESAW

Shift the whole pattern.

The last character matches but one of the m-1 characters does not match and the last character is not uniques

Shift the pattern until the last occurence of the character.

Horspool - Preprocessing

• The number of allowed shifts depends on the character type only.

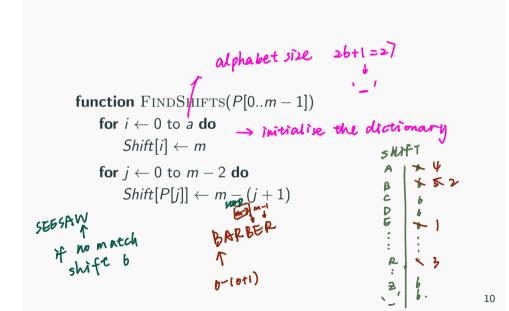
Horspool - Preprocessing

- The number of allowed shifts depends on the character type only.
- Horspool builds a <u>dictionary</u> with all characters in the alphabet and their corresponding allowed skips for a pattern.

Horspool - Preprocessing

- The number of allowed shifts depends on the character type only.
- Horspool builds a <u>dictionary</u> with all characters in the alphabet and their corresponding allowed skips for a pattern. RARBER

$\label{eq:horspool} \textbf{Horspool} \textbf{ - FindShifts}$



Horspool - Algorithm

```
function Horspool(P[0..m-1], T[0..n-1])
                       FINDSHIFTS(P) \rightarrow get the shift table
                       i \leftarrow m-1 \rightarrow look at last charater
                        while i < n do
                                            k \leftarrow 0 gend of patter mismatch while k < m and P[m-1-k] = T[i-k] do
                                                         k \leftarrow k + 1
                                               if k = m then
                                                                                                                                                                                                                                                                         return i - m + 1

    Start of the match
    ■
    Start of the match
    ■
    ■
    Start of the match
    ■
    ■
    Start of the match
    ■
    ■
    ■
    Start of the match
    Start of the match
    ■
    Start of the match
    Start of the matches
    Start 
                                               else
                                                                     i \leftarrow i + Shift[T[i]]  \triangleright Slide the pattern along
                        return -1
```

Horspool - **Properties**

- Worst-case still $O(m \times n)$.
- For random strings, it's linear and faster in practice compared to the brute force version.

Other String Search algorithms

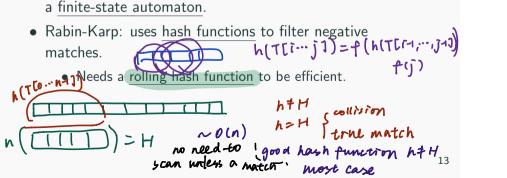
• Boyer-Moore: extends Horspool to allow shifts based on suffixes.

Other String Search algorithms

- Boyer-Moore: extends Horspool to allow shifts based on suffixes.
- Knuth-Morris-Pratt: also preprocess the patten but builds a finite-state automaton.

Other String Search algorithms

- Boyer-Moore: extends Horspool to allow shifts based on suffixes.
- Knuth-Morris-Pratt: also preprocess the patten but builds a finite-state automaton.



• String search algorithms can be sped up by input enhancement.

- String search algorithms can be sped up by input enhancement.
- Allocates extra memory to preprocess the pattern.

- String search algorithms can be sped up by input enhancement.
- Allocates extra memory to preprocess the pattern.
- Horspool uses a dictionary of shifts.

- String search algorithms can be sped up by input enhancement.
- Allocates extra memory to preprocess the pattern.
- Horspool uses a dictionary of shifts.
- The Boyer-Moore extension is one of the most used string searching algorithms, for instance in the "grep" Linx-command line tool.

- String search algorithms can be sped up by input enhancement.
- Allocates extra memory to preprocess the pattern.
- Horspool uses a dictionary of shifts.
- The Boyer-Moore extension is one of the most used string searching algorithms, for instance in the "grep" Linx command line tool.

Next week: trade memory for speed by storing intermediate solutions.