

COMP261 Lecture 17

Marcus Freen

String Search 1 of 2



Victoria
UNIVERSITY OF WELLINGTON

*Te Whare Wānanga
o te Ūpoko o te Ika a Māui*



CAPITAL CITY UNIVERSITY

String search

“Given a string S and a text T ,
look for an occurrence of S as a substring of T ”

- Which one? (first, all...)
- What do I do when I find it?
- If found, return index of first character of S in T ;
otherwise return -1 (or some other index outside of T).
- What would you expect the cost to be?

String search

- Find the string "**vtewfvtxqwfcsrdzcaj**" in this text:

qwerxcvvtewfzxcfasfedrsadfsdacfasdrtvtewqwertcsvte
wfvtxqwfcsrdzfeceeaeszxcvtsafsersdxzcvtedfaevsadv
tewfvtxqwfcszsvzxgvtasfvtcasrfvtewqtrwtravtewfxtrac
wrtrdtgfdvxvvsbdgfstqtretydfxvzccadawqeewtertgvb
vczfafsvtewfvtxqwfcszsgfsdfdxvzvzvtvsgfsgtfwt6fqwt
qwrclfxtvtewfwtqwfzvwqgtfvtqfwcxetwfazreqresdqxrdqc
fwqdxvgfewcvtwefxvtrfczrqesxqecaqrzfzvtqwxvbwyegcbe
bcwtfexvtfwxcrqxeqdcqzrwdfvtwxefvctyvtewfwefxqtfxc
qcdzrqxesrzqxrrcwqtfxtewfcwverygcviewytxvqewtcxzdc
qwfxtewfvtxqwfcsrdzcajwfcsxtqwefdvetwqfvxdtqfwvq

String search - some variations

- Just check whether it's there, returning Boolean.
- Find first/last/any occurrence of **S** in **T**.
- Find all occurrences of **S** in **T**.
 - What if occurrences overlap?
- Find occurrence(s) as a whole word/anywhere?
- Find occurrences within lines/allow occurrences to extend across line breaks?
- Assume random data? English text? Other data?

String search

- In Java, we can do this by using:
 - `T.indexOf(S);`
 - `T.lastIndexOf(S);`
 - `T.contains(S);`
- But we'd like to know if these are good choices – or if we can do better.
- Let's start with a simple algorithm, and see how we can improve upon it.

Brute force approach

- string: $S = \text{ananaba}$
- text: $T = \text{bannabanabananaban}$
- Look for S , starting at $T[0]$:
 ananaba
 $\text{bannabanabananaban}$
- Look for S , starting at $T[1]$:
 ananaba
 $\text{bannabanabananaban}$
- Look for S , starting at $T[2]$:
 ananaba
 $\text{bannabanabananaban}$
- Etc. till found, or none left.

Brute force algorithm

- Basic idea:
Look for S in T,
starting at positions T[0], T[1],
- What is last position in T we need to consider?

```
for k ← 0 to T.length() - S.length()  
    if T.substring(k, S.length()).equals(S) then return k  
return -1
```

- What is the cost?

Brute force algorithm

Can we improve?

```
for k ← 0 to T.length() - S.length()
  if T.substring(k, S.length()).equals(S) then return k
return -1
```

- First, some very simple “improvements”:
 - Don’t call **length** methods in the loop.
Avoid cost of method call (compiler may inline it).
 - Don’t call **substring** method in the loop.
Don’t need to copy the substring to a new string to compare with S.

Brute force algorithm

Assigning $m \leftarrow S.length()$ and $n \leftarrow T.length()$ first:

- **for** $k \leftarrow 0$ to $n-m$
 $found \leftarrow true$
 for $i \leftarrow 0$ to $m-1$
 if $S[i] \neq T[k+i]$ **then** $found \leftarrow false,$
 break
 if $found$ **then** **return** k
 return -1
- Okay what is the cost?

Brute force algorithm: cost

- $S = s_0 \quad s_1 \quad s_2 \quad s_3 \quad \dots$ (m of these)
 $T = t_0 \quad t_1 \quad t_2 \quad t_3 \quad t_4 \quad t_5 \quad t_6 \quad t_7 \quad \dots$ (n of these)
- What is best/worst/expected cost?
- What if text is random? English?
- What case gives best/worst cost (for any m and n)?
 - How many positions in T need to be considered?
 - How many characters need to be considered at each position?

Brute force algorithm: best cost

- $S = s_0 \quad s_1 \quad s_2 \quad s_3 \quad \dots$ (m of these)
 $T = t_0 \quad t_1 \quad t_2 \quad t_3 \quad t_4 \quad t_5 \quad t_6 \quad t_7 \quad \dots$ (n of these)
- Suppose s_0 doesn't occur in T.
 - s_0 will be compared to t_0, t_1, \dots
 - So cost will be?
- Suppose S is a prefix of T.
 - Will compare s_0 with t_0 , s_1 with t_1, \dots
 - So cost is?

Brute force algorithm: worst cost

- $S = s_0 \quad s_1 \quad s_2 \quad s_3 \quad \dots$ (m of these)
 $T = t_0 \quad t_1 \quad t_2 \quad t_3 \quad t_4 \quad t_5 \quad t_6 \quad t_7 \quad \dots$ (n of these)

What case will force the algorithm to do the most comparisons?

- *Hint 1:* Want S not in T, so it tries the maximum number of positions.
- *Hint 2:* At each position, want algorithm to do the **most possible comparisons before failing**.
 $\rightarrow \rightarrow$ Fail on the last character in S!

What inputs would do this?

- what about

$S = \text{aaaaab}$

$T = \text{aaaaaaaaaaaaaaaaaaaaaaaaa}$

What is the cost?

Would this ever happen with English text?

What sort of data then?

String search: can we do better?

- ideally, we'd have an algorithm that never needs to re-trace its steps in the long string. Can we check each letter just once?

fail

- abcdh????????????????????

abcdefg

- having got to a fail point, where should we check next?
 - jump ahead, and re-start at the fail point?
- this could speed up search a lot!
 - is it “safe”?

String search: can we do better?

fail

- `anan#????????????????????`
`anangfg`
 - having got to a fail point, where should we check next?
 - jump ahead, and re-start at... where?
- `anan#????????????????????`
`ananafg`
 - what about now?
 - unsafe to jump straight to the fail point
- Key idea of KMP algorithm: Use characters in partial match to determine where to start next match attempt.

String search: Example

- T = abc_abc_dab_abc_dabcdabde
S = abc_dabd

- T = abc_abc_dab_abc_dabcdabde
S = a_bcdabd

- T = abc_abc_dab_abc_dabcdabde
S = abc_dab_d

- T = abc_abc_dab_abc_dabcdabde
S = abc_dabd

continued
next slide

String search: Example

- T = abc_abc_dab_abc_dabcdabde
S = abcdabd

- T = abc_abc_dab_abc_dabcdabde
S = abcdabd

- T = abc_abc_dab_abc_dabcdabde
S = abcdabd

- T = abc_abc_dab_abc_dabcdabde
S = abcdabd

Knuth-Morris-Pratt (KMP) algorithm

- The “Knuth” here is Donald Knuth –
https://en.wikipedia.org/wiki/Donald_Knuth

After a mismatch, advance to the earliest place where search string could possibly match.

- never has to re-check a character

How far can we advance safely?

- Use a table based on the search string.
- Let $M[0..m-1]$ be a table showing how far to back up the search if a prefix of S has been matched.

String search

- Simple search

- Slide the window by 1

- $t = t + 1;$

ab~~cd~~mndsjhhhsjgrjgslagf
 ab~~bdd~~efg

- KMP

- Slide the window faster

- $t = t + s - M[s]$

ananfd~~f~~joi~~j~~toi~~i~~inkjjkjgghfj
 anan~~g~~bgba

- Never re-check the matched characters

- If there a “suffix == prefix”?

- No, skip these characters

- » $M[s] = 0$

- Yes, reuse, no need to recheck these characters

- » $M[s]$ is the length of the “reusable” suffix

Knuth Morris Pratt

input: string $S[0 \dots m-1]$, text $T[0 \dots n-1]$, partial match table $M[0 \dots m-1]$

output: the position in T at which S is found, or -1 if not present

variables: $k \leftarrow 0$ *start of current match in T*
 $i \leftarrow 0$ *position of current character in S*

while $k + i < n$

if $S[i] = T[k + i]$ **then** *// match*
 $i \leftarrow i + 1$
 if $i = m$ **then return** k *// found S*
 else if $M[i] = -1$ **then** *// mismatch, no self overlap*
 $k \leftarrow k + i + 1, i \leftarrow 0$
 else *// mismatch, with self overlap*
 $k \leftarrow k + i - M[i]$ *// match position jumps forward*
 $i \leftarrow M[i]$

return -1 *// failed to find S*