

Knuth-Morris-Pratt deterministic string matching

Text: H O C U S P O C U S A B R A C A D A B R A
 Pattern: A B R A C A D A B R A

A B R A C A D A B R A
 A B R A C A D A B R A
 A B R A C A D A B R A
 A B R A C A D A B R A

$T[1..n]$
 $P[1..m]$

Once we've matched $T[i]$ to some $P[j]$, we shouldn't need to compare $T[i]$ to anything else

The next shift is smallest s such that $T[s..i-1]$ are a ^{shorter} prefix of the pattern.

But what is this?

for now, magic

P[i]	A	B	R	A	C	A	D	A	B	R	A
fail[i]	0	1	1	1	2	1	2	1	2	3	4

Failure function for the string ABRACADABRA

```

KNUTHMORRISPRATT(T[1..n], P[1..m]):
    j ← 1
    for i ← 1 to n
        while j > 0 and T[i] ≠ P[j]
            j ← fail[j]
        if j = m    <<Found it!>>
            return i - m + 1
        j ← j + 1
    return NONE
    
```

Running time: Crudely $O(nm)$


increment $i \leq n$
 # increment $j \leq n$
 # decrease $j \leq \text{\#increment } j \leq n$ | $O(n)$ time

Prefix: 

proper = not the whole string

Suffix: 

Border = both a ^{nonempty} proper prefix and a suffix of $P[1..i-1]$

ABACABADABACABA


Borders are nested!

Border of a string w is either
longest border of w
or a border of a border of w

$P[1..fail[j]-1]$ is the longest proper prefix of $P[1..j-1]$
that is also a suffix of $T[1..i-1]$.

= longest border of $P[1..j-1]$

$P[1..fail[j]-1]$ is the longest proper prefix of $P[1..j-1]$
that is also a suffix of $P[1..j-1]$.

$P[i]$	A	B	R	A	C	A	D	A	B	R	A
$fail[i]$	0	1	1	1	2	1	2	1	2	3	4

Failure function for the string **ABRACADABRA**

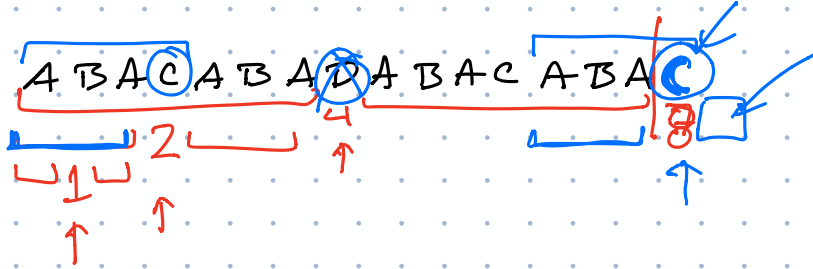
Naive $O(m^3)$
DP $O(m^2)$

COMPUTEFAILURE($P[1..m]$):

```

j ← 0
for i ← 1 to m
  fail[i] ← j (*)
  while j > 0 and P[i] ≠ P[j]
    j ← fail[j]
  j ← j + 1
  
```

→ $O(m)$ time



$$fail^c(i-1) = \begin{cases} fail(fail^{c-1}(i-1)) & \text{if } c > 0 \\ i-1 & \text{if } c = 0 \end{cases}$$

$$fail[i] = \begin{cases} 0 & \text{if } i = 0, \\ \max_{c \geq 1} \{ fail^c[i-1] + 1 \mid P[i-1] = P[fail^c[i-1]] \} & \text{otherwise.} \end{cases}$$

KNUTHMORRISPRATT($T[1..n], P[1..m]$):

```
 $j \leftarrow 1$   
for  $i \leftarrow 1$  to  $n$   
  while  $j > 0$  and  $T[i] \neq P[j]$   
     $j \leftarrow fail[j]$   
  if  $j = m$     «Found it!»  
    return  $i - m + 1$   
   $j \leftarrow j + 1$   
return NONE
```

COMPUTEFAILURE($P[1..m]$):

```
 $j \leftarrow 0$   
for  $i \leftarrow 1$  to  $m$   
   $fail[i] \leftarrow j$     (*)  
  while  $j > 0$  and  $P[i] \neq P[j]$   
     $j \leftarrow fail[j]$   
   $j \leftarrow j + 1$ 
```

Handwritten annotations:
- A red arrow points from the first 'A' (index 1) to the first 'A' (index 10).
- A blue arrow points from the 'A' at index 10 to the 'A' at index 15.
- The string is: A B A C A B A D A B A C A B A

Fail: X 0 1 1 2 1 2 3 4 1 2 3 4 5 6 7