# STRING MATCHING

DATA STRUCTURES AND ALGORITHMS (INFO-F413)
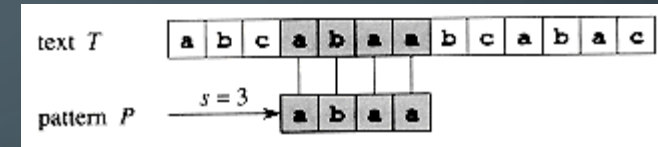
CHARLOTTE NACHTEGAEL

# YOUR MISSION, IF YOU ACCEPT IT

Your favorite agent 007 called you with an important task : find the indexes of the secret word in the text to use the numbers to defuse an destructive bomb menacing the world !

The fate of the world rests on your shoulders….
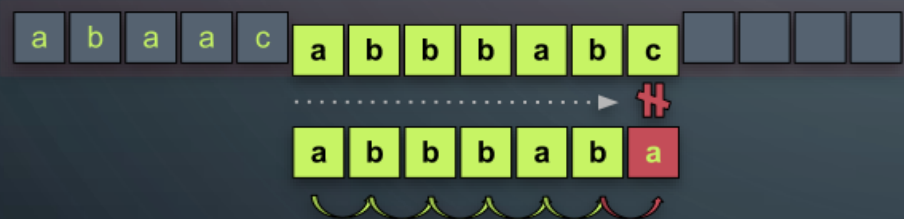
# FORMALLY SPEAKING

- Text is an array T[1…n]

- Pattern/Secret word is an array P[1…m]



-  P occurs with shift s in text T (or, equivalently, that pattern P occurs beginning at position s + 1 in text T) if $0 \leq s \leq n - m$ and T[s + 1 . . s + m] = P[1 . . m] (that is, if T[s + i] = P[i], for $1 \leq i \leq m$).

- The string-matching problem is the problem of finding all valid shifts with which a given pattern P occurs in a given text T.

# FIRST IDEA : LET'S DO THIS ONE-BY-ONE
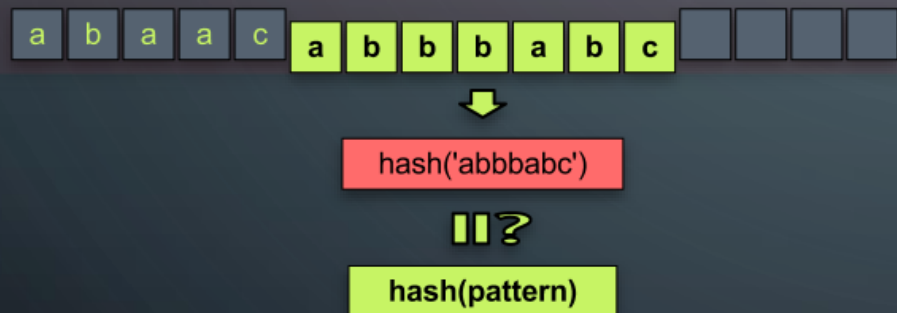
**Brute-force matching**

| a | b | a | a | c | a | b | b | b | a | b | c |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| a | b | b | b | a | b | a |
|---|---|---|---|---|---|---|

- Checks the condition $P[1 .. m] = T[s + 1 .. s + m]$ for each of the n - m + 1 possible values of s

- $\Theta((n-m+1)m)$ !

➔ Too long, the world would explose before you finish !

# SECOND IDEA : LET'S DO HASHING !

**Rabin-Karp**

| a | b | a | a | c | a | b | b | b | a | b | c | | | | |

↓

hash('abbbabc')

‖?

hash(pattern)

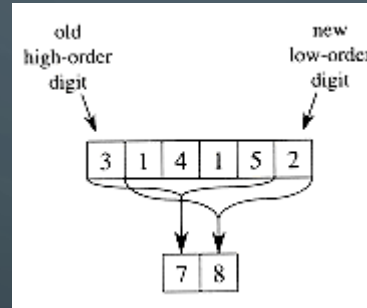- Introduction of a powerful technique : the *hash function*

- Calculate a number H for the pattern and each substring of T of length m

➔ Better than one-by-one as you compare two integers and not two strings character by character !

# ROLLING HASH

With d, the size of alphabet and q, a large prime number

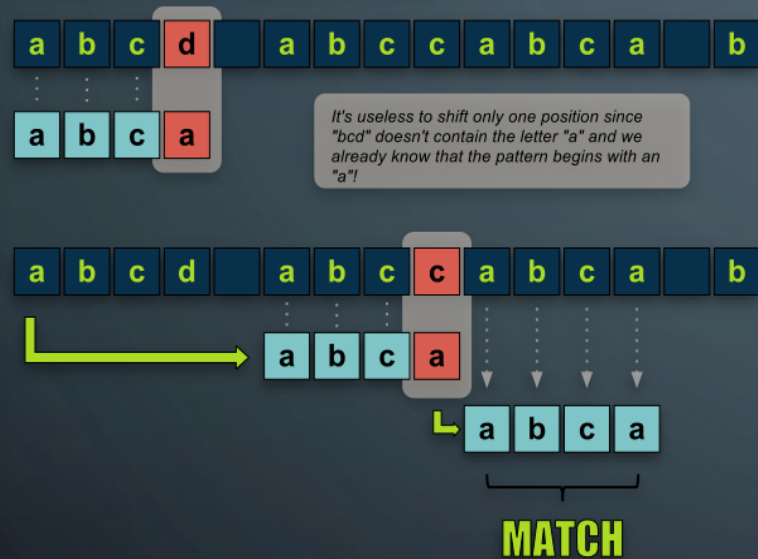$$H(T) = (T[0]d^{m-1} + T[1]d^{m-2} + \cdots T[m-1]d^0) \ \% \ q$$



$$H(T_{s+1}) = (d(T_s - T[s+1]d^{m-1}\%q) + T[s+m+1]) \ \% \ q$$

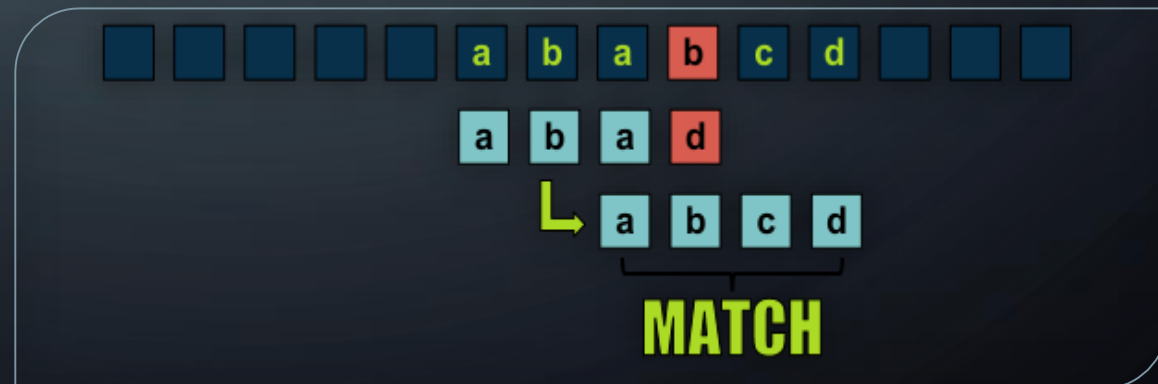➔ Preprocessing O(m) + Looking at the text O(n) * rolling hash to look O(1)

➔ O(m + n)

# THIRD IDEA : LET'S JUMP FROM ONE INTERESTING PLACE TO ANOTHER



Morris-Pratt String Searching

- You look for the prefix of your pattern within your pattern

- Shift directly to the next beginning of pattern

- If no next prefix, shift further

# BUILDING THE PREFIX TABLE



COMPUTE-PREFIX-FUNCTION(P)

m ← length[P]

π[1] ← 0

k ← 0

for q ← 2 to m

    do while k > 0 and P[k + 1] ≠ P[q]

        do k ← π[k]

    if P[k + 1] = P[q]

        then k ← k + 1

    π[q] ← k

return π

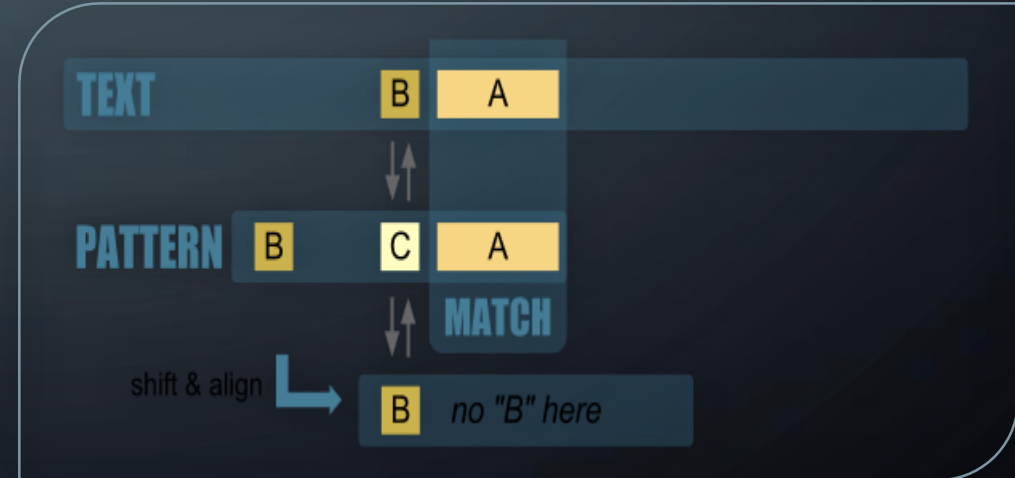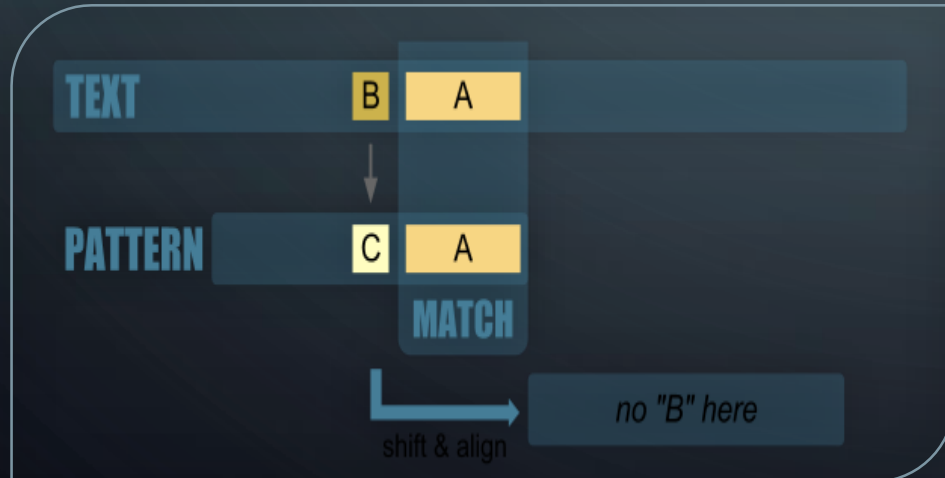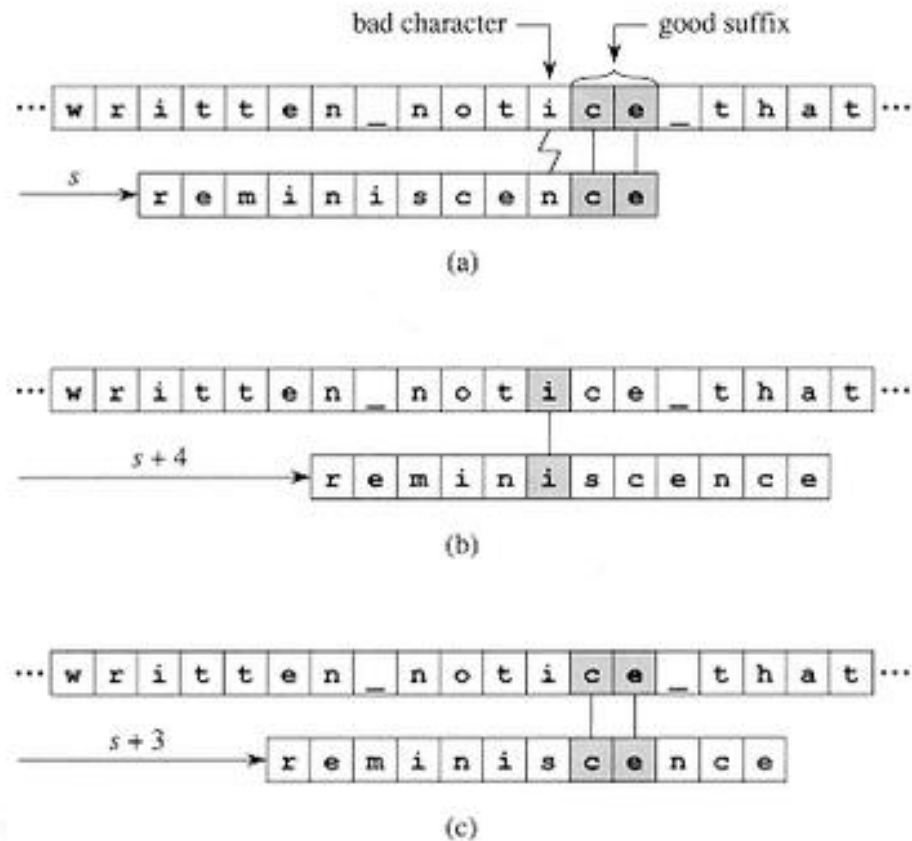➔ Preprocessing O(m) + Looking in the text O(n) ➔ O(n + m)

# FOURTH IDEA : LET'S BEGIN WITH THE END



- The letters of the pattern are compared from right to left !

- Two kinds of shifts :
  - Good suffix shift
  - Bad character shift
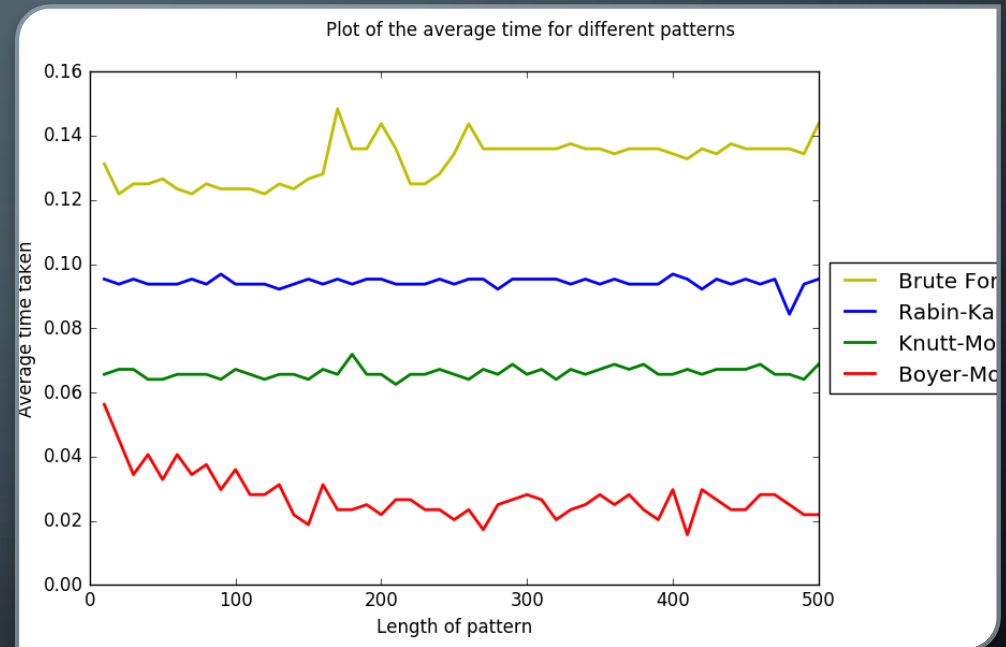
# THE BAD CHARACTER & THE GOOD SUFFIX



- Bad character : shift to the last occurrence of the bad character in the pattern

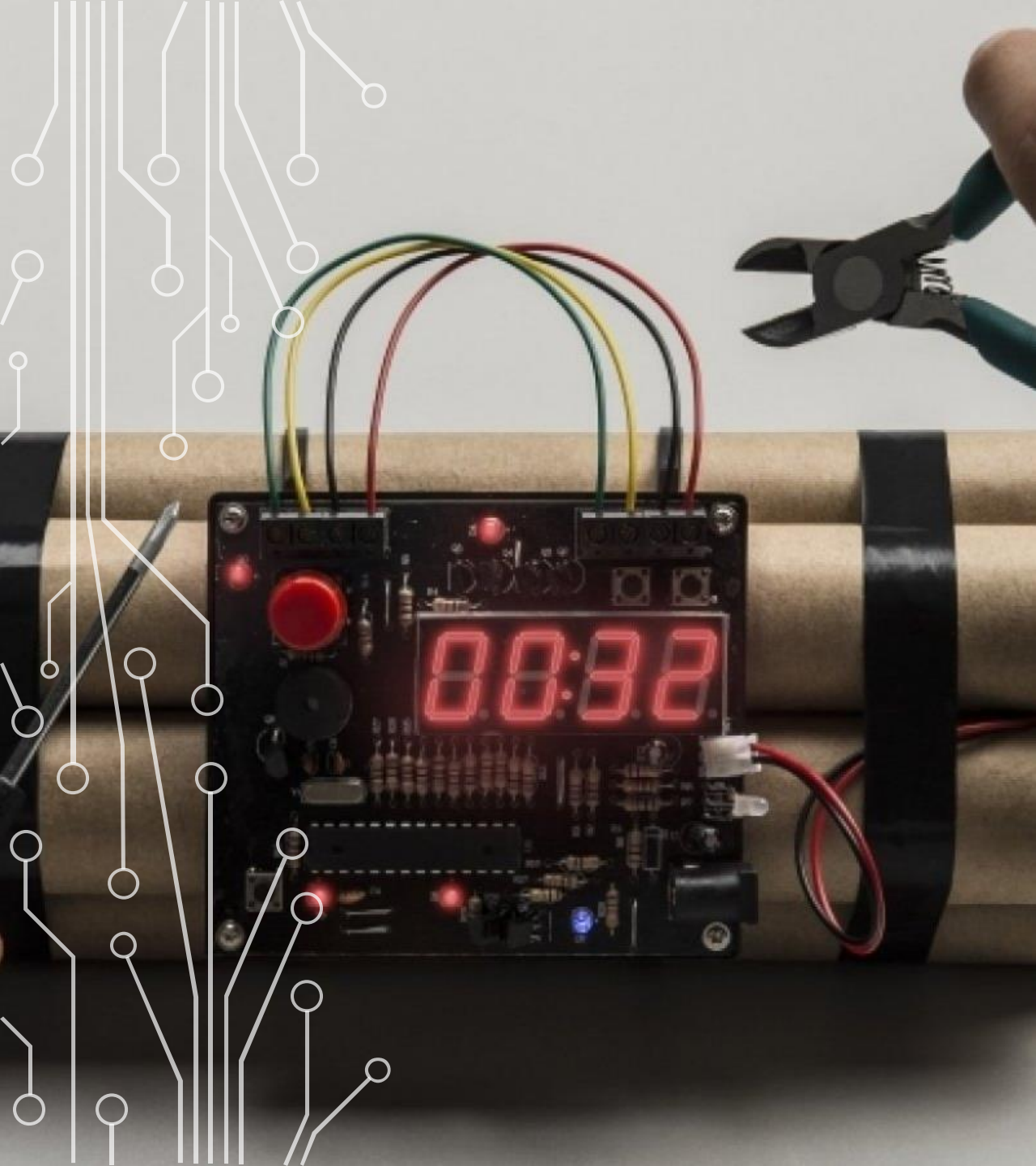- Good suffix : shift to the next pattern found with the suffix

➔ Choose the biggest shift !

➔ Preprocessing O(m) and O(m + size alphabet) + O(m) time to validate each shift s
➔ $O((n-m+1)m)$ but in practise : $\Omega(n/m)$ and $O(nm)$

# WHAT'S THE BEST IDEA ?

- How many words must we find ?

- What is the size of your alphabet ?

- What is the size of your pattern ?

- What is the size of your text ?



Plot of the average time for different patterns

Thank you for listening !

Now, it's your turn to save the world !

But before that, any questions ?