

# ALGORITHMES DE RECHERCHE DE MOTIFS (SUR LES MOTS)

## INTRODUCTION

### Problématique

La problématique est la suivante : on a une chaîne de caractère  $C$  : 0110011000101010110... et un motif  $M = 010$ . On recherche les occurrences de  $M$  dans la chaîne  $C$  :

$$C = 0110011000\mathbf{0101010110} \dots$$

### Notations

On appellera  $\Sigma$  dans la suite un alphabet fini (ensemble fini de caractères).

On appellera  $\Sigma^k$  l'ensemble des mots de longueur  $k$  sur  $\Sigma$ .

On appellera  $\Sigma^*$  :

$$\Sigma^* = \bigcup_{k \geq 0} \Sigma^k.$$

Soit  $w \in \Sigma^*$ , on appellera  $\text{len}(w)$  la longueur de  $w$ . Soit  $i \leq \text{len}(w)$ , alors  $w[i]$  sera le  $i$ -ième caractère de  $w$ .  $w[i : j] = w[i]w[i+1] \dots w[j-1]$ .

## 1. ALGORITHME NAÏF

### 1.1. Présentation

On se donne une chaîne  $C$  et un motif  $M$ .

On teste l'égalité de  $M$  avec chaque sous-mot  $C[i : i + \text{len}(M)]$ .

**DÉFINITION 1.1.0.1.** —

Lorsque  $M = C[i : i + \text{len}(M)]$  on dit que  $M$  apparaît avec décalage  $i$  dans la chaîne  $C$ .

```

1 function rechercheMotifNaif(C,M)
2     n = len(C)
3     m = len(M)
4     S = []
5     for i from 0 to n - m do
6         if M[0:m] = C[i : i + m] then
7             S = S + [i]
8     return S

```

## 1.2. Coût de l'algorithme naïf

On cherche à compter le nombre de comparaisons de caractères,  $C(n)$ .

Chaque test  $M = C[i : i + m]$  a un coût au plus de  $m$  comparaisons. De plus, il y a  $n - m + 1$  tests d'égalité donc :

$$C(n) = O(m(n - m + 1)).$$

De plus la borne est atteinte pour :  $M = a^m$  avec  $a \in \Sigma$  et  $C = a^n$  avec  $n \geq m$ .

## 2. ALGORITHME DE KARP-RABIN

R. KARP, M. RABIN : 1987.

C'est un algorithme qui se base sur le naïf mais enrichi par une technique de *hachage*.

Dans la suite,  $\Sigma = \{0, 1, \dots, 9\}$ . On peut donc représenter des mots par des entiers. Le mot  $M \in \Sigma^*$  peut être représenté par l'entier :

$$m = \sum_{k=0}^{\text{len}(M)} M[k]10^{\text{len}(M)-k}.$$

*Coût.* — Le coût de la représentation est un  $O(m)$  <sup>(1§)</sup>.

*Exemple.* — Par exemple, si  $C = '23451'$  et  $\text{len}(M) = 4$ . On a  $C[0 : 4] = C[0]C[1]C[2]C[3] = '2345'$  et sa représentation est  $c_0 = 2345$ . On veut obtenir l'entier  $c_1$  représentant  $C[1 : 5]$  à partir de  $c_0$  et  $C[4] : c_1 = 10(c_0 - 10^3 \times 2) + 1 = 10(c_0 - 10^3 C[0]) + C[4]$ .

Ainsi, on peut calculer  $10^{m-1}$  en  $O(\log m)$  opérations (mais  $O(m)$  suffit). On obtient  $c_{i+1}$  à partir de  $c_i$  par :

$$c_{i+1} = d \times (c_0 - d^{m-1} \times C[i]) + C[i + m],$$

où  $d = |\Sigma|$ . Donc on obtient  $c_{i+1}$  à partir de  $c_i$  en un nombre constant d'opérations.

---

<sup>1§</sup>. Voir : schéma de HORNER.

*Conclusion.* — L'ensemble des  $c_0, c_1, \dots, c_{n-m+1}$  peut être obtenu en  $O(m+n-m+1) = O(n)$ .