

Algorithmes sur les mots (séquences)

Introduction

Algorithmes sur les mots (textes, séquences, chaînes de caractères)

Nombreuses applications :

- bases de données bibliographiques
- bioinformatique (séquences de biomolécules)
- sécurité informatique

Livres :

- M.Crochemore, C.Hancart, T.Lecroq, Algorithmique du texte, Vuibert 2001

<http://igm.univ-mlv.fr/~mac/CHL/CHL.html>

Recherche d'un motif dans un texte

$T[1..n]$ texte

$P[1..m]$ motif

Tâche : localiser toutes les occurrences de P dans T
(variantes : vérifier si P apparaît dans T , compte le nombre d'occurrences)

Algorithme naïf : $O(n \cdot m)$

$T = \text{aaaaaa...aa} = a^n$

$P = \text{aa...ab} = a^{m-1}b$

Recherche d'un motif dans un texte

$T = a b a b a a a b a b a a b a \dots$

$P = a b a a b$

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
P = a b a a b

705

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
P = a b a a b

706

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
P = a b a a b

707

Recherche d'un motif dans un texte

UMLV ©



T = a b a **b** a a a b a b a a b a ...
P = a b a a b

708

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
P = a b a a b
 a b a a b

709

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
P = a b a a b
 a b a a b

710

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
P = a b a a b
 a b a a b

711

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a **a** b a b a a b a ...
P = a b a a b
 a b a a b

712

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
P = a b a a b
 a b a a b
 a b a a b

713

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
P = a b a a b
 a b a a b
 a b a a b

714

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
P = a b a a b
 a b a a b
 a b a a b

715

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a **b** a a b a ...
P = a b a a b
 a b a a b
 a b a a b

716

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
 P = a b a a b
 a b a a b
 a b a a b
 a b a a b

717

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
 P = a b a a b
 a b a a b
 a b a a b
 a b a a b

718

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b a b a a b a ...
 P = a b a a b
 a b a a b
 a b a a b
 a b a a b

719

Recherche d'un motif dans un texte

UMLV ©



T = a b a b a a a b **a b a a b** a ...
 P = a b a a b
 a b a a b
 a b a a b
 a b a a b

720

Recherche d'un motif dans un texte

UMLV ©

↓

T = a b a b a a a b a b a a b a ...
 P = a b a a b
 a b a a b
 a b a a b
 a b a a b
 a b a a b

721

Recherche d'un motif dans un texte

UMLV ©

↓

T = a b a b a a a b a b a a b a ...
 P = a b a a b
 a b a a b
 a b a a b
 a b a a b
 a b a a b

722

Recherche d'un motif dans un texte

UMLV ©

Quelques observations :

- à chaque pas on progresse d'une lettre dans le texte
- l'état de la recherche est défini par les positions dans le texte et dans le motif
- le décalage du motif est défini en fonction de la *position d'échec* dans le motif et la lettre du texte qui a provoqué l'échec

T = ... * * * * * * * * ...
 = = = ≠
 P = * * * * *
 * * * * *

723

Recherche d'un motif dans un texte

UMLV ©

Quelques observations :

- à chaque pas on progresse d'une lettre dans le texte
- l'état de la recherche est défini par les positions dans le texte et dans le motif
- le décalage du motif est défini en fonction de la *position d'échec* dans le motif et la lettre du texte qui a provoqué l'échec

T = ... * * * * * * * * ...
 = = = ≠
 P = * * * * *
 * * * * *

⇒ Automate fini !

724

Recherche d'un motif dans un texte

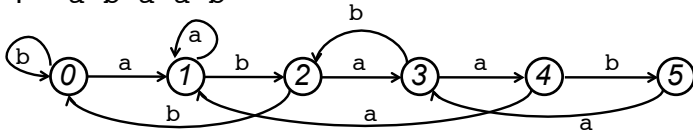
UMLV ©

$P \Rightarrow$ automate fini

état \equiv préfixe de P (position) $q \equiv P_q = P[1..q]$

$$\sigma(q,a) = \begin{cases} q+1, & \text{si } a=P[q+1] \\ \max\{i \mid P_i \text{ est un suffixe de } P_q a\} & \text{sinon} \end{cases}$$

$P = a \ b \ a \ a \ b$



Invariant: P_q est le préfixe maximal de P qui est un suffixe de la partie lue de T

725

Recherche d'un motif dans un texte

UMLV ©

Algorithme qui en résulte :

1. *Prétraitement* : calculer l'automate $O(m \cdot |A|)$
2. *Recherche* : faire marcher l'automate sur le texte $O(n)$

But : obtenir un algorithme ne dépendant pas de la taille d'alphabet

726

Algorithme de Knuth-Morris-Pratt

UMLV ©

Algorithme de Knuth-Morris-Pratt

UMLV ©

$T = \dots a \ b \ a \ b \ * \ * \ * \ * \ \dots$

$= = = \neq$

$P = \quad \quad b \ a \ b \ a \ b$

Quels sont les décalages possibles du motif ?

Fonction d'échec (de suppléance) :

$f(q) = \max\{k \mid k < q \text{ et } P_k = P[1..k] \text{ est un suffixe de } P_q\}$

$P = a \ b \ a \ b \ a \ c \ a$

$q \quad 0 \ 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7$

$f(q) \ -1 \ 0 \ 0 \ 1 \ 2 \ 3 \ 0 \ 1$

727

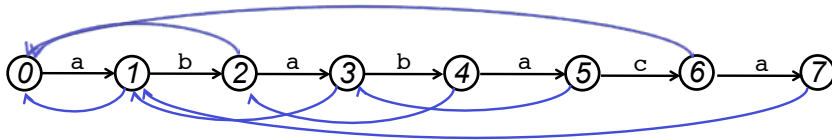
728

Algorithme de Knuth-Morris-Pratt

UMLV ©

P = a b a b a c a
 q 0 1 2 3 4 5 6 7
 f(q) -1 0 0 1 2 3 0 1

Automate de Knuth-Morris-Pratt



729

Algorithme de Knuth-Morris-Pratt

UMLV ©

Une fois la fonction f calculée ...

```

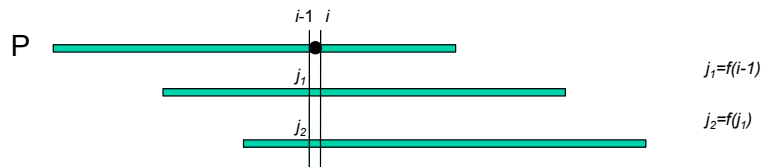
KMP(T[1..n], f)
  j=0 /* pointeur dans P */
  for i=1 to n do
    while j ≥ 0 and P[j+1] ≠ T[i] do
      j=f(j) endwhile
    j=j+1
    if j==m then
      output(occurrence de P à la position (i-m))
      j=f(j) endif
  endfor
    
```

730

Algorithme de Knuth-Morris-Pratt

UMLV ©

Comment calculer la fonction d'échec



Décaler le motif par rapport à lui-même jusqu'à ce que $P[j_q+1]=P[i]$

731

Algorithme de Knuth-Morris-Pratt

UMLV ©

Calcul de la fonction d'échec

```

FE(P[1..m])
  f[0]=-1
  f[1]=0
  k=0
  for j=2 to m do
    while k ≥ 0 and P[k+1] ≠ P[j] do
      k=f(k) endwhile
    k=k+1
    f(j)=k
  endfor
    
```

732

Algorithme de Knuth-Morris-Pratt

UMLV ©

Version optimisée (KMP vs MP)

```
T = ... a b a * * * * * ...
      = = ≠
P =   a b a b a c a
      a b a b a c a   ← décalage inutile
```

$h(3)=0$ h fonction d'échec optimisée

q	0	1	2	3	4	5	6	7
$f(q)$	-1	0	0	1	2	3	0	1
$h(q)$	-1	0	0	0	0	3	0	1

733

Algorithme de Knuth-Morris-Pratt

UMLV ©

Calcul de la fonction d'échec *optimisée* h

```
FE(P[1..m])
  h[0]=-1
  h[1]=0
  k=0
  for j=2 to m do
    while k≥0 and P[k+1]≠P[j] do
      k=h(k) endwhile
    f(j)=k → if P[j] ≠P[k] then h(j)=h(k) else h(j)=k
  endfor
```

734

Algorithme de Knuth-Morris-Pratt

UMLV ©

Différence avec l'approche automate : on peut « piétiner » sur une position de texte

au plus $\log_{\phi} m$ décalages sur une position, où $\phi=(1+\sqrt{5})/2$ nombre d'or

Complexité amortie

$O(n)$ pour la recherche (KMP)

$O(m)$ pour le pré-traitement (FE)

Historique : Morris-Pratt (1969), Knuth (1970)

735

Algorithme de Aho-Corasick

UMLV ©

736

Algorithme de Aho-Corasick

UMLV ©

Les idées de KMP peuvent se généraliser pour la *recherche simultanée de plusieurs motifs* ⇒ Algorithme de Aho-Corasick (1974)

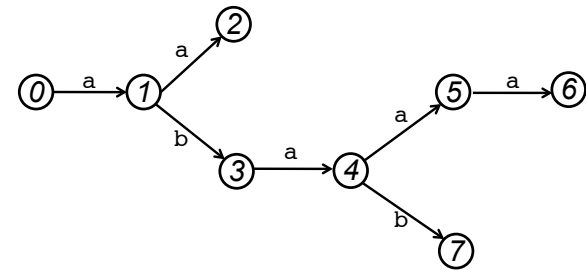
737

Algorithme de Aho-Corasick

UMLV ©

$S = \{aa, abaaa, abab\}$

arbre digital (trie)

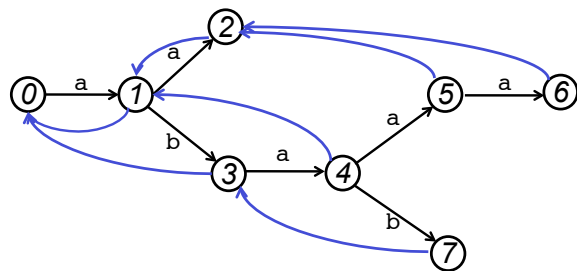


Algorithme de Aho-Corasick

UMLV ©

$S = \{aa, abaaa, abab\}$

arbre digital (trie) + fonction d'échec

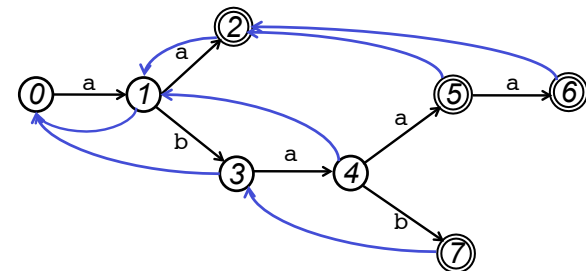


Algorithme de Aho-Corasick

UMLV ©

$S = \{aa, abaaa, abab\}$

arbre digital (trie) + fonction d'échec + états finaux

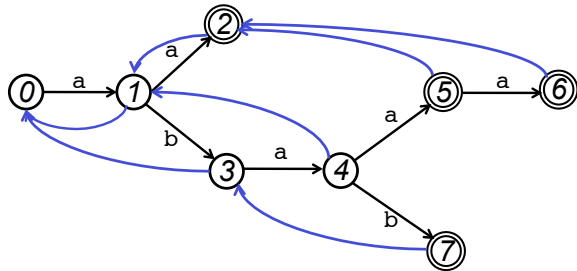


Algorithme de Aho-Corasick

UMLV ©

$S = \{aa, abaaa, abab\}$

arbre digital (trie) + fonction d'échec + états finaux =
Automate de Aho-Corasick (AC)



L'automate de AC peut être construit en temps $O(m)$ où m est la taille **totale** des mots de S

Algorithme de Karp-Rabin (1987)

UMLV ©

742

Algorithme de Karp-Rabin

UMLV ©

- soit $A = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$
- chaque mot $P[1..m]$ peut être encodé par un nombre:
 $p = P[1] \cdot 10^{m-1} + P[2] \cdot 10^{m-2} + \dots + P[m-1] \cdot 10 + P[m]$
- p peut être calculé en temps $O(m)$ avec le schéma de Horner :
 $p = P[m] + 10 \cdot (P[m-1] + 10 \cdot (P[m-2] + \dots + 10 \cdot (P[2] + 10 \cdot P[1]) \dots))$
- idée générale :
 - comparer p successivement avec les codages t_i de $T[i..i+m-1]$, pour $i = 1..n-m+1$
 - le passage de t_i à t_{i+1} se fait en temps constant :
 $t_{i+1} = 10 \cdot (t_i - 10^{m-1} \cdot T[i]) + T[i+m]$ (en supposant que 10^{m-1} est pré-calculé)

743

Algorithme de Karp-Rabin

UMLV ©

- soit $A = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$
- chaque mot $P[1..m]$ peut être encodé par un nombre:
 $p = P[1] \cdot 10^{m-1} + P[2] \cdot 10^{m-2} + \dots + P[m-1] \cdot 10 + P[m]$
- p peut être calculé en temps $O(m)$ avec le schéma de Horner :
 $p = P[m] + 10 \cdot (P[m-1] + 10 \cdot (P[m-2] + \dots + 10 \cdot (P[2] + 10 \cdot P[1]) \dots))$
- idée générale :
 - comparer p successivement avec les codages t_i de $T[i..i+m-1]$, pour $i = 1..n-m+1$
 - le passage de t_i à t_{i+1} se fait en temps constant :
 $t_{i+1} = 10 \cdot (t_i - 10^{m-1} \cdot T[i]) + T[i+m]$ (en supposant que 10^{m-1} est pré-calculé)

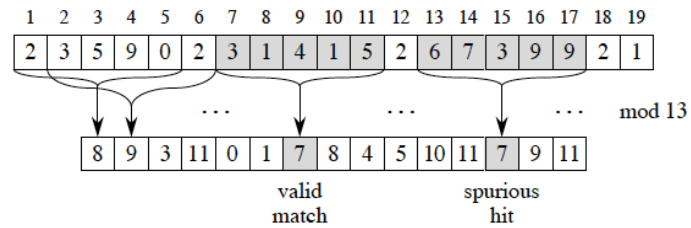
- **Problème** : les nombres manipulés peuvent être très grands
- **Parade** : manipuler les nombres modulo q

744

Algorithme de Karp-Rabin

UMLV ©

- on a alors
 $t_{i+1} = 10 \cdot (t_i - h \cdot T[i]) + T[i+m]$ modulo q , où $h = 10^{m-1}$ modulo q
- **attention, on peut avoir des « faux positifs » !**

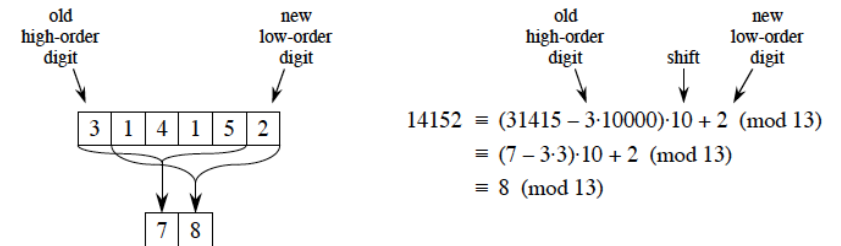


745

Algorithme de Karp-Rabin

UMLV ©

- passage de t_i à t_{i+1} (illustration) :



- dès qu'on trouve une fenêtre candidate ($t_i = p$), on vérifie en comparant P et T[i..i+m-1] caractère par caractère

746

Algorithme de Karp-Rabin

UMLV ©

- l'algorithme nécessite un pré-traitement en $O(m)$ et marche en temps $O(mn)$ dans le pire des cas
- en moyenne il marche en temps $O(n + m \cdot (v+n/q))$ où v est le nombre d'occurrences du motif
- très facile à implémenter, plus efficace que l'algorithme naïf

747