

Exemples d'applications de l'algorithme de KMP

801

Application 1 : périodicité

Période d'un mot $T[1..n]$: entier p tel que $T[i]=T[i+p]$ pour tout $i \in [1..n-p]$

Exemples :

abababa : périodes 2,4,6

bbabb : périodes 3,4

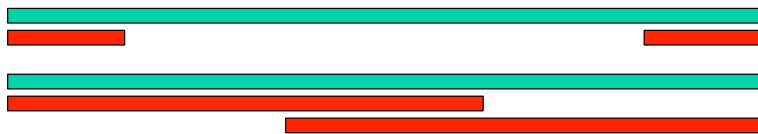
bbabaa : pas de période

Problème : étant donné un mot T , trouver sa période *minimale*

802

Application 1 : périodicité

Un **bord** d'un mot T : facteur qui apparaît en préfixe et en suffixe



Lemme : p est une période de $T[1..n]$ ssi T possède un bord de longueur $n-p$

La fonction d'échec (non-optimisée) de l'algorithme KMP calcule le bord maximal pour chaque préfixe $T[1..i]$

803

Application 1 : périodicité

⇒ calcul de la période minimale en temps $O(n)$ à l'aide de l'algorithme de calcul de la fonction d'échec

	T = 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	période minimale = 7-1=6

	T = a b a a b a a b									
q	0	1	2	3	4	5	6	7	8	
f(q)	-1	0	0	1	1	2	3	4	5	période minimale = 8-5=3

804

Arbre des suffixes

UMLV ©

Trie des suffixes (arbre des suffixes non-compacté): trie stockant tous les suffixes du texte

T = a b a b b b

suffixes de T :

```

a b a b b b
  b a b b b
    a b b b
      b b b
        b b
          b
  
```

809

Arbre des suffixes

UMLV ©

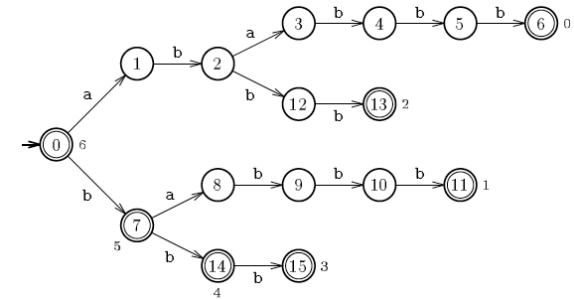
Trie des suffixes (arbre des suffixes non-compacté): trie stockant tous les suffixes du texte

T = a b a b b b

suffixes de T :

```

a b a b b b
  b a b b b
    a b b b
      b b b
        b b
          b
  
```



810

Arbre des suffixes

UMLV ©

Observation : il est facile de localiser un motif à l'aide d'un trie des suffixes : « épeler » le motif en partant de la racine ; si réussi, rapporter tous les nœuds terminaux se trouvant dans le sous-arbre correspondant

811

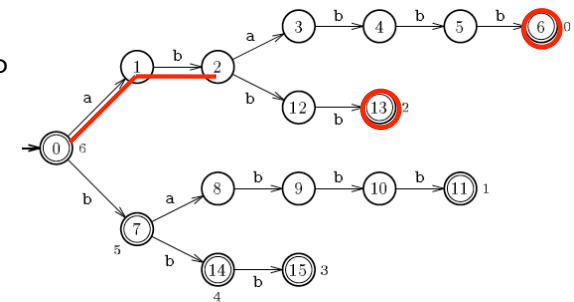
Arbre des suffixes

UMLV ©

Observation : il est facile de localiser un motif à l'aide d'un trie des suffixes : « épeler » le motif en partant de la racine ; si réussi, rapporter tous les nœuds terminaux se trouvant dans le sous-arbre correspondant

T = a b a b b b

Exemple :
localisation du motif P = a b



812

Arbre des suffixes

UMLV ©

Problème : le trie des suffixes peut avoir la taille $O(n^2)$ (n la longueur de texte)

Exemple : le trie des suffixes pour $a^k b^k$ a $\Theta(k^2)$ nœuds

813

Arbre des suffixes

UMLV ©

Problème : le trie des suffixes peut avoir la taille $O(n^2)$ (n la longueur de texte)

Exemple : le trie des suffixes pour $a^k b^k$ a $\Theta(k^2)$ nœuds

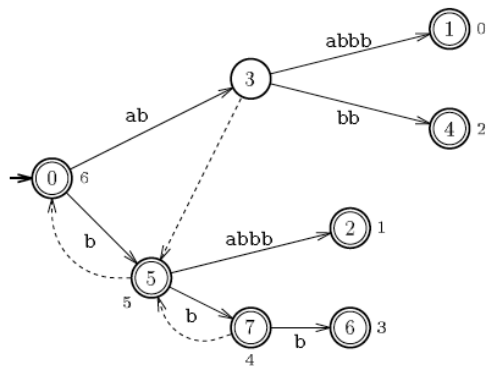
Idée : compacter le trie des suffixes : « compresser » chaque branche sans bifurcation et sans nœuds terminaux en un seul arc

814

Arbre des suffixes

UMLV ©

trie des suffixes compacté de $a b a b b b$:

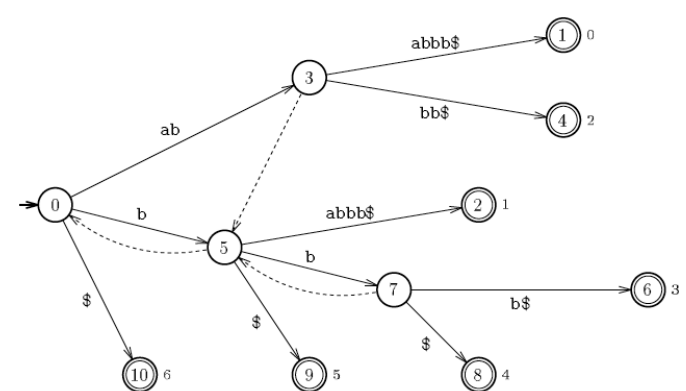


815

Arbre des suffixes

UMLV ©

trie des suffixes compacté de $a b a b b b \$$:



816

Arbre des suffixes

UMLV ©

Théorème : le trie des suffixes compacté ne peut pas avoir plus de $2n$ nœuds

Preuve : il ne peut pas y avoir plus de n feuilles (une feuille par suffixe). Comme chaque nœud interne de l'arbre est branchant ou terminal (ou les deux), il n'y a au plus n nœuds internes branchant.

817

Arbre des suffixes

UMLV ©

Théorème : le trie des suffixes compacté ne peut pas avoir plus de $2n$ nœuds

Preuve : il ne peut pas y avoir plus de n feuilles (une feuille par suffixe). Comme chaque nœud interne de l'arbre est branchant ou terminal (ou les deux), il n'y a au plus n nœuds internes branchant.

Problème : la représentation n'est pas linéaire car la taille totale des étiquettes reste quadratique dans le pire cas

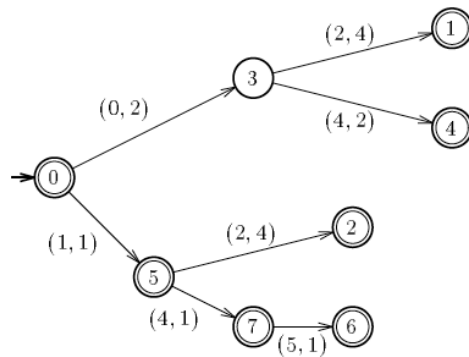
Remède : remplacer chaque étiquette par un couple (début, longueur) \Rightarrow arbre des suffixes

818

Arbre des suffixes

UMLV ©

arbre des suffixes de a b a b b b :



819

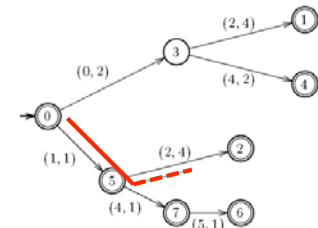
Arbre des suffixes

UMLV ©

- occupe $O(n)$ mémoire
- peut être construit en temps $O(n)$
 - Weiner (1973) (parcours du texte de droite à gauche)
 - McCreight (1976)
 - Ukkonen (1992) algorithme **on-line**
- motif $P[1..m]$ peut être localisé en temps $O(m)$

$T = a b a b b b$

$P = b a b$



820

Arbre des suffixes : applications

UMLV ©

- avec un pré-traitement en $O(n)$ les requêtes suivantes peuvent être répondues en $O(m)$:
 - première occurrence de P
 - dernière occurrence de P
 - nombre d'occurrences de P
- le plus long facteur répété dans T peut être trouvé en $O(n)$
- le plus long facteur commun à deux mots peut être trouvé en $O(n_1+n_2)$

821

Structures pour index.
Automate des suffixes
(*DAWG: Directed Acyclic Word Graph*)

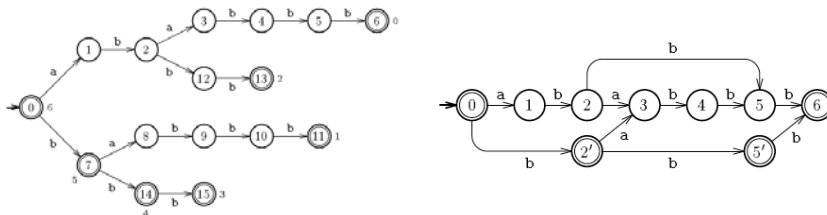
UMLV ©

822

Automate des suffixes

UMLV ©

Autre façon de compacter le trie des suffixes



Automate des suffixes : automate minimal acceptant tous les suffixes de T. Contient $\leq 2n-2$ états et $\leq 3n-4$ transitions

Peut être construit en temps $O(n)$ [Blumer et al 1983, Crochemore 1984]

823

Structures pour index.
Table des suffixes
(*Suffix Array*)

UMLV ©

824

Table des suffixes

UMLV ©

A T C A C A T C A T C A
 0 1 2 3 4 5 6 7 8 9 10 11

825

Table des suffixes

UMLV ©

A T C A C A T C A T C A
 0 1 2 3 4 5 6 7 8 9 10 11

	pos
ATCACATCATCA	0
TCACATCATCA	1
CACATCATCA	2
ACATCATCA	3
CATCATCA	4
ATCATCA	5
TCATCA	6
CATCA	7
ATCA	8
TCA	9
CA	10
A	11

826

Table des suffixes

UMLV ©

A T C A C A T C A T C A
 0 1 2 3 4 5 6 7 8 9 10 11

pos	rang	pos
ATCACATCATCA	0 A	11
TCACATCATCA	1 ACATCATCA	3
CACATCATCA	2 ATCA	8
ACATCATCA	3 ATCACATCATCA	0
CATCATCA	4 ATCATCA	5
ATCATCA	5 CA	10
TCATCA	6 CACATCATCA	2
CATCA	7 CATCA	7
ATCA	8 CATCATCA	4
TCA	9 TCA	9
CA	10 TCACATCATCA	1
A	11 TCATCA	6

827

Table des suffixes

UMLV ©

A T C A C A T C A T C A
 0 1 2 3 4 5 6 7 8 9 10 11

pos	rang	pos
ATCACATCATCA	0 A	11
TCACATCATCA	1 ACATCATCA	3
CACATCATCA	2 ATCA	8
ACATCATCA	3 ATCACATCATCA	0
CATCATCA	4 ATCATCA	5
ATCATCA	5 CA	10
TCATCA	6 CACATCATCA	2
CATCA	7 CATCA	7
ATCA	8 CATCATCA	4
TCA	9 TCA	9
CA	10 TCACATCATCA	1
A	11 TCATCA	6

828

Table des suffixes

A T C A C A T C A T C A
 0 1 2 3 4 5 6 7 8 9 10 11

pos	rang	pos
ATCACATCATCA 0	0 A 11	} A T C
TCACATCATCA 1	1 ACATCATCA 3	
CACATCATCA 2	2 ATCA 8	
ACATCATCA 3	3 ATCACATCATCA 0	
CATCATCA 4	4 ATCATCA 5	
ATCATCA 5	5 CA 10	
TCATCA 6	6 CACATCATCA 2	
CATCA 7	7 CATCA 7	
ATCA 8	8 CATCATCA 4	
TCA 9	9 TCA 9	
CA 10	10 TCACATCATCA 1	
A 11	11 TCATCA 6	

Table des suffixes

- Recherche de motif en temps $O(m \cdot \log(n))$, peut être optimisé en $O(m + \log(n))$
- Construction en temps $O(n \cdot \log(n))$, peut être optimisé en $O(n)$ (2003)