

Introduction à l'algorithmique de texte

Chaine de caractères } de longueur n sur un alphabet Σ .

Croquer dans cette chaine les occurrences

- d'un motif
- d'un mot
- d'un ensemble de mots

Applications

- Extraction d'information
- Recherche dans un éditeur
- Génomique, bioinformatique
- Correction orthographique
- ...

Def. Un préfixe ^(strict) d'un mot $u \in \Sigma^*$ est un mot $v \in \Sigma^*$ t.q. il existe un mot $w \in \Sigma^*$ avec $u = vw$.

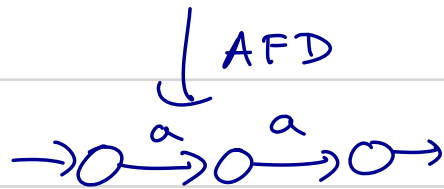
Un suffixe d'un mot $u \in \Sigma^*$ est un mot $v \in \Sigma^*$ t.q. il existe un mot $w \in \Sigma^*$ avec $u = wv$.

Mot }
 Chaîne } de taille n dans laquelle on veut rechercher
 les motifs.

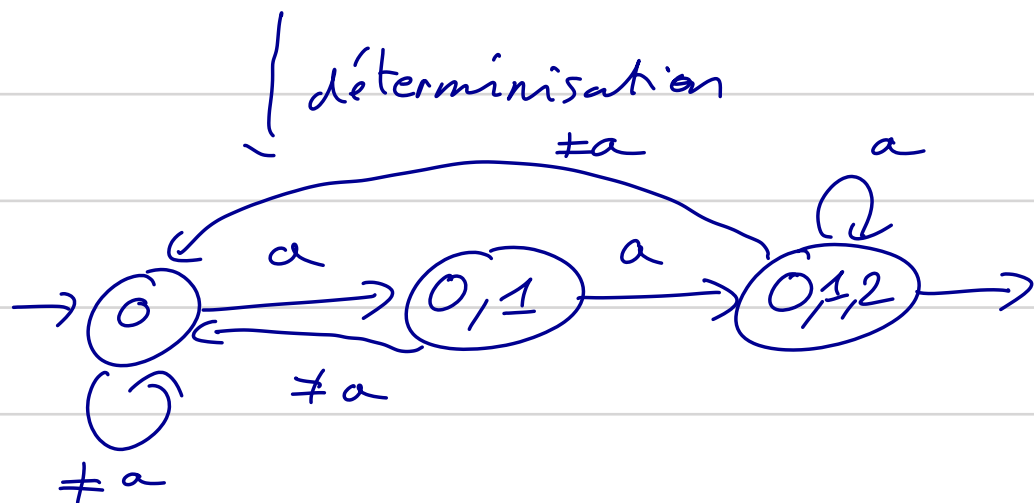
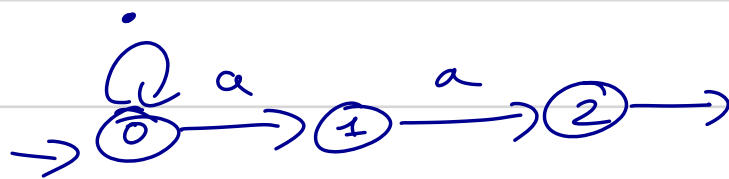
Motif	Recherche	Technique	Complexité de compilation	Complexité de recherche
e. x. e. de taille k	1 occurrence ou toutes les occurrences	AFD pour $.^*e$ (une occurrence à chaque passage dans un état final)	$O(\beta^k)$ pour $\beta > 2$	$O(n)$
chaîne de taille k	1 occ. ou toutes les occ.	naïf	O	$O(n \times k)$
chaîne de taille k	1 occ. ou toutes les occ.	Knuth-Morris-Pratt (KMP)	$O(k)$	$O(n)$
ensemble de m chaînes de taille k	correspondance exacte	Trie (ou Arbre PATRICIA)	$O(m \times k)$	$O(n)$
ensemble de m chaînes de taille k	1 occ. ou toutes les occ.	$m \times$ KMP	$O(m \times k)$	$O(m \times n)$
=====	=====	Aho-Corasick (Trie + KMP, automates à repli arborescents)	$O(m \times k)$	$O(n)$ + taille des résultats)

Ex. pom . * e

Rechercher toutes les occurrences de "aa".



AFND pour . * aa :



$a \underline{b} a \underline{b} a \underline{a} \underline{a} \underline{a} \underline{b} a \underline{a} \underline{b}$

← états finaux
 ≡ fins de motif

Ex. algo naïf

Motif

abcabd

Chaîne

abacabcabcabd

ab—

— a—

— abcab—

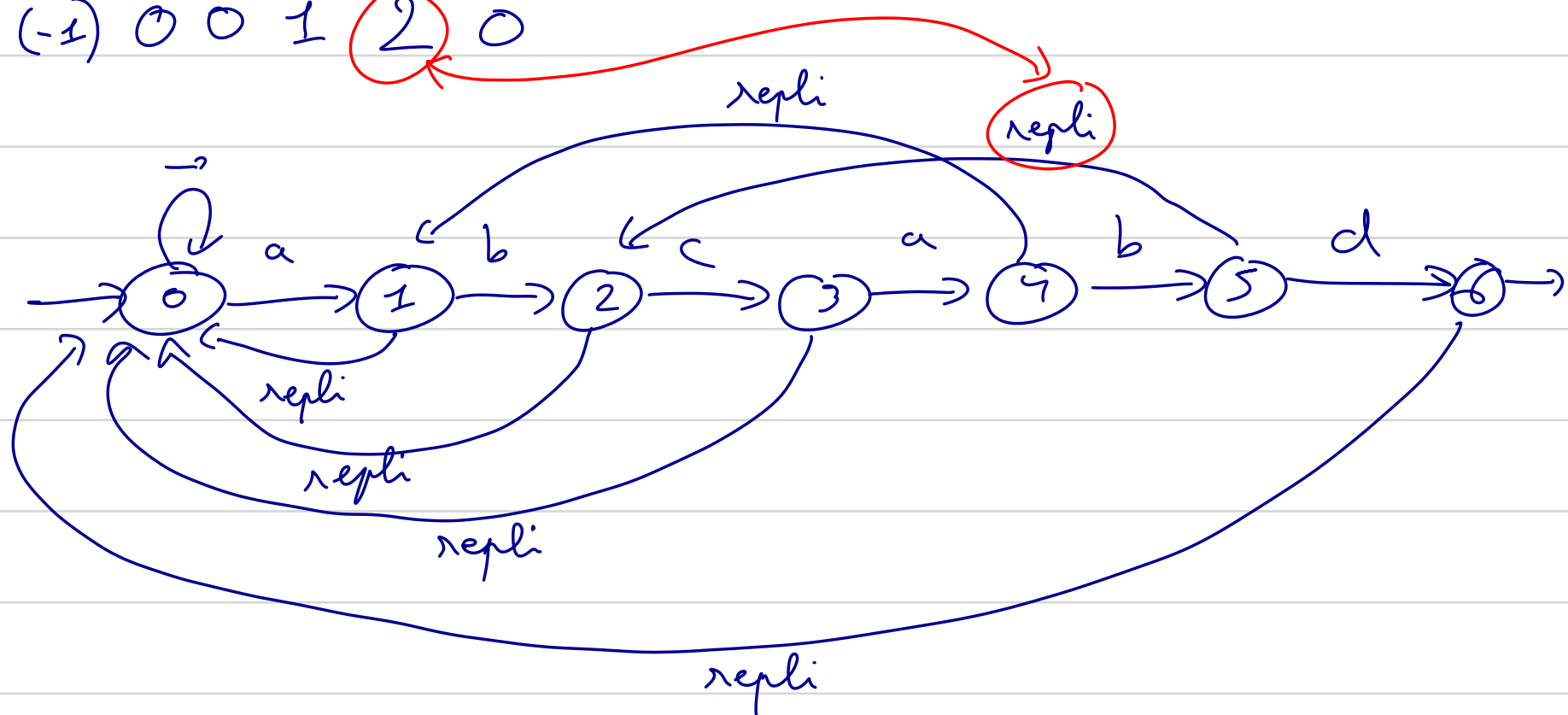
— abcabd

KMP

Compilation : On se souvient, dans le motif, de la taille du plus long préfixe strict qui est un suffixe de chaque préfixe du motif.

abcabd
 (-1) 0 0 1 2 0

1	2	3	4	5	6
a	b	c	a	b	d
(-1)	0	0	1	2	0



a b a c a b c a b c a b d

1 2
0 1
0 0

0 1 2 3 4 5

2 3 4 5 6 : match

KMP

Complexité de reconnaissance : $O(n + \# \text{replis})$

$\leq n$: nombre de fois où on va de droite à gauche \leq nombre de fois où on va de gauche à droite.

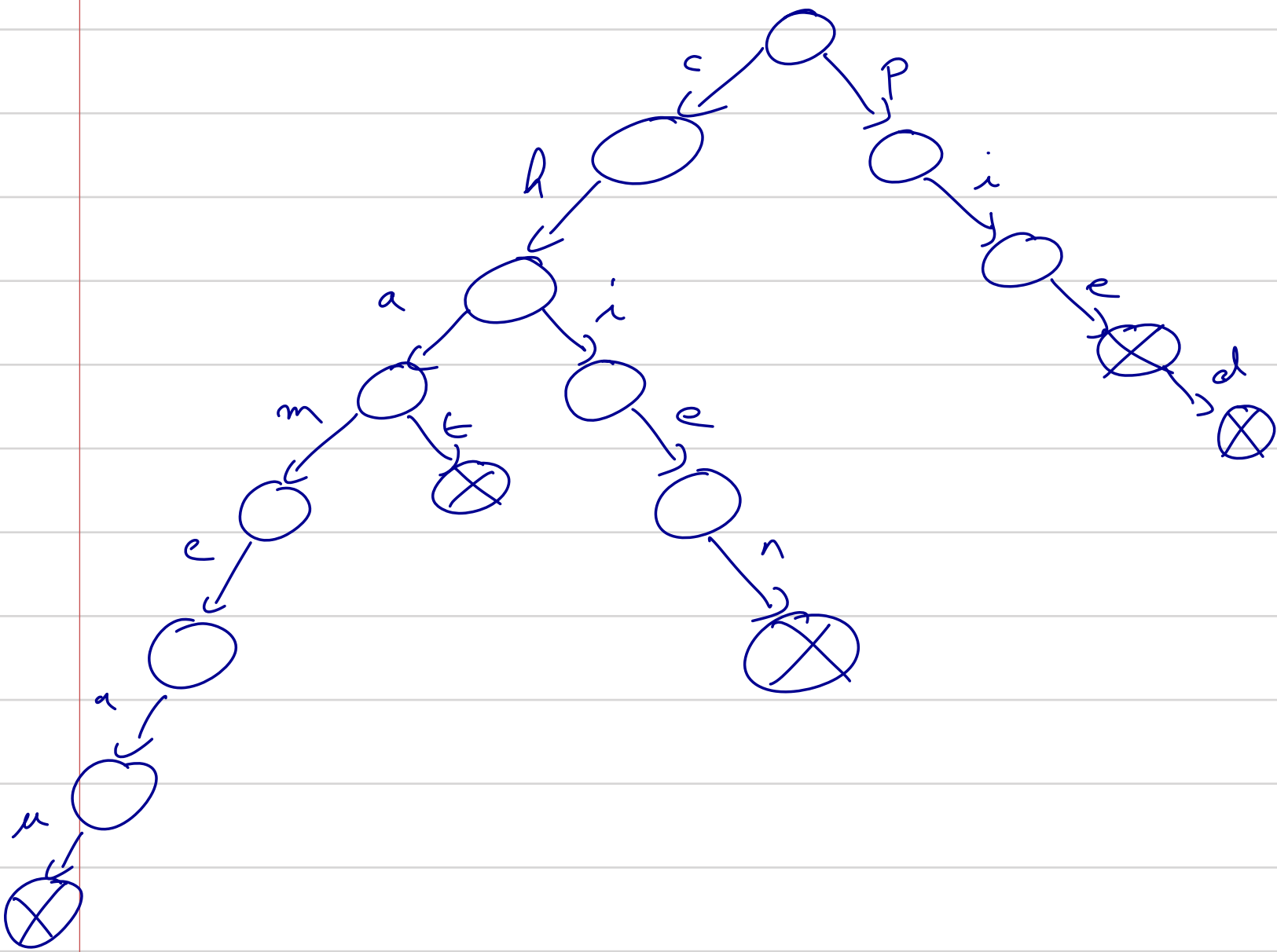
Compilation : $O(k)$ par programmation dynamique

1	2	3	4	5	6
a	b	c	a	b	d
(-1)	0	X	1	2	X
		0			X
					0

Trie ("retrieval")

Arbre dont chaque nœud correspond aux préfixes des mots d'un ensemble, chaque lien enfant - parent correspond à un symbole, et où les nœuds des mots de l'ensemble sont marqués.

{ chien, chat, chameau, pie, pied } m = 5
k = 7



Arbres radix ou PATRICIA

Raffinement des Trie dans lequel les arêtes enfant - parent peuvent être étiquetées par des sous-chaînes

