

# Recherche exacte de motif

Algorithmes naïf, MP et KMP

Sèverine Bérard



ISE-M – FDS, Université de Montpellier



## Bord et période (encore)

2

- Soit  $u$  un mot non vide,  $p$  un entier tel que  $0 < p \leq |u|$
- $p$  est une **période** de  $u$  si une de ces conditions est satisfaite :
  1.  $u[i] = u[i + p]$ , pour  $1 \leq i \leq |u| - p$
  2.  $u$  est un préfixe d'un mot  $y^k$ ,  $k > 0$ ,  $|y| = p$
  3.  $u = yw = wz$ , pour des mots  $y$ ,  $z$  et  $w$  avec  $|y| = |z| = p$   
Le mot  $w$  est appelé **bord** de  $u$
- $period(u)$  est la plus petite période de  $u$  (peut être  $|u|$ )
- $border(u)$  est le plus long bord de  $u$  (peut être  $\epsilon$ )

4	abacaba
8	aba
10	a
11	$\epsilon$

Périodes et bords du mot *abacabacaba* de longueur 11

## Plan

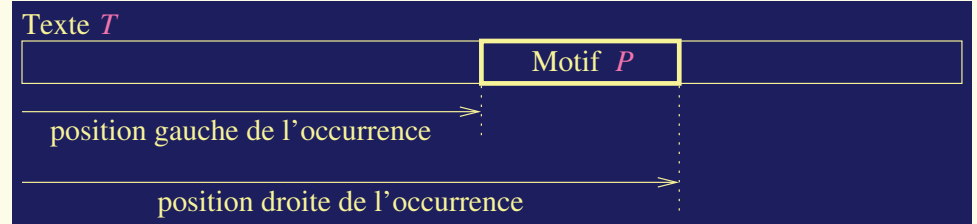
1

- Introduction
- Algorithme naïf
- Algorithme MP
- Algorithme KMP
- Références

## Recherche de motif exact

3

Trouver toutes les occurrences d'un motif  $P$  de longueur  $m$  à l'intérieur d'un texte  $T$  de longueur  $n$



- **Motif** : un mot  $P$  de longueur  $m$

t a t a

- **Texte** : un mot  $T$  de longueur  $n$

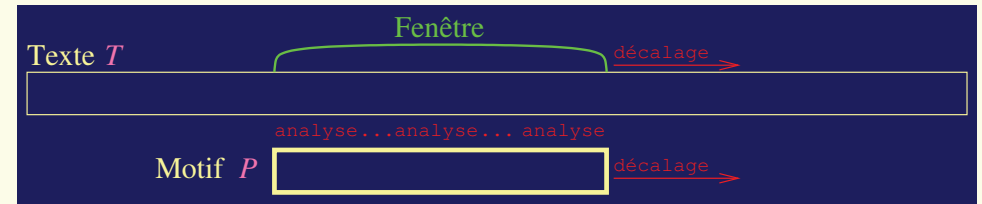
a g g c t c a c g t a t a t a t g c g t t a t a a t

- **Occurrences** :

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26
a	g	g	c	t	c	a	c	g	t	a	t	a	t	a	t	g	c	g	t	t	a	t	a	a	t
									t a t a																
										t a t a															

Le motif  $P$  apparaît aux positions 10, 12 et 21 dans le texte  $T$

- **Opérations élémentaires** : comparaison de symboles ( $=, \neq$ )



### Algorithme 1 : Mécanisme « analyse et décalage »

**Données** : Deux chaînes  $T$  et  $P$  de longueurs respectives  $n$  et  $m$ .

Placer la **Fenêtre** au début du texte ;

**tant que** la **Fenêtre** est sur le texte **faire**

**analyse** : **si** **Fenêtre** = **Motif** **alors** le rapporter ;

**décalage** : déplacer la **Fenêtre** vers la droite et  
                 mémoriser des informations à utiliser durant les prochaines  
                 analyses et décalages ;

- Introduction
- Algorithme naïf
- Algorithme MP
- Algorithme KMP
- Références

- Principes : Pas de mémorisation, décalage d'une seule position
- Complexité :  $O(m \times n)$

### Algorithme 2 : Recherche naïve

**Données** : Deux chaînes  $T$  et  $P$  de longueurs respectives  $n$  et  $m$ .

$pos := 1$  ;

**tant que**  $pos \leq n - m + 1$  **faire**

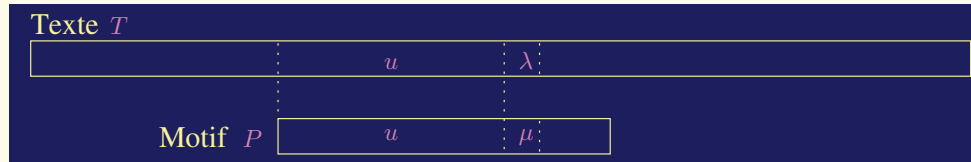
$i := 1$  ;

**tant que**  $(i \leq m \text{ et } P[i] = T[pos + i - 1])$  **faire**  $i := i + 1$  ;

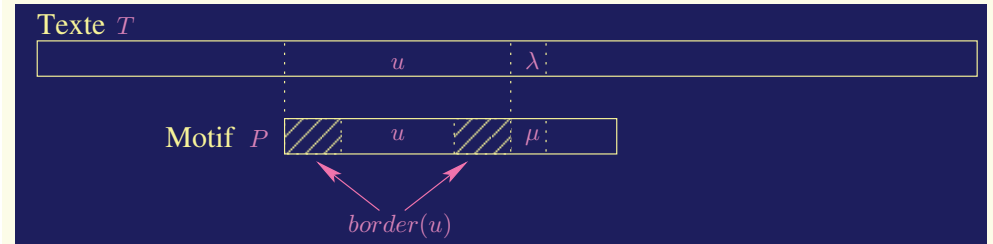
**si**  $i = m + 1$  **alors**

        Écrire("P apparaît à la position ",  $pos$ ) ;

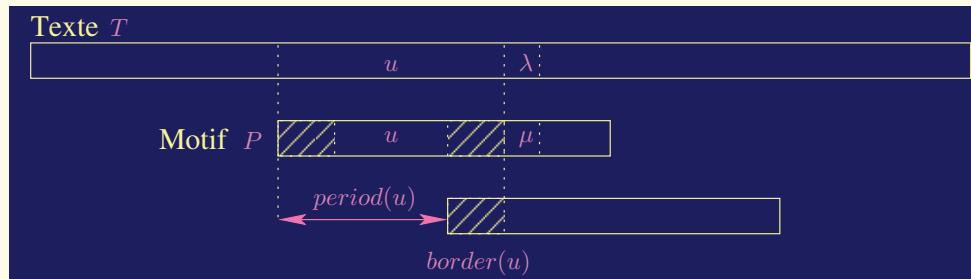
$pos := pos + 1$  ;



- Situation d'échec : identité des caractères de  $u$  mais  $\lambda \neq \mu$
- $u$  est un préfixe de  $P$



- Situation d'échec : identité des caractères de  $u$  mais  $\lambda \neq \mu$
- $u$  est un préfixe de  $P$



- Situation d'échec : identité des caractères de  $u$  mais  $\lambda \neq \mu$
- $u$  est un préfixe de  $P$
- $|u| - |border(u)| = period(u)$

- Introduction
- Algorithme naïf
- Algorithme MP
- Algorithme KMP
- Références

- Longueur du décalage  $\geq 1$  : taille de la période  $period(u)$
- Mémorisation des bords

---

**Algorithme 3 : Mécanisme « analyse et décalage » amélioré**


---

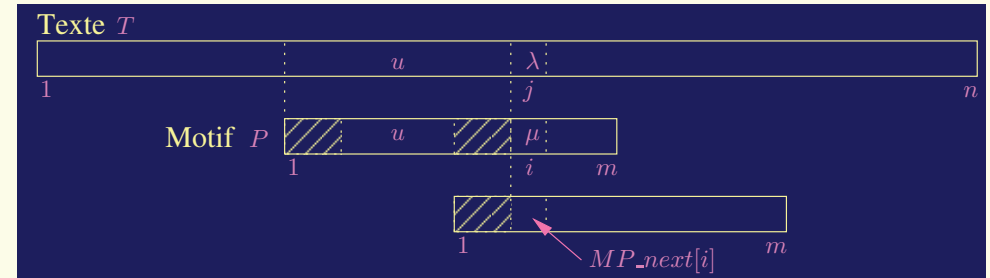
**Données** : Deux chaînes  $T$  et  $P$  de longueurs respectives  $n$  et  $m$ .

Placer la Fenêtre au début du texte ;

**tant que** la Fenêtre est sur le texte **faire**

$u$  := plus long préfixe commun entre la Fenêtre et le Motif ;  
**si**  $u = \text{Motif}$  **alors** le rapporter ;  
 Décaler la Fenêtre de  $period(u)$  places vers la droite ;  
 Mémoriser  $border(u)$  ;

---




---

**Algorithme 4 : Algorithme MP**


---

**Données** : Deux chaînes  $T$  et  $P$  de longueurs respectives  $n$  et  $m$ .

$i := 1$  ;  $j := 1$  ;

**tant que**  $j \leq n$  **faire**

**tant que** ( $i > 0$  et  $P[i] \neq T[j]$ ) **faire**  
      $i := MP\_next[i]$  ;  
 $i := i + 1$  ;  $j := j + 1$  ;  
**si**  $i = m + 1$  **alors**  
     Écrire ("P apparaît à la position ",  $j - i + 1$ ) ;  
      $i := MP\_next[i]$  ;

---

- Remarque : un bord d'un bord de  $u$  est un bord de  $u$
- On calcule un tableau  $BORD$ , tel que  $BORD[i] = |border(P[1..i])|$

---

**Algorithme 5 : Calcule  $BORD$** 


---

**Données** : Un mot  $P$  de longueur  $m$

$BORD[0] := -1$  ;

**pour** ( $i$  de 1 à  $m$ ) **faire**

$j := BORD[i - 1]$  ;  
**tant que** ( $j \geq 0$  et  $P[i] \neq P[j + 1]$ ) **faire**  $j := BORD[j]$  ;  
 $BORD[i] := j + 1$  ;

---

- $i$  parcourt les longueurs des préfixes de manière croissante
- $j$  parcourt les longueurs des bords de manière décroissante

Remarque :

$MP\_next[i] = BORD[i - 1] + 1$ , pour  $i$  de 1 à  $m + 1$

⇒ Quasiment même algorithme

#### Algorithme 6 : Calcule *MP\_next*

**Données** : Un mot  $P$  de longueur  $m$

$MP\_next[1] := 0$  ;  $j := 0$  ;

**pour** ( $i$  de 1 à  $m$ ) **faire**

/\* À chaque entrée de boucle on a  $j := MP\_next[i]$  \*/

**tant que** ( $j > 0$  **et**  $P[i] \neq P[j]$ ) **faire**  $j := MP\_next[j]$  ;

$j := j + 1$  ;

$MP\_next[i + 1] := j$  ;

- $i$  parcourt les longueurs des préfixes de manière croissante
- $j - 1$  parcourt les longueurs des bords de manière décroissante

- Introduction
- Algorithme naïf
- Algorithme MP
- Algorithme KMP
- Références

#### Complexité :

- Pré-traitement :  $O(m)$  en temps et en espace
- Phase de recherche :  $O(n + m)$  en temps et en espace

- Une fois le pré-traitement effectué sur un motif, on peut le rechercher dans autant de textes que souhaités

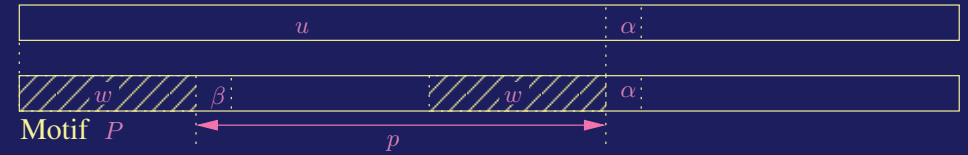
#### Intérêts :

- Phase de recherche plus rapide que l'algorithme naïf
- Très facile à implémenter

- Inconvénients** : besoin de temps et d'espace supplémentaire ( $O(m)$ )

Peut-on faire mieux ? OUI!

Motif  $P$



- Notions de bords stricts ( $w$ ) et de périodes interrompues ( $p$ )
- Modifie seulement la phase de pré-traitement

#### Algorithme 7 : Mécanisme de MP amélioré

**Données** : Deux chaînes  $T$  et  $P$  de longueurs respectives  $n$  et  $m$ .

Placer la Fenêtre au début du texte ;

**tant que** la Fenêtre est sur le texte **faire**

$u :=$  plus long préfixe commun entre la Fenêtre et le Motif ;

**si**  $u = \text{Motif}$  **alors** le rapporter ;

Décaler la Fenêtre de  $\text{period\_int}(u)$  places vers la droite ;

Mémoriser  $\text{border\_strict}(u)$  ;

- Soit  $P$  un mot fixé et  $u$  un préfixe non vide de  $P$
- $w$  est un **bord strict** de  $u$  si à la fois :
  1.  $w$  est un bord de  $u$
  2.  $w\beta$  est un préfixe de  $P$  mais pas  $u\beta$
- $p$  est une **période interrompue** de  $u$  si  $p = |u| - |w|$  avec  $w$  un bord strict de  $u$
- Exemple : soit  $u = abacabacaba$  un préfixe de  $P = abacabacabacc$

10	$a$
11	$\epsilon$

Périodes interrompues et bords strict de  $abacabacaba$

## Conclusion sur l'algorithme KMP (1977)

- **Complexité** : idem MP pour le pire des cas
  1. Pré-traitement :  $O(m)$  en temps et en espace
  2. Phase de recherche :  $O(n + m)$  en temps et en espace
- Une fois le pré-traitement effectué sur un motif, on peut le rechercher dans autant de textes que souhaités
- **Intérêts** :
  1. Phase de recherche optimisée par rapport à MP
  2. Aussi facile à implémenter

Peut-on encore faire mieux ? **OUI!**

→ Recherche de droite à gauche (cf. Boyer-Moore)

## Pré-traitement de l'algorithme KMP : $KMP\_next$

- $k = MP\_next[i]$
- $KMP\_next[i] = \begin{cases} k & \text{si } P[i] \neq P[k] \text{ ou si } i = m + 1 \\ KMP\_next[k] & \text{si } P[i] = P[k] \end{cases}$

### Algorithme 8 : Calcule $KMP\_next$

**Données** : Un mot  $P$  de longueur  $m$

$KMP\_next[1] := 0$  ;  $j := 0$  ;

**pour** ( $i$  de 1 à  $m$ ) **faire**

**tant que** ( $j > 0$  **et**  $P[i] \neq P[j]$ ) **faire**  $j := KMP\_next[j]$  ;

$j := j + 1$  ;

**si** ( $i = m$  **ou**  $P[i + 1] \neq P[j]$ ) **alors**

$KMP\_next[i + 1] := j$  ;

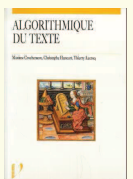
**sinon**

$KMP\_next[i + 1] := KMP\_next[j]$

## Références

- Ce cours a été construit à partir du cours de MAXIME CROCHEMORE et THIERRY LECROQ : **Text searching**
- Pour aller plus loin :

1. M. Crochemore, C. Hancart et T. Lecroq, Algorithmique du texte, Vuibert, 2001, 347 pages. ISBN 2-7117-8628-5. Disponible en version pdf : <http://www-igm.univ-mlv.fr/~mac/CHL/CHL-2011.pdf>



2. T. Cormen, C. Leiserson, R. Rivest et C. Stein, Algorithmique, Dunod, 2010 - 3e édition, 1296 pages. ISBN : 9782100545261

~ 15 exemplaires disponibles à la BU (2e et 3e édition)

