

SORBONNE UNIVERSITÉ

FACULTÉ DE SCIENCES ET INGÉNIERIE MASTER II STL  
DAAR

---

# Clone de egrep unix

---

*Étudiants*

Lyes LOUNIS

Arezki GALOU

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Algorithmes de recherches</b>	<b>2</b>
2.1	Methode Aho Ulman . . . . .	2
2.1.1	Transformation en automate avec $\varepsilon$ transitions . . . . .	2
2.1.2	Transformation en automate sans $\varepsilon$ transitions . . . . .	3
2.1.3	Transformation en automate deterministe . . . . .	4
2.1.4	Structure finale . . . . .	4
2.2	Méthode Knuth-Morris-Prat . . . . .	4
<b>3</b>	<b>Teste</b>	<b>5</b>
3.0.1	Aho Ulman VS KMP . . . . .	5
3.0.2	Aho Ulman vs Egrep linux . . . . .	6
<b>4</b>	<b>Conclusion</b>	<b>7</b>

## 1 Introduction

La commande `egrep` de Unix permet de chercher un motif dans un texte et afficher toutes les lignes contenant ce motif.

le motif peut être soit une chaîne de caractère simple ("Babylone") ou une expression régulière : une représentation d'un langage algébrique reconnaissable par un automate fini déterministe.

Le but du projet est d'implémenter différents algorithmes afin de reproduire `egrep`, ici pour les expressions régulières on considère seulement les symboles suivants :

- concatenation : "ab" .
- étoile : "c\*" la lettre c répéter zéro ou plusieurs fois ("cccc ...." ou "").
- alternative : "motif1" | "motif2" soit le motif1 ou le motif2

Dans notre implémentation si le motif se réduit à une chaîne de caractère simple (pas une `regE`) on réalise la recherche avec l'algorithme de : Knuth-Morris-Pratt.

Dans le cas où le motif à chercher est une `regE` on réalise la recherche avec l'algorithme de : Aho Ulman

## 2 Algorithmes de recherches

### 2.1 Méthode Aho Ulman

On utilise cette méthode afin de créer un automate déterministe qu'on utilisera lors de la recherche du motif dans un texte, pour cela on passe par les étapes suivantes :

- On transforme la `regEx` en un arbre syntaxique .
- À partir de cet arbre construire un automate avec  $\varepsilon$ -transitions .
- On enlève les  $\varepsilon$ -transitions
- Enfin construire notre automate déterministe .

#### 2.1.1 Transformation en automate avec $\varepsilon$ transitions

On parcourt l'arbre syntaxique et on applique cet algorithme :  
Si le nœud est une **feuille** (caractère) : On construit deux nœuds (initial et final) avec une transition de ce caractère depuis le nœud initial vers le nœud final

.Si le noeud est une **Alternative** : on construit deux automates récursivement avec ces deux fils, on construit un noeud initial, on ajoute deux  $\varepsilon$  transitions vers les deux nouveaux automates depuis le noeud initial, on ajoute  $\varepsilon$  transition depuis les deux noeuds finaux des nouveaux automates vers un nouveau noeud final .

.Si le noeud est une **Concaténation** : On construit deux automates avec ces deux fils, on ajoute une  $\varepsilon$ -transition depuis le noeud final du premier fils vers le noeud initial de l'autre automate

.Si le noeud est une **Etoile** : On construit un automate avec son fils , on crée deux noeuds (initial ni et final nf), on ajoute deux  $\varepsilon$  transitions depuis l'initial(ni) vers le final(nf) et vers le noeud initial du nouvel automate (fils), on ajoute  $\varepsilon$  transition depuis le noeud final du nouvel automate vers le noeud final (nf) .

### 2.1.2 Transformation en automate sans $\varepsilon$ transitions

Dans cette étape on parcourt l'automate avec  $\varepsilon$  transitions et on construit entièrement un nouvel automate de cette manière :

**algo** (courant : node, dernierImportant : node)

On garde une référence du dernier noeud important visité (un noeud qui a une transition entrante autre que  $\varepsilon$  plus le noeud initial), au début on initialise avec le noeud initial(dernierImportant = initial)

Pour chaque noeud visité on parcourt ses transitions :

Soit une transition  $\alpha$  vers un noeud n :

- Si  $\alpha = \varepsilon$  : on fait un appel récursif (on visite n)
- Sinon : Dans ce cas le noeud n est un noeud important alors :
  - Créer un nouvel noeud (newNode) si ce n'est déjà créé qui représente ce noeud important .
  - On ajoute une transition (dernierImportant,  $\alpha$ , newNode).
  - Enfin mettre à jour dernierImportant = newNode et visiter n.

**les états acceptables** : Si le noeud qu'on visite (courant) est un noeud acceptable alors dernierImportant est acceptable aussi .

### 2.1.3 Transformation en automate déterministe

Dans cette étape on parcourt l'automate sans  $\varepsilon$  transitions et on construit entièrement un nouvel automate déterministe de cette manière :

**Principe** : Pour chaque nouvel noeud crée on garde une référence vers l'ensemble des noeuds regroupé .

**initialisation** : Créer un nouvel noeud  $n$ , qui a comme ensemble le noeud initial de notre automate sans  $\epsilon$  .

**Parcours** : Pour chaque noeud visité on parcourt son ensemble de noeud regroupé : et pour chaque noeud  $n$  dans cet ensemble on parcourt ses transitions sortantes pour construire une map (caractere : ensemble de noeud) .

On parcourt la map et pour chaque ensemble de noeuds faire :

- Si on a déjà créer un noeud contenant cet ensemble alors ne rien faire .
- Sinon créer un nouvel noeud (`newNode`) contenant cet ensemble, puis on ajoute une transitions (noeud courant, clé, `newNode`), enfin on visite `newNode` .

construire un nouvel noeud contenant cet ensemble de noeud (On vérifie d'abord si le noeud existe pas déjà )

**Les états acceptables** : Un état est acceptable si au moins un état dans son ensemble regroupé est un état acceptable

### 2.1.4 Structure finale

Après avoir construit notre automate déterministe, on le transforme en une matrice (nombre de noeud x 256 )

$\text{matrice}[i][\alpha] = k$  : transition depuis l'état  $i$  vers l'état  $j$  avec  $\alpha$

## 2.2 Méthode Knuth-Morris-Prat

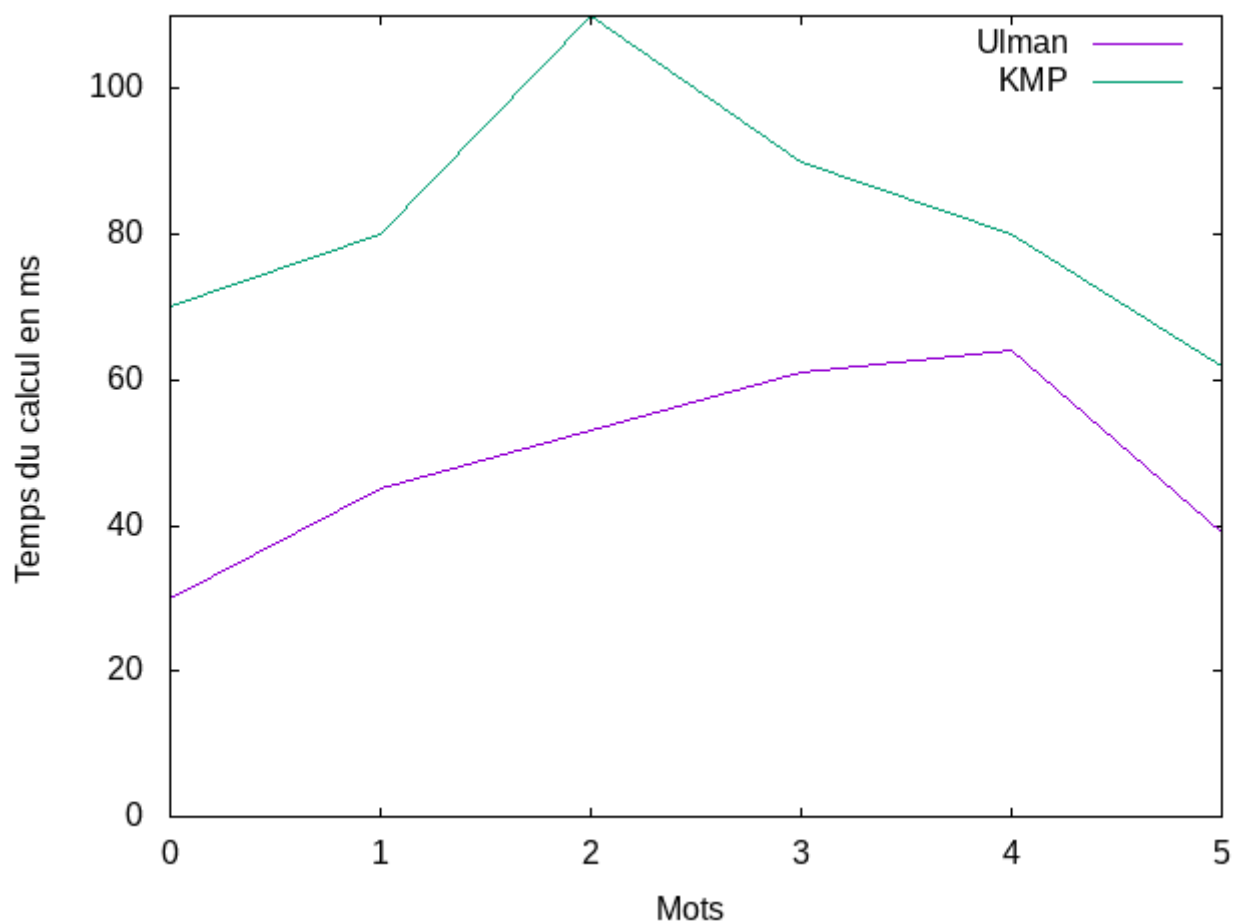
L'algorithme KMP est conçu pour améliorer la recherche naïve d'une sous chaîne de caractère, la particularité de cet algorithme réside dans son pré-traitement de la sous chaîne, ainsi on a 2 étapes à effectuer :

- **Calcule de la retenue** : on calcule un tableau de taille égale à la sous chaîne, chaque case indique la prochaine position
- **La recherche** : on se servant du tableau **retenue** on évite de revenir à chaque échec à la position  $i+1$  .

## 3 Teste

### 3.0.1 Aho Ulman VS KMP

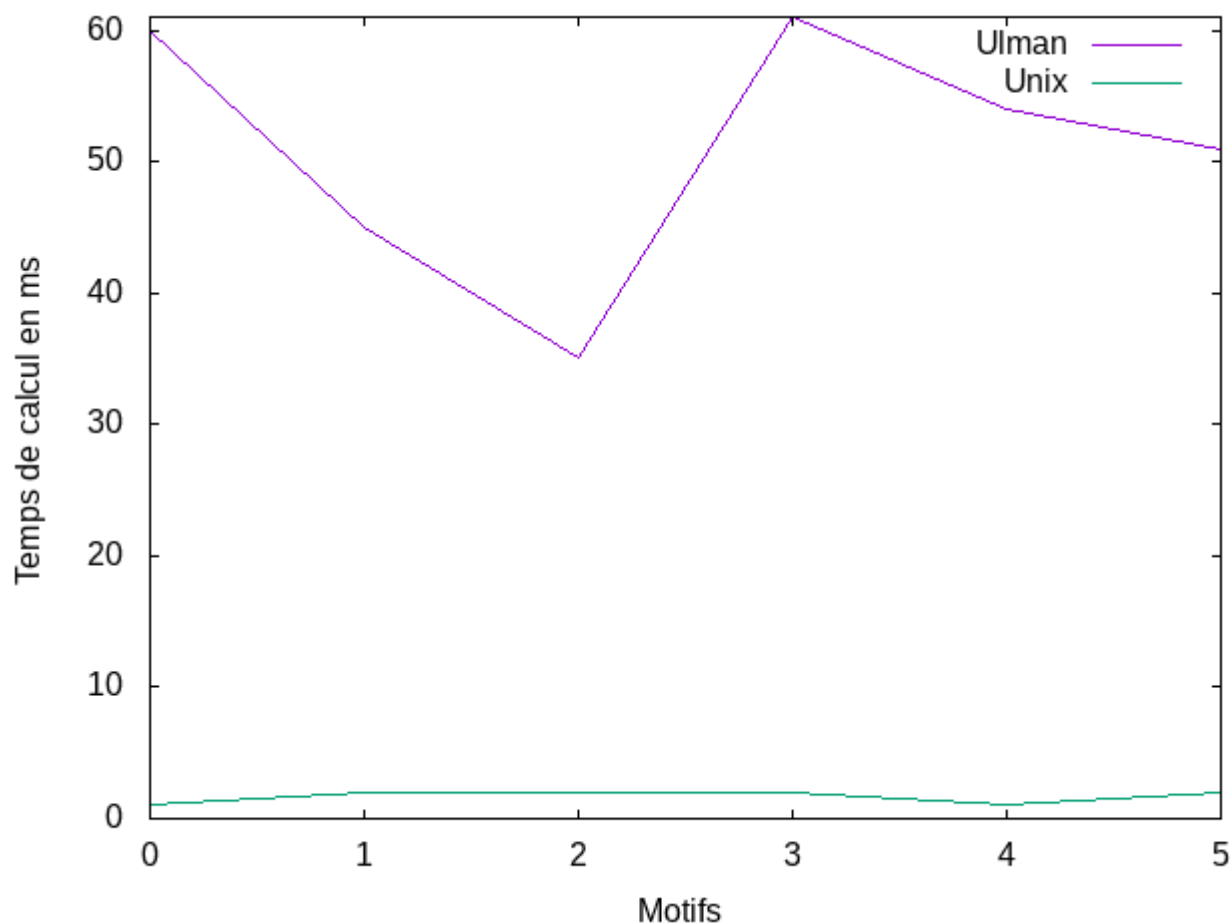
On compare le temps d'exécutions entre les algorithmes Aho Ulman et KMP avec différents mots .



On remarque que la méthode Ulman est plus efficace que KMP, le temps de construction et de recherche est beaucoup moins important dans Ulman

### 3.0.2 Aho Ulman vs Egrep linux

On compare le temp d'exécution entre les algorithmes Aho Ulman et egrep de Unix avec différents motifs



On remarque que egrep de Unix est beaucoup plus efficace, le temps de construction de notre automate est important ce qui augmente le temps de calcul .

## 4 Conclusion

Unix reste plus rapide que nos deux implementation, par contre la méthode de Ulman est plus efficace que celle de KMP, donc on prefere s'en passer de kmp, et utiliser la méthode de Ulman meme dans le cas d'une chaine caractere simple .