

## TP Arbre des suffixes

S  verine B  rard - ISE-M, Facult   des Sciences, Universit   de Montpellier

Severine.Berard@umontpellier.fr

Le but de ce TP est de programmer l'algorithme de construction de l'arbre des suffixes quadratique vu en cours.

Un canevas de programmation de ce TP en C++ est disponible sur Moodle : `Canevas_TP_2019_ArbreSuf.cc`.    utiliser sauf bonne raison contraire. La ligne pour compiler est `g++ Canevas_TP_2019_ArbreSuf.cc -o NomExe`.

Soit  $N_i$  l'arbre qui encode tous les suffixes de 1     $i$ , pour passer de  $N_i$      $N_{i+1}$  :

1. Partir de la racine de  $N_i$
2. Trouver le plus long chemin depuis la racine qui correspond    un pr  fixe de  $T[i+1..n]$   
*   un seul chemin possible car aucun suffixe de  $T$  n'est pr  fixe d'un autre suffixe, on est alors soit    un n  ud  $w$  ( $\neq$  feuille) ou au milieu d'une ar  te*
3. Ins  rer   ventuellement un nouveau n  ud et cr  er un nouvel arc  
     tiquet   avec les derniers caract  res de  $T[i+1..n]$  non mis en correspondance  
   finissant sur une nouvelle feuille   tiquet  e  $i+1$

---

**Algorithme 1** : Approche na  ve

---

**Donn  es** : Le texte  $T$  de lg  $n$

$\mathcal{T} := \text{ArbreVide}$  ;

**pour** ( $i$  de 1     $n$ ) **faire**

    Ins  rer l'arc  $T[i..n]$  dans  $\mathcal{T}$  ;

---

Il vous faudra r  fl  chir aux points suivants :

-    Comment coder l'alphabet ?
-    Quelle structure choisir pour le codage de l'arbre ? O   stocker les   tiquettes des arcs ?
-    Comment parcourir l'arbre ?
-    ...

Dans le canevas, une structure est propos  e avec les fonctions d'affichages et d'exportation vers le format `.dot`. Mais vous pouvez tout    fait cr  er votre propre structure. Pour visualiser vos arbres, une fois le fichier dot g  n  r  , vous n'avez plus qu'   taper en ligne de commande : `dot -Tpdf monfic.dot -o monfic.pdf`.

**Remarque importante** : Nous rappelons que nous consid  rons que la num  rotation des cha  nes de caract  res commence    1. Or dans les langages de programmation la premi  re lettre d'une cha  ne est stock  e    la position 0. Pour   viter des manipulations d'indices fastidieuses, nous rangerons donc la premi  re lettre d'une cha  ne    la position 1 en ajoutant un caract  re arbitraire devant nos cha  nes.

On veut un seul fichier de code avec toutes les fonctions    l'int  rieur. Les langages autoris  s sont C, C++ et java avec une forte pr  f  rence pour C++. Ce fichier sera bien s  r propre et comment  .