



# **Rapport de Projet**

**Rendu 3**

**Algorithmique des Arbres**

Gaël Paquette

Hervé Nguyen

22 avril 2022

L2 Informatique

# Table de matières

Table de matières	2
Introduction	3
Module ArbreBK	3
Recherche	4
Évaluation des performances	6
Structuration du rendu	8
Organisation du travail	9
Comment compiler les programmes ?	9
Comment utiliser les programmes ?	10
Problèmes rencontrés	11
Addendum	12
Conclusion	13

# Introduction

Ce dernier rendu consiste à développer une implémentation des Arbres BK. Cette structure de donnée étant plus adaptée dans le cadre de la recherche de corrections. Un Arbre BK sera alors représenté sous la forme d'un arbre binaire avec des chaînages « Fils Gauche » et « Frère Droit » pour simuler un Arbre n-aire.

Pour ce second rendu , nous avons repris notre rendu 2 et puis nous avons ajouté un nouveau module ArbreBK.

Ce module ArbreBK permettra d'employer une correction orthographique en se passant de la méthode brute force en recherchant (via Levenshtein) avec les ATR.

## Module ArbreBK

Ce module a une dépendance avec le module Levenshtein pour les calculs de valeurs de noeuds et aussi une dépendance avec le module Listes pour la fonction de recherche dans un Arbre BK qui renvoie une liste.

Comme pour les ATR, la subtilité ici est autour de la recherche et l'insertion. L'insertion consiste à calculer au fur et à mesure la distance de Levenshtein

entre le mot à insérer et le mot du noeud père, après un calcul on cherche à insérer au fils gauche du père actuel.

Soit le fils gauche n'existe pas donc on crée un noeud sur cette position sinon, on se déplace suivant le fils gauche si  $d(\text{Mot à ajouter}, \text{Mot du père}) == d(\text{Mot du noeud actuel}, \text{Mot du père})$ .

Dans le cas contraire on se déplace suivant le frère droit si  $d(\text{Mot à ajouter}, \text{Mot du père}) != d(\text{Mot du noeud actuel}, \text{Mot du père})$ .

On répète ces instructions jusqu'à que l'on se retrouve à un noeud vide ou bien si le mot à insérer est déjà présent.

L'affichage sur la console consiste simplement à bien reconnaître les profondeurs de chaque noeud et espacer correctement les noeud. Un frère droit a la même profondeur que le noeud courant, tandis que son fils gauche est à une profondeur « + 1 » à celui-ci.

Un affichage de l'Arbre BK sur pdf avec dot est aussi disponible si le système où le programme tourne a le paquet graphviz installé.

## Recherche

La recherche (de correction) est dans ce rendu très important car c'est un des points qui permettent au programme d'avoir une complexité bien moindre que le parcours naïf avec des ATR.

On s'est servi d'une propriété du cours :

**Propriété :**

Soit  $v_i$  la valuation du sous-arbre  $A_i$  vérifiant  $|d - v_i| > S_0$ .

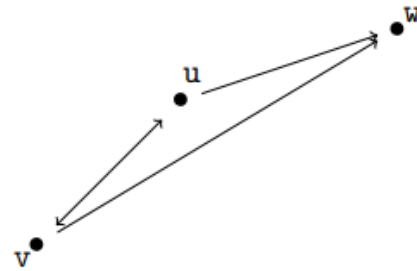
Alors, le sous-arbre  $A_i$  ne contient pas de mots au plus à distance  $S_0$  de  $m$ .

Qui provient de cette inégalité triangulaire :

**Propriété : Seconde inégalité triangulaire**

Soit  $u$ ,  $v$  et  $w$  trois chaînes de caractères.  
Alors, on a :

$$|\mathcal{L}(u, v) - \mathcal{L}(u, w)| \leq \mathcal{L}(v, w) . \quad (2)$$



Cela nous permet de ne pas parcourir des sous-arbres avec des mots qui ont peu d'intérêt (Distance de Levenshtein trop important avec le mot recherché).

Nous avons pris peu de liberté et nous avons implémenté la recherche décrite par le cours :

**Principe détaillé :**

`recherche(ArbreBK A, char * m, Liste * corrections, int * seuil) :`

- Si l'arbre est vide, on ne fait rien
- Sinon, soit  $u$  la racine de l'arbre
  - $d = \mathcal{L}(u, m)$
  - Si  $d = *seuil$ , ajouter  $r$  en tête de la liste chaînée `*corrections`
  - Si  $d < *seuil$  :  
vider `*corrections` et y ajouter ensuite  $r$   
modifier `*seuil` en le fixant à  $d$
  - Pour chaque enfant  $e$  de  $r$ , relié par une arête de valuation  $v_i$  :  
Si  $|v_i - d| \leq *seuil$  :  
    `recherche(e, m, corrections, *seuil)`

Jusque là nous avons décrits un algorithme de recherche de correction, mais nous avons également écrit une fonction de recherche qui renvoie 0 ou 1 selon la présence d'un mot dans l'Arbre BK. Cette fonction de recherche reprends le même parcours de l'insertion.

# Évaluation des performances

Nous avons temporairement modifié nos programmes du précédent rendu et de celui-ci pour évaluer les performances entre l'Arbre BK et l'ATR.

Nous avons mesuré le temps d'exécution des deux programmes avec `gettimeofday()` de `<sys/time.h>` pour calculer l'intervalle au début et à la fin des exécutions. Il y a eu 2 test, le premier avec un petit texte puis le deuxième avec un texte plus volumineux.

Pour le premier test, nous avons utilisé le texte `a_corriger_1.txt` avec le dictionnaire `dico_3.dico` :

- Le programme `correcteur_1` a en moyenne 110ms en temps d'exécution.
- Le programme `correcteur_2` a en moyenne 180ms en temps d'exécution.

Le programme utilisant les ATR s'est exécuté 70ms plus rapidement que celui avec les Arbre BK, cependant ce text était assez court donc on ne peut pas encore juger si les ATR sont supérieur au niveau de la complexité.

Pour le second test, nous avons utilisé le texte `a_corriger_2.txt` avec le dictionnaire `dico_3.dico`, ce texte est un extrait du Chapitre VI « Jean Valjean » Des Misérables de Victor Hugo (accessible sur ce [lien](#)) :

- Le programme `correcteur_1` a en moyenne 2100ms en temps d'exécution.
- Le programme `correcteur_2` a en moyenne 900ms en temps d'exécution.

On remarque alors que les ArbresBK commencent à avoir un avantage vis-à-vis des ATR. Avec le programme avec les ATR qui s'exécute 2 fois plus lentement que les ArbresBK.

Nous avons triplé la taille de ce même texte (copier-coller) pour confirmer l'avantage en complexité de correcteur\_2 et nous avons eu :

- correcteur\_1 avec 6000ms en temps d'exécution.
- correcteur\_2 avec 2000ms en temps d'exécution.

On peut alors conjecturer avec une confiance relativement haute que les ArbresBK sont beaucoup plus adaptés que les ATR en terme de complexité pour les corrections orthographiques, sur ce dernier essai correcteur\_2 s'est exécuté avec 1/3 du temps d'exécution de correcteur\_1. L'écart semble bien se creuser selon la taille du texte.

# Structuration du rendu

Le rendu contient les fichiers en-têtes des modules dans le sous répertoire « ../include »

Les fichiers sources sont dans le sous répertoire « ../src ».

Le makefile, les fichiers régénérables, le rapport se trouvent au répertoire courant du rendu.

Voici la liste des modules :

- ATR (Permet d'avoir une implémentation d'un dictionnaire)
- Listes (Implémentation de listes chaînées simples)
- Levenshtein (Permet de calculer des distances de Levenshtein)
- ArbreBK (Permet de manipuler des Arbres BK)

Voici les fichiers « main » :

- correcteur\_0.c (Fichier main du rendu 1 pour le programme correcteur\_0)
- correcteur\_1.c (Fichier main du rendu 2 pour le programme correcteur\_1, affiche des suggestions de correction en plus du comportement habituel de correcteur\_0)
- correcteur\_2.c (Fichier main du rendu 3 pour le programme correcteur\_2)

Deux fichiers textes sont fournis :

- a\_corriger\_0.txt
- a\_corriger\_1.txt
- a\_corriger\_2.txt



Trois fichiers dictionnaires sont fournis :

- dico\_1.dico
- dico\_2.dico
- dico\_3.dico

## Organisation du travail

Les communications et échanges de fichiers ont été sur Discord.

## Comment compiler les programmes ?

Nous avons fourni un makefile qui permet de compiler les deux programmes de ce rendu.

Il suffit de faire :

`make`

Pour supprimer les fichiers objets :

`make clean`

Pour supprimer la totalité des fichiers régénérables :

`make mrproper`

## Comment utiliser les programmes ?

Le premier programme est correcteur\_0.

Celui-ci permet d'afficher les erreurs contenu dans le fichier texte passé en argument selon le dictionnaire aussi passé en argument.

Pour le lancer, il suffit de faire :

```
./correcteur_0 [fichier_texte] [fichier_dictionnaire]
```

Le second affiche les erreurs et aussi des suggestions de correction pour chaque mot erronés.

Pour le lancer il suffit de faire :

```
./correcteur_1 [fichier_texte] [fichier_dictionnaire]
```

Le troisième affiche les erreurs et aussi des suggestions de correction pour chaque mot erronés mais en utilisant des ArbresBK pour faire des recherches approximatives.

Pour lancer la correction il suffit de faire :

```
./correcteur_2 [fichier_texte] [fichier_dictionnaire]
```

Il est possible d'afficher l'arbre BK avec deux modes,

Pour afficher en pdf avec evince + dot (graphviz) :

```
./correcteur_2 -g [fichier_dictionnaire]
```

Pour afficher sur le terminal :

```
./correcteur_2 -a [fichier_dictionnaire]
```

## Problèmes rencontrés

Dans ce rendu, nous avons rencontrés de nombreux problème comme par exemple des erreurs de segmentations, la plupart sont dû à des mauvaises manipulation de pointeurs. Nous avons en général, utilisé lldb / gdb pour retrouver les lignes des erreurs de segmentations.

Un autre problème rencontré, était que le programme pouvait avoir un segmentation fault sur MacOS ARM alors que le programme fonctionnait sur WSL2 (Ubuntu).

Le problème était visiblement lié au fait que l'on faisait deux fopen sur un même fichier sur deux pointeurs différents, nous avons estimé que le comportement sur les manipulations de pointeurs de fichiers entre les deux systèmes sont différents. Pour corriger cela, nous avons simplement gardé un seul pointeur et utilisé rewind() pour remettre à 0 la tête de lecture sur le fichier.

# Addendum

Nous avons renommé les fichiers Main\_0.c / Main\_1.c / Main\_2.c vers correcteur\_0.c / correcteur\_1.c / correcteur\_2.c pour mieux suivre le sujet contrairement aux deux derniers rendus.

Dans le sujet, la détection de mots erronés pour ce 3ème programme aurait pu se faire avec un ATR ou bien un Arbre BK. Étant donné que c'était assez ambiguë nous avons décidé que ce 3ème programme (correcteur\_2) utilisera exclusivement des Arbres BK pour l'implémentation de dictionnaires.

Nous avons choisi le standard C99 seulement parce-que nous avons voulu utiliser des tableaux statique de taille défini par variable. Sinon, hormis cela le programme n'a aucun warning en compilant avec -ansi (C89).

# Conclusion

Comparé aux autres rendu, celui-ci n'était pas nécessairement plus difficile. Il était nécessaire de d'implémenter correctement les algorithmes décrits en cours et de ne pas se perdre dans les utilisations de pointeurs.

Il y a eu un peu d'ambiguïtés dans le sujet pour cette 3 étape mais ce n'était pas très gênant.

L'intérêt des structures de données « complexe » dans les programme était très apparent ici. L'objectif étant de limiter un maximum la complexité des algorithmes utilisés.