

Ce TP se déroule sur deux séances (quatre heures). Les étudiants qui se sentent en difficulté devraient réaliser la section 1, uniquement. Les autres peuvent réaliser la section 2, après s'être assurés qu'ils savent effectivement faire la première partie.

1 Implantation d'un module d'ABR

On veut réaliser un module dédié aux ABR, qui permette de faire fonctionner le programme principal suivant. Les valeurs sont des entiers. Vous pouvez vous inspirer de l'extrait de code disponible sur le site du cours mais il serait préférable que vous refassiez tout vous-même.

```
#include "ABR.h"
#include <stdio.h>

int main ()
{
    struct ABR* racine;
    int z;
    racine = NULL;
    scanf ("%d", &z);
    while (z != -1)
    {
        racine = ajouter_ABR (x, racine);
        afficher_ABR (racine);
        scanf ("%d", &z);
    }
    printf ("la hauteur de l'ABR est %d\n", hauteur_ABR (racine));
    printf ("le nombre de noeuds de l'ABR est %d\n", nombre_noeuds_ABR (racine));

    clear_ABR (racine);
    return 0;
}
```

Question 1. Combien de fichiers doit-on écrire ?

Question 2. On souhaite compiler séparément ce qui peut l'être. Quelles seront les commandes de compilation nécessaires?

Question 3. Écrire le fichier d'entête. Spécifier la structure de données, dans un commentaire, placé dans le fichier.

Question 4. Écrire le fichier source. Lors des essais, commentez, dans le programme principal, les appels aux fonctions que vous n'avez pas encore réalisées.

Question 5. Écrire une fonction qui imprime toutes les valeurs des nœuds, en les indentant en fonction de la profondeur du nœud, dans l'arbre (la valeur de la racine en première colonne, celles des fils de la racine en colonne 4, celles des petits-fils en colonne 8, etc.).

Question 6. Écrire une fonction qui permette de visualiser un ABR avec dot.

Question 7. Visualiser un ABR obtenu en insérant plusieurs fois la même valeur. Que constate-t-on? Corriger ce qui doit l'être.

2 Les codages de Huffman

Le TP porte sur les codages de Huffman, qui constituent une technique de compression de données très utilisée. Dans ce TP, la donnée à compresser est un texte. Du point de vue des structures de données, ce TP nous amène à utiliser des arbres binaires et des files de priorité.

L'idée consiste à coder chaque caractère c d'un texte par une suite de bits, qui dépend du nombre d'occurrences de c dans le texte. Plus le caractère est fréquent, plus la suite de bits est courte. Prenons pour exemple, le texte : « exemple de codage de Huffman », formé de 29 caractères. Le codage de Huffman qui lui est associé est présenté sous la forme d'un arbre binaire, Figure 1.

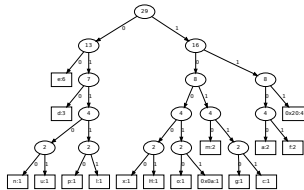


FIGURE 1 – Arbre binaire représentant un codage de Huffman. Chaque feuille est étiquetée par un caractère et son nombre d'occurrences. L'espace et le retour chariot sont représentés par leur code ASCII en hexadécimal. Chaque nœud intérieur est étiqueté par le nombre d'occurrences total des feuilles du sous-arbre dont il est la racine. Les arcs vers les fils gauches sont étiquetés 0; les arcs vers les fils droits sont étiquetés 1.

Pour obtenir sa liste de bits qui code un caractère *c*, il suffit de suivre le chemin qui part de la racine vers la feuille et d'écrire les étiquettes des arcs suivis. Par exemple, le caractère « e », qui apparaît 6 fois dans le texte, est codé par la suite de deux bits « 00 ». Le caractère « B » qui n'apparaît qu'une fois, est codé par la suite de cinq bits « 10001 ». Le texte complet est codé par la concaténation des codages des caractères. Il commence donc par « 001000000 », c'est-à-dire « exe ». Au total, la chaîne se situe sur 107 bits (14 octets). C'est deux fois plus court que les 7 × 29 = 203 bits utilisés par le codage ASCII. Ce codage a de nombreuses propriétés, très intéressantes. Voir [1, chapitre 16.3].

Construction de l'arbre

On se donne une file de priorité F d'arbres de Huffman. Un arbre de Huffman H_1 est plus prioritaire qu'un arbre de Huffman H_2 si le nombre d'occurrences qui étiquette la racine de H_1 est inférieur à celui de H_2 . On présente le présent, caractère par caractère. Pour chaque caractère c , deux cas de figure se présentent : si c est lu pour la première fois, on crée un nouvel arbre de Huffman (une feuille) étiquetée par c et le nombre d'occurrences 1, qu'on enfila dans F ; si c a déjà été lu, on incrémente le nombre d'occurrences de la feuille qui lui correspond (le caractère est forcément présent dans la file F , sous la forme d'une feuille) et on restructure la file, puisque la priorité de c a baissé.

À la fin de la lecture du texte, on a donc une file de priorité, ne comportant que des feuilles. Pour former l'arbre H , il suffit alors d'appliquer l'algorithme de la Figure 2. Une trace d'exécution est donnée Figure 3.

while la file F contient deux arbres ou plus de

$$G := \text{défiler } (F)$$
$$D := \text{défiler}(F)$$

$N :=$ un nouveau nœud avec fils gauche G , fils droit D (opération de *fusion*)

le nombre d'occurrences qui étiquette N doit être égal à la somme de

nombre d'occurrences qui étiquettent G et D

enfiler N dans F

end do

$$H := \text{défiler}(F)$$

FIGURE 2 – Algorithme de construction de l'arbre de Huffman. Initialement, la file de priorité contient les feuilles correspondant aux caractères lus. À la fin, la file ne contient qu'un arbre : l'arbre de Huffman.

Travail à faire

Question 8. On veut écrire un algorithme qui construise l'arbre de Huffman d'un texte et qui imprime le nombre de bits nécessaire au codage de Huffman de ce texte. Pour cela, on demande de mettre au point deux structures de données : une pour les arbres de Huffman et une pour les files de priorité.

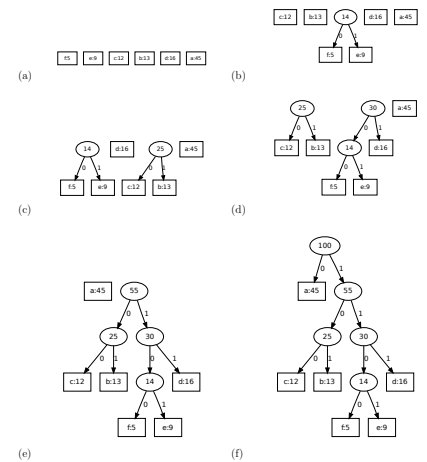


FIGURE 3 – États successifs de la file de priorité, lors de l'exécution de l'algorithme de la Figure 2, sur un exemple. Les éléments de la file sont des arbres. Initialement, la file contient 6 feuilles correspondant aux caractères lus. À la fin, la file ne contient qu'un arbre : l'arbre de Huffman de l'exemple.

Votre programme devrait donc se répartir sur cinq fichiers : deux fichiers d'entête et deux fichiers source correspondant aux deux structures, ainsi qu'un programme principal.

Bien spécifier la structure d'arbre de Huffman.

Question 9. Votre programme fini, comparer le taux de compression que vous obtiendriez avec celui de l'utilitaire `gzip`. Que constatez-vous ? Expliquez rapidement pourquoi en effectuant une recherche sur internet.

Question 10. On souhaite maintenant mettre au point une structure de données permettant d'imprimer les suites de bits, correspondant au codage d'un texte, sur la sortie standard. On ne peut imprimer ces séquences de bits que par paquets de huit, sous la forme d'un caractère. Quelle structure de données vous paraît la plus appropriée ? Spécifiez-la.

Question 11. Déterminer le codage d'un caractère dans un arbre de Huffman n'est pas complètement immédiat. Quelle solution proposez-vous ?

Compression à la volée

A priori, compresser un texte se fait en deux temps : dans un premier temps, on compte les occurrences des caractères ; dans un deuxième temps, on construit le codage et on compresse le texte.

Peut-on compresser un texte à la volée (sans le lire deux fois de suite) avec le codage de Huffman ? Oui. On lit le texte caractère par caractère. Pour chaque caractère c , deux cas de figure se présentent : si c est lu pour la première fois, on l'imprime « en clair » sur la sortie standard, puis on l'insère dans F ; si c a déjà été lu, on construit l'arbre de Huffman correspondant à l'état courant de la file F , on imprime la séquence de bits correspondant à c dans cet arbre, puis on met à jour F en incrémentant le nombre d'occurrences de c et en abaissant sa priorité. À chaque caractère lu, un nouvel arbre de Huffman est créé.

Comment l'algorithme de décodage fait-il pour distinguer les séquences de bits correspondant à un caractère écrit « en clair » des séquences de bits du codage de Huffman ? Il suffit de se donner un caractère spécial (appelons-le NYT) avec un nombre d'occurrences 0, qui n'appartient pas au texte¹, de l'insérer dans F en début d'algorithme, et d'imprimer le codage de Huffman de NYT juste avant les huit bits de c .

Pour pouvoir décoder le texte, l'algorithme de décodage est obligé de mimer le comportement de l'algorithme de codage et de faire évoluer, lui aussi, l'arbre de Huffman du texte, à chaque fois qu'un caractère est lu. Il peut ainsi repérer la séquence de bits codant NYT et lire les huit bits suivants, qui donnent un nouveau caractère « en clair ».

Une variante de cette méthode est connue sous le nom d'algorithme de Vitter [2].

Question 12. Implantez cet algorithme.

1. C'est facile en UTF-8.

Références

- [1] Thomas Cormen, Charles Leiserson, Ronald Rivest, and Clifford Stein. *Introduction à l'algorithmique*. Dunod, Paris, 2ème édition, 2002.
- [2] Jeffrey Scott Vitter. Design and Analysis of Dynamic Huffman Codes. *Journal of the ACM*, 34(4) :825–845, 1987.