

Professor: Fernando V. Paulovich (paulovic at icmc.usp.br)
Aluno PAE: Felipe S. L. G. Duarte (fgduarte at icmc.usp.br)

Trabalho 04: Compactação de Huffman

1 Prazos, Especificações e Observações importantes:

- O trabalho descrito a seguir é individual e não será tolerado qualquer tipo de plágio ou cópia em partes ou totalidade do código. Caso seja detectada alguma irregularidade, os envolvidos serão chamados para conversar com o professor responsável pela disciplina e os trabalhos serão zerados.
- A entrega deverá ser feita única e exclusivamente por meio do *Sistema Run Codes* no endereço eletrônico <https://run.codes/> até o dia **02 de Dezembro de 2016 as 23 horas e 59 minutos**. Sejam responsáveis com o prazo final para entrega, o *Run Codes* está programado para não aceitar submissões após este horário e não será aceito entrega fora do sistema.
- A interpretação desta descrição faz parte do trabalho. Leia a descrição do trabalho com atenção e várias vezes, anotando os pontos principais e as possíveis formas de resolver o problema. Comece a trabalhar o quanto antes para não ficar dúvidas e você consiga entregar o trabalho a tempo.
- O trabalho deverá ser submetido em formato zip/Makefile contendo o arquivo principal, o Makefile e possíveis TAD's implementadas (pares .h .c). Atente para o fato que todos os arquivos enviados deverão conter um cabeçalho contendo nome e número USP.
- A implementação é livre, crie quantos TAD's julgarem necessárias. Importante notar que a modularização do código bem como a forma como as TAD's foram criadas serão levados em consideração na atribuição final da nota.
- Compile o programa em um sistema operacional linux antes de submeter ao *Run Codes*. **Trabalho em que o comando make e make run não funcione, não será corrigido e consequentemente receberá nota zero.** Caso não tenha o sistema operacional linux instalado em sua máquina, aconselha-se a utilização de uma máquina virtual. Como instalar o sistema operacional linux em uma máquina virtual está descrito no seguinte tutorial: <http://maistutoriais.com/2012/01/instalar-maquina-virtual-usando-o-virtual-box/>.
- Referencie, com um comentário no próprio código, qualquer algoritmo ou trecho de código retirado da internet. Código copiado sem a devida referencia é considerado plágio que por sua vez é crime.
- Nenhuma saída do trabalho será em arquivo, ou seja, o resultado dos algoritmos deverá ser exibidos utilizando o stdout `printf`.

2 Restrições:

- O trabalho deverá obrigatoriamente utilizar o arquivo `trab4.c` sem que se tenha qualquer alteração até mesmo no nome do arquivo.
- Este trabalho deverá ser construído utilizando de um Heap Dinâmico. Caso sua implementação esteja correta e utilize Heap estático o seu trabalho receberá somente metade dos pontos (5 pontos). Ainda, caso seu trabalho não utilize qualquer forma de Heap, sua nota máxima é 3
- A implementação do Heap Dinâmico deverá ter a mesma complexidade computacional do Heap estático.

3 Huffman

A codificação de Huffman é um método de compressão que usa as probabilidades de ocorrência dos símbolos no conjunto de dados a ser comprimido para determinar códigos de tamanho variável para cada símbolo.

Para atribuir aos caracteres mais frequentes os códigos binários de menor comprimento, constrói-se uma árvore tomando como base as probabilidades de ocorrência de cada símbolo. Nesta árvore as folhas representam os símbolos presentes nos dados, associados com suas respectivas probabilidades de ocorrência. Os nós intermediários representam a soma das probabilidades de ocorrência de todos os símbolos presentes em suas ramificações e a raiz representa a soma da probabilidade de todos os símbolos no conjunto de dados. O processo se inicia pela junção dos dois símbolos de menor probabilidade, que são então unidos em um nó ao qual é atribuída a soma de suas probabilidades. Este novo nó é então tratado como se fosse uma folha da árvore, isto é, um dos símbolos do alfabeto, e comparado com os demais de acordo com sua probabilidade. O processo se repete até que todos os símbolos estejam unidos sob o nó raiz.

3.1 Proposta

3.1.1 Árvore de frequência

Para determinar o mapa de bit por letra é importante calcular a frequência com que os caracteres aparecem no texto. Em particular, em nosso trabalho, iremos trabalhar com os caracteres de a-z, A-Z, 0-9 e o caractere espaço ‘ ’.

Ao determinar a frequência dos termos, monta-se a árvore de Huffman em que nós a esquerda recebem bit zero e nós a direita bit 1. Assim, neste trabalho, deverá ser construído um Heap Dinâmico baseando-se nas probabilidades de ocorrência de cada símbolo em que as operações tenham a mesma complexidade computacional da implementação de um Heap Estático. Importante ressaltar que esta não é a única forma de implementar o algoritmo de Huffman. Sendo assim, caso sua implementação esteja correta e utilize Heap estático o seu trabalho receberá somente metade dos pontos (5 pontos). Ainda, caso seu trabalho não utilize qualquer forma de Heap, sua nota máxima é 3.

Uma vez que determinou-se quais os dois caracteres vão ser unidos em uma nova árvore, o novo nó terá como frequência a somas dos nós que lhe deram origem e o seu termo é a concatenação dos termos dos nós filhos.

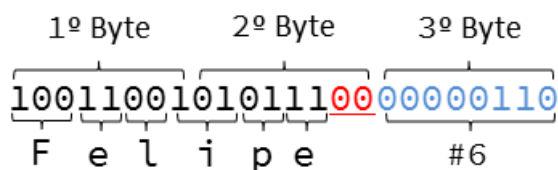
3.1.2 Compactação de Huffman

O sistema aqui implementado deverá ser capaz de compactar um texto plano, de no máximo 10000 (Dez mil) caracteres válidos, utilizando o algoritmo de Huffman. Uma vez que a árvore de frequência foi gerada (Heap Dinâmico), é necessário criar uma tabela que converta cada caractere do texto em um código binário único. Uma importante característica do processo de construção da tabela é que os caracteres com maior frequência recebem o menor número de bits o que garante no final um texto representado com o menor número de bits possível no algoritmo de Huffman.

Essa tabela de símbolo/bits tem também a função de ajudar na descompactação da informação recebida o que torna necessário, ao compactar o texto original, que antes tenhamos a tabela de bits de cada termo. Uma possível abstração para a tabela de símbolo/bits para o texto “Felipe” seria:

F - 100
e - 11
i - 101
l - 00
p - 01

A compactação é feita substituindo-se os caracteres pelos bits correspondentes de acordo com a tabela já calculada. A cada 8 bits concatenados, gera-se 1 byte do texto compactado.



Um problema comum neste tipo de compactação é quando o número total de bits não é múltiplo de 8, ou seja, como somente é possível salvar informações com a unidade de bytes, o ultimo byte do texto pode não ser totalmente utilizado. Com o intuito de resolver tal problema, existem varias formas de informar o fim do arquivo ou o número total de bits utilizado. Em nosso trabalho, utilizaremos o ultimo byte para indicar quantos bits foram usados no ultimo byte do texto compactado.

Na imagem anterior, o texto ‘Felipe’ utilizou somente 6 bits do último byte, desta forma os últimos dois bits foram completados com zero o ultimo byte traz o número 6 que indica quantos bits foram utilizados no byte final do texto.

Segue alguns links que podem ajudar na elaboração do algoritmo de Huffman e na transformação em byte do texto:

- <http://www.programminglogic.com/implementing-huffman-coding-in-c/>
- <http://michael.dipperstein.com/huffman/>

3.1.3 Descompactação de Huffman

A descompactação deverá fazer o serviço contrario da compactação. Este comando deverá ser capaz de processar um documento compactado e com a ajuda da tabela retornar a mensagem original. O resultado da descompactação do texto ‘???’ é o texto original ‘Felipe’.

4 trab4.c

```
1
3 #include <stdio.h>
4 #include <stdlib.h>
5 #include <string.h>
6
7 int main(int argc, char *argv[]) {
8
9     int i;
10    char *text, *textCompac, *textDescompac;
11    text = (char *)malloc(10001 * sizeof(char));
12    TABELA_HUFFMAN *tabela = criarTabelaHuffman(); // cria a tabela de huffman
13
14    // le o texto que sera compactado
15    for(i = 0; (text[i] = getchar()) != '\n' && i < 10000; ++i);
16    text[i] = '\0';
17
18    // compacta o texto e modifica a tabela
19    int sizeCompac = compactarHuffman(text, &textCompac, &tabela);
20    // descompacta o texto baseado na tabela criada
21    textDescompac = descompactarHuffman(textCompac, sizeCompac, tabela);
22
23    // imprime o tamanho do texto original
24    printf("%d\n", (int)strlen(text));
25    // imprime o tamanho do texto compactado
26    printf("%d\n", sizeCompac);
27    // imprime o texto descompactado
28    printf("%s\n", textDescompac);
29
30    // libera a memoria alocada para a tabela de Huffman
31    liberarTabelaHuffman(&tabela);
32    free(text);
33    free(textCompac);
34    free(textDescompac);
35
36    return 0;
37 }
38
39
```

Listing 1: trab4.c source code