

Laborator 13 - Cod Huffman

13.2. Coduri Huffman

Codurile Huffman reprezintă o modalitate optimă de arhivare a unui text dat. Ideea de bază este că fiecare caracter va fi înlocuit cu un cod format din biți de 0 și 1 astfel încât un caracter cu frecvență mare de apariție să aibă un cod cât mai mic posibil față de unul cu frecvență mică. Codurile asociate sunt coduri cu prefix (nici un cod valid nu reprezintă prefixul unui alt cod valid). În codarea Huffman textul dat este analizat pentru a se calcula frecvența de apariție a fiecărui caracter. Se construiește apoi un arbore binar format din noduri frunze (asociate caracterelor) și noduri interne, astfel: se porneste cu n noduri libere (frunze), unde n reprezintă numărul de caractere distincte din text; fiecare nod conține un caracter și frecvența sa de apariție; se creează apoi un nod nou având ca fii 2 noduri cu frecvență minimă de apariție; se continuă analog cu mulțimea de noduri din care s-au eliminat cele două noduri folosite și la care s-a adăugat noul nod creat; procesul se încheie atunci când se rămâne cu un singur nod, acesta devenind și rădăcina unui arbore binar Huffman. Codul unui caracter se calculează parcurgând lanțul de la rădăcina la nodul frunză care îl conține și considerând 0 dacă se merge pe direcția fiului stâng, respectiv 1 pentru fiul drept. Pentru a se dezarhiva un text codat Huffman trebuie să se știe care sunt codurile caracterelor (date de arborele Huffman).

Exemplul 13.2. Pentru textul "*this is an example of a huffman tree*" se calculează următoarele frecvențe de apariție și se creează arborele Huffman de mai jos

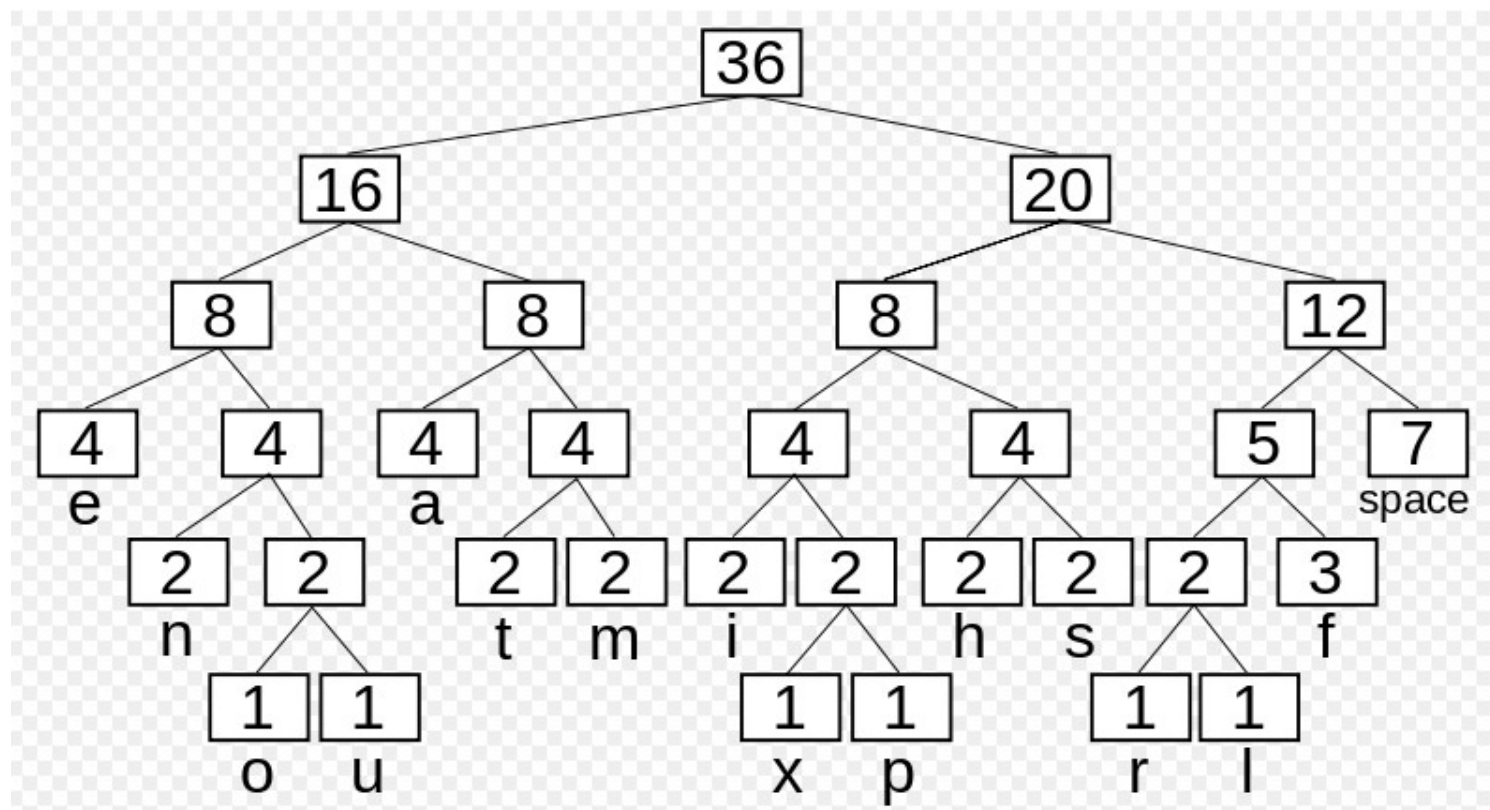
(vezi http://commons.wikimedia.org/wiki/File:Huffman_tree.svg)

Textul este astfel codat și rezultă arhivarea

"01101010100010111111000101111101000101110001001001001111001111001000111001101101111010111010001111101110101110100010111011010001011011011000000000" pe 135 de biți, reprezentând o compresie de 53% față de numărul de 288 de biți folosiți pentru memorarea codurilor ASCII corespunzătoare caracterelor din text (36 de caractere \times 8 biți de caracter), parcurgându-se arborele de la nodul rădăcină la caracter, astfel: 0 pentru subramură stângă, respectiv 1 pentru dreaptă, fără nodul rădăcină.

Char	Freq	Code	Char	Freq	Code	Char	Freq	Code
space	7	111	m	2	0111	p	1	10011
a	4	010	n	2	0010	r	1	11000
e	4	000	s	2	1011	u	1	00111
f	3	1101	t	2	0110	x	1	10010
h	2	1010	l	1	11001			
i	2	1000	o	1	00110			

Arborele rezultat este:

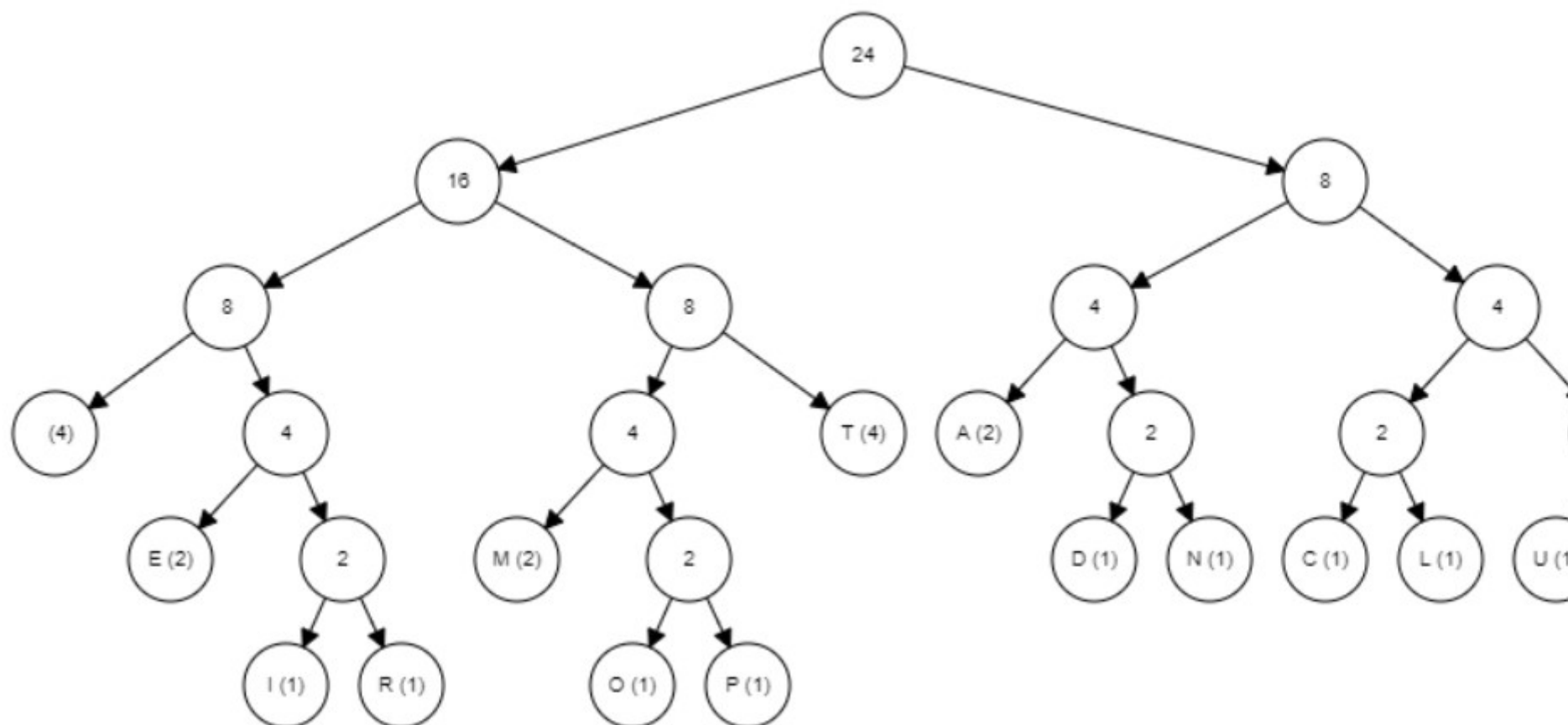


Exercitiul 13.1

Construiti arborele binar asociat codului Huffman pentru textul: *text de comprimat*

Exercitiul 13.2

Pentru arborele Huffman de mai jos:



Determinati forma decompresata a secventei: 11000101001000101100111001100100100011

Exercitiul 13.3 Fisierul [date.in](#) contine un text format numai din literele mari ale alfabetului latin.

1. Scrieti un program C prin care calculati frecventa de aparitie a fiecarui caracter din text continut de fisier

Exercitiul 13.4.*. Construiti apoi codul Huffman si arborele asociat acestuia.

