

# **STRUCTURI DE DATE**

## **TEMA 2**

**Data publicare:** 19 Aprilie 2015

**Deadline:** 10 Mai 2015

Sevastian Emma  
**Facultatea de Automatică și Calculatoare**

## 1. Contextul cerinței

**1.1 Compresia** este procesul de minimizare a spațiului ocupat sau a timpului necesar transmiterii unei anumite cantități de informație.

Termenul de comprimare a datelor a apărut în contextul în care se manifesta o necesitate evidentă de a atinge rate mari de transfer în rețele sau de a stoca o cantitate cât mai mare de informații folosind cât mai puțin spațiu.

Ca o formă primitivă a compresiei de date putem considera prescurtările din viața de zi cu zi. Compresia de date fără pierdere, prezentă în programele de arhivare, în sistemele de transmisie a datelor, a evoluat de-a lungul timpului pornind de la algoritmi simpli (suprimarea zerourilor, codarea pe șiruri) și ajungând la algoritmii complecși folosiți în prezent.

## 1.2 Codificare Huffman

**Codificarea Huffman** se bazează pe frecvența de apariție a simbolurilor. Ideea algoritmului este de a atașa coduri de lungime mică simbolurilor cu frecvență mare, iar celor cu frecvență mică, coduri de lungime mai mare.

**Algoritmul de codificare Huffman** presupune că simbolurile ce doresc a fi codificate fac parte dintr-o mulțime finită și au asociată o probabilitate de apariție. Pentru exemplul de față vom considera cazul șirurilor de caractere și vom explica algoritmul pe baza acestui exemplu.

“ana are mere”

mulțimea simbolurilor  $S = \{ 'a', 'n', ' ', 'r', 'm', 'e' \}$

lungimea șirului inițial = 12

Asociem fiecărui caracter probabilitatea de apariție calculată ca: numărul de repetiții a caracterului respectiv în șirul inițial / lungimea șirului.

Simbol	Probabilitate
'a'	1/4
'n'	1/12
' '	1/6
'r'	1/6
'e'	1/4
'm'	1/12

Bazându-ne pe frecvența de apariție a fiecărui caracter vom crea noua codificare asociată fiecărui caracter.

## Crearea arborelui Huffman

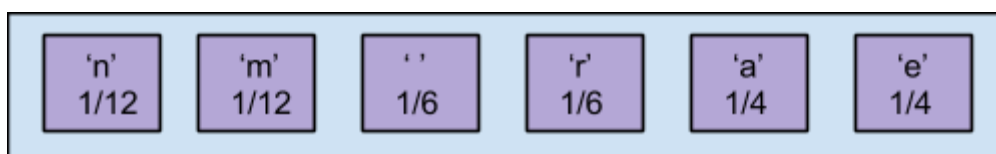
Arborele Huffman este un arbore binar ce conține în fiecare frunză un simbol/caracter alături de probabilitatea sa de apariție. Nodurile interne conțin doar un număr egal cu suma probabilităților de apariție ale descendenților direcți. Plecând de la aceste reguli, rădăcina va avea asociat numărul 1. Suma probabilităților de apariție ale tuturor simbolurilor/caracterelor este 1.

### Pași pentru construirea arborelui:

1. se construiește câte un nod frunză pentru fiecare simbol/caracter. Acest nod va conține simbolul și probabilitatea de apariție asociată.



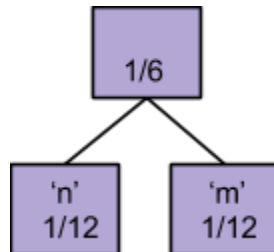
2. nodurile create anterior se adaugă într-o coadă cu priorități, sortată crescător după probabilități, astfel încât primele care vor fi extrase vor fi cele cu probabilitatea cea mai mică.



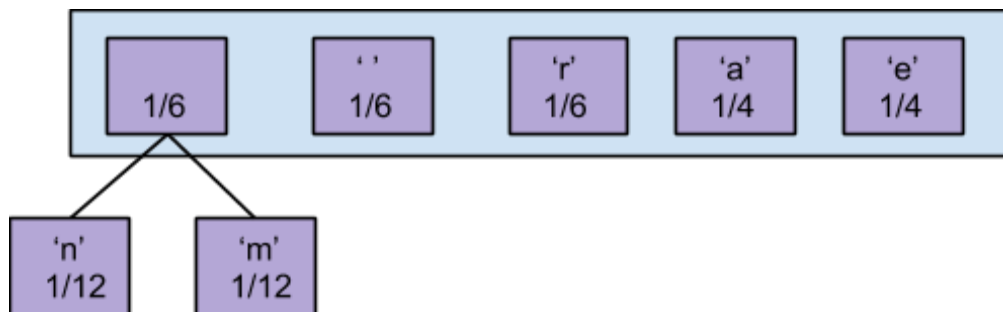
3. cât timp coada conține cel puțin 2 elemente:

**3.1** se extrag primele 2 elemente, cele cu frecvența de apariție cea mai mică, și se creează un nod nou, cu probabilitatea egală cu suma celor două noduri.

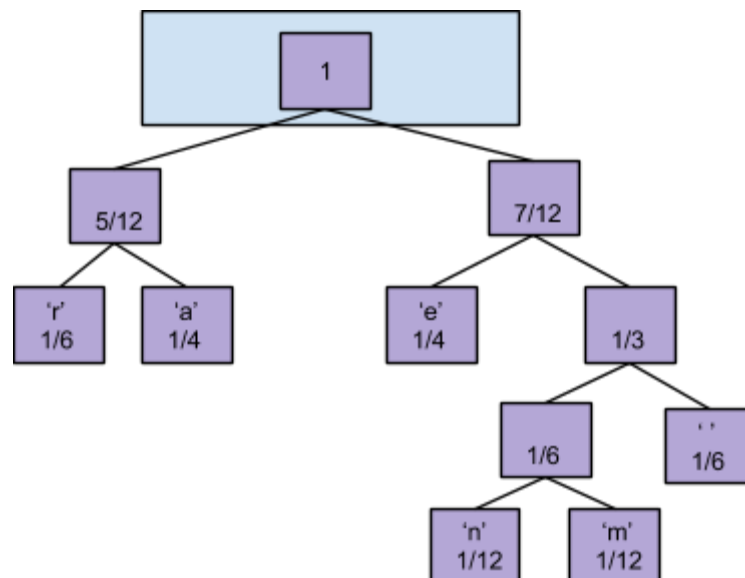
**3.2** noul nod creat va conține cele 2 noduri ca și descendent stâng, respectiv drept.



**3.3** noul nod se inserează înapoi în coadă, cu păstrarea ordinii crescătoare după probabilități

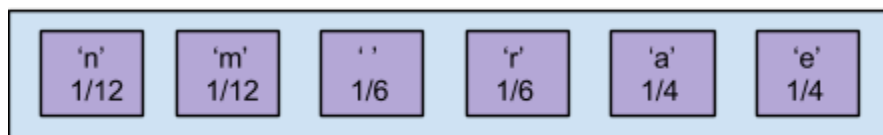


4. la final coada va conține un singur element: rădăcina arborelui Huffman

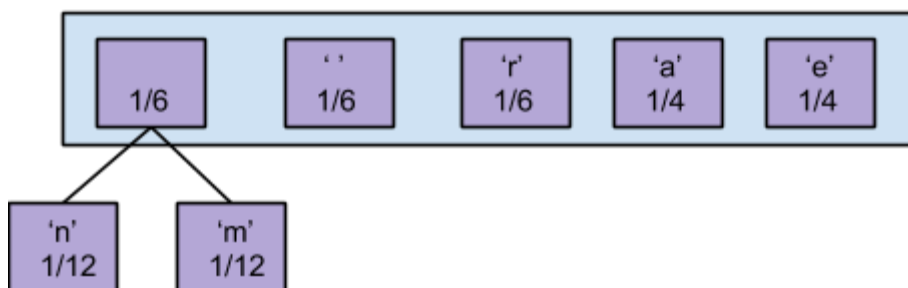


## Procesul de construcție al arborelui ilustrat pe exemplul de mai sus

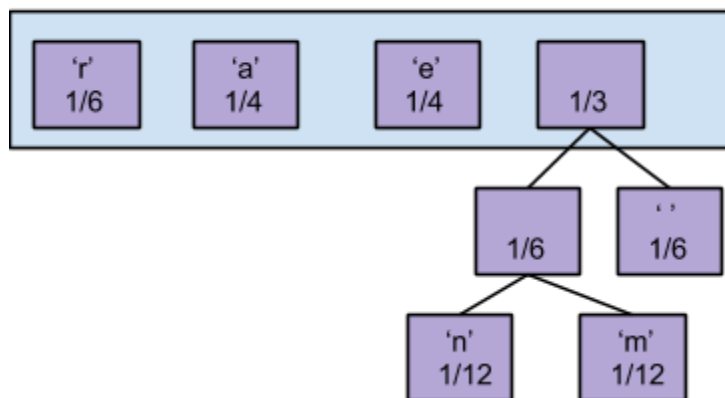
Pas 1



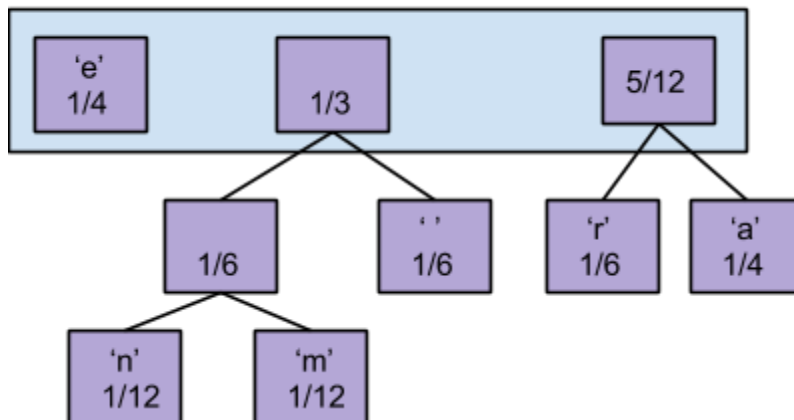
Pas 2



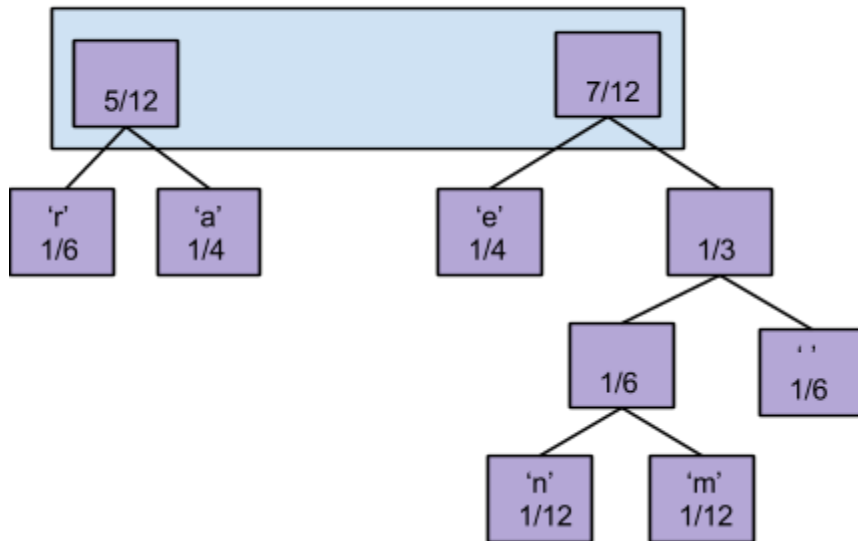
Pas 3



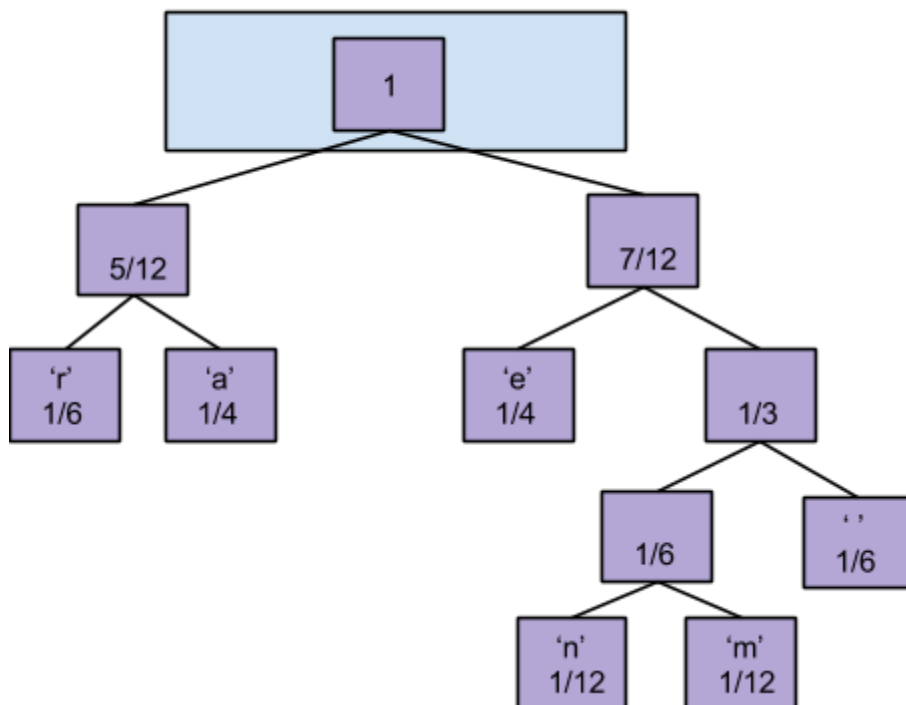
Pas 4



Pas 5

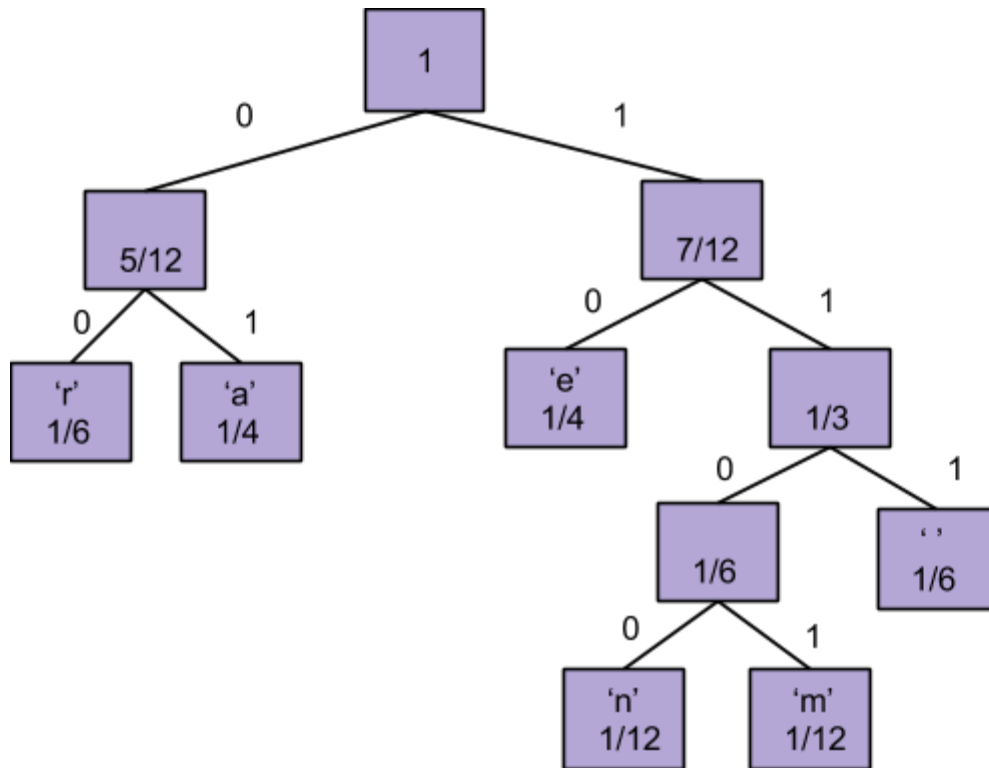


Pas 6



## Codificarea caracterelor

Plecând de la rădăcina arborelui Huffman se poate genera codul fiecărui simbol/caracter astfel: la fiecare pas al parcurgerii, pentru o alegere a nodului stâng se adaugă un 0, pentru o alegere a nodului drept, se adaugă un 1.



Pe baza arborelui creat se generează următoarea codificare a caracterelor:

Simbol	Codificare
'r'	00
'a'	01
'e'	10
' '	111
'n'	1100
'm'	1101

După cum se poate observa, caracterele care au o frecvență mai mare în șirul inițial au o lungime mai mică a codului asociat.

Șirul inițial devine:

'a'	'n'	'a'	' '	'a'	'r'	'e'	' '	'm'	'e'	'r'	'e'
01	1100	01	111	01	00	10	111	1101	10	00	10

Codul ASCII asociat șirului inițial are nevoie de  $12 * 8 = 96$  biți pentru reprezentare, caracterele ASCII fiind reprezentate pe câte un octet (1 octet = 8 biți). Pentru reținerea noului șir este nevoie de 30 de biți. Pentru a putea fi reprezentat în memoria calculatorului însă, este nevoie de un număr de biți divizibil cu 8, astfel se vor adăuga 2 biți ajungându-se la 32.

## 2. Cerință și punctaj

Cerința temei este de a realiza codificarea unui fișier folosind codificarea Huffman. Se consideră că fișierele folosesc encodare extended ASCII (256 de caractere). Detalii la [\[1\]](#), [\[2\]](#).

### Format fișier comprimat

1. Fișierul comprimat va conține la început detaliile caracterelor care au apărut în fișierul inițial astfel:

- `numar_caractere` de tip `uint32_t` ce specifică numărul de caractere din fișier
- `numar_noduri` de tip `uint16_t` ce specifică numărul de noduri ale arborelui Huffman creat
- `numar_noduri` structuri de tipul

```
typedef struct TagHuffmanNode {
    unsigned char value;
    int16_t left;
    int16_t right;
} __attribute__((__packed__)) TagHuffmanNode;
```

Detalii despre `__attribute__` găsiți la [\[3\]](#).



Fiecare element de tip `TagHuffmanNode` va reprezenta un nod al arborelui creat.

**Detalii câmpuri structură:**

- **value:** reprezintă caracterul asociat nodului respectiv (relevant doar în cazul frunzelor)
- **left:** reprezintă indicele în vectorul de structuri a nodului stâng
- **right:** reprezintă indicele în vectorul de structuri a nodului drept

**Detalii noduri:**

- dacă **nodul este frunză** el nu are descendenți, caz în care `left` și `right` vor fi setați la -1, iar `value` va avea valoarea caracterului din nodul respectiv.
- dacă **nodul este unul interior**, el are descendenți, și nu are un caracter asociat. `left` și `right` vor reprezenta indicii nodului stâng, respectiv drept din vectorul de structuri.

**Numerotarea indicilor începe de la 0, pe prima poziție reținându-se rădăcina arborelui Huffman.**

2. Urmează datele efective, cu simbolurile scrise codificate în ordinea în care au fost citite din fișierul inițial (în urma aplicării arborelui Huffman asupra simbolurilor, se obține un șir de biți). În cazul în care lungimea codificării nu este divizibilă cu 8, se va adăuga un padding de biți de 0 până la cea mai apropiată valoare divizibilă cu 8.

**Codificarea se va face pe fiecare octet de la bitul cel mai semnificativ spre cel nesemnificativ.**

**Bonus:** implementarea decompresării.

**Detalii despre pașii urmăți în cazul comprimării și al decompresării îi găsiți în secțiunea Hints.**

**Punctaj:**

- 80 puncte pentru implementarea corectă și completă a comprimării
- 10 puncte coding-style
- 10 puncte README
- 20 puncte pentru implementarea corectă și completă a decompresării

Programul vostru va primi la intrare numele fișierului de intrare și cel de ieșire și o opțiune în felul următor:

```
./huffman -c | -d nume_fisier_intrare nume_fisier_iesire
```

unde

-c indică faptul că programul va realiza o compresie

-d indică faptul că programul va realiza o decompresie

nume\_fisier\_intrare reprezintă numele fișierul sursă, asupra căruia se va opera cu compresia sau cu decompresia

nume\_fisier\_iesire reprezintă fișierul destinație, în care se va scrie, în funcție de comanda dată, fișierul comprimat sau decomprimat.

### 3. Format arhivă și testare

**Nu se vor folosi biblioteci cu cozi, arbori deja implementate.**

**Se vor depuncta pierderile de memorie.**

Temele vor fi testate automat pe vmchecker. Acesta suportă temele rezolvate în C sau C++.

Arhiva cu rezolvarea temei trebuie să fie .zip și să conțină:

- fișierul / fișierele sursă
- Makefile
- README - va conține descrierea soluției. Absența sau descrierea sumară poate duce la depunctari de până la 1 punct.

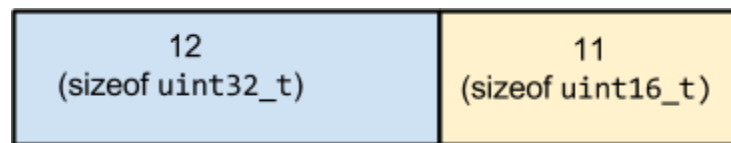
Fișierul pentru make trebuie denumit obligatoriu Makefile și să conțină (cel puțin) următoarele reguli:

- build, care va compila sursele și va obține **executabilul care se va numi huffman**
- clean

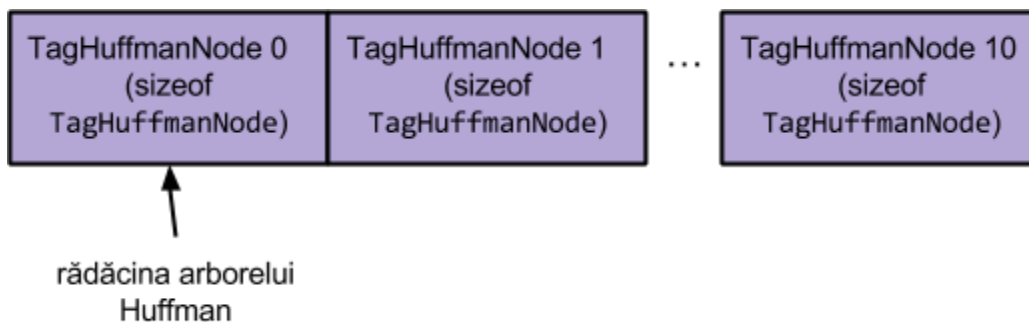
**Temele se vor puncta doar pe vmchecker.**

## 4. Hints și cerințe

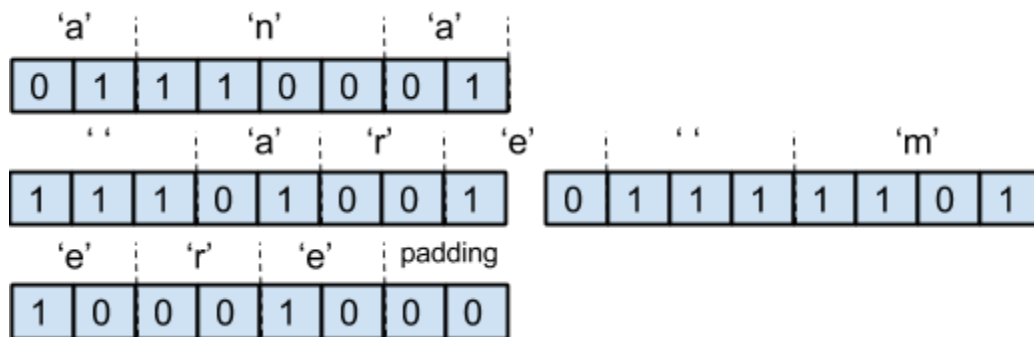
1. Se garantează că fișierele de intrare furnizate pentru compresie sau decompresie sunt corecte.
2. Se garantează că fișierele de intrare au cel puțin 2 caractere distincte.
3. Pași pentru realizarea compresiei:
  - 3.1 se calculează frecvența de apariție a fiecărui caracter  
se reține numărul total de caractere din fișier  
frecvența de apariție a unui caracter este data de  $n_a(c) / n_t$  unde  
 $n_a(c)$  = numărul de apariții a caracterului  $c$ , iar  $n_t$  = numărul total de caractere din fișier.
  - 3.2 arborele Huffman se creează urmând specificațiile prezentate mai sus.
  - 3.3 folosind arborele creat, vom asocia fiecărui caracter un șir de biți.
  - 3.4 scrierea în fișier
4. Exemplu format fișier comprimat ilustrat pentru exemplul de mai sus:
  - inițial fișierul va conține numărul de caractere din fișier (12) urmat de numărul de noduri ale arborelui Huffman creat (11)



- urmat de descrierea arborelui Huffman (pentru acest exemplu 11 structuri de tipul TagHuffmanNode, prima reprezentând rădăcina arborelui Huffman)



- urmat de codificarea fișierul de intrare
  - pe fiecare octet, scrierea codificării se face de la bitul cel mai semnificativ spre cel nesemnificativ.
  - la final, în cazul în care numărul de biți ocupați de codificare nu este divizibil cu 8, se adaugă un padding de 0. Pentru acest exemplul, codificarea șirului (ana are mere) ocupă 30 de biți, adăugând un padding de 0 până la 32 de biți (4 octeți).



## 5. Pași pentru realizarea decompresiei:

**5.1** citirea datelor necesare creării arborelui Huffman

**5.2** parcurgând restul fișierului (asociat fișierului inițial) și ținând cont de codificarea fiecărui caracter, se reface fișierul inițial.

## 5. Precizări

**Tema este individuală**

**Deadline hard pe 10.05.2015, 23:55**

**Nu uitați să eliberați memoria. Pentru a vă testa programul ca să fiți siguri că nu are pierderi de memorie puteți folosi utilitarul numit Valgrind [4].**

## 6. Documentație

[1] <http://www.ascii-code.com/>

[2] [http://en.wikipedia.org/wiki/Extended\\_ASCII](http://en.wikipedia.org/wiki/Extended_ASCII)

[3] <https://gcc.gnu.org/onlinedocs/gcc/Type-Attributes.html>

[4] <http://techblog.rosedu.org/valgrind-introduction.html>