

# Laborator PS

## Algoritmi de compresie a datelor

Lucrare 1: Introducere în Compresia Datelor

Lucrare 2: Algoritmul Shannon-Fano

Lucrare 3: Algoritmul Huffman static

**Lucrare 4: Algoritmul Huffman dinamic**

Prof. Dan Ștefănoiu <[dan.stefanoiu@acse.pub.ro](mailto:dan.stefanoiu@acse.pub.ro)>

Ș.I. Alexandru Dumitrașcu <[alexandru.dumitrascu@acse.pub.ro](mailto:alexandru.dumitrascu@acse.pub.ro)>

## A. Objective

Implementarea Algoritmului de compresie dinamică al lui Huffman. Testarea performanței acestuia, prin comparație cu cele ale Algoritmilor Shannon-Fano, Huffman static și ale programelor de uz general **WinZIP** și **WinRAR**.

## B. Suport teoretic

Performanțele algoritmilor de compresie bazați pe modele statistice de tip static (neadaptiv) sunt limitate. Pentru a le ameliora, în loc să se efectueze construcția modelului statistic la sfârșitul citirii întregului set de date, aceasta se realizează „din mers”, adică în cursul citirii datelor. Au fost elaborate, în acest fel, metode de compresie de tip *adaptiv*, unde modelul statistic suferă o permanentă îmbunătățire, în funcție de noile date citite de pe fluxul de intrare.

În cazul utilizării alfabetului static de ordin 0, arborele Huffman este construit în mod identic, atât în faza de compresie, cât și în cea de decompresie, pe baza setului de contoare asociate simbolilor din alfabet. Alfabetul și contoarele trebuie transmise alături de setul comprimat de date, pentru ca decriptarea să fie exactă.

O primă idee de a construi arborele binar în manieră adaptivă este următoarea: odată cu fiecare simbol citit/decriptat se reconstruiește complet arborele. Această manieră de construcție este, însă, extrem de ineficientă, deoarece efortul de construcție de la pasul anterior se pierde integral. Ca urmare, algoritmul adaptiv astfel proiectat este cronofag, ceea ce face inoportună utilizarea sa într-o aplicație de timp real. Rezultă că strategia de construcție trebuie să se bazeze pe conceptul de *reactualizare*. Aceasta presupune

conservarea unei cât mai mari părți din efortul de construcție a arborelui binar, efort efectuat în cursul citirilor/decriptărilor anterioare. Mai precis, este de dorit ca arborele existent înainte de citirea/decriptarea simbolului curent să fie modificat ulterior, în funcție de informația nouă adusă de acest simbol. Astfel:

- dacă simbolul respectiv apare pentru prima dată, atunci arborelui i se va adăuga un nou nod în zona inferioară (opusă rădăcinii); un fiu al acestui nod este ocupat de noul simbol;
- dacă simbolul există deja în arbore, atunci nodul ocupat de el trebuie reactualizat prin:
  - creșterea cu o unitate a contorului simbolului (a ponderii nodului ocupat de acesta);
  - rearanjarea ramurilor arborelui în funcție de noua structură de contoare, prin schimbări reciproce de locuri între noduri.

În cazul adaptiv, deoarece nu se pot cunoaște dinainte dimensiunile setului de date și ale alfabetului corespondent, nu se mai poate adopta aceeași manieră de construcție. Totuși, se poate păstra ideea care stă la baza relației de ordine dintre noduri, pentru a o redefini, cu scopul utilizării ei în cazul construcției adaptive a arborelui. Această nouă relație de ordine este sintetizată de *proprietatea de fraternitate*.

Prin definiție, un arbore binar cu noduri ponderate posedă *proprietatea de fraternitate* dacă:

- fiii aceluiași nod părinte (numite și *(noduri) frați*) sunt indexați consecutiv;
- ponderea fiecărui nod descrește sau rămâne constantă odată cu creșterea indexului asociat;
- când două noduri ocupă același nivel ierarhic, fiii nodului de indice mai mic au indici inferiori indicilor fiilor celuilalt nod.

Tot prin definiție, un arbore binar cu noduri ponderate având proprietățile:

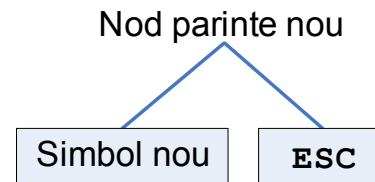
- fraternitate;
- ponderea unui nod părinte este egală cu suma ponderilor fiilor săi direcți;

se numește *arbore Huffman*. Prin Metoda Huffman de tip adaptiv, este construit și reactualizat un arbore Huffman asociat setului de date procesat. Fiecare nod al unui arbore Huffman este caracterizat de

3 parametri: ponderea asociată, poziția sa în cadrul arborelui și tipul său. La reprezentarea în memorie a unui astfel de arbore, acești parametri trebuie luați în considerare, construcția, reactualizarea și parcurgerea arborelui depinzând în mod esențial de ei.

### + Inițializarea arborelui Huffman

Necunoașterea a priori a dimensiunilor setului de date și alfabetului de ordin 0 asociat impune adăugarea la alfabet a 2 simboluri virtuale, diferiți de oricare dintre simbolii ce ar putea apărea în setul de date. Primul dintre ei se numește *simbol de sfârșit* și este notat prin **EOS** (*end of string*). Rolul său este de a marca sfârșitul setului de date (în varianta originală sau comprimată). (Acest simbol poate fi utilizat și în varianta neadaptivă a algoritmului.) Simbolul **EOS** este reprezentat inițial pe doi octeți, având adesea valoarea **256=0x100**. Al doilea simbol virtual se numește *simbol de evitare* și este notat prin **ESC** (*escape*). Rolul său este extrem de important atât în faza de inițializare a arborelui, cât și în cea de reactualizare, în cazul în care un nou simbol (care nu a mai apărut pe fluxul de intrare) este citit/decriptat. Mai precis, **ESC** va reprezenta fratele noului simbol, adică fiul pereche al nodului părinte ce trebuie creat cu ocazia apariției noului simbol, ca în **Figura 3.1**.



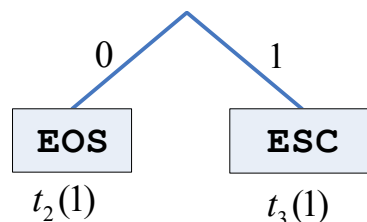
**Figura 4.1. Adăugarea unui nou simbol în arborele Huffman dinamic.**

Codul său original (inițial) este reprezentat tot pe 2 octeți, având valoarea uzuală de **257=0x101**. În cursul compresiei, cei 2 simboluri virtuale primesc coduri diferite de cele originale, în funcție de poziția lor în arbore. Simbolul de evitare (**ESC**) poate apărea de mai multe ori în șirul de date utile comprimate (ori de câte ori codul care îi urmează are 8 biți și descrie un nou simbol original, diferit de toți cei anteriori), în timp ce simbolul de

sfârșit (**EOS**) apare o singură dată (pe ultima poziție a șirului de coduri, marcând sfârșitul acestuia). Datorită acestui fapt, lor li se atribuie câte un contor fix, nemodificabil în cursul compresiei sau al decompresiei, de valoare egală cu 1.

În aceste condiții, arborele inițial (de la care se începe construcția) are doar 3 noduri: rădăcina și două frunze ocupate de simbolii virtuali **EOS** și **ESC** (ca în [Figura 4.2](#)). Alfabetul inițial conține de asemenea doar cei doi simbolii virtuali.

$$\mathcal{A}^{0,2} = \{\mathbf{EOS}, \mathbf{ESC}\}$$



*Figura 4.2. Arborele Huffman dinamic inițial.*

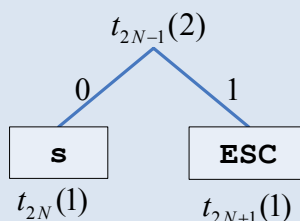
Arborele inițial posedă proprietatea de fraternitate, fiind un arbore Huffman unic determinat. Mai mult, așa cum vom vedea, ultimul nod (de indice maxim) va fi întotdeauna ocupat de simbolul de evitare **ESC**. Noul cod al simbolului de evitare este 1, în această fază inițială. Ordinea celor 2 simbolii virtuali este extrem de importantă aici, deși ea nu corespunde întru totul intuiției (paradoxal, nu simbolul de sfârșit este ultimul din arbore, ci simbolul de evitare). Dacă ei ar figura în ordine inversă, atunci, prin reactualizare, indicele atribuit simbolului de sfârșit va fi mereu modificat, operație evitată în cazul ordinii precizate. De altfel, în [Figura 4.2](#) este respectată ordinea lexicografică pe același nivel ierarhic, așa cum impune definiția arborelui Huffman.

## + Reactualizarea arborelui Huffman în etapa de compresie

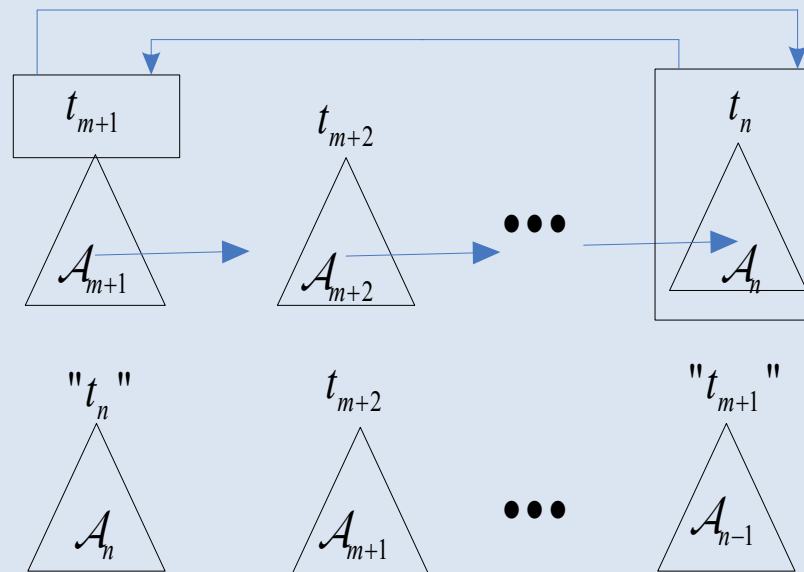
Etapa de reactualizare a arborelui Huffman în faza de compresie este cea mai dificilă și complexă. Aceasta deoarece, chiar dacă este conservat efortul de construcție de la pasul anterior, cei doi arbori Huffman succesivi construiți recursiv pot avea structuri extrem de diferite prin informația transportată de simbolul curent citit/decriptat. Procedura de reactualizare este descrisă de **Algoritmul 4.1**.

### Algoritmul 4.1. Reactualizarea arborelui Huffman pentru compresie.

1. Dacă noul simbol nu se regăsește în alfabetul curent (reprezentat de frunzele arborelui), atunci el trebuie adăugat alfabetului după strategia care urmează:
  - a) Nodul frunză ocupat de simbolul de evitare devine nod intermediar,  $t_{2N-1}$ , iar ponderea sa devine egală cu 2.
  - b) Noul nod intermediar,  $t_{2N-1}$ , devine nod părinte pentru două noduri fii:  $t_{2N}$  (ocupat de simbolul  $s$ ), cu ponderea  $\mathcal{N}(t_{2N})$  (căci  $s$  apare pentru prima dată) și  $t_{2N+1}$  (ocupat de simbolul **ESC**), cu aceeași pondere ( $\mathcal{N}(t_{2N+1}) = \mathcal{N}(t_{2N})$ ).
  - c) Deoarece ponderea nodului  $t_{2N-1}$  s-a modificat, este posibil ca proprietatea de fraternitate a arborelui să se fi pierdut. De aceea, ea trebuie refăcută atât prin mutări de subarbori cât și prin reactualizări ale tuturor contoarelor nodurilor noului arbore. Figura următoare ilustrează mai clar modificarea locală adusă de simbolul  $s$  în acest caz:



2. Dacă noul simbol există deja în alfabetul curent, atunci se apelează la procedura de mai jos:
  - a) Contorul simbolului  $s$  trebuie incrementat cu o unitate. Dacă nodul frunză ocupat de  $s$  este  $t_n$  (cu  $n \in \overline{2, 2N-2}$ ), atunci noul contor asociat va fi egal cu  $\mathcal{N}(t_n) + 1$ .
  - b) Schimbarea ponderii nodului  $t_n$  poate conduce la pierderea proprietății de fraternitate a arborelui.
3. Refacerea proprietății de fraternitate și reactualizarea ponderilor întregului arbore binar. În arborele binar se află un *nod critic*  $t_n$  (cu  $n \in \overline{2, 2N-2}$ ), a cărui pondere a crescut cu o unitate. Pentru refacerea proprietății de fraternitate, trebuie găsit noul loc pe care îl va ocupa nodul critic  $t_n$  împreună cu toți descendenții săi (dacă există) și să se incrementeze cu o unitate ponderea noului său părinte. Acesta va deveni, la rândul său, noul nod critic și astfel se va relua procedura pînă cînd se va ajunge la rădăcină. Reindexarea nodului critic  $t_n$  și a altor noduri în scopul recâștigării fraternității se realizează astfel:
  - a) Plecând de la indicele  $n$ , se caută, în ordinea descrescătoare a indicilor, primul nod cu pondere cel puțin egală valorii  $\mathcal{N}(t_n) + 1$ . Fie acesta nodul  $t_m$ .
  - b) Dacă  $m = n - 1$ , nu mai este necesară mutarea lui  $t_n$  și a descendenților săi, ci se incrementează ponderea părintelui lui  $t_n$ , care devine la rândul său nod critic și algoritmul se reia de la pasul 3.
  - c) Dacă  $m < n - 1$ , între cele două noduri există cel puțin un nod cu pondere strict inferioară ponderii nodului critic. Pentru refacerea proprietății de fraternitate este necesară interschimbarea pozițiilor nodurilor  $t_{m+1}$  și  $t_n$ , împreună cu subarborii lor. Astfel:
    - dacă  $m = n - 2$ , între  $t_m$  și  $t_n$  există un singur nod,  $t_{n-1}$ ;
    - dacă  $m < n - 2$ , între  $t_m$  și  $t_n$  există cel puțin două noduri intermediare ( $t_{m+1}, t_{m+2}, \dots, t_{n-1}$ ); în acest caz, interschimbarea pozițiilor lui  $t_{m+1}$  și  $t_n$  se realizează astfel (conform figurii care urmează):



- $t_n$  se mută împreună cu arborele lui descendent în noua poziție,  $t_{m+1}$ ;
- arborele descendent al lui  $t_{n-1}$  devine noul arbore descendent al lui  $t_{m+1}$ ;
- subarborii nodurilor intermediare își schimbă părintele direct cu o poziție spre stânga;
- în final, nodul  $t_{m+2}$  primește arborele descendent al nodului  $t_{m+1}$ ;
- dacă printre nodurile intermediare se întâlnesc frunze, ele vor fi sărite;
- după interschimbare, ponderea nodului  $t_{m+1}$  se incrementează și astfel devine noul nod critic, algoritmul reluându-se de la pasul 3.

Mutarea arborilor descendenți ai nodurilor în cursul reactualizării *nu necesită* reindexarea nodurilor componente. Se poate demonstra că această proprietate (care simplifică foarte mult algoritmul) este corectă, ea datorându-se proprietății de fraternitate.



## Scalarea contoarelor în cursul reactualizării

În cazul unor șiruri lungi de date destinate compresiei/decriptării, este posibil ca operația de incrementare a contoarelor asociate nodurilor arborelui Huffman să conducă la depășirea numărului maxim ce poate fi reprezentat în memorie. Acest neajuns se poate remedia printr-o (re)scalare a contoarelor asociate nodurilor. Scalarea constă în două operații:

1. Împărțirea la 2 a tuturor contoarelor asociate frunzelor (cu păstrarea intactă a contoarelor unitare).
2. Refacerea integrală a structurii arborelui HUFFMAN, în manieră statică, respectând proprietatea de fraternitate.

Împărțirea la 2 poate conduce la schimbarea fundamentală a formei arborelui, deoarece această operație poate degrada frecvența reală de apariție a simbolilor pe fluxul de intrare.

Să considerăm arborele Huffman din stânga [Figurii 4.3](#), surprins înainte de scalarea contoarelor sale. După scalare, noul arbore nu mai verifică proprietatea de fraternitate. De aceea, el trebuie reactualizat. Se poate constata ușor că forma finală corectă (ilustrată în dreapta [Figurii 4.3](#)) este extrem de diferită de cea a arborelui inițial, existent înaintea rescalării.

Comparând cei doi arbori Huffman din acest exemplu, se constată că a avut loc o transformare importantă a lungimilor codurilor asociate frunzelor. Astfel, lungimea maximă a codului, de 4 biți pentru simbolii  $s_4$  și **ESC**, a scăzut cu o unitate, lungimile de cod ale simbolilor  $s_1$ ,  $s_2$  și **EOS** s-au conservat, iar lungimea de cod a simbolului  $s_3$  a crescut cu o unitate.

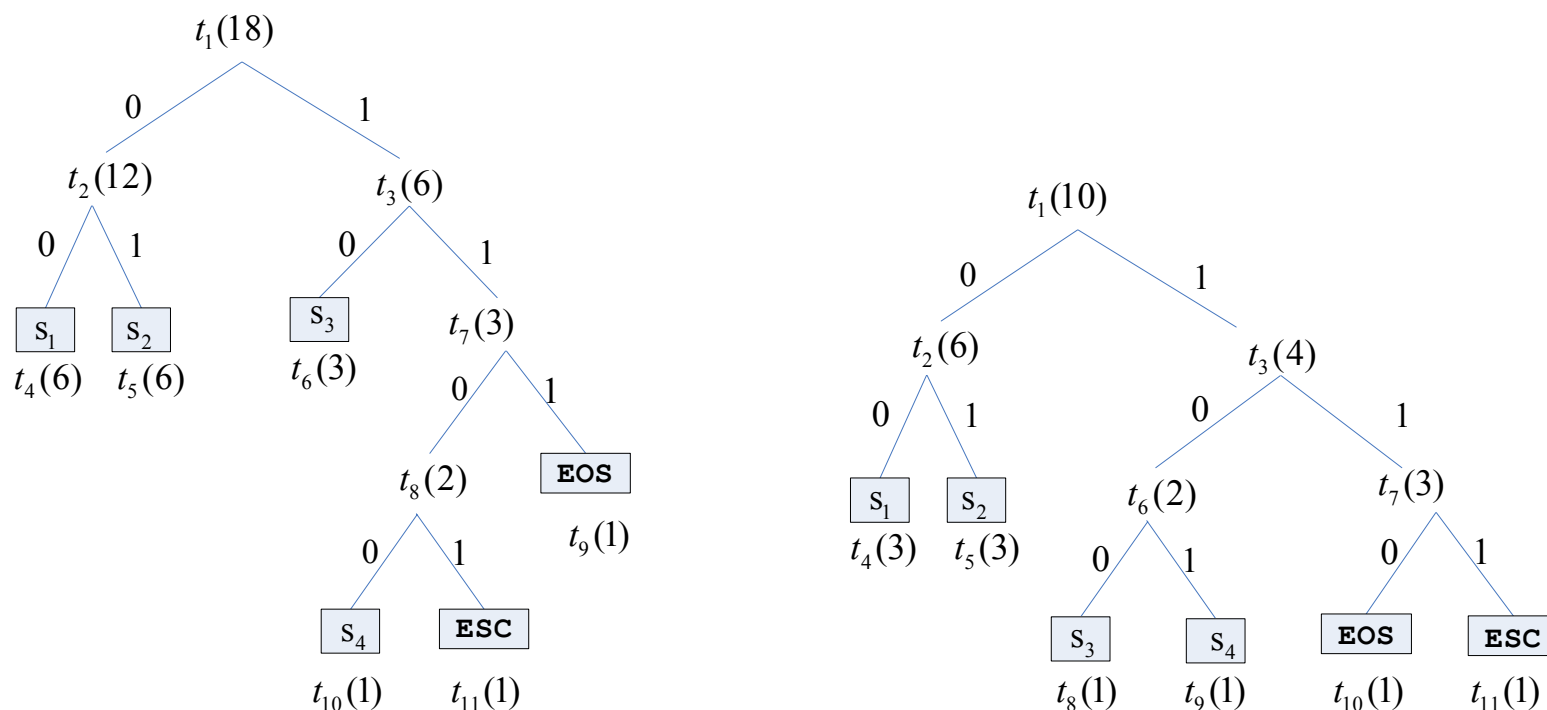


Figura 4.3. Un exemplu de arbore Huffman dinamic înainte de (stânga) și după (dreapta) scalarea contoarelor aferente.

Așadar, se poate considera că scalarea contoarelor a adus un câștig în materie de compresie (codul cel mai lung are acum doar 3 biți), deși arborele Huffman final este foarte diferit de cel inițial. Totodată, aceste modificări ale lungimilor de cod pun în evidență și un mecanism autocompensator de eroare (unele lungimi scad, altele cresc).

Așa cum reiese și din exemplul prezentat, scalarea contoarelor este o operație necesară pentru a evita blocarea algoritmului. Realizarea ei presupune, însă, stabilirea unui criteriu care să definească precis momentul când trebuie declanșată scalarea. Obiectivul urmărit este dublu: prevenirea depășirii capacității de reprezentare a contoarelor și evitarea emisiei de coduri a căror lungime depășește o valoare dată (chiar dacă nu a fost depășită capacitatea de reprezentare a contoarelor).

## Reactualizarea arborelui Huffman în etapa de decompresie

Spre deosebire de construcția adaptivă a unui cod, decriptarea sa adaptivă este o operație mult mai simplă, ea presupunând doar o parcurgere descendentă a arborelui. Vom prezenta în continuare algoritmul asociat acestei operații.

### Algoritmul 4.2. *Reactualizarea arborelui Huffman pentru decompresie.*

1. Se citește fluxul de intrare la nivel de bit (în mod normal, se citește câte un octet, în cadrul căruia se analizează fiecare bit). În paralel, se construiește și se reactualizează arborele Huffman întocmai ca în faza de compresie. Inițial, arborele are forma din **Figura 4.2**. Primul simbol decriptat va fi **EOS** (cod 0), dacă fluxul de intrare este vid, sau **EOS** (cod 1), în caz contrar. Simbolul de evitare va fi urmat de codul original (pe 8 biți) al primului simbol valid din setul de date.
2. Se decriptează simbolul curent de pe fluxul de intrare, prin parcurgerea arborelui Huffman existent, începând cu rădăcina și terminând cu una dintre frunzele sale. Parcurgerea arborelui este ghidată de succesiunea biților fluxului de intrare:
  - a) Dacă simbolul decriptat este **EOS**, atunci se citește următorul octet de pe fluxul de intrare. Acesta va fi codul original al simbolului curent.
  - b) Dacă simbolul decriptat este diferit de **EOS**, atunci simbolul a mai apărut pe fluxul de intrare, frunza ocupată de el în arborele Huffman oferind codul original (pe 8 biți) al acestuia.În ambele situații, codul original al simbolului trebuie înscris pe fluxul de ieșire. Dacă ultimul simbol decriptat este **EOS**, algoritmul se încheie. În caz contrar, se trece la pasul următor.
3. Se reactualizează arborele HUFFMAN utilizând **Algoritmul 4.1** nemodificat, deoarece, în această fază, nu mai este necesară o tabelă de coduri. Algoritmul se reia de la pasul 2.

Exemplul care urmează are rolul de a pune în evidență câteva aspecte importante pentru algoritmul de generare și manevrare a codurilor (mai ales în faza de compresie). Înainte de a-l prezenta, sunt necesare două remarci:

- ori de câte ori simbolul citit nu se regăsește printre frunzele arborelui binar, sunt emise două coduri succesive: primul este cel al simbolului de evitare **ESC** (cod preluat din arborele Huffman înainte de reactualizare), iar al doilea este codul original pe 8 biți (nemodificat) al simbolului;
- odată cu mutarea unui nod în altă poziție (în cursul reactualizării arborelui), este posibil ca unul sau mai mulți arbori descendenți să fie, la rândul lor, mutați în alte poziții.

Reactualizarea codurilor frunzelor constă în reactualizarea părții mai semnificative a vechiului cod, deoarece partea mai puțin semnificativă nu este afectată de mutare.

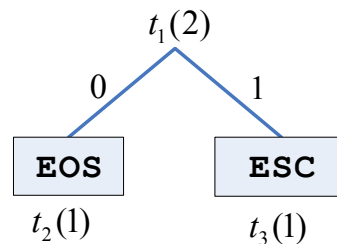
### Exemplu

Se consideră setul de date format din 4 simboluri, primii 3 fiind diferiți:

$\mathcal{D}$ : ADIA

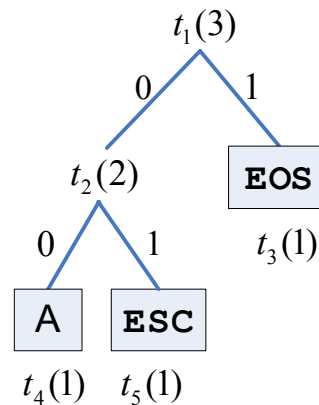
Se dorește compresia adaptivă de tip Huffman a acestuia, urmată de decompresia adaptivă.

Pentru **compresie**, se pleacă de la arborele Huffman inițial din **Figura 4.2**, redat mai jos pentru ușurință.

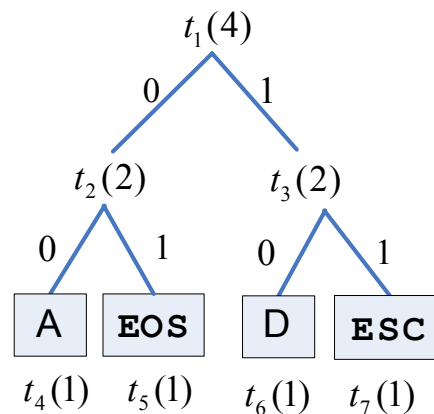


Evoluția ulterioară a acestuia este redată în continuare.

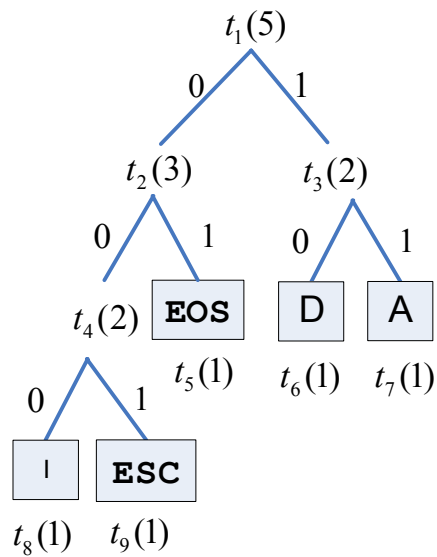
Se citește simbolul 'A' și se emit două coduri (care însumează 9 biți): 1 (**ESC**) și 01000001 (A). Arborele aferent devine:



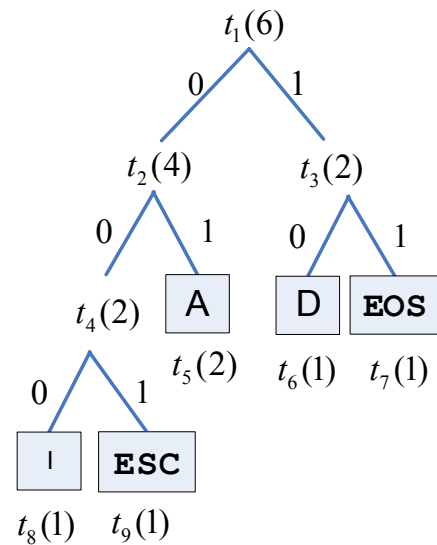
Se citește simbolul 'D' și se emit două coduri (care însumează 10 biți): 01 (**ESC**) și 01000100 (D). Aceasta conduce la modificarea arborelui ca mai jos:



Se citește simbolul 'I' și se emit două coduri (care însumează tot 10 biți): 11 (**ESC**) și 01001001 (I). Arborele se extinde către nivelul următor:



În fine, se citește simbolul 'A' și se emite codul (de numai 2 biți): 11. Acesta înlocuiește codul original pe 8 biți al simbolului 'A'. Arborele Huffman final arată astfel:



Șirul original de simbolii este reprezentat pe cel puțin 41 de biți (32 de biți ai caracterelor, plus cel puțin 9 biți care indică sfârșitul mesajului). Fluxul de date comprimate, care va fi utilizat în procedura de decompresie este constituit din următoarele coduri, numărând 33 de biți:

1 01000001 01 01000100 11 01001001 11 11

Deși setul de date este extrem de scurt, Algoritmul Huffman dinamic reușește să realizeze compresia (pe seama faptului că unul dintre simbolii se repetă). Această caracteristică este generală (cu condiția să existe repetiții de simbolii) și situează algoritmul în rândul celor mai eficiente rețete de compresie a datelor.

Pentru a realiza **decompresia (decriptarea)** șirului de biți de mai sus, se procedează astfel:

- Se citește bitul de valoare 1 și se decriptează simbolul **ESC**, folosind arborele Huffman inițial. Codul său original nu este emis, însă, pe fluxul de ieșire. Operația care se execută în continuare este citirea următorului octet de pe fluxul de intrare: 01000001 (litera 'A'), care este înscris nemodificat pe fluxul de ieșire.
- Se reactualizează arborele Huffman și se resetează parcurgerea sa din rădăcină.
- Se citesc biții succesivi de valori 01 (conform noului arbore) și se constată că simbolul decriptat este tot **ESC**. În consecință, următorul octet al fluxului de intrare, 01000100 (litera 'D'), este înscris nemodificat pe fluxul de ieșire.
- Urmează o nouă reactualizare a arborelui.
- Cum următorii 2 biți citați, 11, conduc tot la simbolul de evitare, următorul octet citit, 01001001 (litera 'I') este și el descărcat direct pe fluxul de ieșire.
- Se reactualizează arborele, care va avea acum 5 frunze.

- Se reia citirea fluxului de intrare. După codul 11, se constată că, folosind arborele Huffman curent, s-a ajuns la nodul frunză ocupat de litera 'A'. Codul original al acesteia (01000001) este preluat din informația de nod a arborelui și descărcat pe fluxul de ieșire.
- O nouă reactualizare a arborelui binar nu aduce noi frunze acestuia, ci o organizare diferită.
- Ultimii 2 biți citați, aceiași 11, conduc la decriptarea simbolului **EOS**, care arată că decompresia a luat sfârșit și provoacă închiderea fluxului de ieșire.

Din acest exemplu, se observă clar cum un același cod (11) poate avea semnificații total diferite la momente de timp diferite (**ESC**, apoi 'A' și, în fine, **EOS**). Cu toate acestea, decriptarea este exactă, deoarece structura arborelui Huffman este precisă, unic determinată de fiecare simbol decriptat și identică în etapele de compresie-decompresie.

## C. Sarcini de lucru

### Tema 1 (Implementarea algoritmilor Huffman dinamici)

- Scriveți rutinele corespunzătoare **Algoritmilor 4.1, 4.2** și metodelor de compresie-decompresie din Metoda Huffman statică.
- Verificați corectitudinea lor cu ajutorul unui program de test aplicat exemplului din secțiunile precedente, adică setului de date:

***D*** : IT\_IS\_BETTER\_LATER\_THAN\_NEVER.

(textul trebuie să fie reconstituit perfect).

- Evaluați factorul de compresie realizat în acest exemplu (ținând cont și de informația auxiliară). Este capabilă metoda să realizeze compresia?



## Tema 2 (Testarea algoritmilor Huffman dinamici)

- a. Scrieți un program care realizează compresia și decompresia celor 12 fișiere din corpusul de date descris în Lucrarea 1, cu ajutorul rutinelor de la tema precedentă.
- b. Calculați norma diferenței dintre setul inițial de date și cel decomprimat. Pentru fiecare set de date, indicați factorul de compresie aferent (incluzând informația auxiliară) și evaluați duratele de compresie, respectiv decompresie. Comentați rezultatul obținut, cu referire la capacitatea de compresie a metodei pentru diferite tipuri de fișiere.
- c. Comparați factorii de compresie de la punctul precedent cu cei obținuți prin utilizarea Metodelor Shannon-Fano, Huffman static și a celor două programe de uz general **WinZIP** și **WinRAR**. (Utilizați funcția **fread** pentru a evalua exact numărul de biți ai arhivelor produse de aceste utilitare.) Comentați rezultatele obținute, cu referire la compromisul dintre complexitatea algoritmilor de compresie și performanțele acestora.