

## MODULUL 9 - TABELE DE REPARTIZARE

Așa cum s-a menționat în primul capitol al acestui curs, structurile de date au fost dezvoltate cu scopul memorării și manipulării *eficiente* a unei mulțimi dinamice de date. Această manipulare eficientă depinde desigur semnificativ de scopul urmărit și de tipurile de operații care se doresc a fi efectuate.

Una dintre principalele operații pe un set de date este căutarea sau *accesul prin cheie*. Am văzut în capitolele precedente că problema căutării poate fi rezolvată în complexitate logaritmică prin utilizarea arborilor de căutare echilibrați.

Există însă și structuri de date care permit accesul prin cheie în timp constant, atunci când nu sunt necesare operații de sortare și păstrare ordonată a setului de date. Aceste structuri sunt *tabelele de repartizare*.

O tabelă de repartizare - *hash table* - este de fapt o generalizare a noțiunii de vector - *array*. Operațiile de bază sunt inserția, căutarea și eventual ștergerea. Tabelele de repartizare pot fi utilizate cu succes pentru implementarea de dicționare, în care căutarea este cea mai frecventă operație. O altă utilizare a tabelor de repartizare este în cadrul compilatoarelor pentru implementarea tabeli de identificatori sau pentru cuvintele cheie. Vom vedea pe parcurs că, într-o tabelă de repartizate operațiile au în medie complexitatea  $O(1)$ . O tabelă de repartizare se memorează într-un tablou liniar.

### 1 Tabele cu adresare directă

Considerăm  $U = \{0, 1, \dots, m-1\}$  universul cheilor posibile,  $m$  relativ mic și  $T[0, \dots, m-1]$  tabloul în care se memorează elementele din tabelă. Elementul cu cheia  $k$  se plasează în  $T$  pe poziția  $T[k]$ . Dacă în tabelă nu există cheia  $k$ , atunci  $T[k] = NULL$ .

**Observație:** o astfel de tabelă de adresare directă este potrivită, atunci când universul  $U$  al cheilor este relativ redus și numărul de elemente memorate în tabelă este comparabil cu  $m$ .

### 2 Tabele de repartizare - *Hash Tables*

Pentru situații în care numărul de elemente stocate este mult mai mic decât universul cheilor se recomandă înlocuirea tabeli cu adresare directă printr-o tabelă de repartizare.

În cazul unei tabele de repartizare  $T[0, \dots, m-1]$  accesul la un element se face prin intermediul unei funcții de repartizare - *hash function* -  $h : U \rightarrow \{0, 1, \dots, m-1\}$ , în care de obicei  $|U|$  este semnificativ mai mare decât  $m$ .

Elementul cu cheia  $k$  va fi plasat în tabelă pe poziția  $T[h(k)]$ . Spunem că, elementul cu cheia  $k$  este repartizat pe poziția  $h(k)$  (*hashes to slot  $h(k)$* ).

### Exemplu de funcție de repartizare:

$$h : U \rightarrow \{0, 1, \dots, m-1\}, h(k) = k \bmod m$$

Considerând  $m = 11$  și cheile  $\{3, 12, 15, 17\}$  atunci:  $h(3) = 3$ ,  $h(12) = 1$ ,  $h(15) = 4$ ,  $h(17) = 6$  (fig. 2).

*	12	*	3	15	*	17	*	*	*	*
0	1	2	3	4	5	6	7	8	9	10

Figure 1: Tabela de repartizare cu funcția de repartizare  $h(k) = k \bmod 11$  în care au fost plasate cheile 3, 12, 15 și 17. Au fost marcate cu \* pozițiile neocupate din tabelă.

### Coliziunea

În cazul unei funcții de repartizare  $h$  este posibil ca pentru două chei diferite  $k_1$  și  $k_2$  să se obțină  $h(k_1) = h(k_2)$ , adică ambele elemente ar fi repartizate pe aceeași poziție. O astfel de situație se numește *coliziune*. În exemplul de mai sus  $h(24) = h(35) = 2$ .

Se pune problema rezolvării coliziunilor. Ideal ar fi, evitarea completă a acestora. Din faptul că se presupune că  $|U| > m$ , prin tipul de repartizare descris mai sus, acest lucru nu este posibil. Totuși, prin construirea atentă a funcției de repartizare, poate fi redusă semnificativ probabilitatea unei coliziuni. Vor fi discutate pe parcursul acestui curs câteva metode de construcție a unei funcții de repartizare, care să îndeplinească acest lucru.

Rezolvarea problemei coliziunilor se poate realiza dacă în loc de a stoca într-o poziție a tabelii de repartizare un singur element, se păstrează o listă de elemente (*bucket*). Lista de la poziția  $h(k)$  va conține toate elementele care au această valoare de repartizare.

**Exemplu:** Considerăm  $h(k) = k \bmod 11$  și cheile  $\{23, 34, 78, 15, 37, 42, 45\}$ .  $h(23) = 1$ ,  $h(34) = 1$ ,  $h(78) = 1$ ,  $h(15) = 4$ ,  $h(37) = 4$ ,  $h(42) = 9$ ,  $h(45) = 1$

### Operațiile uzuale:

```
TAB_INSERT( $T, x$ )
    Insereaza  $x$  în capul listei  $T[h(x.cheie)]$ 
RETURN
```

```

TAB_CAUT( $T, k$ )
    Cauta elementul cu cheia  $k$  in lista  $T[h(k)]$ 
RETURN

TAB_STERG( $T, x$ )
    sterge  $x$  din lista  $T[h(x.cheie)]$ 
RETURN

```

O tabelă de repartizare cu dimensiunea  $m = 11$  care utilizează liste înlanțuite și în care au fost inserate cheile 33, 35, 55, 46, 12, 2, 7, 10, 11 este reprezentată în figura 2.

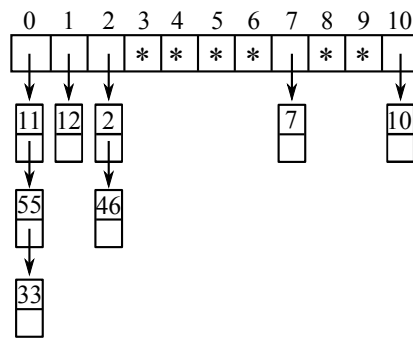


Figure 2: Tabela de repartizare cu funcția de repartizare  $h(k) = k \bmod 11$  și utilizând liste înlanțuite în care au fost plasate succesiv cheile 33, 35, 55, 46, 12, 2, 7, 10 și 11. Au fost marcate cu \* pozițiile neocupate din tabelă.

### Complexitate:

- funcția de inserție are evident complexitatea  $O(1)$ ,
- la funcția de căutare complexitatea depinde de lungimea listelor înlanțuite
- funcția de ștergere are complexitatea  $O(1)$  dacă lista e dublu înlanțuită, altfel depinde de lungimea listei.

### Analiza complexității

Considerăm tabela de repartizare  $T$  cu  $m$  poziții, în care se stochează  $n$  elemente, folosind liste înlanțuite pentru rezolvarea coliziunilor. Atunci numărul mediu de elemente stocate într-o listă este  $\alpha = n/m$ , unde  $\alpha$  poate fi mai mic sau mai mare decât 1.

Cea mai defavorabilă situație este aceea când toate cele  $n$  elemente se repartizează prin  $h$  pe aceeași poziție, adică sunt stocate în aceeași listă. În această situație complexitatea la căutare este  $O(n)$ . Pentru a evita astfel de situații este necesară alegerea unei funcții de repartizare potrivită. Modul de construcție a unor funcții de repartizare care să permită o repartizare cât mai uniformă va fi discutată ulterior.

Repartizarea uniformă presupune ca probabilitatea ca un anumit element să fie de repartizat pe oricare dintre pozițiile din tabelă este aceeași.

Considerând o funcție de repartizare care repartizează cheile uniform în tabela  $T$  și pentru care calculul lui  $h(k)$  este  $O(1)$ , se demonstrează faptul că funcția de căutare a unei chei este  $O(1 + \alpha)$  (vezi bibliografia recomandată).

Dacă  $n$  este proporțional cu  $m$ , atunci  $n = am$ , deci complexitatea la căutare este  $O(1 + am/m) = O(1 + a) = O(1)$ .

### 3 Metode de repartizare

Din discuția complexității rezultă că, o funcție de repartizare bună permite o repartizare aproape uniformă în tabelă. Pentru acest lucru poziția obținută la repartizare ar trebui să fie independentă de orice **pattern** care ar putea fi prezent în setul de date.

În plus, funcțiile de repartizare au ca mulțime a valorilor poziții în tabela de repartizare, deci elemente din mulțimea numerelor naturale. În cazul în care cheile sunt de exemplu șiruri de caractere - de exemplu **hash tables** utilizate pentru tabela de simboluri din faza de analiză lexicală a unui compilator - este necesară stabilirea unei metode de interpretare a acestora ca numere naturale, trebuie deci determinată o funcție de la universul cheilor către mulțimea numerelor naturale.

În continuare sunt prezentate câteva metode pentru construcția funcțiilor de repartizare.

#### 3.1 Metoda diviziunii

Funcția de repartizare  $h : U \rightarrow \{0, 1, \dots, m-1\}$ ,  $h(k) = k \bmod m$

**Exemplu:**  $h(k) = k \bmod 11$ .

**Observație:** o alegere bună pentru  $m$  este un număr prim nu prea apropiat de o putere a lui 2. În plus acest număr trebuie ales depinzând de numărul de elemente care se estimează că vor fi introduse în tabelă, precum și de factorul de încărcare dorit.

**Exemplu:** dacă  $n = 3000$  și numărul mediu de elemente per poziție este considerat 4 atunci,  $\alpha = n/m \Rightarrow m = n/\alpha = 3000/4 = 750$ . Pot considera  $m = 751$ , care este un număr prim nu prea apropiat de o putere a lui 2.

#### 3.2 Metoda multiplicării

În acest caz funcția de repartizare este de tipul:

$$h(k) = \lfloor m(kA - \lfloor kA \rfloor) \rfloor, 0 < A < 1$$

unde  $\lfloor x \rfloor$  reprezintă partea întreagă a lui  $x$ .

În cazul acestor funcții modul de alegere a lui  $m$  nu influențează modul de repartizare. Se demonstrează (Knuth) faptul că, o alegere bună pentru  $A$  este

$$A = (\sqrt{5} - 1) / 2 \approx 0.618033$$

### 3.3 Repartizarea universală

Dacă repartizarea în tabelă este realizată printr-o funcție dată, există posibilitatea ca, pentru anumite seturi de date, toate cheile să fie repartizate pe aceeași poziție, ajungându-se din nou la complexitatea  $O(n)$  la căutare.

Acest lucru poate fi evitat, dacă în loc să se folosească o singură funcție de repartizare, se utilizează o mulțime  $H$  de funcții de repartizare. La fiecare execuție se selectează în mod aleator din  $H$  una dintre funcțiile de repartizare, care va fi apoi utilizată.

Utilizare repartizării universale are ca efect faptul că, pentru același set de date de intrare, execuții diferite ale programului produc în general tabele diferite, astfel putându-se evita în orice situație cel mai defavorabil caz, eventual prin reluare a repartizării - **rehashing**.

**Definiție:** O mulțime finită de funcții de repartizare  $H$ , care repartizează cheile dintr-un univers  $U$  într-o tabelă  $T[0, \dots, m-1]$  se numește universală, dacă pentru orice două chei  $k$  și  $l$ ,  $k \neq l$ , numărul de funcții din  $H$  pentru care  $h(k) = h(l)$  este cel mult  $|H|/m$ .

Acest lucru asigură faptul că, pentru o funcție aleasă în mod aleator din  $H$ , probabilitatea unei coliziuni între  $k$  și  $l$  este  $1/m$ . (Evident, probabilitatea alegerii oricărei funcții din  $H$  este  $1/|H|$ . Probabilitatea alegerii unei funcții care repartizează  $k$  și  $l$  pe aceeași poziție este  $(|H|/m) * 1/|H| = 1/m$ ).

## Construcție a unei clase universale de funcții de repartizare

Se alege un număr prim  $p$  relativ mare,  $p \gg m$ . Pentru fiecare  $a \in \{1, \dots, p-1\}$  și  $b \in \{0, 1, \dots, p-1\}$  se definește funcția de repartizare:

$$h_{ab}(k) = ((ak + b) \bmod p) \bmod m$$

Mulțimea de astfel de funcții de repartizare este

$$H_{pm} = \{h_{ab} | a \in \{1, \dots, p-1\} \text{ și } b \in \{0, 1, \dots, p-1\}\}$$

Se poate demonstra că mulțimea astfel definită este universală (vezi bibliografia recomandată).

## 3.4 Adresare deschisă

În cazul adresării deschise, fiecare poziție a tabelului de repartizare  $T$  conține un singur element. Pentru rezolvarea coliziunilor la inserție nu se utilizează liste înlanțuite, ci se testează diferite poziții, obținute pe baza funcției de repartizare, până când se determină o poziție liberă, pe care poate fi inserat elementul dorit.

La căutare de asemenea se testează diferite poziții până când găsește elementul căutat sau se ajunge la o poziție neocupată.

Faptul că nu se utilizează liste înlanțuite, ci fiecare poziție conține cel mult un element, are ca urmare posibilitatea umplerii tabelului. În schimb se evită utilizarea pointerilor, iar memoria, care altfel ar fi fost utilizată pentru listele înlanțuite, poate fi utilizată pentru o tabelă de dimensiune mai mare.

Modul de testare a pozițiilor în tabelă în cazul adresării deschise nu este liniar ci depinde de cheia inserată și de numărul de testări efectuate până la momentul curent.

O funcție de repartizare pentru repartizare deschisă este de forma:

$$h : U \rightarrow \{0, 1, \dots, m-1\} \times \{0, 1, \dots, m-1\}$$

$h(k, t)$  = poziția de repartizare a cheii  $k$  după  $t$  teste.

Inserția unei chei într-o tabelă  $T[0, \dots, m-1]$

```
TAD_INSERT( $T, k$ )
   $i = 0$ 
  repeta
     $j = h(k, i)$ 
    dacă  $T[j] = NULL$  atunci
       $T[j] = k$ 
```

```

        RETURN  $j$ 
    sfarsit daca
     $i = i + 1$ 
    pana cand  $i = m$ 
    scrie("tabela este plina")
    RETURN -1

```

Algoritmul de căutare testează pentru o anumită cheie  $k$  aceeași secvență de poziții ca și algoritmul de inserție al cheii  $k$ .

```

TAD_CAUT( $T, k$ )
     $i = 0$ 
    repeta
         $j = h(k, i)$ 
        daca  $T[j] = k$  atunci
            RETURN  $j$ 
         $i = i + 1$ 
    pana cand  $i = m$  sau  $T[j] = NULL$ 
    RETURN -1

```

**Observație:** Operația de ștergere este problematică în cadrul adresării deschise. Dacă o cheie  $k$  este ștersă prin marcarea poziției cu  $NULL$  atunci o cheie  $p$ , inserată după  $k$ , dar care la inserție/ștergere presupune testarea poziției pe care s-a aflat  $k$ , nu va mai fi găsită prin algoritmul de mai sus  $\Rightarrow$  este nevoie de un algoritm de complexitate mai mare la căutare.

O soluție a acestei probleme este marcarea pozițiilor șterse cu un marcaj special de poziție ștersă, dar tot presupune complexitate crescută pentru număr mare de poziții șterse.

Atunci când este nevoie de operația de ștergere, se recomandă folosirea de tabele de repartizare cu liste înlănțuite.

În ceea ce privește alegerea unei funcții de repartizare pentru repartizarea cu adresare deschisă, aceasta poate fi realizată în diferite moduri. Vor fi descrise mai jos trei metode de construcție a unei astfel de funcții.

### (1) Funcții de repartizare cu testare liniară

Se consideră o funcție de repartizare obișnuită:  $h_1 : U \rightarrow \{0, 1, \dots, m-1\}$ . Se definește:  $h : U \times \{0, 1, \dots, m-1\} \rightarrow \{0, 1, \dots, m-1\}$

$$h(k, i) = (h_1(k) + i) \bmod m$$

**Exemplu:** considerăm  $m = 11$ ,  $h_1(k) = k \bmod 11$  și

$$h(k, i) = ((k \bmod 11) + i) \bmod 11.$$

Atunci cheile 22, 33, 45, 59, 67, 13, 71 vor fi repartizate după cum urmează:

$$h(22, 0) = (0 + 0) \bmod 11 = 0$$

$$h(33, 0) = (0 + 0) \bmod 11 = 0, \text{ este deja ocupat, deci testarea continuă}$$

$$h(33, 1) = (0 + 1) \bmod 11 = 1$$

$$h(45, 0) = (1 + 0) \bmod 11 = 1, \text{ este deja ocupat, deci testarea continuă}$$

$$h(45, 1) = (1 + 1) \bmod 11 = 2$$

$$h(59, 0) = (4 + 0) \bmod 11 = 4$$

$$h(67, 0) = (1 + 0) \bmod 11 = 1, \text{ este deja ocupat, deci testarea continuă}$$

$$h(67, 1) = (1 + 1) \bmod 11 = 2, \text{ este deja ocupat, deci testarea continuă}$$

$$h(67, 2) = (1 + 2) \bmod 11 = 3,$$

$$h(13, 0) = (2 + 0) \bmod 11 = 2, \text{ este deja ocupat, deci testarea continuă}$$

$$h(13, 1) = (2 + 1) \bmod 11 = 3, \text{ este deja ocupat, deci testarea continuă}$$

$$h(13, 2) = (2 + 2) \bmod 11 = 4$$

$$h(71, 0) = (5 + 0) \bmod 11 = 5$$

**Observație:** Dezavantajul testării liniare este acela că, pe măsură ce se adaugă elemente, cresc secvențele de poziții succesive ocupate - **primary clustering**, ceea ce duce la creșterea complexității de inserție/căutare.

Un mod de îmbunătățire a acestei situații este utilizarea funcțiilor de repartizare cu testare pătratică descrise în continuare.

## (2) Funcțiile de repartizare cu testare pătratică

Se consideră o funcție de repartizare obișnuită:  $h_1 : U \rightarrow \{0, 1, \dots, m-1\}$ . Se definește:  $h : U \times \{0, 1, \dots, m-1\} \rightarrow \{0, 1, \dots, m-1\}$

$$h(k, i) = (h_1(k) + c_1 i + c_2 i^2) \bmod m$$

$h_1$ , funcție de repartizare obișnuită,  $c_1$  și  $c_2$  constante întregi pozitive auxiliare.

Prin acest mod de repartizare, se îmbunătățește ușor performanța. Problema este că, pentru două elemente  $k_1$  și  $k_2$  distincte, pentru care  $h(k_1, 0) = h(k_2, 0)$ , oricare ar fi  $i$ ,  $h(k_1, i) = h(k_2, i)$ . Acest lucru are ca urmare faptul că, dacă se inserează numeroase elemente care au aceeași repartizare inițială, crește complexitatea inserției. Rezultă o formă oarecum atenuată de *clustering*, numită *secondary clustering*.

Una dintre cele mai bune metode de repartizare deschisă se realizează cu ajutorul **dublei repartizări**.

## (3) Funcții de repartiție cu dublă repartizare

Se consideră două funcții de repartizare obișnuite distincte  $h_1(k)$  și  $h_2(k)$ . Se construiește funcția cu dublă repartizare:

$$h(k, i) = (h_1(k) + i h_2(k)) \bmod m$$



**Observație:** Se recomandă ca funcția  $h_2(k)$  să producă valori care sunt prime față de  $m$ . Acest lucru se poate obține dacă

- $m$  putere a lui 2 și  $h_2(k)$  produce întotdeauna un număr impar
- $m$  prim și  $h_2(k)$  produce numere naturale din  $\{0, 1, \dots, m-1\}$ , de exemplu considerăm funcțiile de repartiție

$$h_1(k) = k \bmod m$$

$$h_2(k) = 1 + (k \bmod m_1)$$

cu  $m$  prim și  $m_1$  astfel încât  $m_1$  ceva mai mic decât  $m$  - de exemplu  $m_1 = m - 1$ .

**Exemplu:**  $h_1(k) = k \bmod 11$ ,  $h_2(k) = 1 + (k \bmod 10)$ . Atunci cheile 22, 33, 45, 59, 67, 13, 71 vor fi repartizate după cum urmează:

$$h(22, 0) = (22 \bmod 11) \bmod 11 = 0$$

$$h(33, 0) = (33 \bmod 11) \bmod 11 = 0, \text{ este deja ocupat, deci continuă testarea}$$

$$h(33, 1) = ((33 \bmod 11) + (1 + 33 \bmod 10)) \bmod 11 = (0 + 4) \bmod 11 = 4$$

$$h(45, 0) = (45 \bmod 11) \bmod 11 = 1$$

$$h(59, 0) = (59 \bmod 11) \bmod 11 = 4, \text{ este deja ocupat, deci continuă testarea}$$

$$h(59, 1) = (59 \bmod 11 + 1 + 59 \bmod 10) \bmod 11 = (4 + 10) \bmod 11 = 3 \quad h(67, 0) = (67 \bmod 11) \bmod 11 = 1, \text{ este deja ocupat, deci continuă testarea}$$

$$h(67, 1) = ((67 \bmod 11) + (1 + 67 \bmod 10)) \bmod 11 = 9$$

$$h(13, 0) = (13 \bmod 11) \bmod 11 = 2$$

$$h(71, 0) = (71 \bmod 11) \bmod 11 = 6$$

Se observă că sunt necesare mai puține încercări decât în cazul repartizării cu testare liniară.

**Complexitate:** considerând factorul de încărcare al tabeli de repartizare  $\alpha = n/m < 1$ , unde  $n$  = numărul de poziții ocupate și  $m$  = dimensiunea tabeli atunci:

- numărul de testări în cazul unei căutări fără succes este în cel mai defavorabil caz  $1/(1 - \alpha)$  (presupunând repartizare uniformă)  $\Rightarrow$  căutarea fără succes are complexitate  $O(1)$  pentru  $\alpha$  constant;
- numărul de testări în cazul unei căutări cu succes este în cel mai defavorabil caz  $\frac{1}{\alpha} \ln \frac{1}{1-\alpha}$  (presupunând repartizare uniformă)  $\Rightarrow$  căutarea cu succes are complexitate  $O(1)$  pentru  $\alpha$  constant

**Aplicații:**

- Compilatoare - pentru păstrarea tabeli de identificatori se recomandă utilizarea unei tabeli de repartizare, deoarece este necesar acces rapid la informație.
- Criptografie - pentru securitatea parolilor se utilizează *hashing* criptografic. Acest tip de repartizare este mai dificil de realizat și presupune îndeplinirea anumitor criterii de securitate, care depășesc însă limitele acestui curs.

## 4 Repartizarea cheilor de tip șir de caractere

O funcție de repartizare ideală va produce pentru orice două chei inserate poziții de inserție diferite. Există situații, atunci când sunt cunoscute toate cheile posibile (de exemplu în cazul unui dicționar) să fie construite funcții de repartizare ideale în care să se evite orice coliziune. Acest lucru se numește *perfect hashing* și există documentație în internet. În cazul general acest lucru este practic imposibil, astfel încât se recomandă construirea de funcții de repartizare care să producă o repartizare uniformă a cheilor în tabelă, rezultând o probabilitate mică de coliziune. Pentru chei numere întregi au fost prezentate câteva criterii de construcție a acestor funcții. În plus, o funcție de repartizare bună trebuie să utilizeze toate părțile unei chei pentru generarea valorii de repartizare. De exemplu, dacă se utilizează *string*-uri pentru chei, nu se vor considera doar primele 5 caractere pentru generarea valorii de repartizare sau dacă se folosesc chei numere întregi cu mai multe cifre, nu se vor utiliza doar anumite cifre ale acestora.

Se spune despre o funcție de repartizare că obține o avalanșă - *achieves avalanche* - dacă și doar prin modificarea unui singur bit al unei chei se obține o valoare de repartizare complet diferită de cea anterioară. Acest efect ajută la o distribuire uniformă a cheilor în tabelă și la o reducere a probabilității coliziunilor. Conceptul de *avalanche* este derivat de la *hashing*-ul criptografic.

Este extrem de dificil de elaborat o funcție de repartizare bună. De aceea se recomandă utilizarea unor funcții deja existente, despre care se știe că produc rezultate bune. În continuare sunt prezentate câteva funcții de repartizare cunoscute împreună cu anumite proprietăți ale lor. Sunt considerate chei de tip *string*, dar pot fi utilizate și pentru chei de tip întreg, în loc de câte un caracter considerând câte un byte al cheii.