

# Tabele de dispersie

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*. În cazul mulțimilor sortate, căutarea binară permite găsirea unei anumite chei în complexitate logaritmică. Vom vedea în capitolele ce tratează arborii binari de căutare, cum poate fi rezolvată în complexitate logaritmică problema căutării în cazul mulțimilor dinamice de date, deci a mulțimilor care se modifică prin inserții și ștergeri.

În capitolul de față sunt prezentate tabele de dispersie - *hash tables*, *hash maps* - structuri de date care permit accesul prin cheie în timp constant. O astfel de structură de date este eficientă, atunci când nu sunt necesare operații de sortare și păstrare ordonată a setului de date.

O tabelă de dispersie este de fapt o generalizare a noțiunii de vector - *array*. Operațiile de bază sunt inserția, căutarea și eventual ștergerea. Tabelele de dispersie 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 tabelii de identificatori sau pentru cuvintele cheie. Vom vedea pe parcurs că, într-o tabelă de dispersie operațiile au în medie complexitatea  $O(1)$ . O tabelă de dispersie se memorează într-un *array*.

## 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  $T[k] = NIL$ . O astfel de tabelă este reprezentată grafic în figura 1.

Algoritmii pentru operațiile de căutare, inserție și ștergere sunt extrem de simpli.

---

**Algoritm:** AD-CAUT

---

**Intrare:** tabela  $T$  și cheia  $k$   
RETURN  $T[k]$

---

---

**Algoritm:** AD-INSEREAZA

---

**Intrare:** tabela  $T$  și elementul  $x$   
 $T[x.cheie] \leftarrow x$

---

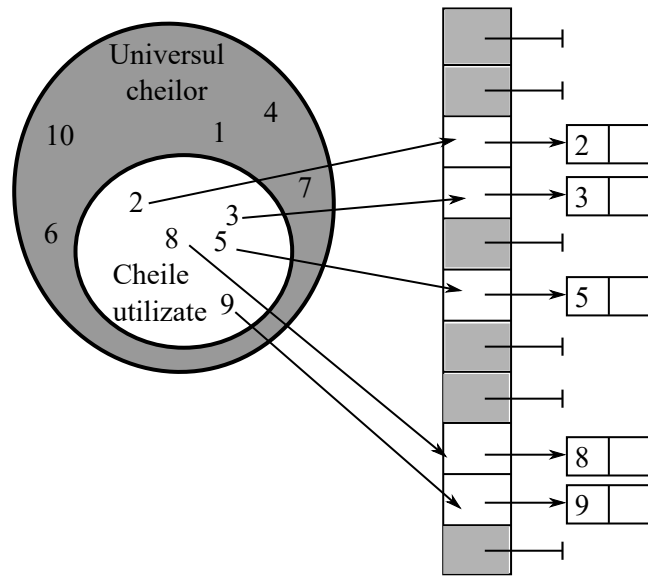


Figura 1: Tabelă cu adresare directă.

---

**Algoritm:** AD-STERGE

---

**Intrare:** tabela  $T$  și elementul  $x$

$T[x.cheie] \leftarrow NIL$

---

**Complexitate:** evident complexitatea este  $O(1)$  pentru toate cele trei operații.

**Observații:**

- se presupune faptul că fiecare element  $x$  are o cheie unică, diferită de cheile tuturor celorlalte elemente din tabelă. Practic un element este identificat în mod unic prin cheie.
- o tabelă cu adresare directă este potrivită, atunci când universul  $U$  al cheilor este relativ redus și numărul de elemente memorate în tabelă este comparabil cu  $|U| = m$ .
- de obicei elementele dintr-o tabelă cu adresare directă sau dintr-o tabelă de dispersie au un câmp pentru cheie, după care se face căutarea, și un câmp pentru informația de interes.

## 2 Tabele de dispersie - *Hash Tables*

Pentru situații în care numărul de elemente stocate într-o tabelă este mult mai mic decât universul cheilor, de exemplu atunci când orice număr natural poate fi cheie, tabelele cu adresare directă devin practic inutilizabile, datorită cantității de memorie necesare, care poate fi astfel nelimitată. În această situație se înlocuiește tabela cu adresare directă printr-o tabelă de dispersie.

Pentru memorarea unei tabele de dispersie se utilizează de asemenea un *array*  $T$  de dimensiune  $m$ :  $T[0, \dots, m-1]$ . Accesul la un element se face prin intermediul unei funcții de dispersie (repartizare) - *hash function* -  $h : U \rightarrow \{0, 1, \dots, \dim - 1\}$ ,  $U \gg m$ .

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  $h(k)$* ). În figura 2 este reprezentată grafic o tabelă de dispersie.

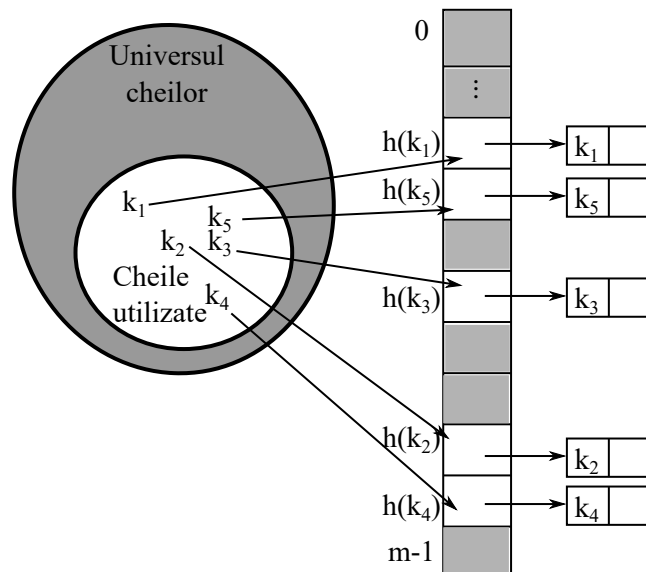


Figura 2: Tabelă de dispersie.

## 2.1 Problema coliziunilor

### 2.1.1 Exemplu de funcție de repartizare:

O funcție foarte simplă pentru repartizarea numerelor naturale într-o tabelă de dispersie este dată prin funcția *modulo* ( $\text{mod}$ ), care reprezintă restul împărțirii la un număr natural.

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

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

### Coliziunea

În cazul unei funcții de dispersie  $h$  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  $h(28) = h(17) = 6$  (fig. 4).

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

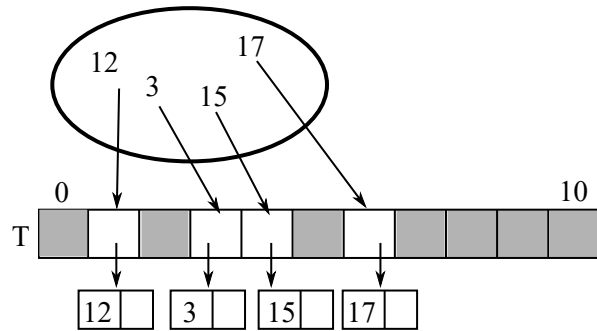


Figura 3: Tabela de repartizare cu funcția de dispersie  $h(k) = k \bmod 11$  în care au fost plasate cheile 3, 12, 15 și 17.

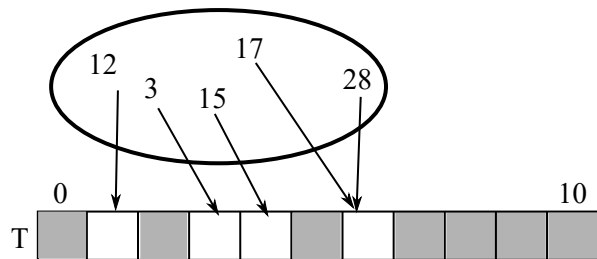


Figura 4: Coliziune între cheile 28 și 17.

lucru nu este posibil. Totuși, prin construirea atentă a funcției de dispersie, poate fi redusă semnificativ probabilitatea unei coliziuni. Vor fi discutate pe parcursul acestui curs câteva metode de construcție a unei funcții de dispersie, care să îndeplinească acest lucru.

În practică rezolvarea coliziunilor pentru tabele de dispersie poate fi făcută de exemplu cu ajutorul listelor înlănțuite. Astfel, în loc de a stoca într-o poziție a tabelului de repartizare un singur element, se păstrează o listă de elemente (ideea de *bucket*). Lista de la poziția  $k$  va conține toate elementele a cărei cheie este repartizată pe poziție.

**Exemplu:** Se consideră o tabelă de dispersie  $T$  cu dimensiunea  $m = 11$  care utilizează liste înlănțuite și în care se introduc elementele cu cheile 33, 35, 55, 46, 12, 2, 7, 10, 11. Rezultatul este reprezentat în figura 5. Inserția în liste s-a realizat la începutul listelor.

O altă metodă de rezolvare a coliziunilor este *adresare deschisă*, metodă prezentată la sfârșitul capitolului.

## 2.2 Operațiile uzuale

În continuare sunt prezentați algoritmi pentru operațiile de căutare, inserare și ștergere a unui element dintr-o tabelă de dispersie ce utilizează liste înlănțuite pentru rezolvarea coliziunilor. Se consideră că fiecare element  $x$  din tabelă este o structură cu

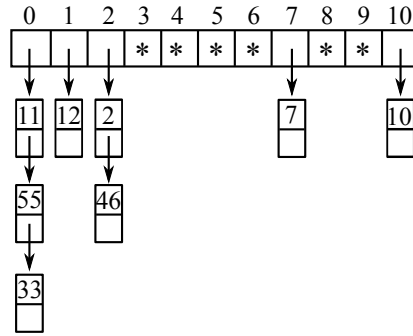


Figura 5: Tabela de dispersie de dimensiune  $m = 11$  în care au fost plasate succesiv cheile 33, 35, 55, 46, 12, 2, 7, 10 și 11. Au fost marcate cu \* pozițiile neocupate din tabelă.

două câmpuri  $x.cheie$  pentru cheie și  $x.val$  pentru informația de interes. În algoritmi se consideră notația TDL pentru Tabelă de Dispersie cu Liste înlanțuite.

---

**Algoritm:** TDL-INSEREAZA

---

**Intrare:** tabela  $T$  de dimensiune  $m$  cu funcția de dispersie  $h$  și elementul  $x$   
 Insereaza  $x$  in capul listei  $T[h(x.cheie)]$

---



---

**Algoritm:** TDL-CAUTA

---

**Intrare:** tabela  $T$  de dimensiune  $m$  cu funcția de dispersie  $h$  și cheia  $k$   
**Iesire:** elementul cu cheia  $x$  sau NIL în caz de eșec  
 Cauta elementul  $x$  cu cheia  $k$  in lista  $T[h(k)]$   
 return  $x$

---



---

**Algoritm:** TDL-STERGE

---

**Intrare:** tabela  $T$  de dimensiune  $m$  cu funcția de dispersie  $h$  și elementul  $x$   
 Sterge  $x$  din lista  $T[h(x.cheie)]$

---

**Complexitate:**

- funcția de inserție are evident complexitatea  $O(1)$ ,
- la funcția de căutare complexitatea depinde de lungimea listelor înlanțuite
- funcția de ștergere, care are ca parametru un nod  $x$  deja identificat în lista înlanțuită, are complexitatea  $O(1)$  dacă lista e dublu înlanțuită, altfel depinde de lungimea listei. Dacă întâi trebuie căutat elementul cu cheia  $k$ , apoi șters, se adaugă complexitatea căutării.

## 2.3 Analiza complexității

Considerăm tabela de repartizare  $T$  cu  $m$  poziții, în care rezolvarea coliziunilor se realizează prin înlanțuire. Presupunem că în  $T$  se stochează  $n$  elemente. Atunci numărul mediu de elemente stocate într-o listă este  $\alpha = n/m$ , unde  $\alpha$  poate fi mai mic sau mai mare decât 1.  $\alpha$  se numește *factor de încărcare* al tablei.

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.

Notăm cu  $n_j$  lungimea listei de pe poziția  $T[j]$ ,  $j = \overline{0, m-1}$ . Observăm că  $n_0 + n_1 + \dots + n_{m-1} = n$ , iar valoarea medie sau valoarea **estimată** (vezi Cormen) pentru lungimea oricărei liste  $n_j$  este dată prin:  $E[n_j] = n/m = \alpha$ .

Considerând o funcție de dispersie 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  $\Theta(1 + \alpha)$  cu ajutorul următoarelor două teoreme.

**Teorema 2.1** *Într-o tabelă de dispersie care tratează coliziunile prin înlănțuire, o căutare fără succes are complexitatea medie de timp  $\Theta(1 + \alpha)$ . (citată Cormen)*

**Demonstrație:** Presupunem că se caută în  $T$  cheia  $k$ , care nu există. Datorită faptului că s-a presupus o repartizare uniformă cu  $h$ , probabilitatea repartizării cheii  $k$  pe oricare dintre poziții este aceeași. Timpul necesar pentru a constata că această cheie nu este în  $T$  este timpul necesar de parcurgere a listei  $T[h(k)]$ . În medie acest timp depinde de lungimea medie a listei  $h(k)$ , deci de  $E[n_h(k)] = \alpha$ . Dat fiind că se ia ține cont și timpul necesar pentru calculul valorii  $h(k)$  se obține o complexitate de  $\Theta(1 + \alpha)$ .

**Teorema 2.2** *Într-o tabelă de dispersie care tratează coliziunile prin înlănțuire, o căutare cu succes are complexitatea medie de timp  $\Theta(1 + \alpha)$ . (citată Cormen)*

**Demonstrație** Situația căutării cu succes diferă de cea a căutării fără succes prin faptul că, timpul de căutare al elementului  $x$  cu cheia  $k$  din tabelă depinde de numărul de elemente aflate în lista  $T[h(k)]$  înaintea lui  $x$ . Modul de calcul al acestui timp mediu presupune ceva cunoștințe de probabilități și depășește cadrul acestui curs. O descriere detaliată poate fi găsită în bibliografie (Cormen)

Dacă numărul de elemente stocate în tabelă,  $n$ , este proporțional cu dimensiunea  $m$  a tablei, adică  $n = am$ , atunci complexitatea la căutare devine  $\Theta(1 + am/m) = \Theta(1 + a) = \Theta(1)$ . Cu alte cuvinte, dacă funcția de dispersie este aleasă în mod potrivit, astfel încât să ofere o repartizare uniformă și dacă dimensiunea tablei este aleasă proporțional cu numărul estimat de elemente ce vor fi introduse în tabelă, timpul de căutare a unei chei devine în medie constant.

### 3 Metode de repartizare

Din discuția complexității rezultă că, o funcție de dispersie 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 dispersie au ca mulțime a valorilor poziții în tabela de dispersie, deci elemente din mulțimea numerelor naturale. În cazul în care cheile sunt de exemplu șiruri de caractere - situație frecventă de altfel - 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. La finalul capitolului vor fi prezentate câteva funcții de dispersie pentru *string*-uri.

În continuare sunt prezentate câteva metode pentru construcția funcțiilor de dispersie pentru numere naturale.

### 3.1 Metoda diviziunii

Funcția de dispersie în acest caz este dată prin:

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

**Alegerea dimensiunii tablei.** O alegere bună pentru dimensiunea tablei  $m$  este un număr prim nu prea apropiat de o putere a lui 2.

**Justificare.** Dacă s-ar alege  $m = 2^p$  atunci de fapt împărțirea la  $m$  ar reprezenta o *shift*-are a biților cheii, ceea ce ar crea o dependență de *patter*-uri în structura cheilor, lucru ce trebuie evitat.

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,  $\alpha = n/m \Rightarrow m = n/\alpha = 3000/4 = 750$ . Se poate 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 dispersie 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 Dispersie 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  $\Theta(n)$  la căutare.

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

Utilizarea dispersiei 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 - **reshing**.

**Definiție:** (Cormen) O mulțime finită de funcții de dispersie  $H$ , 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 aleatoriu din  $H$ , probabilitatea unei coliziuni între  $k$  și  $l$  este cel mult  $1/m$ . (Evident, probabilitatea alegerii oricărei funcții din  $H$  este  $1/|H|$ . Probabilitatea alegerii unei funcții care repartizează  $k$  și  $l$  pe aceeași poziție este cel mult  $(|H|/m) * 1/|H| = 1/m$ ).

### 3.3.1 Construcție a unei clase universale de funcții de dispersie

Se alege un număr prim  $p$  suficient de mare, astfel încât să fie mai mare decât orice cheie posibilă. Evident  $p \gg m$ . Pentru fiecare  $a \in \{1, \dots, p-1\}$  și  $b \in \{0, 1, \dots, p-1\}$  se definește funcția de dispersie:

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

Mulțimea de astfel de funcții de dispersie este

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

cu  $|H_{pm}| = p(p-1)$ .

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 dispersie  $T$  conține un singur element. Pentru rezolvarea coliziunilor la inserție nu se utilizează liste înlănțuite, ci se testează diferite poziții, obținute pe baza funcției de dispersie, 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 înlănțuite, ci fiecare poziție conține cel mult un element, are ca urmare posibilitatea umplerii tabelului. În schimb se evită pointerii, iar memoria, care altfel ar fi fost necesară pentru listele înlănțuite, 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.



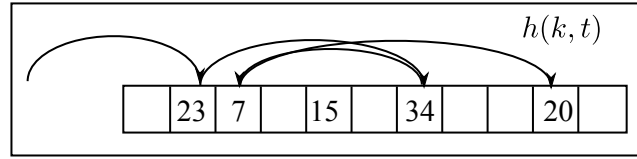


Figura 6: Căutarea elementului cu cheia 20. Acesta este găsit după 3 teste.

Forma generală a unei funcții de dispersie în cazul adresării deschise este de forma:

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

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

În mod evident, pentru a permite testarea tuturor pozițiilor posibile și a nu reveni periodic la aceleași poziții, testate anterior, funcția  $f(k, t)$  trebuie să producă pentru valorile lui  $t \in \{0, \dots, m - 1\}$  o permutare a mulțimii pozițiilor  $\{0, \dots, m - 1\}$ . Căutarea unei chei într-o astfel de tabelă este ilustrată în figura 6.

În continuare sunt prezentați algoritmi pentru inserarea și căutare într-o tabelă cu adresare deschisă.

---

**Algoritm:** TAD-INSEREAZA

---

**Intrare:** tabela  $T$  cu  $m$  elemente, elementul cu  $x$  cheia  $k$

$t \leftarrow 0$  //numarul de teste deja efectuate

**repetă**

$j \leftarrow h(k, t)$  //determina pozitia pentru testare

**dacă**  $T[j] = NIL$  **atunci**

$T[j] \leftarrow x$

**return**  $j$

**sfarsit\_dacă**

$t \leftarrow t + 1$  //incrementarea numarului de teste

**pana\_cand**  $t = m$ ;

**return** -1 //terminare cu esec

---

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

---

**Algoritm:** TAD-CAUTA

---

**Intrare:** tabela  $T$  cu  $m$  elemente, elementul cu  $x$  cheia  $k$

$t \leftarrow 0$

**repetă**

$j \leftarrow h(k, t)$

**dacă**  $T[j].cheie = k$  **atunci**

**return**  $j$

**sfarsit\_dacă**

$t \leftarrow t + 1$

**pana\_cand**  $t \neq m$  și  $T[j] \neq NIL$ ;

**return** -1

---

**Operația de ștergere** este problematică în cadrul adresării deschise. Dacă o cheie  $k$  este ștersă prin marcarea poziției cu  $NIL$  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.

O soluție a acestei probleme este marcarea pozițiilor șterse cu un marcaj special de poziție ștersă, acest lucru presupune modificarea algoritmilor de căutare și inserare și în plus duce la o complexitate crescută a operației de căutare pentru număr mare de poziții șterse. Pe măsură ce crește numărul de poziții marcate ca șterse, crește timpul de căutare și scade performanța. Acest lucru se numește *contaminare*. Singura posibilitate de rezolvare a acestei probleme este prin *rehashing*.

În ceea ce privește alegerea unei funcții de dispersie pentru adresarea deschisă, aceasta poate fi realizată în diferite moduri. Vor fi descrise mai jos trei metode de construcție a unei astfel de funcții. Toate cele trei metode garantează faptul că secvența  $\{h(k, 0), h(k, 1), \dots, h(k, m-1)\}$  este o permutare a mulțimii  $\{0, 1, \dots, m-1\}$  pentru orice cheie  $k$ .

### 3.4.1 Funcții de dispersie cu testare liniară

Se consideră o funcție de dispersie auxiliară, de tipul celor descrise la începutul capitolului:  $h_1 : U \rightarrow \{0, 1, \dots, m-1\}$ . Se definește:  $h : U \times \{0, 1, \dots, m-1\} \rightarrow \{0, 1, \dots, m-1\}$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

$$h(13, 2) = (2 + 2) \bmod 11 = 4 \text{ este deja ocupat, deci testarea continuă}$$

$$h(71, 0) = (2 + 3) \bmod 11 = 5$$

Inserarea cheilor din acest exemplu este ilustrată în figura 7.

**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.

	Tabela inițială	22	33	45	59	67	13
0	█	22	22	22	22	22	22
1	█	█	33	33	33	33	33
2	█	█	█	45	45	45	45
3	█	█	█	█	67	67	67
4	█	█	█	█	59	59	59
5	█	█	█	█	█	13	13
6	█	█	█	█	█	█	█
7	█	█	█	█	█	█	█
8	█	█	█	█	█	█	█
9	█	█	█	█	█	█	█
10	█	█	█	█	█	█	█

Figura 7: Inserarea cheilor 22, 33, 45, 59, 67, 13 într-o tabelă de dispersie de dimensiune  $m = 11$  folosind testarea liniară. Pozițiile testate pentru fiecare element sunt marcate cu săgeți.

### 3.4.2 Funcțiile de dispersie cu testare pătratică

Se consideră o funcție de dispersie auxiliară de tipul celor descrise la începutul capitolului:  $h_1 : U \rightarrow \{0, 1, \dots, m-1\}$ . Se definește:  $h : U \times \{0, 1, \dots, m-1\} \rightarrow \{0, 1, \dots, m-1\}$

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

$c_1$  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  $t$ ,  $h(k_1, t) = h(k_2, t)$ . 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*. (citată)

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

### 3.4.3 Adresarea deschisă cu dublă repartizare

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

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

**Observații:**

- poziția inițială testată este  $h(k, 0) = h_1(k)$ , deci nu depinde de  $h_2$ , pentru orice cheie  $k$
- poziția  $h(k, t)$  este la distanța  $h_2(k)$ , față de poziția  $h(k, t-1)$ , pentru orice  $t > 0$

- pentru ca secvența de poziții  $\{h(k, 0), h(k, 1), \dots, h(k, m - 1)\}$  să fie o permutare a pozițiilor  $\{0, 1, \dots, m - 1\}$  trebuie ca funcția  $h_2(k)$  să producă valori care sunt prime față de  $m$ . Dacă există o cheie  $k$ , pentru care  $h_2(k)$  și  $m$  au un divizor comun  $1 < d < m$ , atunci cel mult după  $d$  teste se repetă din nou testarea poziției inițiale. Adică vor rămâne în tabelă poziții netestate!!

Satisfacerea condiției discutate mai sus, 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 repartitie

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

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

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

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

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

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

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

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

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

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

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

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

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

Procesul de inserție a acestor chei este ilustrat în figura 8. Se observă că sunt necesare semnificativ mai puține încercări decât în cazul repartizării cu testare liniară.

### 3.4.4 Analiza complexității

Considerăm o tabelă de dispersie cu factorul de încărcare  $\alpha = n/m < 1$ , unde  $n$  = numărul de poziții ocupate și  $m$  = dimensiunea tablei atunci. De asemenea presupunem că funcția de dispersie asigură repartizarea uniformă. Acest lucru presupune că probabilitatea ca secvența de poziții testate pentru o anumită cheie  $k$  să fie oricare dintre cele  $m!$  permutări ale mulțimii  $\{0, 1, \dots, m - 1\}$  este aceeași. Atunci complexitatea căutării unei chei este determinată pe baza următoarelor două teoreme.

**Teorema 3.1** *Pentru o tabelă de dispersie cu adresare deschisă, cu factorul de încărcare  $\alpha < 1$ , numărul mediu teste necesare pentru căutarea fără succes este cel mult  $1/(1 - \alpha)$ , presupunând o distribuție uniformă.*

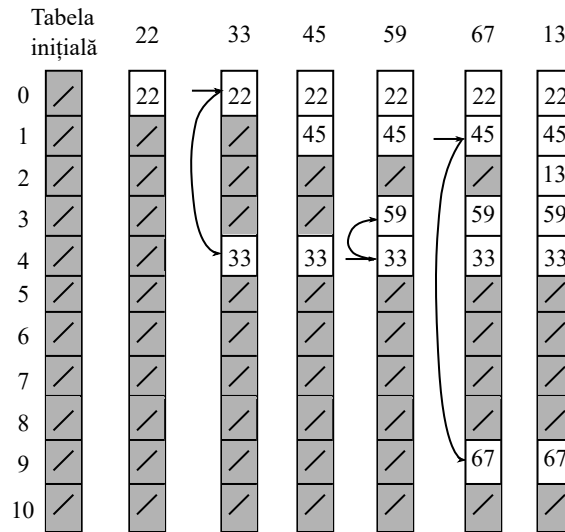


Figura 8: Inserarea cheilor 22, 33, 45, 59, 67, 13 într-o tabelă de dispersie de dimensiune  $m = 11$  folosind testarea dubla repartizare. Pozițiile testate pentru fiecare element sunt marcate cu săgeți.

**Teorema 3.2** Pentru o tabelă de dispersie cu adresare deschisă, cu factorul de încărcare  $\alpha < 1$ , numărul mediu teste necesare pentru căutarea cu succes este cel mult  $\frac{1}{\alpha} \log \frac{1}{1-\alpha}$ , presupunând o repartizare uniformă.

Ambele teoreme sunt demonstrate cu ajutorul relațiilor probabilistice în (Cormen).

### Observații

- Testarea liniară și testarea pătratică nu garantează distribuție uniformă datorită fenomenului de *clustering*.
- Dubla repartizare se apropie cel mai mult de o distribuție uniformă.
- În practică, dacă factorul de încărcare devine foarte mare eficiența inserției, care presupune căutarea unei poziții neocupate, deci de fapt căutarea fără succes a cheii elementului ce se înserează, scade semnificativ. Astfel încât, atunci când  $\alpha$  se apropie de 80% se recomandă redimensionarea tabelii și *rehashing*. De exemplu, în limbajul **Python**, care implementează dicționarele folosind adresarea deschisă, se pornește cu o tabelă de dimensiune 8 și atunci când s-a umplut 75% din tabelă se redimensionează. (<https://adamgold.github.io/posts/python-hash-tables-under-the-hood/#dictionary-resize>)

## 3.5 Perfect hashing

O funcție de distribuție 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 dispersie ideale în care să se evite

orice coliziune. Acest lucru se numește *perfect hashing*. În (Cormen) este propus un algoritm de dispersie perfect, care are la bază tot ideea de *bucket-uri* similar cu rezolvarea coliziunilor prin înlanțuire, doar că, în loc de liste înlanțuite se utilizează tabele de dispersie secundare.

În cazul general acest lucru este practic imposibil, astfel încât se recomandă construirea de funcții de dispersie 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 dispersie bună trebuie să utilizeze toate părțile unei chei pentru generarea valorii de dispersie. De exemplu, dacă se utilizează *string-uri* pentru chei, nu se vor considera doar primele 5 caractere pentru generarea valorii de dispersie sau dacă se folosesc chei numere întregi cu mai multe cifre, nu se vor utiliza doar anumite cifre ale acestora.

## 4 Tratarea cheilor de tip șir de caractere

În partea teoretică prezentată până acum s-a considerat doar distribuirea de chei numere naturale. În practică există chei de diferite tipuri, frecvent caractere sau șiruri de caractere. Acestor chei trebuie mai întâi să li se asocieze un număr natural.

În primul rând putem considera pentru fiecare caracter ASCII codul său numeric. Problema devine, cum combinăm codurile unei secvențe de caractere într-un număr natural, astfel încât să nu depășim valoarea maximă admisă și să evităm coliziunile care se produc, dacă pentru două șiruri de caractere se generează aceeași valoare numerică.

### 4.1 Funcții elementare de transformare a unui *string* într-un număr natural

Cele mai simple metode de a transforma o cheie de tip *string* într-o cheie numerică sunt *hashing-ul* aditiv și XOR-*hash*, descrise mai jos.

**Hashing - aditiv** - presupune adunarea codurilor numerice corespunzătoare tuturor caracterelor din șir. Acest tip de calcul al unei chei numerice este extrem de prost, deoarece generează coliziuni pentru toate permutările unui șir de caractere. În practică este deci inutilizabil.

**XOR-*hash*** - presupune aplicarea iterativă a operației de XOR asupra caracterelor din șir. Practic se pornește cu valoarea  $h\_val = 0$ , care pentru fiecare caracter  $sir[i]$  din șir se modifică prin instrucțiunea

$$h\_val \leftarrow h\_val \text{ XOR } sir[i]$$

Nici această metodă nu dă rezultate bune în ceea ce privește coliziunile, dar operația XOR poate fi utilizată în cazul unor funcții de repartizare mai sofisticate.

## 4.2 Funcții polinomiale de dispersie

Un tip funcție de dispersie frecvent utilizată pentru a asocia unui *string* un număr natural este un polinom de grad  $L - 1$  unde  $L$  este lungimea șirului considerat:

$$h(sir) = sir[0] + sir[1] * p + sir[2] * p^2 + \dots + sir[L - 1] * p^{L-1} = \sum_{i=0}^{L-1} sir[i] * p^i$$

În cazul unor șiruri de caractere lungi există riscul de a obține valori extrem de mari, care să excedă maximul admis pentru valorile numerice. Din acest motiv, adesea valoarea obținută se consideră modulo  $M$ , unde  $M$  este suficient de mare, ca să poată permite obținerea tuturor numerelor naturale din intervalul posibil ( $[0, \text{MAX\_INT}]$ ). Formula devine astfel:

$$h(sir) = \left( \sum_{i=0}^{L-1} sir[i] * p^i \right) \bmod M$$

În [?] se folosesc metoda lui Horner și regulile pentru restul împărțirii la un număr natural pentru a calcula iterativ formulă, astfel

$$h_n = 0, h_i = (h_{i+1} * p + sir[i]) \bmod M, i = n - 1, \dots, 0$$

Pentru a avea o probabilitate cât mai redusă de coliziuni, este esențială alegerea parametrilor  $p$  și  $M$ . În [?] este discutat modul de alegere al acestora. În diferite exemple practice se utilizează pentru  $p$  un număr prim, de exemplu 31 sau 33 (Bernstein hash), iar pentru  $M$  poate fi considerat un număr prim mai mic decât  $2^{64}$ , o valoare întâlnită în literatura online fiind  $10^9 + 7$ .

## 4.3 Algoritmul FNV

Unul dintre cei mai cunoscuți algoritmi de asociere a unei valori naturale pentru un *string*, algoritmul FNV, utilizează o metodă similară cu cea a polinoamelor, doar că în loc de adunare între puteri se folosește operația de XOR. Algoritmul general este dat de:

---

**Algorithm:** FNV-Hash

---

**Intrare:** cheia  $k$   
 $hash \leftarrow FNV\_Basis$   
**pentru** fiecare byte  $b$  din șir **executa**  
     $hash \leftarrow hash \times FNV\_prime$   
     $hash \leftarrow hash \text{ XOR } b$   
**sfarsit\_for**  
**return** hash

---

Constantele  $FNV\_Prime$ , și  $FNV\_Basis$  depind de numărul de biți pe care este reprezentată valoarea de dispersie rezultată. În [?] este explicat modul de alegere al acestor constante, iar pentru valori pe 32 de biți se consideră:

$$FNV\_Prime = 16777619, FNV\_Basis = 2166136261$$

**Observație** Evident, în cazul tuturor acestor funcții care calculează o valoare naturală asociată unui **string**, pentru inserarea într-o tabelă de repartizare valoarea obținută trebuie considerată modulo dimensiunea tabelii.