

---

## 12 Tabele de dispersie

Multe aplicații necesită o mulțime dinamică pentru care să se aplice numai operațiile specifice pentru dicționare INSEREAZĂ, CAUTĂ și ȘTERGE. De exemplu, un compilator pentru un limbaj de programare întreține o tabelă de simboluri, în care cheile elementelor sunt șiruri arbitrare de caractere ce corespund identificatorilor din limbaj. O tabelă de dispersie este o structură eficientă de date pentru implementarea dicționarelor. Deși căutarea unui element într-o tabelă de dispersie poate necesita la fel de mult timp ca și căutarea unui element într-o listă înlanțuită – un timp  $\Theta(n)$  în cazul cel mai defavorabil – în practică, dispersia funcționează extrem de bine. Pe baza unor ipoteze rezonabile, timpul preconizat pentru căutarea unui element într-o tabelă de dispersie este  $O(1)$ .

O tabelă de dispersie este o generalizare a noțiunii mai simple de tablou. Adresarea directă într-un tablou folosește abilitatea noastră de a examina o poziție arbitrară în tablou într-un timp  $O(1)$ . Secțiunea 12.1 discută în detaliu despre adresarea directă. Adresarea directă este aplicabilă în cazul în care ne putem permite să alocăm un tablou care are câte o poziție pentru fiecare cheie posibilă.

Când numărul cheilor memorate efectiv este relativ mic față de numărul total de chei posibile, tabelele de dispersie devin o alternativă eficientă la adresarea directă într-un tablou, deoarece o tabelă de dispersie folosește în mod normal un tablou de mărime proporțională cu numărul de chei memorate efectiv. În loc să folosim direct cheia ca indice în tablou, indicele este *calculat* pe baza cheii. În secțiunea 12.2 sunt prezentate ideile principale, iar în secțiunea 12.3 se descrie cum pot fi calculați indicii din tablou pe baza cheilor, folosind funcții de dispersie. Sunt prezentate și analizate diferite variații ale temei de bază; ideea de bază este că dispersia reprezintă o tehnică extrem de eficientă și practică; operațiile de bază pentru dicționare necesită, în medie, doar un timp  $O(1)$ .

---

### 12.1. Tabele cu adresare directă

Adresarea directă este o tehnică simplă care funcționează bine atunci când universul  $U$  al cheilor este rezonabil de mic. Să presupunem că o aplicație necesită o mulțime dinamică în care fiecare element are o cheie aleasă dintr-un univers  $U = \{0, 1, \dots, m-1\}$ , unde  $m$  nu este foarte mare. Vom presupune că nu există două elemente având aceeași cheie.

Pentru a reprezenta mulțimea dinamică folosim un tablou sau o **tabelă cu adresare directă**  $T[0..m-1]$ , în care fiecare poziție sau **locație** corespunde unei chei din universul  $U$ . Figura 12.1 ilustrează această abordare; locația  $k$  referă elementul având cheia  $k$  din mulțime. Dacă mulțimea nu conține un element cu cheia  $k$ , atunci  $T[k] = \text{NIL}$ .

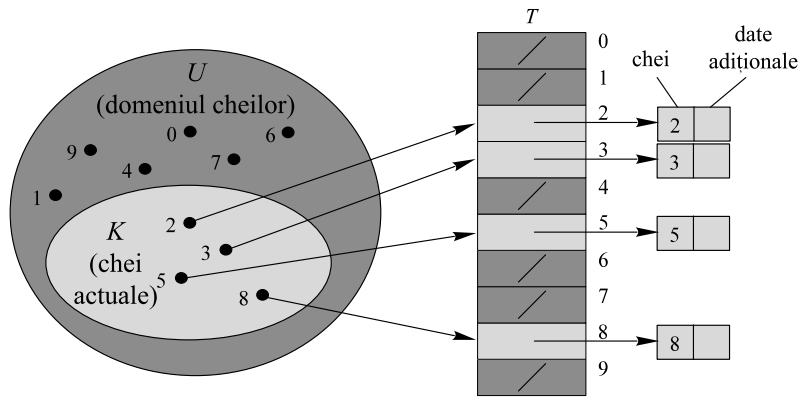
Operațiile pentru dicționare sunt ușor de implementat.

ADRESARE-DIRECTĂ-CAUTĂ( $T, k$ )

**returnează**  $T[k]$

ADRESARE-DIRECTĂ-INSEREAZĂ( $T, x$ )

$T[\text{cheie}[x]] \leftarrow x$



**Figura 12.1** Implementarea unei mulțimi dinamice printr-un tablou cu adresare directă  $T$ . Fiecare cheie din universul  $U = \{0, 1, \dots, 9\}$  corespunde unui indice în tablou. Mulțimea  $K = \{2, 3, 5, 8\}$  a cheilor efective determină locațiile din tablou care conțin pointeri către elemente. Celelalte locații, hașurate mai întunecat, conțin NIL.

ADRESARE-DIRECTĂ-ȘTERGE( $T, x$ )

$T[\text{cheie}[x]] \leftarrow \text{NIL}$

Fiecare dintre aceste operații este rapidă: este necesar doar un timp  $O(1)$ .

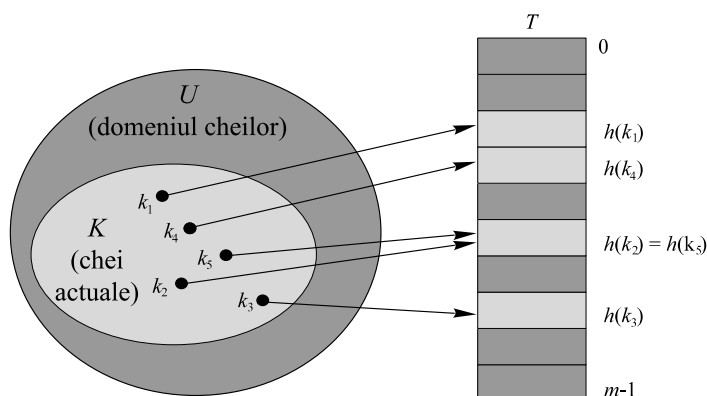
Pentru unele aplicații, elementele din mulțimea dinamică pot fi memorate chiar în tabela cu adresare directă. Aceasta înseamnă că în loc să memorăm cheia elementului și datele adiționale într-un obiect extern tabelii cu adresare directă, printr-un pointer dintr-o poziție din tabelă către obiect, putem să memorăm obiectul în locația respectivă, economisind astfel spațiu. Mai mult, deseori nu este necesar să memorăm câmpul cheie al obiectului, deoarece dacă avem indicele unui obiect în tabelă avem și cheia sa. Oricum, dacă cheile nu sunt memorate, trebuie să existe o modalitate de a afirma că locația este goală.

## Exerciții

**12.1-1** Se consideră o mulțime dinamică  $S$  care este reprezentată printr-o tabelă  $T$  cu adresare directă, de lungime  $m$ . Descrieți o procedură care găsește elementul maxim din  $S$ . Care este performanța procedurii pentru cazul cel mai defavorabil?

**12.1-2** Un **vector de biți** este un tablou simplu de biți (cu 0 și 1). Un vector de biți de lungime  $m$  ocupă mult mai puțin spațiu decât un tablou de  $m$  pointeri. Descrieți cum se poate folosi un vector de biți pentru a reprezenta o mulțime dinamică având elemente distincte, fără date adiționale. Operațiile pentru dicționare ar trebui să funcționeze într-un timp  $O(1)$ .

**12.1-3** Sugerați cum se poate implementa o tabelă cu adresare directă în care cheile elementelor memorate nu sunt neapărat distincte și elementele pot avea date adiționale. Toate cele trei operații pentru dicționare (INSEREAZĂ, ȘTERGE și CAUTĂ) ar trebui să se execute într-un timp  $O(1)$ . (Nu uitați că ȘTERGE are ca argument un pointer la obiectul ce va fi șters și nu o cheie.)



**Figura 12.2** Folosirea unei funcții de dispersie  $h$  pentru a transforma chei în poziții din tabela de dispersie. Cheile  $k_2$  și  $k_5$  se transformă în aceeași poziție, deci sunt în coliziune.

**12.1-4** ★ Dorim să implementăm un dicționar folosind adresarea directă pe un tablou *urîa*. La început, intrările în tablou pot conține date ne semnificative, iar inițializarea întregului tablou este nepractică datorită mărimei sale. Descrieți o schemă pentru implementarea unui dicționar cu adresare directă printr-un tablou urîș. Fiecare obiect memorat ar trebui să utilizeze un spațiu  $O(1)$ ; operațiile CAUTĂ, INSEREAZĂ și ȘTERGE ar trebui să funcționeze fiecare într-un timp  $O(1)$ ; inițializarea structurii de date ar trebui să necesite un timp  $O(1)$ . (*Indica ie:* Folosiți o stivă suplimentară, a cărei mărime să fie numărul de chei efectiv memorate în dicționar, pentru a determina dacă o intrare dată în tabloul urîș este validă sau nu.)

## 12.2. Tabele de dispersie

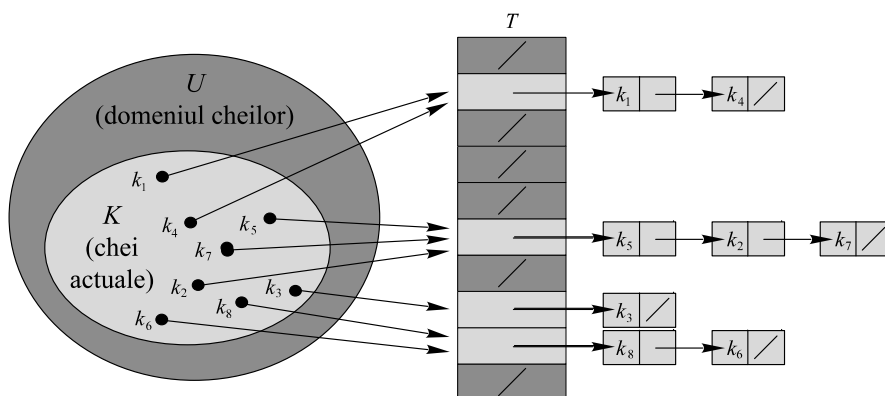
Dificultatea în adresarea directă este evidentă: dacă universul  $U$  este mare, memorarea tabelului  $T$  poate fi nepractică, sau chiar imposibilă, dată fiind memoria disponibilă a unui calculator uzual. Mai mult, mulțimea  $K$  a cheilor *efectiv memorate* poate fi atât de mică relativ la  $U$ , încât majoritatea spațiului alocat pentru  $T$  ar fi irosit.

Când mulțimea  $K$  a cheilor memorate într-un dicționar este mult mai mică decât universul  $U$  al cheilor posibile, o tabelă de dispersie necesită un spațiu de memorie mult mai mic decât o tabelă cu adresare directă. Mai exact, cerințele de memorare pot fi reduse la  $\Theta(|K|)$ , chiar și atunci când căutarea unui element în tabela de dispersie necesită tot un timp  $O(1)$ . (Singurul punct slab este că această margine este stabilită pentru *timpul mediu*, în timp ce pentru adresarea directă ea era valabilă pentru *cazul cel mai defavorabil*.)

Prin adresare directă, un element având cheia  $k$  este memorat în locația  $k$ . Prin dispersie, acest element este memorat în locația  $h(k)$ ; aceasta înseamnă că o **funcție de dispersie**  $h$  este folosită pentru a calcula locația pe baza cheii  $k$ . În acest caz,  $h$  transformă universul  $U$  al cheilor în locații ale unei **tabele de dispersie**  $T[0..m-1]$ :

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

Vom spune că un element cu cheia  $k$  **se dispersează** în locația  $h(k)$ ; de asemenea vom spune că



**Figura 12.3** Rezolvarea coliziunii prin înlanțuire. Fiecare locație din tabela de dispersie  $T[j]$  conține o listă înlanțuită a tuturor cheilor a căror valoare de dispersie este  $j$ . De exemplu,  $h(k_1) = h(k_4)$  și  $h(k_5) = h(k_2) = h(k_7)$ .

$h(k)$  este **valoarea de dispersie** a cheii  $k$ . Figura 12.2 ilustrează ideea de bază. Scopul funcțiilor de dispersie este de a reduce domeniul indicilor tabloului care trebuie manipulați. În loc de  $|U|$  valori, va trebui să manipulăm doar  $m$  valori. Cerințele de memorare sunt reduse corespunzător.

Dezavantajul acestei idei frumoase este că două chei se pot dispersa în aceeași locație, adică se produce o **coliziune**. Din fericire, există tehnici eficiente pentru rezolvarea conflictelor create de coliziuni.

Desigur, soluția ideală ar fi să nu existe coliziuni. Putem încerca să atingem acest scop printr-o alegere potrivită a funcției de dispersie  $h$ . O idee este de a-l face pe  $h$  să pară “aleator”, evitând astfel coliziunile sau cel puțin minimizând numărul lor. Termenul de “dispersie”, care evocă imaginea unei fărâmișări și amestecări aleatoare, captează spiritul acestei abordări. (Bineînțeles, o funcție de dispersie  $h$  trebuie să fie deterministă, în sensul că o intrare dată  $k$  trebuie să producă întotdeauna aceeași ieșire  $h(k)$ .) Deoarece  $|U| > m$ , există cu siguranță două chei care să aibe aceeași valoare de dispersie; de aceea, evitarea totală a coliziunilor este imposibilă. Prin urmare, deși o funcție de dispersie “aleatoare”, bine proiectată, poate minimiza numărul coliziunilor, avem în continuare nevoie de o metodă pentru rezolvarea coliziunilor ce apar.

Restul acestei secțiuni prezintă cea mai simplă tehnică de rezolvare a coliziunilor, numită înlanțuire. În secțiunea 12.4 se introduce o metodă alternativă de rezolvare a coliziunilor, numită adresare deschisă.

## Rezolvarea coliziunii prin înlanțuire

Prin **înlanțuire** punem toate elementele ce se dispersează în aceeași locație, într-o listă înlanțuită, după cum se arată în figura 12.3. Locația  $j$  conține un pointer către capul listei tuturor elementelor care se dispersează în locația  $j$ ; dacă nu există astfel de elemente, locația  $j$  conține NIL.

Operațiile pentru dicționare sunt ușor de implementat pe o tabelă de dispersie, în cazul în care coliziunile sunt rezolvate prin înlanțuire.

DISPERSIE-CU-ÎNLĂȚUIRE-INSEREAZĂ( $T, x$ )  
inserează  $x$  în capul listei  $T[h(\text{cheie}[x])]$

DISPERSIE-CU-ÎNLĂȚUIRE-CAUTĂ( $T, k$ )  
caută un element cu cheia  $k$  în lista  $T[h(k)]$

DISPERSIE-CU-ÎNLĂȚUIRE-ȘTERGE( $T, x$ ) indexDispersie-Cu-Inlantuire-Sterge@DISPERSIE-CU-ÎNLĂȚUIRE-ȘTERGE  
șterge  $x$  din lista  $T[h(\text{cheie}[x])]$

Timpul de execuție pentru inserare în cazul cel mai defavorabil este  $O(1)$ . Pentru căutare, timpul de execuție în cazul cel mai defavorabil este proporțional cu lungimea listei; vom analiza îndeaproape această situație în cele ce urmează. Ștergerea unui element  $x$  poate fi realizată într-un timp  $O(1)$  dacă listele sunt dublu înlănțuite. (Dacă listele sunt simplu înlănțuite, atunci trebuie întâi să-l găsim pe  $x$  în lista  $T[h(\text{cheie}[x])]$ , astfel încât legătura *next* a predecesorului lui  $x$  să fie modificată corespunzător ca să-l “ocolească” pe  $x$ ; în acest caz, ștergerea și căutarea au în esență același timp de execuție.)

## Analiza dispersiei cu înlănțuire

Cât de bine funcționează dispersia prin înlănțuire? În particular, cât durează căutarea unui element având o cheie dată?

Fiind dată o tabelă de dispersie  $T$  cu  $m$  locații ce memorează  $n$  elemente, vom defini **factorul de încărcare**  $\alpha$  pentru  $T$  prin  $n/m$ , raport care reprezintă numărul mediu de elemente memorate într-o înlănțuire. Analiza noastră se bazează pe  $\alpha$ ; adică, ne imaginăm  $\alpha$  ca fiind fix în timp ce  $n$  și  $m$  tind la infinit. (Se observă că  $\alpha$  poate fi mai mic, egal sau mai mare decât 1.)

Comportamentul, în cazul cel mai defavorabil, al dispersiei prin înlănțuire este slabă: toate cele  $n$  chei se dispersează în aceeași locație, creând o listă de lungime  $n$ . Timpul de căutare pentru cazul cel mai defavorabil este astfel  $\Theta(n)$  plus timpul pentru calculul funcției de dispersie—cu nimic mai bun decât în cazul în care am fi folosit o listă înlănțuită a tuturor elementelor. Este clar că tabelele de dispersie nu sunt folosite pentru performanța lor în cazul cel mai defavorabil.

Performanța dispersiei în cazul mediu depinde de cât de bine distribuie (în medie) funcția de dispersie  $h$  mulțimea cheilor ce trebuie memorate în cele  $m$  locații. În secțiunea 12.3 se discută aceste probleme, dar deocamdată vom presupune că orice element se poate dispersa în oricare din cele  $m$  locații cu aceeași probabilitate, independent de locul în care s-au dispersat celelalte elemente. Vom numi această ipoteză de lucru **dispersie uniformă simplă**.

Presupunem că valoarea de dispersie  $h(k)$  poate fi calculată într-un timp  $O(1)$ , astfel încât timpul necesar pentru a căuta un element având cheia  $k$  depinde liniar de lungimea listei  $T[h(k)]$ . Lăsând de o parte timpul  $O(1)$  necesar pentru calculul funcției de dispersie și pentru accesul la locația  $h(k)$ , să luăm în considerare numărul mediu de elemente examinate de algoritmul de căutare, adică, numărul de elemente din lista  $T[h(k)]$  care sunt verificate pentru a vedea dacă cheia lor este egală cu  $k$ . Vom considera două cazuri. În primul caz, căutarea este fără succes: nici un element din tabelă nu are cheia  $k$ . În al doilea caz, căutarea găsește cu succes un element având cheia  $k$ .

**Teorema 12.1** Într-o tabelă de dispersie în care coliziunile sunt rezolvate prin înlănțuire, o căutare fără succes necesită, în medie, un timp  $\Theta(1 + \alpha)$ , în ipoteza dispersiei uniforme simple.

**Demonstrație.** În ipoteza dispersiei uniforme simple, orice cheie  $k$  se poate dispersa cu aceeași probabilitate în oricare din cele  $m$  locații. Astfel, timpul mediu pentru o căutare fără succes, pentru o cheie  $k$ , este timpul mediu pentru căutarea până la coada uneia din cele  $m$  liste. Deci, numărul mediu de elemente examinate într-o căutare fără succes este  $\alpha$  și timpul total necesar (incluzând timpul pentru calculul lui  $h(k)$ ) este  $\Theta(1 + \alpha)$ . ■

**Teorema 12.2** Într-o tabelă de dispersie în care coliziunile sunt rezolvate prin înlănțuire, o căutare cu succes necesită, în medie, un timp  $\Theta(1 + \alpha)$ , în ipoteza dispersiei uniforme simple.

**Demonstrație.** Presupunem că cheia care este căutată poate fi, cu aceeași probabilitate, oricare din cele  $n$  chei memorate în tabelă. De asemenea, presupunem că procedura DISPERSIE-CU-ÎNLĂNȚUIRE-INSEREAZĂ inserează un nou element la sfârșitul unei liste și nu la începutul ei. (Conform exercițiului 12.2-3, timpul mediu pentru o căutare cu succes este același indiferent dacă noile elemente se inserează la începutul sau la sfârșitul listei.) Numărul mediu de elemente examinate într-o căutare cu succes este cu 1 mai mare decât numărul de elemente examinate în cazul în care elementul căutat a fost inserat (deoarece fiecare nou element se adaugă la sfârșitul listei). În concluzie, pentru a afla numărul mediu de elemente examinate, vom considera media celor  $n$  obiecte din tabelă, ca fiind 1 plus lungimea presupusă a listei în care se adaugă al  $i$ -lea element. Lungimea medie a listei respective este  $(i - 1)/m$  și deci, numărul mediu de elemente examinate într-o căutare cu succes este:

$$\frac{1}{n} \sum_{i=1}^n \left(1 + \frac{i-1}{m}\right) = 1 + \frac{1}{nm} \sum_{i=1}^n (i-1) = 1 + \left(\frac{1}{nm}\right) \left(\frac{(n-1)n}{2}\right) = 1 + \frac{\alpha}{2} - \frac{1}{2m}.$$

Deci, timpul total necesar unei căutări cu succes (incluzând timpul pentru calculul funcției de dispersie) este  $\Theta(2 + \alpha/2 - 1/2m) = \Theta(1 + \alpha)$ . ■

Ce semnificație are această analiză? Dacă numărul de locații din tabela de dispersie este cel puțin proporțional cu numărul de elemente din tabelă, avem  $n = O(m)$  și, prin urmare,  $\alpha = n/m = O(m)/m = O(1)$ . Deci, căutarea necesită, în medie, un timp constant. Deoarece inserarea necesită un timp  $O(1)$ , în cazul cel mai defavorabil (vezi exercițiul 12.2-3) iar ștergerea necesită un timp  $O(1)$  în cazul cel mai defavorabil, când listele sunt dublu înlănțuite, toate operațiile pentru dicționare pot fi efectuate în medie într-un timp  $O(1)$ .

## Exerciții

**12.2-1** Presupunând că folosim o funcție de dispersie aleatoare  $h$  pentru a dispersa  $n$  chei distincte într-un tablou  $T$  de dimensiune  $m$ , care este numărul mediu de coliziuni? Mai exact, care este cardinalul probabil al mulțimii  $\{(x, y) : h(x) = h(y)\}$ ?

**12.2-2** Ilustrați inserarea cheilor 5, 28, 19, 15, 20, 33, 12, 17, 10 într-o tabelă de dispersie cu coliziunile rezolvate prin înlănțuire. Tabela are 9 locații, iar funcția de dispersie este  $h(k) = k \bmod 9$ .

**12.2-3** Argumentați că timpul mediu pentru o căutare cu succes prin înlănțuire este același indiferent dacă noile elemente se inserează la începutul, respectiv la sfârșitul unei liste. (*Indica ie:* Arătați că timpul mediu pentru o căutare cu succes este același pentru *oricare* două ordine ale unei liste.)

**12.2-4** Profesorul Marley emite ipoteza că se poate îmbunătăți esențial performanța dacă modificăm schema de înlănțuire astfel încât fiecare listă să fie păstrată ordonată. În ce mod afectează modificarea profesorului timpul de execuție pentru căutări cu succes, căutări fără succes, inserări și ștergeri?

**12.2-5** Sugați cum poate fi alocat și dealocat spațiul de gestionare pentru elemente în cadrul tabelii de dispersie prin legarea tuturor locațiilor nefolosite într-o listă liberă. Presupunem că o locație poate memora un “indicator” (engl. *flag*) împreună cu un element și un pointer, sau cu doi pointeri. Toate operațiile pentru dicționare pentru lista liberă trebuie să funcționeze într-un timp  $O(1)$ . Este necesar ca lista liberă să fie dublu înlănțuită sau este suficient să fie simplu înlănțuită?

**12.2-6** Arătați că dacă  $|U| > nm$ , există o submulțime a lui  $U$  de mărime  $n$  ce conține chei care se dispersează toate în aceeași locație, astfel încât timpul de căutare, pentru dispersie cu înlănțuire, în cazul cel mai defavorabil, este  $\Theta(n)$ .

## 12.3. Funcții de dispersie

În această secțiune vom discuta unele aspecte legate de construirea de funcții de dispersie bune, apoi vom prezenta trei scheme pentru crearea lor: dispersia prin diviziune, dispersia prin înmulțire și dispersia universală.

### Ce înseamnă o funcție de dispersie bună?

O funcție de dispersie bună satisface (aproximativ) ipoteza dispersiei uniforme simple: fiecare cheie se poate dispersa cu aceeași probabilitate în oricare dintre cele  $m$  locații. Mai formal, să presupunem că fiecare cheie este aleasă independent din  $U$  conform unei distribuții de probabilitate  $P$ ; deși,  $P(k)$  este probabilitatea de a fi aleasă cheia  $k$ . Atunci ipoteza dispersiei uniforme simple constă în

$$\sum_{k:h(k)=j} P(k) = \frac{1}{m} \text{ pentru } j = 0, 1, \dots, m-1. \quad (12.1)$$

Din păcate, în general, nu este posibilă verificarea acestei condiții, deoarece de regulă distribuția  $P$  nu este cunoscută.

Uneori (relativ rar) distribuția  $P$  este cunoscută. De exemplu, presupunem că se cunoaște faptul că cheile sunt numere reale aleatoare  $k$ , independent și uniform distribuite în intervalul  $0 \leq k < 1$ . În acest caz se poate arăta că funcția de dispersie

$$h(k) = \lfloor km \rfloor$$

satisface relația (12.1).

În practică, se pot folosi tehnici euristice pentru a crea funcții de dispersie care par a se comporta bine. Informația calitativă despre  $P$  este uneori utilă în procesul de construcție al lor. De exemplu, să considerăm o tabelă de simboluri a unui compilator, în care cheile sunt șiruri de caractere arbitrare, reprezentând identificatori dintr-un program. O situație des întâlnită este aceea în care simboluri foarte asemănătoare, ca **pt** și **pts**, apar în același program. O funcție de dispersie bună va minimiza posibilitatea ca asemenea variații să se disperseze în aceeași locație.

O abordare uzuală este de a obține valoarea de dispersie într-un mod care se vrea independent de orice model sau șablon ce poate exista între date. De exemplu, “metoda diviziunii” (discutată în detaliu mai jos) calculează valoarea de dispersie ca fiind restul împărțirii cheii la un număr prim specificat. În afara cazului în care numărul este într-un fel dependent de șabloanele din distribuția de probabilitate  $P$ , această metodă dă rezultate bune.

În final, să remarcăm faptul că unele aplicații ale funcțiilor de dispersie pot impune proprietăți mai tari decât cele asigurate de dispersia uniformă simplă. De exemplu, am putea dori ca unele chei care sunt “apropiate” într-un anumit sens să producă valori de dispersie care să fie total diferite. (Această proprietate este dorită în special când folosim verificarea liniară, definită în secțiunea 12.4.)

## Interpretarea cheilor ca numere naturale

Majoritatea funcțiilor de dispersie presupun universul cheilor din mulțimea  $\mathbb{N} = \{0, 1, 2, \dots\}$  a numerelor naturale. Astfel, dacă cheile nu sunt numere naturale, trebuie găsită o modalitate pentru a le interpreta ca numere naturale. De exemplu, o cheie care este un șir de caractere poate fi interpretată ca un întreg într-o bază de numerație aleasă convenabil. Prin urmare, identificatorul **pt** poate fi interpretat ca o pereche de numere zecimale (112,116), pentru că **p** = 112 și **t** = 116 în mulțimea codurilor ASCII; atunci **pt** exprimat ca un întreg în baza 128 devine  $(112 \cdot 128) + 116 = 14452$ . În mod obișnuit, în orice aplicație se poate crea direct o astfel de metodă simplă de a interpreta fiecare cheie ca un număr natural (posibil mare). În continuare, vom presupune că avem chei numere naturale.

### 12.3.1. Metoda diviziunii

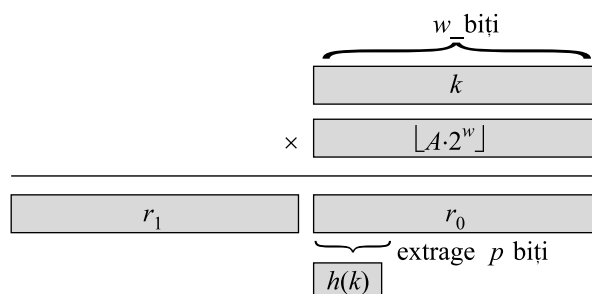
Prin *metoda diviziunii* pentru crearea funcțiilor de dispersie, transformăm o cheie  $k$  într-una din cele  $m$  locații considerând restul împărțirii lui  $k$  la  $m$ . Prin urmare, funcția de dispersie este

$$h(k) = k \bmod m.$$

De exemplu, dacă tabela de dispersie are dimensiunea  $m = 12$  și cheia este  $k = 100$ , atunci  $h(k) = 4$ . Deoarece necesită doar o singură operație de împărțire, dispersia prin diviziune este rapidă.

Când folosim metoda diviziunii, de obicei evităm anumite valori ale lui  $m$ . De exemplu,  $m$  nu ar trebui să fie o putere a lui 2, pentru că dacă  $m = 2^p$ , atunci  $h(k)$  reprezintă doar primii  $p$  biți ai lui  $k$ . În afara cazului în care se cunoaște apriori că distribuția de probabilitate pe chei produce cu aceeași probabilitate oricare dintre șabloanele primilor  $p$  biți, este mai bine să construim funcția de dispersie ca dependentă de toți biții cheii. Puterile lui 10 ar trebui evitate dacă aplicația lucrează cu numere zecimale ca și chei, pentru că în acest caz funcția de dispersie





**Figura 12.4** Metoda înmulțirii pentru dispersie. Reprezentarea pe  $w$  biți a cheii  $k$  este înmulțită cu valoarea pe  $w$  biți  $\lfloor A \cdot 2^w \rfloor$ , unde  $0 < A < 1$  este o constantă aleasă convenabil. Primii  $p$  biți ai celei de a doua jumătăți a produsului, de lungime  $w$  biți, formează valoarea de dispersie  $h(k)$  dorită.

nu depinde de toate cifrele zecimale ale lui  $k$ . În final, se poate arăta că dacă  $m = 2^p - 1$  și  $k$  este un șir de caractere interpretat în baza  $2^p$ , două șiruri de caractere care sunt identice, exceptând o transpoziție a două caractere adiacente, se vor dispersa în aceeași valoare.

Valori bune pentru  $m$  sunt numerele prime nu prea apropiate de puterile exacte ale lui 2. De exemplu, să presupunem că dorim să alocăm o tabelă de dispersie, cu coliziunile rezolvate prin înlanțuire, pentru a reține aproximativ  $n = 2000$  șiruri de caractere, unde un caracter are 8 biți. Nu ne deranjează să examinăm în medie 3 elemente într-o căutare fără succes, deci vom alocă o tabelă de dispersie de dimensiune  $m = 701$ . Numărul 701 este ales pentru că este un număr prim apropiat de  $2000/3$ , dar nu este apropiat de nici o putere a lui 2. Tratatând fiecare cheie  $k$  ca un întreg, funcția noastră de dispersie va fi

$$h(k) = k \bmod 701.$$

Ca o măsură de precauție, putem să verificăm cât de uniform distribuie această funcție de dispersie mulțimi de chei în locații, unde cheile sunt alese din date “reale”.

### 12.3.2. Metoda înmulțirii

**Metoda înmulțirii** pentru crearea de funcții de dispersie operează în doi pași. În primul pas, înmulțim cheia  $k$  cu o constantă  $A$  din intervalul  $0 < A < 1$  și extragem partea fracționară a lui  $kA$ . Apoi, înmulțim această valoare cu  $m$  și considerăm partea întreagă inferioară a rezultatului. Pe scurt, funcția de dispersie este

$$h(k) = \lfloor m(kA \bmod 1) \rfloor,$$

unde “ $kA \bmod 1$ ” înseamnă partea fracționară a lui  $kA$ , adică,  $kA - \lfloor kA \rfloor$ .

Un avantaj al metodei înmulțirii este că valoarea lui  $m$  nu este critică. De obicei o alegem ca fiind o putere a lui 2 –  $m = 2^p$  pentru un întreg  $p$  – pentru că atunci putem implementa funcția pe majoritatea calculatoarelor după cum urmează. Să presupunem că lungimea cuvântului mașinii este de  $w$  biți și că  $k$  încapă pe un singur cuvânt. Referindu-ne la figura 12.4, înmulțim întâi  $k$  cu întregul pe  $w$  biți  $\lfloor A \cdot 2^w \rfloor$ . Rezultatul este o valoare pe  $2w$  biți  $r_1 2^w + r_0$ , unde  $r_1$  este primul cuvânt al produsului și  $r_0$  este al doilea cuvânt al produsului. Valoarea de dispersie de  $p$  biți dorită constă din cei mai semnificativi  $p$  biți ai lui  $r_0$ .

Deși această metodă funcționează cu orice valoare a constantei  $A$ , ea lucrează mai bine cu anumite valori decât cu altele. Alegerea optimă depinde de caracteristicile datelor care sunt dispersate. Knuth [123] discută alegerea lui  $A$  în detaliu și sugerează că

$$A = (\sqrt{5} - 1)/2 \approx 0.6180339887... \quad (12.2)$$

se va comporta relativ bine.

De exemplu, dacă avem  $k = 123456$ ,  $m = 10000$  și  $A$  din relația (12.2), atunci

$$\begin{aligned} h(k) &= \lfloor 10000 \cdot (123456 \cdot 0.61803... \bmod 1) \rfloor = \lfloor 10000 \cdot (76300.0041151... \bmod 1) \rfloor = \\ &= \lfloor 10000 \cdot 0.0041151... \rfloor = \lfloor 41.151... \rfloor = 41. \end{aligned}$$

### 12.3.3. Dispersia universală

Dacă cheile ce se dispersează sunt alese de un adversar malițios, atunci el poate alege  $n$  chei care toate se vor dispersa în aceeași locație, rezultând un timp mediu de acces  $\Theta(n)$ . În cazul cel mai defavorabil, orice funcție de dispersie fixată este vulnerabilă la acest tip de comportament; singura modalitate eficientă de a îmbunătăți această situație este de a alege funcția de dispersie în mod *aleator* astfel încât alegerea să fie *independent* de cheile care se vor memora. Această abordare, numită **dispersie universală**, produce în medie o performanță bună, indiferent de ce chei sunt alese de către adversar.

Ideea care stă la baza dispersiei universale este de a selecta funcția de dispersie în mod aleator în momentul execuției dintr-o clasă de funcții construită cu atenție. Ca și în cazul sortării rapide, randomizarea asigură că nici o intrare nu va produce comportamentul în cazul cel mai defavorabil. Datorită randomizării, algoritmul se poate comporta diferit la fiecare execuție, chiar pentru o aceeași intrare. Această abordare asigură o performanță bună pentru cazul mediu, indiferent de cheile care sunt date ca intrare. Revenind la exemplul cu tabela de simboluri a unui compilator, descoperim că alegerea identificatorilor de către programator nu poate cauza, din punct de vedere al consistenței, performanțe slabe la dispersie. Performanța slabă apare doar atunci când compilatorul alege o funcție de dispersie aleatoare care provoacă o dispersie slabă a mulțimii identificatorilor, dar probabilitatea ca această situație să apară este mică și este aceeași pentru orice mulțime de identificatori de aceeași mărime.

Fie  $\mathcal{H}$  o colecție finită de funcții de dispersie care transformă un univers dat  $U$  al cheilor, în domeniul  $\{0, 1, \dots, m-1\}$ . O astfel de colecție se numește **universală** dacă pentru fiecare pereche de chei distincte  $x, y \in U$ , numărul de funcții de dispersie  $h \in \mathcal{H}$  pentru care  $h(x) = h(y)$  este exact  $|\mathcal{H}|/m$ . Cu alte cuvinte, cu o funcție de dispersie aleasă aleator din  $\mathcal{H}$ , șansa unei coliziuni între  $x$  și  $y$  când  $x \neq y$  este exact  $1/m$ , care este exact șansa unei coliziuni dacă  $h(x)$  și  $h(y)$  sunt alese aleator din mulțimea  $\{0, 1, \dots, m-1\}$ .

Următoarea teoremă arată că o clasă universală de funcții de dispersie dă un comportament bun în cazul mediu.

**Teorema 12.3** Dacă  $h$  este aleasă dintr-o colecție universală de funcții de dispersie și este folosită pentru a dispersa  $n$  chei într-o tabelă de dimensiune  $m$ , unde  $n \leq m$ , numărul mediu de coliziuni în care este implicată o cheie particulară  $x$ , este mai mic decât 1.

**Demonstrație.** Pentru fiecare pereche  $y, z$  de chei distincte, fie  $c_{yz}$  variabila aleatoare care are valoarea 1 dacă  $h(y) = h(z)$  (adică, dacă  $y$  și  $z$  sunt în coliziune folosind  $h$ ) și 0 în caz contrar.

Deoarece, prin definiție, o singură pereche de chei sunt în coliziune cu probabilitatea  $1/m$ , avem

$$E[c_{yz}] = 1/m.$$

Fie  $C_x$  numărul total de coliziuni în care este implicată cheia  $x$  dintr-o tabelă de dispersie  $T$  de dimensiune  $m$  conținând  $n$  chei. Din relația (6.24) obținem

$$E[C_x] = \sum_{\substack{y \in T \\ y \neq x}} E[c_{xy}] = \frac{n-1}{m}.$$

Deoarece  $n \leq m$ , obținem  $E[C_x] < 1$ . ■

Dar cât de ușor este să construim o clasă universală de funcții de dispersie? Este relativ ușor, după cum se poate demonstra folosind unele cunoștințe din teoria numerelor. Să alegem dimensiunea  $m$  a tabelii noastre ca fiind număr prim (ca și în metoda diviziunii). Descompunem o cheie  $x$  în  $r+1$  octeți (de exemplu, caractere sau subșiruri binare de caractere de lungime fixă), astfel încât  $x = \langle x_0, x_1, \dots, x_r \rangle$ ; singura cerință este ca valoarea maximă a unui octet să fie mai mică decât  $m$ . Notăm prin  $a = \langle a_0, a_1, \dots, a_r \rangle$  un șir de  $r+1$  elemente alese aleator din mulțimea  $\{0, 1, \dots, m-1\}$ . Definim o funcție de dispersie corespunzătoare  $h_a \in \mathcal{H}$  astfel:

$$h_a(x) = \sum_{i=0}^r a_i x_i \pmod{m}. \quad (12.3)$$

Cu această definiție,

$$\mathcal{H} = \bigcup_a \{h_a\}. \quad (12.4)$$

are  $m^{r+1}$  elemente.

**Teorema 12.4** Clasa  $\mathcal{H}$  definită de relațiile (12.3) și (12.4) este o clasă universală de funcții de dispersie.

**Demonstrație.** Considerăm o pereche oarecare de chei distincte  $x$  și  $y$ . Presupunem că  $x_0 \neq y_0$ . (O argumentație similară se poate face pentru diferența dintre oricare alte poziții ale octeților.) Pentru orice valori fixe ale lui  $a_1, a_2, \dots, a_r$  există exact o valoare a lui  $a_0$  care satisface ecuația  $h(x) = h(y)$ ; acest  $a_0$  este soluția ecuației

$$a_0(x_0 - y_0) \equiv - \sum_{i=1}^r a_i(x_i - y_i) \pmod{m}.$$

Pentru a observa această proprietate să remarcăm faptul că deoarece  $m$  este prim, cantitatea nenulă  $x_0 - y_0$  are un invers față de înmulțirea modulo  $m$  și astfel există o soluție unică pentru  $a_0$  modulo  $m$ . (Vezi secțiunea 33.4.) Prin urmare, fiecare pereche de chei  $x$  și  $y$  este în coliziune pentru exact  $m^r$  valori ale lui  $a$ , deoarece ea este în coliziune exact o dată pentru fiecare valoare posibilă a lui  $\langle a_1, a_2, \dots, a_r \rangle$  (de exemplu, pentru valoarea unică a lui  $a_0$  observată mai sus). Deoarece există  $m^{r+1}$  valori posibile pentru secvența  $a$ , cheile  $x$  și  $y$  sunt în coliziune exact cu probabilitatea  $m^r/m^{r+1} = 1/m$ . În concluzie,  $\mathcal{H}$  este universală. ■

## Exerciții

**12.3-1** Presupunem că dorim să efectuăm o căutare într-o listă înlănțuită de lungime  $n$ , în care fiecare element conține o cheie  $k$  împreună cu o valoare de dispersie  $h(k)$ . Fiecare cheie este un șir lung de caractere. Cum putem profita de valorile de dispersie când căutăm în listă un element cu o cheie dată?

**12.3-2** Presupunem că un șir de  $r$  caractere este dispersat pe  $m$  poziții, fiind considerat ca număr în baza 128, căruia i se aplică metoda diviziunii. Numărul  $m$  este ușor de reprezentat în memorie ca un cuvânt pe 32 de biți, dar șirul de  $r$  caractere, tratat ca un număr în baza 128, necesită multe cuvinte. Cum putem aplica metoda diviziunii pentru a calcula valoarea de dispersie a șirului de caractere fără a folosi mai mult decât un număr suplimentar constant de cuvinte de memorie în afară de șirul propriu-zis?

**12.3-3** Se consideră o versiune a metodei diviziunii în care  $h(k) = k \bmod m$ , unde  $m = 2^p - 1$  și  $k$  este un șir de caractere interpretat în baza  $2^p$ . Arătați că dacă șirul de caractere  $x$  poate fi obținut din șirul de caractere  $y$  prin permutări de caractere, atunci  $x$  și  $y$  se dispersează în aceeași valoare. Dați un exemplu de aplicație în care această proprietate nu este dorită pentru o funcție de dispersie.

**12.3-4** Se consideră o tabelă de dispersie de dimensiune  $m = 1000$  și funcția de dispersie  $h(k) = \lfloor m(kA \bmod 1) \rfloor$  pentru  $A = (\sqrt{5} - 1)/2$ . Calculați locațiile în care se pun cheile 61, 62, 63, 64, și 65.

**12.3-5** Arătați că dacă punem restricția ca fiecare componentă  $a_i$  din  $a$  din relația (12.3) să fie nenulă, atunci mulțimea  $\mathcal{H} = \{h_a\}$  definită ca în relația (12.4) nu este universală. (*Indica ie:* Se consideră cheile  $x = 0$  și  $y = 1$ .)

---

## 12.4. Adresarea deschisă

Prin **adresare deschisă**, toate elementele sunt memorate în interiorul tabelii de dispersie. Prin urmare, fiecare intrare în tabelă conține fie un element al mulțimii dinamice, fie NIL. Când căutăm un element, vom examina sistematic locațiile tabelii fie până când este găsit elementul dorit, fie până când este clar că elementul nu este în tabelă. Nu există liste sau elemente memorate în afara tabelii, așa cum se întâmplă în cazul înlănțuirii. Prin urmare, prin adresare deschisă, tabelă de dispersie se poate “umple” astfel încât nici o inserare nu mai este posibilă; factorul de încărcare  $\alpha$  nu poate depăși niciodată valoarea 1.

Desigur, am putea memora listele înlănțuite pentru a crea legăturile în cadrul tabelii de dispersie, în locațiile altfel neutilizate ale acesteia (vezi exercițiul 12.2-5), dar avantajul adresării deschise este că evită total folosirea pointerilor. Secvența de locații care se examinează nu se determină folosind pointerii, ci se *calculează*. Spațiul de memorie suplimentar, eliberat prin faptul că nu memorăm pointerii, oferă tabelii de dispersie un număr mai mare de locații pentru același spațiu de memorie, putând rezulta coliziuni mai puține și acces mai rapid.

Pentru a realiza inserarea folosind adresarea deschisă, examinăm succesiv sau **verificăm** tabela de dispersie până când găsim o locație liberă în care să punem cheia. În loc să fie fixat în ordinea  $0, 1, \dots, m - 1$  (care necesită un timp de căutare  $\Theta(n)$ ), șirul de poziții examinate

depinde de cheia ce se inserează. Pentru a determina ce locații sunt verificate, extindem funcția de dispersie astfel încât ea să includă și numărul de verificare (începând de la 0) ca un al doilea argument. Astfel, funcția de dispersie devine

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

În cadrul adresării deschise, cerem ca pentru fiecare cheie  $k$ , **secvența de verificare**

$$\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$$

să fie o permutare a lui  $\langle 0, 1, \dots, m-1 \rangle$ , astfel încât fiecare poziție din tabela de dispersie să fie considerată ca o eventuală locație pentru o nouă cheie pe măsură ce se umple tabela. În algoritmul următor, presupunem că elementele din tabela de dispersie  $T$  sunt chei, fără informații adiționale; cheia  $k$  este identică cu elementul care conține cheia  $k$ . Fiecare locație conține fie o cheie, fie NIL (dacă locația este liberă).

DISPERSIE-INSEREAZĂ( $T, k$ )

```

1:  $i \leftarrow 0$ 
2: repetă
3:    $j \leftarrow h(k, i)$ 
4:   dacă  $T[j] = \text{NIL}$  atunci
5:      $T[j] \leftarrow k$ 
6:     returnează  $j$ 
7:   altfel
8:      $i \leftarrow i + 1$ 
9:   până când  $i = m$ 
10: eroare “depășire tabela de dispersie”
```

Algoritmul pentru căutarea unei chei  $k$  examinează aceeași secvență de locații pe care o folosește și algoritmul de inserare atunci când s-a inserat cheia  $k$ . Prin urmare, căutarea se poate termina (fără succes) când se întâlnește o locație liberă, pentru că  $k$  ar fi fost inserat în acea locație și nu mai încolo în secvența de verificare. (Se observă că această argumentație presupune că o dată introduse, cheile nu mai sunt șterse din tabela de dispersie.) Procedura DISPERSIE-CAUTĂ are ca intrare o tabelă de dispersie  $T$  și o cheie  $k$  și returnează  $j$  dacă locația  $j$  conține cheia  $k$  sau NIL dacă cheia  $k$  nu există în tabela  $T$ .

DISPERSIE-CAUTĂ( $T, k$ )

```

1:  $i \leftarrow 0$ 
2: repetă
3:    $j \leftarrow h(k, i)$ 
4:   dacă  $T[j] = k$  atunci
5:     returnează  $j$ 
6:    $i \leftarrow i + 1$ 
7: până când  $T[j] = \text{NIL}$  sau  $i = m$ 
8: returnează NIL
```

Ștergerea dintr-o tabelă de dispersie cu adresare deschisă este dificilă. Când ștergem o cheie dintr-o locație  $i$ , nu putem marca pur și simplu acea locație ca fiind liberă memorând în ea

valoarea NIL. Procedând astfel, va fi imposibil să accesăm orice cheie  $k$  a cărei inserare a verificat locația  $i$  și a găsit-o ocupată. O soluție este de a marca locația, memorând în ea valoarea specială ȘTERS în loc de NIL. Apoi vom modifica procedura DISPERSIE-INSEREAZĂ astfel încât să trateze astfel de locații ca și când ar fi libere, astfel încât o nouă cheie să poată fi inserată. Nu este necesară nici o modificare în DISPERSIE-CAUTĂ, deoarece va trece peste valorile ȘTERS în timpul căutării. Procedând astfel, timpii de căutare nu mai depind de factorul de încărcare  $\alpha$  și din acest motiv în cazul în care cheile trebuie șterse, se preferă, în mod uzual, înlănțuirea ca tehnică de rezolvare a coliziunilor.

În analiza noastră, folosim ipoteza **dispersiei uniforme**: presupunem că fiecare cheie considerată poate avea, cu aceeași probabilitate, oricare dintre cele  $m!$  permutări ale mulțimii  $\{0, 1, \dots, m-1\}$  ca secvență de verificare. Dispersia uniformă generalizează noțiunea de dispersie uniformă simplă, definită anterior, pentru situația în care funcția de dispersie produce nu doar un singur număr, ci o întreagă secvență de verificare. Oricum, dispersia uniformă reală este dificil de implementat și în practică sunt folosite aproximări convenabile (ca dispersia dublă, definită mai târziu).

În mod obișnuit sunt utilizate trei tehnici pentru a calcula secvențele de verificare necesare adresării deschise: verificarea liniară, verificarea pătratică și dispersia dublă. Toate aceste tehnici garantează că  $\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$  este o permutare a lui  $\langle 0, 1, \dots, m-1 \rangle$  pentru fiecare cheie  $k$ . Nici una din aceste tehnici nu satisface condiția dispersiei uniforme, pentru că nici una dintre ele nu e capabilă să genereze mai mult de  $m^2$  secvențe de verificare diferite (în loc de  $m!$  cât cere dispersia uniformă). Dispersia dublă are cel mai mare număr de secvențe de verificare și, după cum era de așteptat, dă cele mai bune rezultate.

## Verificarea liniară

Fiind dată o funcție de dispersie ordinară  $h' : U \rightarrow \{0, 1, \dots, m-1\}$ , metoda **verificării liniare** folosește funcția de dispersie

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

pentru  $i = 0, 1, \dots, m-1$ . Fiind dată o cheie  $k$ , prima locație examinată (verificată) este  $T[h'(k)]$ . Apoi examinăm locația  $T[h'(k) + 1]$  și așa mai departe până la locația  $T[m-1]$ . Apoi continuăm circular cu locațiile  $T[0], T[1], \dots$ , până când, în final, verificăm locația  $T[h'(k) - 1]$ . Deoarece poziția de start a verificării determină întreaga secvență de verificare liniară, prin care sunt folosite doar  $m$  secvențe de verificări distincte.

Verificarea liniară este ușor de implementat, dar apare o problemă cunoscută sub numele de **grupare primară**. În acest caz se formează șiruri lungi de locații ocupate, crescând timpul mediu de căutare. De exemplu, dacă avem  $n = m/2$  chei într-o tabelă, în care fiecare locație de indice par este ocupată și fiecare locație de indice impar este liberă, atunci căutarea fără succes necesită, 1,5 verificări în cazul mediu. Dacă primele  $n = m/2$  locații sunt cele ocupate, numărul mediu de verificări crește la aproximativ  $n/4 = m/8$ . Grupările au probabilitate mare de apariție, pentru că dacă o locație liberă este precedată de  $i$  locații ocupate, atunci probabilitatea ca locația liberă să fie următoarea completată este  $(i+1)/m$ , comparativ cu probabilitatea de  $1/m$  pentru cazul în care locația precedentă ar fi fost liberă. Prin urmare, șirurile de locații ocupate tind să se lungească și verificarea liniară nu este o aproximare foarte bună a dispersiei uniforme.

## Verificarea pătratică

**Verificarea pătratică** folosește o funcție de dispersie de forma

$$h(k, i) = (h'(k) + c_1 i + c_2 i^2) \bmod m, \quad (12.5)$$

unde (ca și în cazul verificării liniare)  $h'$  este o funcție de dispersie auxiliară,  $c_1$  și  $c_2 \neq 0$  sunt constante auxiliare și  $i = 0, 1, \dots, m-1$ . Poziția verificată inițial este  $T[h'(k)]$ ; următoarele poziții examinate sunt decalate cu cantități ce depind într-o manieră pătratică de numărul de verificare  $i$ . Această metodă lucrează mult mai bine decât verificarea liniară, dar pentru a folosi integral tabela de dispersie valorile lui  $c_1, c_2$  și  $m$  trebuie determinate corespunzător. Problema 12-4 ilustrează o modalitate de determinare a acestor parametri. De asemenea, dacă două chei au aceeași poziție de start a verificării, atunci secvențele lor de verificare coincid, pentru că  $h(k_1, 0) = h(k_2, 0)$  implică  $h(k_1, i) = h(k_2, i)$ . Această situație conduce la o formă mai ușoară de grupare, numită **grupare secundară**. Ca și în cazul verificării liniare, verificarea inițială determină întreaga secvență, deci sunt folosite doar  $m$  secvențe de verificare distincte.

## Dispersia dublă

dispersia dublă este una dintre cele mai bune metode disponibile pentru adresarea deschisă, deoarece permutările produse au multe din caracteristicile permutărilor alese aleator. **dispersia dublă** folosește o funcție de dispersie de forma

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

unde  $h_1$  și  $h_2$  sunt funcții de dispersie auxiliare. Poziția examinată inițial este  $T[h_1(k)]$ ; pozițiile succesive de verificare sunt decalate față de pozițiile anterioare cu  $h_2(k)$  modulo  $m$ . Astfel, contrar situației de verificare liniară sau pătratică, în acest caz secvența de verificare depinde în două moduri de cheia  $k$ , pentru că fie poziția inițială de verificare, fie decalajul, fie amândouă pot varia. Figura 12.5 ilustrează un exemplu de inserare prin dispersie dublă.

Valoarea  $h_2(k)$  și dimensiunea tabelii de dispersie  $m$  trebuie să fie prime între ele pentru a fi parcursă întreaga tabelă de dispersie. În caz contrar, dacă  $m$  și  $h_2(k)$  au un cel mai mare divizor comun  $d > 1$  pentru o cheie  $k$ , atunci o căutare pentru cheia  $k$  va examina doar a  $(1/d)$ -a parte din tabela de dispersie. (Vezi capitolul 33.) O modalitate convenabilă de a asigura această condiție este de a avea un număr  $m$  ca o putere a lui 2 și de a construi  $h_2$  astfel încât să producă întotdeauna un număr impar. O altă modalitate este de a avea  $m$  prim și de a construi  $h_2$  astfel încât să returneze întotdeauna un întreg pozitiv mai mic decât  $m$ . De exemplu, putem alege  $m$  prim și

$$\begin{aligned} h_1(k) &= k \bmod m, \\ h_2(k) &= 1 + (k \bmod m'), \end{aligned}$$

unde  $m'$  este ales să fie un pic mai mic decât  $m$  (să zicem  $m-1$  sau  $m-2$ ). De exemplu, dacă  $k = 123456$ ,  $m = 701$  și  $m' = 700$ , avem  $h_1(k) = 80$  și  $h_2(k) = 257$ , deci prima verificare se află la poziția 80 și fiecare a 257-a locație (modulo  $m$ ) este examinată până când cheia este găsită sau a fost examinată toată tabela.

dispersia dublă reprezintă o îmbunătățire față de verificarea liniară sau pătratică în sensul că sunt folosite  $\Theta(m^2)$  secvențe de verificare, față de  $\Theta(m)$ , pentru că fiecare pereche posibilă  $(h_1(k), h_2(k))$  produce o secvență de verificare distinctă c si când cheia variază, poziția inițială

0	
1	79
2	
3	
4	69
5	98
6	
7	72
8	
9	14
10	
11	50
12	

**Figura 12.5** Inserarea prin dispersie dublă. Avem o tabelă de dispersie de dimensiune 13, cu  $h_1(k) = k \bmod 13$  și  $h_2(k) = 1 + (k \bmod 11)$ . Deoarece  $14 \equiv 1 \bmod 13$  și  $14 \equiv 3 \bmod 11$ , cheia 14 va fi inserată în locația liberă 9, după ce locațiile 1 și 5 au fost examinate și găsite ca fiind deja ocupate.

a verificării  $h_1(k)$  și decalajul  $h_2(k)$  pot varia independent. Ca rezultat, performanța dispersiei duble apare ca fiind foarte apropiată de performanța schemei “ideale” a dispersiei uniforme.

### Analiza dispersiei cu adresare deschisă

Analiza noastră pentru adresarea deschisă, este, ca și analiza pentru înlănțuire, exprimată în termenii factorului de încărcare  $\alpha$  al tabelii de dispersie, când  $n$  și  $m$  tind spre infinit. Să reamintim că dacă sunt memorate  $n$  elemente într-o tabelă cu  $m$  locații, numărul mediu de elemente pe locație este  $\alpha = n/m$ . Desigur, prin adresare deschisă, avem cel mult un element într-o locație și, prin urmare,  $n \leq m$ , ceea ce implică  $\alpha \leq 1$ .

Presupunem că este folosită dispersia uniformă. În această schemă idealizată, secvența de verificare  $\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$  pentru orice cheie  $k$  poate să apară sub forma oricărei permutări a mulțimii  $\langle 0, 1, \dots, m-1 \rangle$ . Aceasta înseamnă că fiecare secvență de verificare posibilă poate fi folosită ca secvență de verificare pentru o inserare sau o căutare. Desigur, o cheie dată are o secvență de verificare unică, fixă, asociată ei; considerând distribuția de probabilitate pe spațiul cheilor și operația funcției de dispersie pe chei, fiecare secvență de verificare posibilă are aceeași probabilitate de apariție.

Analizăm acum numărul mediu de verificări pentru dispersia cu adresare deschisă, în ipoteza dispersiei uniforme și începem cu o analiză a numărului de verificări ce apar într-o căutare fără succes.

**Teorema 12.5** Fiind dată o tabelă de dispersie cu adresare deschisă cu factorul de încărcare  $\alpha = n/m < 1$ , în ipoteza dispersiei uniforme, numărul mediu de verificări dintr-o căutare fără succes este cel mult  $1/(1 - \alpha)$ .

**Demonstrație.** Într-o căutare fără succes, fiecare verificare, cu excepția ultimei, accesează o



locație ocupată care nu conține cheia dorită, iar ultima locație examinată este liberă. Definim

$$p_i = \Pr \{ \text{exact } i \text{ verificări accesează locațiile ocupate} \}$$

pentru  $i = 0, 1, 2, \dots$ . Pentru  $i > n$ , avem  $p_i = 0$ , pentru că putem găsi cel mult  $n$  locații deja ocupate. Astfel, numărul mediu de verificări este

$$1 + \sum_{i=1}^{\infty} ip_i \quad (12.6)$$

Pentru a evalua relația (12.6), definim

$$q_i = \Pr \{ \text{cel puțin } i \text{ verificări accesează locațiile ocupate} \}$$

pentru  $i = 0, 1, 2, \dots$ . Apoi putem folosi identitatea (6.28):

$$\sum_{i=1}^{\infty} ip_i = \sum_{i=1}^{\infty} q_i .$$

Care este valoarea lui  $q_i$  pentru  $i \geq 1$ ? Probabilitatea ca prima verificare să acceseze o locație ocupată este  $n/m$ ; prin urmare,

$$q_1 = \frac{n}{m} .$$

Prin dispersie uniformă, a doua verificare, în cazul în care este necesară, va fi la una din cele  $m-1$  locații rămase neverificate, din care  $n-1$  sunt ocupate. Facem a doua verificare doar dacă prima verificare accesează o locație ocupată; prin urmare,

$$q_2 = \left( \frac{n}{m} \right) \left( \frac{n-1}{m-1} \right) .$$

În general, a  $i$ -a verificare este necesară doar dacă primele  $i-1$  verificări accesează locații ocupate și locația examinată este cu aceeași probabilitate oricare dintre cele  $m-i+1$  locații rămase, din care  $n-i+1$  sunt ocupate. Deci,

$$q_i = \left( \frac{n}{m} \right) \left( \frac{n-1}{m-1} \right) \dots \left( \frac{n-i+1}{m-i+1} \right) \leq \left( \frac{n}{m} \right)^i = \alpha^i$$

pentru  $i = 1, 2, \dots, n$ , pentru că  $(n-j)/(m-j) \leq n/m$  dacă  $n \leq m$  și  $j \geq 0$ . După  $n$  verificări, toate cele  $n$  locații ocupate au fost vizitate și nu vor fi verificate din nou și astfel  $q_i = 0$  pentru  $i > n$ .

Acum se poate evalua relația (12.6). Presupunând  $\alpha < 1$ , numărul mediu de verificări într-o căutare fără succes este

$$1 + \sum_{i=1}^{\infty} ip_i = 1 + \sum_{i=1}^{\infty} q_i \leq 1 + \alpha + \alpha^2 + \alpha^3 + \dots = \frac{1}{1-\alpha} . \quad (12.7)$$

Relația (12.7) are o interpretare intuitivă: prima verificare se face întotdeauna, a doua verificare este necesară cu o probabilitate aproximativ egală cu  $\alpha$ , a treia verificare este necesară cu o probabilitate aproximativ egală cu  $\alpha^2$  și așa mai departe. ■

Dacă  $\alpha$  este o constantă, teorema 12.5 prognozează faptul că o căutare fără succes necesită un timp  $O(1)$ . De exemplu, dacă tabela de dispersie este pe jumătate plină, numărul mediu de verificări într-o căutare fără succes este cel mult  $1/(1 - .5) = 2$ . Dacă este 90% ocupată, numărul mediu de verificări este cel mult  $1/(1 - .9) = 10$ .

Teorema 12.5 ne dă, aproape imediat, performanța procedurii DISPERSIE-INSEREAZĂ.

**Corolarul 12.6** Inserarea unui element într-o tabelă de dispersie cu adresare deschisă, cu factorul de încărcare  $\alpha$  necesită, în cazul mediu, cel mult  $1/(1 - \alpha)$  verificări, în ipoteza dispersiei uniforme.

**Demonstrație.** Un element este inserat doar dacă există spațiu în tabelă și, deci,  $\alpha < 1$ . Inserarea unei chei necesită o căutare fără succes, urmată de o plasare a cheii în prima locație găsită liberă. Prin urmare, numărul mediu de verificări este cel mult  $1/(1 - \alpha)$ . ■

Calculul numărului mediu de verificări pentru o căutare cu succes necesită un pic mai mult de lucru.

**Teorema 12.7** Fiind dată o tabelă de dispersie cu adresare deschisă, cu factorul de încărcare  $\alpha < 1$ , numărul mediu de verificări într-o căutare cu succes este cel mult

$$\frac{1}{\alpha} \ln \frac{1}{1 - \alpha},$$

în ipoteza dispersiei uniforme și presupunând că fiecare cheie din tabelă este căutată cu aceeași probabilitate.

**Demonstrație.** O căutare pentru o cheie  $k$  urmează aceeași secvență de verificare ce a fost urmată când elementul având cheia  $k$  a fost inserat. Conform corolarului 12.6, dacă  $k$  a fost a  $i + 1$ -a cheie inserată în tabela de dispersie, numărul mediu de verificări făcute într-o căutare a lui  $k$  este cel mult  $1/(1 - i/m) = m/(m - i)$ . Calculând media peste toate cele  $n$  chei din tabela de dispersie, obținem numărul mediu de verificări dintr-o căutare cu succes:

$$\frac{1}{n} \sum_{i=0}^{n-1} \frac{m}{m-i} = \frac{m}{n} \sum_{i=0}^{n-1} \frac{1}{m-i} = \frac{1}{\alpha} (H_m - H_{m-n}),$$

unde  $H_i = \sum_{j=1}^i 1/j$  este al  $i$ -lea număr armonic (definit ca în relația (3.5)). Folosind marginile  $\ln i \leq H_i \leq \ln i + 1$  din relațiile (3.11) și (3.12), obținem

$$\begin{aligned} \frac{1}{\alpha} (H_m - H_{m-n}) &= \frac{1}{\alpha} \sum_{k=m-n+1}^m 1/k \leq \frac{1}{\alpha} \int_{m-n}^m (1/x) dx \\ &= \frac{1}{\alpha} \ln(m/(m-n)) = \frac{1}{\alpha} \ln(1/(1-\alpha)) \end{aligned}$$

ca mărginire pentru numărul mediu de verificări într-o căutare cu succes. ■

Dacă tabela de dispersie este pe jumătate ocupată, numărul mediu de verificări într-o căutare cu succes este mai mic decât 1,387. Dacă tabela de dispersie este 90% ocupată, numărul mediu de verificări este mai mic decât 2,559.

## Exerciții

**12.4-1** Se consideră că se inserează cheile 10, 22, 31, 4, 15, 28, 17, 88, 59 într-o tabelă de dispersie de lungime  $m = 11$  folosind adresarea deschisă cu funcția primară de dispersie  $h'(k) = k \bmod m$ . Ilustrați rezultatul inserării acestor chei folosind verificarea liniară, verificarea pătratică cu  $c_1 = 1$  și  $c_2 = 3$  și folosind dispersia dublă cu  $h_2(k) = 1 + (k \bmod (m - 1))$ .

**12.4-2** Scrieți un algoritm în pseudocod pentru DISPERSIE-ȘTERGE după descrierea din text și modificați DISPERSIE-INSEREAZĂ și DISPERSIE-CAUTĂ pentru a încorpora valoarea specială ȘTERS.

**12.4-3** ★ Să presupunem că folosim dispersia dublă pentru a rezolva coliziunile; mai exact, folosim funcția de dispersie  $h(k, i) = (h_1(k) + ih_2(k)) \bmod m$ . Arătați că secvența de verificare  $\langle h(k, 0), h(k, 1), \dots, h(k, m - 1) \rangle$  este o permutare a secvenței de locații  $\langle 0, 1, \dots, m - 1 \rangle$  dacă și numai dacă  $h_2(k)$  este relativ prim cu  $m$ . (*Indica ie:* Vezi capitolul 33.)

**12.4-4** Se consideră o tabelă de dispersie cu adresare deschisă și cu dispersie uniformă și factorul de încărcare  $\alpha = 1/2$ . Dați margini superioare pentru numărul mediu de verificări într-o căutare fără succes și pentru numărul mediu de verificări într-o căutare cu succes. Care este numărul mediu de verificări într-o căutare cu succes? Repetați calculele pentru factorii de încărcare  $3/4$  și  $7/8$ .

**12.4-5** ★ Să presupunem că inserăm  $n$  chei într-o tabelă de dispersie de dimensiune  $m$  folosind adresarea deschisă și dispersia uniformă. Notăm  $\text{cup}(n, m)$  probabilitatea de a nu apare nici o coliziune. Arătați că  $p(n, m) \leq e^{-n(n-1)/2m}$ . (*Indica ie:* Vezi relația (2.7).) Argumentați că atunci când  $n$  depășește  $\sqrt{m}$ , probabilitatea evitării coliziunilor descrește rapid către zero.

**12.4-6** ★ Se consideră o tabelă de dispersie cu adresare deschisă având factorul de încărcare  $\alpha$ . Determinați acea valoare nenulă a lui  $\alpha$  pentru care numărul mediu de verificări într-o căutare fără succes este de două ori mai mare decât numărul mediu de verificări dintr-o căutare cu succes. Folosiți estimarea  $(1/\alpha)\ln(1/(1-\alpha))$  pentru numărul de verificări necesare într-o căutare cu succes.

---

## Probleme

### 12-1 Marginea celei mai lungi verificări pentru dispersie

O tabelă de dispersie de dimensiune  $m$  este folosită pentru a memora  $n$  obiecte, cu  $n \leq m/2$ . Pentru rezolvarea coliziunilor se folosește adresarea deschisă.

- În ipoteza dispersiei uniforme, arătați că pentru  $i = 1, 2, \dots, n$ , probabilitatea ca a  $i$ -a inserare să necesite strict mai mult de  $k$  verificări este cel mult  $2^{-k}$ .
- Arătați că pentru  $i = 1, 2, \dots, n$ , probabilitatea ca a  $i$ -a inserare să necesite mai mult de  $2 \lg n$  verificări este cel mult  $1/n^2$ .

Fie variabila aleatoare  $X_i$ , semnificând numărul de verificări necesare pentru a  $i$ -a inserare. Ați arătat la punctul (b) că  $\Pr\{X_i > 2 \lg n\} \leq 1/n^2$ . Fie variabila aleatoare  $X = \max_{1 \leq i \leq n} X_i$  care semnifică numărul maxim de verificări necesare oricăreia din cele  $n$  inserări.

c. Arătați că  $\Pr\{X > 2 \lg n\} \leq 1/n$ .

d. Arătați că lungimea medie a celei mai lungi secvențe de verificare este  $E[X] = O(\lg n)$ .

### 12-2 Căutarea într-o mulțime statică

Se cere să se implementeze o mulțime dinamică având  $n$  elemente, în care cheile sunt numere. Mulțimea este statică (nu există operațiile INSEREAZĂ și ȘTERGE) și singura operație necesară este CAUTĂ. Se pune la dispoziție un timp arbitrar pentru a preprocesa cele  $n$  elemente, astfel încât operațiile de căutare să se execute rapid.

a. Arătați că procedura CAUTĂ poate fi implementată astfel încât, în cazul cel mai defavorabil să se execute într-un timp  $O(\lg n)$ , fără a folosi spațiu suplimentar față de cel necesar memorării elementelor.

b. Se consideră implementarea mulțimii printr-o tabelă de dispersie cu adresare deschisă având  $m$  locații care folosește dispersia uniformă. Care este spațiul minim necesar (de lungime  $m - n$ ) suplimentar pentru ca performanța, în cazul mediu a unei căutări fără succes, să fie cel puțin la fel de bună ca marginea de la punctul (a)? Răspunsul trebuie să fie o mărginire asimptotică a lui  $m - n$  în funcție de  $n$ .

### 12-3 Margini pentru mărimea locației la înlănțuire

Să presupunem că avem o tabelă de dispersie cu  $n$  locații, având coliziuni rezolvate prin înlănțuire și să presupunem că se inserează  $n$  chei în tabelă. Fiecare cheie se poate dispersa, cu aceeași probabilitate, în oricare dintre locații. Fie  $M$  numărul maxim de chei într-o locație, după ce toate cheile au fost inserate. Demonstrați că  $E[M]$ , valoarea medie a lui  $M$ , are o margine superioară egală cu  $O(\lg n / \lg \lg n)$ .

a. Fie  $Q_k$  probabilitatea ca un număr de  $k$  chei să se disperseze într-o locație particulară. Argumentați că probabilitatea  $Q_k$  este dată de

$$Q_k = \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k} \binom{n}{k}.$$

b. Fie  $P_k$  probabilitatea ca  $M = k$ , adică probabilitatea ca locația care conține cele mai multe chei să conțină  $k$  chei. Arătați că  $P_k \leq nQ_k$ .

c. Folosiți aproximarea Stirling, din relația (2.11), pentru a arăta că  $Q_k < e^k/k^k$ .

d. Arătați că există o constantă  $c > 1$  astfel încât  $Q_{k_0} < 1/n^3$  pentru  $k_0 = c \lg n / \lg \lg n$ . Deduceți că  $P_k < 1/n^2$  pentru  $k = c \lg n / \lg \lg n$ .

e. Arătați că

$$E[M] \leq \Pr \left\{ M > \frac{c \lg n}{\lg \lg n} \right\} \cdot n + \Pr \left\{ M \leq \frac{c \lg n}{\lg \lg n} \right\} \cdot \frac{c \lg n}{\lg \lg n}.$$

Deduceți că  $E[M] = O(\lg n / \lg \lg n)$ .

**12-4 Verificarea pătratică**

Să presupunem că avem o cheie  $k$ , care trebuie căutată într-o tabelă de dispersie având pozițiile  $0, 1, \dots, m-1$  și să presupunem că avem o funcție de dispersie  $h$  care transformă spațiul cheilor în mulțimea  $\{0, 1, \dots, m-1\}$ . Schema de căutare este următoarea.

1. Se calculează valoarea  $i \leftarrow h(k)$  și se inițializează  $j \leftarrow 0$ .
2. Se verifică poziția  $i$  pentru cheia dorită  $k$ . Dacă se găsește sau dacă poziția este liberă atunci căutarea se termină.
3. Se fac atribuirile  $j \leftarrow (i + j) \bmod m$  și se revine la pasul 2.

Se presupune că  $m$  este o putere a lui 2.

- a. Arătați că această schemă este o instanțiere a schemei generale de “verificare pătratică” evidențiind valorile specifice ale constantelor  $c_1$  și  $c_2$  din relația (12.5).
- b. Demonstrați că acest algoritm examinează fiecare poziție din tabelă în cazul cel mai defavorabil.

**12-5 dispersie  $k$ -universală**

Fie  $\mathcal{H} = \{h\}$  o clasă de funcții de dispersie în care fiecare  $h$  transformă universul  $U$  al cheilor în  $\{0, 1, \dots, m-1\}$ . Spunem că  $\mathcal{H}$  este  **$k$ -universală** dacă pentru fiecare secvență fixă de  $k$  chei distincte  $\langle x_1, x_2, \dots, x_k \rangle$  și pentru fiecare funcție  $h$  aleasă aleator din  $\mathcal{H}$ , secvența  $\langle h(x_1), h(x_2), \dots, h(x_k) \rangle$  este, cu aceeași probabilitate, oricare dintre cele  $m^k$  secvențe de lungime  $k$  având elemente din  $\{0, 1, \dots, m-1\}$ .

- a. Arătați că dacă  $\mathcal{H}$  este 2-universală atunci ea este universală.
- b. Arătați că clasa  $\mathcal{H}$  definită în secțiunea 12.3.3 nu este 2-universală.
- c. Arătați că dacă modificăm definiția lui  $\mathcal{H}$  din secțiunea 12.3.3 astfel încât fiecare funcție să conțină un termen constant  $b$ , adică dacă înlocuim  $h(x)$  cu

$$h_{a,b}(x) = \sum_{i=0}^r a_i x_i + b \bmod m,$$

atunci  $\mathcal{H}$  este 2-universală.

---

**Note bibliografice**

Knuth [123] și Gonnet [90] sunt referințe excelente pentru algoritmi de dispersie. Knuth îl creditează pe H. P. Luhn (1953) cu inventarea tabelor de dispersie, împreună cu metoda înlănțuirii pentru rezolvarea coliziunilor. Aproximativ în aceeași perioadă, G. M. Amdahl a avut ideea originală a adresării deschise.