

Laborator 12 – Tabele de dispersie

Marti 16-18

Tabelele de dispersie (hashtables) sunt o structura de date foarte utila daca vrem sa stocam o colectie de elemente si vrem sa verificam in timp cat mai rapid daca un element exista sau nu in multimea noastra. Este o structura de date foarte utila in aplicatii reale. O aplicatie pentru hashtables este un hashmap (un dictionar), in care pe langa faptul ca avem nevoie sa tinem obiectele din multimea noastra, fiecare obiect are si o valoare asociata, astfel incat vom retine perechi de tip (cheie, valoare), avand proprietatea ca trebuie sa avem chei unice. Putem sa ne gandim la un dictionar intre doua limbi, in care dorim sa aflam rapid care este traducerea pentru un cuvint. In cazul unui hashtable simplu, nu ne intereseaza decat sa aflam daca elementul exista sau nu (practic e ca un hashmap in care valoarea nu este relevanta). Diferenta de implementare dintre cele doua este minima, de aceea in limbaje precum C++ sau Java se ofera clase pentru hashmaps.

Un hashtable este de fapt un vector, in care fiecare element din vector stocheaza mai multe obiecte. Aceste elemente se numesc “buckets”. Practic, avem un vector de buckets, fiecare bucket avand mai multe obiecte, care pot fi stocate, de exemplu, intr-o lista.

Exista 3 operatii de baza – insert, remove si find.

Pentru a sti in ce bucket trebuie sa inseram, se foloseste o functie de hash. Aceasta preia un tip de date pe care vrem sa il stocam si intoarce un numar, folosind niste calcule ce tin cont de tipul nostru. De exemplu, pentru tipul de date string/ char*, o functie standard de hash este sa consideram cuvintul scris in baza 26 (daca avem doar litere mici ale alfabetului), si sa obtinem varianta echivalenta in baza 10. Deoarece numarul de bucketuri este limitat, si in general functiile de hash intorc numere foarte mari, se face modulo la fiecare operatie cu numarul total de bucketuri.

Ex: string = ‘abc’

Nr buckets = 50

$\text{hash}(\text{'abc'}) = (1 \cdot 26^2 + 2 \cdot 26 + 3) \% 50$

Exista multe functii de hashing, conditia ca sa functioneze fiind sa distribuie uniform obiectele intr-un numar finit de bucketuri. Daca doua obiecte au aceeasi functie de hash, ele vor trebui inserate in acelasi bucket, cauzand ceea ce se numeste o “coliziune”. Cu cat mai multe coliziuni avem, cu atat operatia de cautare este mai dificila si complexitatea creste. O functie de hash perfecta reuseste sa mapeze intrarile fara nicio coliziune.

Pseudocodul pentru inserare folosind bucketuri de liste ar putea arata in felul urmator:

```
insert(Object o, HashTable h)
{
    bucket = hashfunc(o)
    nod *n = new nod
    n → next = h[bucket]
    h[bucket] = n
    return
}
```

Pentru a cauta un obiect, pur si simplu se parcurge lista corespunzatoare in bucket si se verifica daca gasim obiectul.

```
find(Object o, hashtable h)
{
    bucket = hashfunc(o)
    nod * n = h[bucket]
    while (n!=null)
```

```

        if (n → val == 0)
            return true
        n = n → next
    return false
}

```

Pentru a șterge un element din hashtable, e suficient să ștergem elementul din lista corespunzătoare bucketului lui.

Pentru implementare, puteți folosi aceste structuri pentru hashtable.

```

typedef struct hashtable{
    struct nod ** Table; // vector de nod*, initial toate sunt NULL
    int nr_buckets;
    int (*hfunc) (char*);
    int size; // cate elemente avem in hashtable
} hashtable;

```

```

typedef struct nod{
    char *val;
    struct nod * next;
} nod;

```

Teoretic, complexitatea pentru hashtable este $O(1)$ amortizat, dacă avem puține coliziuni. În general, se ține cont de un procent de umplere – dacă avem mai mult de 75% din numărul de bucketuri, se dublează numărul de bucketuri, iar dacă avem mai puțin de 25% (de exemplu), se împarte numărul de bucketuri la 2. Evident, în acest proces toate elementele trebuie mutate în noul hashtable, dar pentru că aceste operații se fac rar, nu influențează complexitatea totală.

O variantă de hashtable este următoarea – pentru multe coliziuni, în loc să folosim o listă pentru elemente, se folosește un arbore binar de căutare. În acest fel, dacă avem mai multe coliziuni, în loc să căutăm într-o listă în $O(n)$, n fiind numărul de coliziuni pentru același bucket, vom căuta în $O(\log n)$.

Cerinte:

1. Implementați un hashtable folosind liste. Adăugați câteva șiruri în tabel, ștergeți șiruri din tabel și verificați dacă există anumite șiruri, pentru a verifica validitatea tabelului.

2. Modificați implementarea de mai sus, astfel încât nodul unui bucket să nu mai fie un nod de listă, ci un nod de treap. Pentru căutare, inserare și ștergere, se vor folosi funcțiile corespunzătoare din treap implementate în laboratorul 7.