

Cấu trúc dữ liệu

Data Structure

Ts. Nguyễn Đức Thuận
BM Hệ thống Thông Tin



Chương VII

BẢNG BĂM (HASH TABLE)

Bảng băm

■ Nội dung

- Bảng băm
- Định nghĩa hàm băm
- Một số phương pháp xây dựng hàm băm
- Các phương pháp giải quyết đụng độ

Giới thiệu

- Phép băm được đề xuất và hiện thực trên máy tính từ những năm 50 của thế kỷ 20.
- **Ý tưởng:** biến đổi giá trị khóa thành một số (xử lý băm) và sử dụng số này để đánh chỉ số cho bảng dữ liệu: Bảng chỉ mục (bảng băm)

Các phép toán trên bảng băm sẽ giúp hạn chế số lần so sánh

☞ giảm thiểu được thời gian truy xuất.

(Độ phức tạp của các phép toán trên bảng băm thường có bậc là $O(1)$ và không phụ thuộc vào kích thước của bảng băm)

Ví dụ:

Phân 100 học sinh được đánh số tự tự lần lượt vào 10 phòng thi. Tìm phòng thi của học sinh thứ i ?

Bảng băm

- **Bảng băm** (hash table)
- **Định nghĩa**: HashTable là một kiểu dữ liệu danh sách dùng để lưu dữ liệu theo một từ khóa và giá trị của nó.

Address	Key	Values
001	a	Ant
002	b	Book
...
025	y	Yatch
026	z	Zebra

- **Phân loại**
- **Bảng băm đóng** : mỗi khóa ứng với một địa chỉ, thời gian truy xuất là hằng số
- **Bảng băm mở** : một số khóa có cùng địa chỉ, lúc này mỗi mục địa chỉ sẽ là một danh sách liên kết các phần tử có cùng địa chỉ, thời gian truy xuất có thể bị suy giảm đôi chút

Bảng băm

Các khái niệm chính trên cấu trúc bảng băm:

- Khởi tạo (*Initialize*)
- Kiểm tra rỗng (*Empty*)
- Lấy kích thước của bảng băm (*Size*)
- Tìm kiếm (*Search*)
- Thêm mới phần tử (*Insert*)
- Loại bỏ (*Remove*)
- Sao chép (*Copy*)
- Duyệt (*Traverse*)

Thông thường bảng băm được sử dụng khi cần xử lý các bài toán có dữ liệu lớn và được lưu trữ ở bộ nhớ ngoài.

CÁC HÀM C PHỤC VỤ CHO BẢNG BĂM

a. Khai báo cấu trúc bảng băm:

C Code:

```
1. #define NULLKEY -1
2. #define M 100
3. /*
4.  M là số nút có trên bảng băm, du de chưa các nút nhập vào bảng băm
5.  */
6. //khái báo cấu trúc một nút của bảng băm
7. struct node
8. {
9.     int key; //khóa của nút trên bảng băm
10. };
11. //Khái báo bảng băm có M nút
12. struct node hashtable[M];
13. int NODEPTR;
14. /*biến toàn cục chỉ số nút hiện có trên bảng băm*/
```

b. Các tác vụ:

Hàm băm:

C Code:

```
1. int hashfunc(int key)
2. {
3.     return(key% M)
4. }
```

Phép toán khởi tạo (initialize):

Khởi tạo bảng băm.

Gán biến toàn cục N=0.

C Code:

```
1. void initialize( )
2. {
3.     int i;
4.     for(i=0;i<M;i++)
5.         hashtable[i].key=NULLKEY;
6.     N=0;
7.     //số nút hiện có khởi động bảng 0
8. }
```

Phép toán kiểm tra trống (empty):

Kiểm tra bảng băm có trống hay không.

C Code:

```
1. int empty( );
2. {
3.     return(N==0 ? TRUE:FALSE);
4. }
```

Phép toán kiểm tra đầy (full):

Kiểm tra bảng băm đã đầy chưa.

C Code:

```
1. int full( )
2. {
3.     return (N==M-1 ? TRUE: FALSE);
4. }
```

CÁC HÀM C PHỤC VỤ CHO BẢNG BĂM

Phép toán search:

C Code:

```
1. int search(int k)
2. {
3.     int i;
4.     i=hashfunc(k);
5.     while(hashtable[i].key!=k && hashtable[i].key !=NULKEY)
6.     {
7.         //băm lại (theo phương pháp do tuyến tính:fi(key)=f(key)+i) % M
8.         i=i+1;
9.         if(i>=M)
10.            i=i-M;
11.     }
12.     if(hashtable[i].key==k) //tim thay
13.         return(i);
14.     else
15.         //khong tim thay
16.         return(M);
17. }
```

Phép toán insert:

Thêm phần tử có khoá k vào bảng băm.

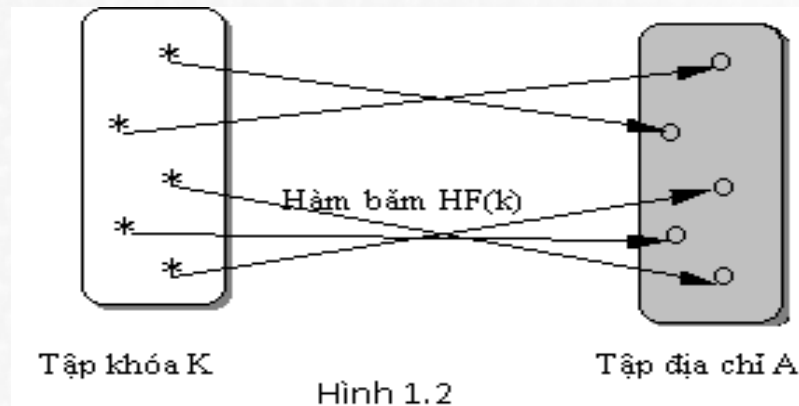
C Code:

```
1. int insert(int k)
2. {
3.     int i, j;
4.     if(full( ))
5.     {
6.         printf("\n Bang bam bi day khong them nut co khoa %d duoc",k);
7.         return;
8.     }
9.     i=hashfunc(k);
10.    while(hashtable[i].key !=NULLKEY)
11.    {
12.        //Băm lại (theo phương pháp do tuyến tính)
13.        i ++;
14.        if(i >M)    i= i-M;
15.    }
16.    hashtable[i].key=k;
17.    N=N+1;
18.    return(i);
19. }
```


Hàm băm

▪ 1. HÀM BẮM (Hash Function)

- *Định nghĩa*: Hàm băm là ánh xạ từ tập giá trị khóa vào một tập các địa chỉ (của bảng băm).



- *Ví dụ* : hàm băm biến đổi khóa dạng chuỗi gồm n kí tự thành 1 địa chỉ (số nguyên)

```
int hashfunc( char *s, int n )  
{  
    int sum = 0;  
    while( n -- ) sum = sum + *s++;  
    return sum % 256; }
```

// Tính địa chỉ của khoá "AB" : **hashfunc("AB",2) → 131**

// Tính địa chỉ của khoá "BA" : **hashfunc("BA",2) → 131**

Hàm băm

- **Khi hàm băm 2 khoá vào cùng 1 địa chỉ thì gọi là đụng độ (Collision)**

Hàm băm tốt thỏa mãn các điều kiện sau:

- ☞ Tính toán nhanh.
- ☞ Các khoá được phân bố đều trong bảng.
- ☞ Ít xảy ra đụng độ .

❖ Một số phương pháp xây dựng hàm băm

a. Hàm băm dạng bảng tra:

Hàm băm có thể tổ chức ở dạng bảng tra (còn gọi là bảng truy xuất) hoặc thông dụng nhất là ở dạng công thức.

Hàm băm

Ví dụ: bảng tra với khóa là bộ chữ cái, bảng băm có 26 địa chỉ từ 0 đến 25. Khóa a ứng với địa chỉ 0, khóa b ứng với địa chỉ 1, ..., z ứng với địa chỉ 25.

Khóa	Địa chỉ	Khóa	Địa chỉ	Khóa	Địa chỉ	Khóa	Địa chỉ
a	1	h	8	o	15	v	22
b	2	i	9	p	16	x	23
c	3	j	10	q	17	y	24
d	4	k	11	r	18	w	25
e	5	l	12	s	19	z	26
f	6	m	13	t	20		
g	7	n	14	u	21		

Hàm băm dạng bảng tra được tổ chức dưới dạng danh sách kê.

Hàm băm

▪ b. Hàm băm sử dụng phương pháp chia

Dùng dư số: $h(k) = k \bmod m$

k là khoá, m là kích thước của bảng.

☞ vấn đề chọn giá trị m

Nếu chọn $m = 2^n$ thông thường không tốt vì $h(k) = k \bmod 2^n$ sẽ chọn cùng n bits thấp của $k \rightarrow$ nên chọn m là nguyên tố (tốt) gần với 2^n

Ví dụ: Ta có tập khoá là các giá trị số gồm 3 chữ số, và vùng nhớ cho bảng địa chỉ có khoảng 100 mục, như vậy ta sẽ lấy hai số cuối của khoá để làm địa chỉ theo phép chia lấy dư cho 100 : chẳng hạn $325 \bmod 100 = 25$

Có nhiều số có cùng địa chỉ. Vì thế, để hàm băm có thể tính địa chỉ khoá ít trùng lặp hơn chọn $m=97$ thay vì $m=100$

(Ví dụ: $325 \bmod 100 = 25$; $125 \bmod 100 = 25$
 $325 \bmod 97 = 34$; $125 \bmod 97 = 28$)

Hàm băm

- c. Hàm băm sử dụng phương pháp nhân

$$h(k) = \lfloor m * (k * A \bmod 1) \rfloor$$

k là khóa, m là kích thước bảng, A là hằng số: $0 < A < 1$

Chọn m và A

Ta thường chọn $m = 2^n$

Theo Knuth chọn $A = 1/2(\sqrt{5}) - 1 \approx 0.618033987$ được xem là tốt.

- Ví dụ: k=123456; m=10000

$$H(k) = \lfloor 10000 (123456 * 0.6180339887 \bmod 1) \rfloor$$

$$H(k) = \lfloor 10000 (76300.0041089472 \bmod 1) \rfloor$$

$$H(k) = \lfloor 10000 (0.0041089472) \rfloor$$

$$H(k) = 41$$

Hàm băm

d. Hàm băm phổ quát (*universal hashing function*)

Định nghĩa: Cho H là một tập hợp hữu hạn các hàm băm: ánh xạ các khóa k từ tập khóa U vào miền giá trị $\{0, 1, 2, \dots, m-1\}$. Tập H là phổ quát nếu:

$\forall f \in H$ và 2 khoá phân biệt k_1, k_2 ta có xác suất: $\Pr\{f(k_1) = f(k_2)\} \leq 1/m$.

Ứng dụng:

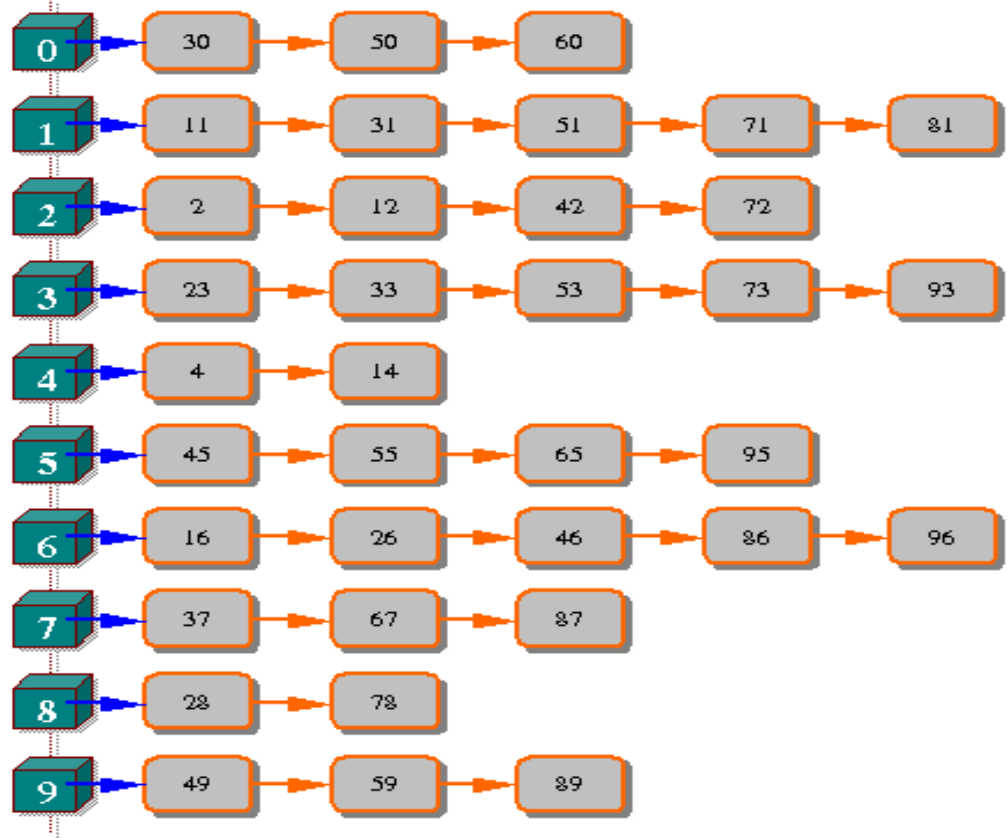
- Khởi tạo một tập các hàm băm H phổ quát
- Chọn $h \in H$ ngẫu nhiên.

Giải quyết đụng độ

- Người ta giải quyết sự đụng độ theo hai phương pháp: *phương pháp nối kết* và *phương pháp băm lại*.

1. Giải quyết sự đụng độ bằng phương pháp nối kết (Chaining Method):

- Các phần tử bị băm vào cùng địa chỉ (các phần tử bị đụng độ) được gom thành một danh sách liên kết (gọi là một bucket).



Giải quyết độ

▪ a. Khai báo cấu trúc bảng băm:

```
#define M 100
struct nodes
{
    int key;
    struct nodes *next };
//khai bao kieu con tro chi nut
typedef struct nodes *NODEPTR;
/*
khai bao mang bucket chua M con tro dau cua M bucket
*/
NODEPTR bucket[M];
```


Giải quyết đụng độ

■ Các phép toán :

- Tính giá trị hàm băm: Giả sử chúng ta chọn hàm băm dạng %: $h(\text{key}) = \text{key} \% M$.
- Phép toán **initbuckets**: khởi tạo các bucket bằng Null.
- Phép toán **emmpybucket(b)**: kiểm tra bucket b có bị rỗng không?
- Phép toán **emmpy**: Kiểm tra bảng băm có rỗng không?
- Phép toán **insert**: Thêm phần tử có khóa k vào bảng băm.
 - + $i = h(k)$
 - + ktra bucket [i]: nếu rỗng => cc o nhỏ cho bucket, gan khoa k
thêm phần tử có khóa k vào ds theo thu tu tang dan.
- Phép toán **remove**: Xóa phần tử có khóa k trong bảng băm.
- Phép toán **clear**: Xóa tất cả các phần tử trong bảng băm.
- Phép toán **traversebucket**: Xử lý tất cả các phần tử trong bucket b.
- Phép toán **traverse**: Xử lý tất cả các phần tử trong bảng băm.
- Phép toán **search**: Tìm kiếm một phần tử trong bảng băm, nếu không tìm thấy hàm này trả về hàm NULL, nếu tìm thấy hàm này trả về địa chỉ của phần tử có khóa k.

ĐỊA CHỈ CÁC HÀM C PHỤC VỤ CHO BẢNG BĂM

- <http://diendan.congdongcviet.com/threads/t3289::ly-thuyet-tim-kiem-bang-bam-hash-table.cpp>

a. Khai báo cấu trúc bảng băm:

#

C Code:

```
1. define M 100
2. struct nodes
3. {
4. int key;
5. struct nodes *next
6. };
```

C Code:

```
1. //khai bao kieu con tro chi nut
2. typedef struct nodes *NODEPTR
3. /*
4. khai bao mang bucket chua M con tro dau cua M bucket
5. */
6. NODEPTR bucket[M];
```

b. Các phép toán:

Hàm băm

Giả sử chúng ta chọn hàm băm dạng %: $f(\text{key}) = \text{key} \% M$.

C Code:

```
1. int hashfunc (int key)
2. {
3. return (key % M);
4. }
```

Phép toán initbuckets:

C Code:

```
1. void initbuckets( )
2. {
3. int b;
4. for (b=0;b<M;b++);
5. bucket[b]=NULL;
6. }
```

Phép toán emmptybucket:

C Code:

```
1. int emptybucket (int b)
2. {
3. return(bucket[b] ==NULL ?TRUE :FALSE);
4. }
```

Phép toán emmpty:

C Code:

```
1. int empty( )
2. {
3. int b;
4. for (b = 0;b<M;b++)
5. if(bucket[b] !=NULL) return(FALSE);
6. return(TRUE);
7. }
```

CÁC HÀM C PHỤC VỤ CHO BẢNG BĂM

Phép toán remove:

C Code:

```
1. void remove ( int k)
2. {
3.     int b;
4.     NODEPTR q, p;
5.     b = hashfunc(k);
6.     p = hashbucket(k);
7.     q=p;
8.     while(p !=NULL && p->key !=k)
9.     {
10.        q=p;
11.        p=p->next;
12.    }
13.    if (p == NULL)
14.        printf("\n không có nút có khóa %d" ,k);
15.    else
16.        if (p == bucket [b])    pop(b);
17.        //Tac vu pop của danh sach lien ket
18.    else
19.        delafter(q);
20.    /*tac vu delafter của danh sach lien ket*/
21. }
```

Phép toán clearbucket:

Xóa tất cả các phần tử trong bucket b.

C Code:

```
1. void clearbucket (int b)
2. {
3.     NODEPTR p,q;
4.     //q là nút trước, p là nút sau
5.     q = NULL;
6.     p = bucket[b];
7.     while(p !=NULL)
8.     {
9.         q = p;
10.        p=p->next;
11.        freenode(q);
12.    }
13.    bucket[b] = NULL; //khởi động lại bucket b
14. }
```

Phép toán clear:

Xóa tất cả các phần tử trong bảng băm.

C Code:

```
1. void clear( )
2. {
3.     int b;
4.     for (b = 0; b<M ; b++)
5.         clearbucket(b);
6. }
```

CÁC HÀM C PHỤC VỤ CHO BẢNG BĂM

Phép toán traversebucket:
Duyệt các phần tử trong bucket b.

C Code:

```
1. void traversebucket (int b)
2. {
3.     NODEPTR p;
4.     p= bucket[b];
5.     while (p !=NULL)
6.     {
7.         printf("%3d", p->key);
8.         p= p->next;
9.     }
10. }
```

Phép toán traverse:
Duyệt toàn bộ bảng băm.

C Code:

```
1. void traverse( )
2. {
3.     int b;
4.     for (b = 0; b<M; b++)
5.     {
6.         printf("\nButket %d:",b);
7.         traversebucket(b);
8.     }
9. }
```

Phép toán search:

Tìm kiếm một phần tử trong bảng băm, nếu không tìm thấy hàm này trả về giá trị NULL,

C Code:

```
1. NODEPTR search(int k)
2. {
3.     NODEPTR p;
4.     int b;
5.     b = hashfunc (k);
6.     p = bucket[b];
7.     while(k > p->key && p !=NULL)
8.         p=p->next;
9.     if (p == NULL || k !=p->key)// không tìm thấy
10.    return(NULL);
11.    else//tìm thấy
12.    // else //tìm thấy
13.    return(p);
14. }
```


Giải quyết đụng độ

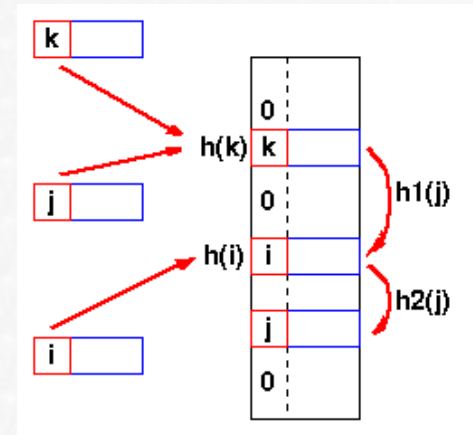
2. Giải quyết sự đụng độ bằng phương pháp băm lại: (Rehash Method)

▪ Phương pháp dò tuyến tính (Linear Probe)

Nếu băm lần đầu bị xung đột thì băm lại lần 1, nếu bị xung đột nữa thì băm lại lần 2,... Quá trình băm lại diễn ra cho đến khi không còn xung đột nữa. Các phép băm lại (*rehash function*) thường sẽ chọn địa chỉ khác cho các phần tử.

$$h_i(\text{key}) = (h(\text{key}) + i) \% M$$

với $h(\text{key})$ là hàm băm chính của bảng băm



Giải quyết đụng độ

▪ Phương pháp dò bậc hai (Quadratic Probing Method)

- Hàm băm lại của phương pháp dò bậc hai là truy xuất các địa chỉ cách bậc 2. Hàm băm lại hàm i được biểu diễn bằng công thức sau:

$$h_i(\text{key}) = (h(\text{key}) + i^2) \% M$$

với $h(\text{key})$ là hàm băm chính của bảng băm.

Nếu đã dò đến cuối bảng thì trở về dò lại từ đầu bảng (?).

Bảng băm với phương pháp dò bậc hai nên chọn số địa chỉ M là số nguyên tố.

CÁC HÀM C PHỤC VỤ CHO BẢNG BĂM

Phép toán search:

Tìm phần tử có khóa k trên bảng băm, nếu không tìm thấy hàm này trả về trị M,

C Code:

```
1. int search(int k)
2. {
3.     int i, d;
4.     i = hashfuns(k);
5.     d = 1;
6.     while(hashtable[i].key!=k&&hashtable[i].key !=NULLKEY)
7.     {
8.         //Băm Lại (theo phương pháp bậc hai)
9.         i = (i+d) % M;
10.    d = d*2;
11.    }
12.    hashtable[i].key =k;
13.    N = N+1;
14.    return(i);
15. }
```

Giải quyết đụng độ

- **Phương pháp băm kép (Double hashing Method):** là một phương pháp băm lại dùng cùng lúc hai hàm băm

Ý tưởng:

- Dùng hai hàm băm $h1(key)$ và $h2(key)$
- *Khi bổ sung 1 khóa* key: $i=h1(key)$ và $j=h2(key)$ sẽ xác định địa chỉ i và j trong khoảng từ 0 đến $M-1$: Khi thêm phần tử có khóa key vào bảng băm
 - Nếu chưa bị xung đột thì thêm phần tử mới tại địa chỉ i này.
 - Nếu bị xung đột thì hàm băm lại lần 1 $h1$ sẽ xét địa chỉ mới $i+j$, nếu lại bị xung đột thì hàm băm lại lần 2 $h2$ sẽ xét địa chỉ $i+2j, \dots$, quá trình cứ thế cho đến khi nào tìm được địa chỉ trống và thêm phần tử vào địa chỉ này.

Giải quyết đụng độ

- *Khi tìm kiếm một phần tử* có khoá key xác định: $i=h_1(\text{key})$ và $j=h_2(\text{key})$
- Xét phần tử tại địa chỉ i , nếu chưa tìm thấy thì xét tiếp phần tử $i+j$, $i+2j$, ..., quá trình cứ thế cho đến khi nào tìm được khoá (trường hợp tìm thấy) hoặc bị rơi vào địa chỉ trống (trường hợp không tìm thấy).
Nếu đã dò đến cuối bảng thì trở về dò lại từ đầu bảng.

Cám ợn đã theo dõi

