

## 第 10 章

### 跳表和散列 (Skip List and Hashing)

#### 本章内容

- 10.1 字典
- 10.2 抽象数据类型
- 10.3 线性表描述
- 10.4 跳表描述
- 10.5 散列表描述
- 10.6 应用

#### 10.1 字典

- 字典 (dictionary) 是一些形如 $(k, v)$ 的数对(元素/记录)所组成的集合。其中 $k$ 是**关键字**,  $v$ 是关键字 $k$ 对应的**值**。任意两个数对, 其关键字各不相同。
- 有关字典的操作有:
  - 确定字典是否为空。
  - 确定字典有多少数对。
  - 在字典中寻找/搜索具有给定关键字的数对。
  - 插入一个数对。
  - 删除一个具有给定关键字的数对。

● **多重字典**: 两个或多个数对有相同的关键字。

#### 10.2 抽象数据类型

抽象数据类型 *Dictionary*{

##### 实例

关键字互不相同的数对(元素)集合

##### 操作

*empty()*: 字典为空时返回true, 否则返回false

*size()*: 返回字典中的数对个数

*find(k)*: 返回关键字为 $k$ 的数对

*insert(p)*: 插入数对 $p$

*erase(k)*: 删除关键字为 $k$ 的数对

}

#### C++抽象类dictionary

```
template <class K, class E>
class dictionary
{
public:
    virtual ~dictionary() {}
    virtual bool empty() const = 0;
        //字典为空时返回true, 否则返回false
    virtual int size() const = 0;
        //返回字典中数对的个数
    virtual pair<const K, E>* find (const K&) const = 0;
        //返回匹配数对中的指针;
    virtual void erase(const K&) = 0;
        //删除匹配的数对
    virtual void insert(const pair<const K, E>&) = 0;
        //在字典中插入一个数对
}
```

#### 访问字典数对

- 随机访问(Random Access)
  - 按照给定的一个关键字来访问字典中的数对。
- 顺序访问 ( Sequential Access)
  - 按照关键字的递增顺序逐个访问字典中的数对。
  - 顺序访问需要操作:
    - *Begin* –用来返回关键字最小的数对
    - *Next* –用来返回下一个数对

## 字典的线性结构描述

### ■ 字典的线性结构描述:

- 线性表描述
- 跳表描述
- 散列表描述

## 10.3 线性表描述

### ■ 字典用有序线性表: $(p_0, p_1, p_2, \dots, p_{n-1})$ , 表示

- 关键字从左到右依次增大

### ■ 使用线性表的描述方法

- 数组描述: sortedArrayList
- 链表描述: sortedChain

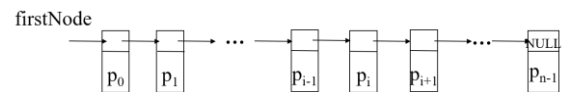
## 数组描述的有序线性表sortedArrayList

### ► 用数组描述的有序线性表类sortedArrayList

- 搜索: 折半搜索, 时间为 $O(\log n)$ 。
- 插入: 首先确认该字典中是否存在相同关键字的数对, 然后进行插入, 插入操作的时间是 $O(n)$ 。
- 删除: 首先要找到欲删除的数对, 然后再进行删除, 删除的时间复杂性为 $O(n)$ 。

## 链表描述的有序线性表sortedChain

### ► 有序链表表示字典



```
template<class K, class E>
struct pairNode{
    pair<const K,E> element;
    pairNode<const K,E> *next;
    .....
};
```

## 类sortedChain

```
template<class K, class E>
class sortedChain
{public:
    .....
    pair<const K,E> * find(const K& theKey) const;
    //返回关键字theKey匹配的数对的指针, 若不存在匹配的数对,
    //则返回NULL
    void insert(const pair<const K,E>& thePair);
    //插入一个数对thePair,覆盖已经存在的数对
    void erase(const K& theKey);
    //删除关键字theKey匹配的数对
protected:
    pairNode<const K,E> *firstNode; //指向链表第一个节点的指针
    int dSize; //表中的数对个数
};
```

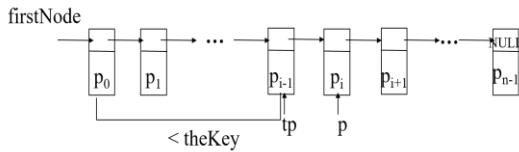
## 方法 'find'

```
template<class E, class K>
pair<const K,E> * find(const K& theKey) const;
{//返回关键字theKey匹配的数对的指针, 若不存在匹配的数对,
//则返回NULL
pairNode<K,E> *currentNode = firstNode;
// 搜索与theKey相匹配的数对
while (currentNode !=NULL &&
        currentNode->element.first < theKey)
    currentNode = currentNode->next;

// 判断是否与theKey匹配
if (currentNode !=NULL &&
    currentNode->element.first == theKey) // 与theKey相匹配
    return &currentNode->element;
return NULL; // 不存在相匹配的数对
}
```

## 方法 ‘insert’

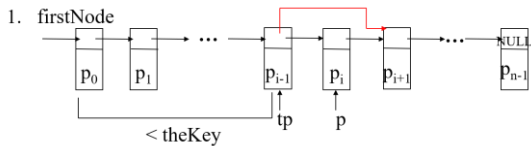
- 在字典中插入一个数对thePair,覆盖已经存在的数对



```
template<class E, class K>
void sortedChain<K,E>::insert(const pair<const K, E>& thePair)
{ // 在字典中插入一个数对thePair,覆盖已经存在的数对
  pairNode<K,E> *p = firstNode, *tp = NULL; // 跟踪p
  // 移动tp,使得thePair可以插在tp的后面
  while (p!=NULL && p->element.first < thePair.first)
    { tp = p; p = p->next; };

  if (p !=NULL&& p->element.first == thePair.first)
    { // 有匹配的数对,覆盖数对中的值
      p->element.second = thePair.second; return; }
  // 若无匹配的数对,则为thePair建立一个新节点
  pairNode<K,E> *newNode = new pairNode<K,E>(thePair, p);
  if (tp==NULL) firstNode = newNode; // p是链首节点
  else tp->next = newNode; // 将新节点插入到tp之后
  dSize++;
  chainNode(const T& element, chainNode<T>* next)
  { this->element = element;
    this->next = next; }
  return;
}
```

## 方法 ‘erase’



```
template<class E, class K>
void sortedChain<K,E>::erase(const K& theKey)
{ // 删除与theKey相匹配的数对
  pairNode<K,E> *p = firstNode,
    *tp = NULL; // 跟踪p

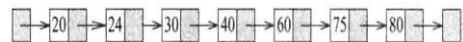
  // 搜索与theKey相匹配的数对
  while (p!=NULL && p->element.first < theKey)
    { tp = p; p = p->next; };

  // 确定是否与theKey匹配
  if (p!=NULL && p->element.first == theKey)
    { // 找到一个相匹配的数对
      if (tp==NULL) firstNode = p->next; // p是链首节点
      else tp->next = p->next;

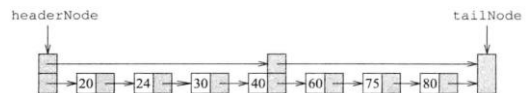
      delete p;
      dSize--; }
}
```

## 10.4 跳表表示

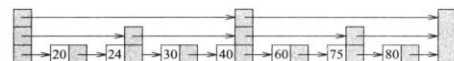
## 理想情况



搜索: 最多比较次数 = n



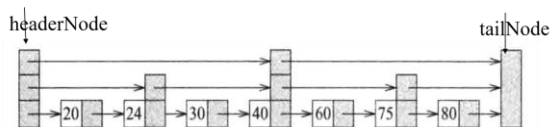
搜索: 最多比较次数 = n/2+1



## 理想情况

- 0 级链表:  $n$  个数对 (元素/记录)
- 1 级链表:  $\frac{1}{2}n$  个数对
- 2 级链表:  $\frac{1}{4}n$  个数对
- .....
- $i$  级链表:  $n/2^i$  个数对
- 一个数对是  $i$  级链表数对  
当且仅当 它属于  $0 \sim i$  级链表  
但不属于  $i+1$  级链表 (若该链表存在)
- $i$  级链表所包含的数对是  $i-1$  级链表所有数对的子集

## 搜索



- 从最高级链表开始, 在每一级链表中搜索, 遇到最后一个小于 theKey 的数对, 跳到下一级链表中搜索, 直到0级链表

## 插入和删除

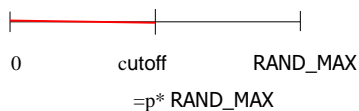
- 当插入和删除时, 要保持跳表的规则结构
- 跳表的规则结构:
- $i$  级链表:  $n/2^i$  个数对
  - $i-1$  级数对属于  $i$  级链表的概率是  $\frac{1}{2}$
  - 一个数对属于  $i$  级链表的概率是  $1/2^i = (1/2)^i$
  - 链的级数是:  $\lfloor \log_2 n \rfloor + 1$

## 插入和删除

- 跳表的规则结构:
  - $i-1$  级链中的数对属于  $i$  级的概率为  $p$
  - 每  $1/p$  个  $i-1$  级链中的数对有一个在  $i$  级
- 插入一个新数对时, 新数对被插入到  $i$  级的概率是  $p^i$
- 链的级数:  $\lfloor \log_{1/p} n \rfloor + 1$

## 级的分配

- 应用一个随机数产生器, 产生0 到RAND\_MAX之间的实数



- 产生的随机数  $\leq$  cutoff ( $= p * RAND\_MAX$ ) 的概率是  $p$
- 分配级数:
 

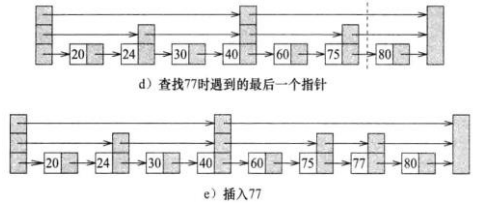
```
int lev = 0
while (rand() <= CutOff) lev++;
```

## 级的分配

- 缺点1: 可能为某些数对分配特别大的级, 从而导致一些数对的级远远超过  $\log_{1/p} N$  ( $N$  为字典中预期的最大数目)
  - 在有  $N$  个数对的跳表中, 级的最大值  $\text{MaxLevel} = \lceil \log_{1/p} N \rceil - 1$ , 以此值作为上限。
- 缺点2: 可能存在下面的情况, 如在插入一个新数对前有三条链, 而在插入之后就有了10条链。这时, 新插入数对的级为9级, .....。
  - 我们可以把新数对的级调整为3。即把新数对的级调整为 **最大级数+1**。

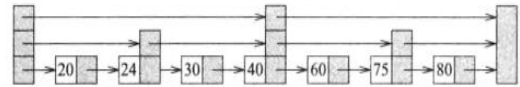
## 插入

- 插入77
- 搜索77



## 删除

- 删除77



- 删除操作，无法控制跳表结构

## 10.5 散列表描述

- 10.5.1 理想散列
- 10.5.2 散列函数和散列表
- 10.5.3 线性探查
- 10.5.4 链式散列

### 10.5.1 理想散列

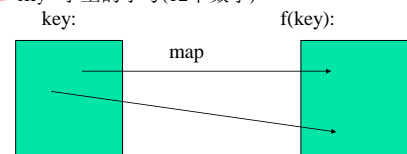
- **散列(Hashing)**:字典的另一种描述方法就是散列.
- **散列方法**:是用一个**散列函数(hash function)**又称**哈希函数**把字典的数对映射到**散列表/哈希表(hash table)**中的具体位置。
- **理想情况**:
  - 数对 $p$ 的关键词为 $k$
  - $f$ 是散列函数
  - 数对 $p$ 在散列表中的位置为 $f(k)$

## 理想散列

- 假定散列表中的每一个位置最多只能存放一个数对(记录)
- 散列表操作:
  - **查找(find)**:计算出 $f(k)$ , 然后看表中 $f(k)$ 处是否有关键字为 $k$ 的数对。
  - **插入(insert)**:计算出 $f(k)$ , 把数对放在 $f(k)$ 位置
  - **删除(erase)**:计算出 $f(k)$ , 把表中 $f(k)$ 位置置为空
- 在理想情况下, 散列表操作时间复杂性:  $\Theta(1)$

## 理想散列

- 例:
  - Key = 学生的学号(12个数字)



关键字的范围:  
202300000000~  
202399999999

$f(\text{key})$ 的范围:  
0~99999999

- 关键字的范围太大, 因此不能像理想散列方法那样建立数组.

## 理想散列

- 当关键字的范围太大，不能用理想方法表示时，可以采用比关键字范围小的散列表。
- ⇒ 散列函数：多个关键字映射到同一位置。

## 10.5.2 散列函数和散列表

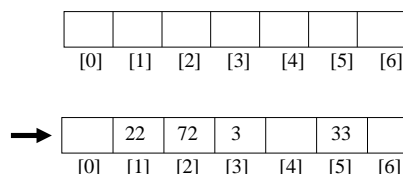
- **桶 (bucket)**: 散列表中的每一个位置
- **起始桶 (home bucket)**: 对关键字为k的数对，位置  $f(k)$  是起始桶
- **散列表的长度或大小**: 桶的数量
- 最常用的散列函数:
- **除余散列 (Hashing by division)**
  - $f(k) = k \% D$  (%为求模操作符)
  - D: 是散列表的大小(即位置数)
  - 散列表的位置索引:  $0 \sim D-1$

## 好的散列函数

- 好的散列函数:
  - 均匀散列函数: 映射到一个桶里的关键字大致相同
- 良好的散列函数: 性能较好的均匀散列函数
- 除余散列:  $f(k) = k \% D$ 
  - D的选择对于散列的性能有着重大的影响(D等于桶的个数b)。
  - 当D为素数或D没有小于20的素数因子时，可以使性能达到最佳。

## 例

- 散列表: `table[0:6]`
- 散列函数:  $f(k) = k \% 7$
- 依次插入关键字: 22, 33, 3, 72



## 冲突和溢出 (overflow)

	22	72	3		33	
[0]	[1]	[2]	[3]	[4]	[5]	[6]

散列函数:  $f(k) = k \% 7$

- 85 放入什么地方?
- $f(85)=1$
- 85 的起始桶已经被另一个数对占用
- 两个不同的关键字起始桶号相同时——**冲突 (collision)**
- 存储桶中若没有空间时就发生**溢出 (overflow)**
- 当每个桶只能存储一个数对时，碰撞和溢出会同时发生。

## 溢出 (overflow)

- 解决溢出的方法:
  - 线性探查
  - 链式散列

## 10.5.2 线性探查

- 当溢出发生时，将元素插入**下一个可用桶**中。在寻找下一个可用桶时，表被视为环形的。

### ■ 例

- 散列表的大小  $b = 11$
- $f(k) = k \% b$
- 在插入 80, 40, 65 后

			80				40			65
0	1	2	3	4	5	6	7	8	9	10

## 线性探查—“插入”

			80				40			65
0	1	2	3	4	5	6	7	8	9	10

$$f(k) = k \% 11$$

- 再依次插入 58, 24, 35

- $f(58)=3$  (冲突);  $f(24)=2$ ;  $f(35)=2$  (冲突)

		24	80	58	35		40			65
0	1	2	3	4	5	6	7	8	9	10

## 线性探查—“搜索”

		24	80	58	35		40			65
0	1	2	3	4	5	6	7	8	9	10

- 例：搜索 24、35、29

### ■ 搜索 (search) 操作

- 首先搜索关键字为  $k$  的起始桶  $f(k)$ ，接着对表中后继桶进行搜索，直到发生以下情况：(1) 存有关键字为  $k$  的桶已找到，即**找到**了要搜索的数对。(2) 到达一个空桶；(3) 又回到起始桶  $f(k)$ 。(2) 和 (3) 说明表中**没有关键字为  $k$  的数对**。

## 线性探查—“删除”

### ■ 删除操作

**删除需要移动若干数对**

- 执行搜索操作，找到关键字为  $k$  的桶。  $f(k) = k \% 11$
- 在完成一次删除操作后，必须保证上述的搜索过程仍能正常进行。不能仅仅把表中相应位置置为空。

98		24	80	58	35		40			65
0	1	2	3	4	5	6	7	8	9	10

- 删除: 58 ( $f(58)=3$ ) ; 然后 搜索: 35 ( $f(35)=2$ )

- 采取办法:

- 从删除位置的下一个桶开始，逐个检查每个桶以确定要移动的数对，直到到达一个空桶或回到删除操作所对应的桶为止。做删除移动时，一定要注意不要把一个数对移到它的起始桶之前，否则对这个数对的查找就可能失败

- 另一个策略：为每个桶增加一个 **NeverUsed** 域。

## 线性探查散列表的类

```
template<class K, class E>
class hashTable <K,E>{
public:
    hashTable(int theDivisor = 11);
    .....//~hashTable(), empty(); size(); output

    pair<const K, E>* find(const K&) const;
    //返回关键字theKey匹配的数对的指针，若不存在匹配的数对，则返回NULL
    void insert(const pair<const K, E>&);
    //在字典中插入一个数对thePair,若存在关键字相同的数对,则覆盖
protected:
    int search(const K&) const;
    pair<const K, E> **table; // 散列表
    int divisor; // 散列函数的除数
    hash<K> hash; //把类型k映射到一个非负整数
    int dSize; // 字典中数对的个数
};
```

## hashTable的构造函数

```
template<class E, class K>
hashTable<K,E>::hashTable(int theDivisor)
{ // 构造函数
    divisor = theDivisor;
    dSize=0;
    // 分配和初始化散列表数组
    table = new pair<const K, E>* [divisor];

    for (int i = 0; i < divisor; i++) //将所有桶置空
        table[i] = NULL;
}
```

### hashTable<K,E>::search

```
template<class E, class K>
int hashTable<K,E>::search(const K& theKey) const
{// 搜索一个公开地址散列表，查找关键字为theKey的数对
// 如果匹配的数对存在，则返回它的位置
// 否则返回关键字为theKey的数对可以插入的位置
int i =(int) hash(theKey) % divisor; // 起始桶
int j = i; // 从起始桶处开始
do {
    if (table[j]==NULL || table[j]->first == theKey) return j;
    j = (j + 1) % divisor; // 下一个桶
} while (j != i); // 是否返回起始桶?
return j; // 表已满
}
```

### hashTable<K,E>:: find

```
template<class E, class K>
pair<const K, E> * hashTable<K,E>::find(const K& theKey) const
{//返回关键字theKey匹配的数对的指针
//若不存在匹配的数对，则返回NULL
//搜索散列表
int b = search(theKey);
//判断table[b]是否是匹配数对
if (table[b] ==NULL || table[b] ->first != theKey)
    return NULL; //未找到匹配数对

return table[b]; //找到匹配数对
}
```

```
template<class E, class K>
```

```
void hashTable<K,E>::insert(const pair<const K, E>& thePair)
{//在字典中插入一个数对thePair,若存在关键字相同的数对,
//则覆盖，若表已满，则抛出异常
int b = search(thePair.first);
// 检查匹配的数对是否存在
//没有匹配的数对，且表不满,则插入数对
if (table[b]==NULL)
    {table[b] = new pair<const K, E> (thePair);
    dSize++; }
// // 有重复的关键字，则覆盖数对中的值
if (table[b]->first == thePair.first)
    table[b]->second=thePair.second;
else throw hashTableFull(); // 表满
}
```

### 性能分析

- 时间复杂性:
- $b$ =散列表中桶的个数
- $f(k)=k \% D$ ;  $b=D$
- 散列表初始化:  $\Theta(b)$
- 最坏情况:
  - 所有 $n$ 个关键字值都在同一个桶中
  - 插入和搜索时间均为:  $\Theta(n)$ ,  $n$ 为表中元素个数
- 散列在最坏情况下的复杂性与线性表在最坏情况下的复杂性是完全相同。

### 性能分析

- 平均性能:
- 搜索因子(loading factor/装填因子、负载因子):
  - $\alpha = n/b$ ;  $0 \leq \alpha \leq 1$
- $S_n$  = 一次成功的搜索中平均搜索桶的个数
- $U_n$  = 一次不成功的搜索中平均搜索桶的个数
- $S_n \approx \frac{1}{2} (1 + 1/(1 - \alpha))$
- $U_n \approx \frac{1}{2} (1 + 1/(1 - \alpha)^2)$

### 性能分析

$\alpha$	$S_n$ (桶数)	$U_n$ (桶数)
0.50	1.5	2.5
0.75	2.5	8.5
0.90	5.5	50.5

$\alpha \downarrow \rightarrow S_n, U_n \downarrow$

$S_n$  = 一次成功的搜索中平均搜索桶的个数

$U_n$  = 一次不成功的搜索中平均搜索桶的个数



## 除余散列中除数D的选择

- $D$ 的选择对于散列的性能有着重大的影响( $D$ 等于桶的个数 $b$ )。

方法 1:

- 确定  $U_n$  和  $S_n$
- 使用  $U_n$  和  $S_n$  的公式, 确定最大的  $\alpha$  值-----  $\alpha_{\max}$

$$\left. \begin{array}{l} \bullet n \\ \bullet \alpha_{\max} \end{array} \right\} \alpha = n/b \rightarrow b_{\min}$$

- 取  $D$  为大于等于  $b_{\min}$  最小的整数, 该整数要么是素数, 要么没有小于20的素数因子。

## 除余散列中除数D的选择

例 10.12:

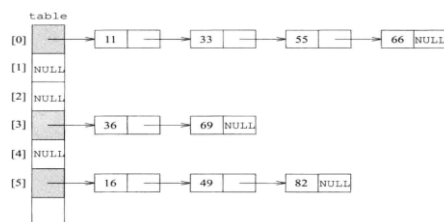
- $n=1000$ ;  $S_n < 4$ ;  $U_n < 50.5$ ;
- $U_n < 50.5 \Rightarrow \alpha \leq 0.9$
- $S_n < 4 \Rightarrow \alpha \leq 6/7$
- $\Rightarrow \alpha \leq \min\{0.9, 6/7\}$ ;  $\alpha_{\max} = 6/7$
- $\alpha = n/b$
- $b_{\min} = n/\alpha_{\max} = 7/6n = 1167$
- $D = 1171$

## 除余散列中除数D的选择

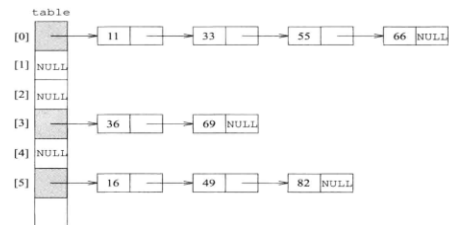
方法 2:

- 根据散列表的最大空间确定  $b_{\max}$  ( $b$ 的最大可能值)
- 取  $D$  为不大于这个最大值的整数, 该整数要么是素数, 要么没有小于20的素数因子。
- 例散列表的最大空间  $b_{\max} = 530$
- 取  $D = 23 * 23 = 529$

## 10.5.4 链式散列



- 设一个桶可以放无限多个数对, 散列表中的桶用一个链表来配置。
- 一个链表上的数对是具有同样起始桶的数对。



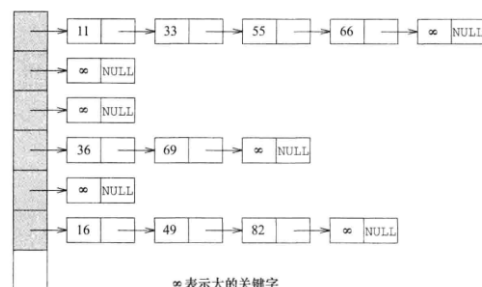
- 查找 (find):
  - 计算  $f(k)$ ; 搜索该桶所对应的链表。
- 插入 (insert):
  - 计算  $f(k)$ ; 搜索; 插入。
- 删除 (erase):
  - 计算  $f(k)$ ; 搜索; 删除。

```
template<class K, class E>
class hashChains
{public:
    .....
    pairNode<const K,E> * find(const K& theKey) const;
    //返回关键字theKey匹配的数对的指针, 若不存在匹配的数对, 则返回NULL
    {return table[hash (& theKey)% divisor].find(theKey);}
    void insert(const pair<const K,E>& thePair);
    //在字典中插入一个数对thePair,若存在关键字相同的数对, 则覆盖
    void erase(const K& theKey);
    //删除关键字theKey匹配的数对
    {table[hash (& theKey)% divisor].erase(theKey);}
protected:
    sortedChain<K,E> *table; // 链表数组
    int divisor; // 位置数
    hash<K> hash; //把类型k映射到一个非负整数
    int dSize; //数对的个数
};
```

## hashChains <K,E>::insert

```
void hashChains <K,E>::insert(const pair<const K,E>& thePair);
{//插入一个数对thePair, 若存在关键字相同的数对,则覆盖
    int homeBucket = (int) hash (thePair.first) % divisor;
    int homeSize = table[homeBucket].size();
    table[homeBucket].insert(thePair);
    if ( table[homeBucket].size() > homeSize )
        dSize++;
}
```

## 一种改进的方法



## 与线性探查比较

- 空间:
  - 线性探查 :  $2b(b \text{ 个指针}) + sn(n \text{ 个数对})$
  - 链式散列 :  $2b(b \text{ 个 firstNode 指针}) + 2b(b \text{ 个 dSize}) + 2n(n \text{ 个 next 指针}) + sn(n \text{ 个数对}) = 4b + (s+2)n$
  - 线性探查空间少

## 与线性探查比较

- 时间复杂性:
  - 线性探查 :
    - 最坏情况:  $\Theta(n)$
    - 平均性能:
      - 一次搜索成功搜索的桶数  $S_n \approx \frac{1}{2} (1 + 1/(1 - \alpha))$
      - 一次搜索不成功搜索的桶数  $U_n \approx \frac{1}{2} (1 + 1/(1 - \alpha)^2)$
  - 链式散列 :
    - 最坏情况:  $O(n)$
    - 平均性能:

## 与线性探查比较

平均性能  $U_n$  :

- 对于有  $i$  个节点的链表, 有  $i+1$  种可能搜索失败
- 一次不成功搜索所要搜索的节点数为:
  - $\frac{1}{i+1} (i + \sum_{j=1}^i j) = \frac{i(i+3)}{2(i+1)}$   $i \geq 1$
  - 0  $i = 0$
- 链表的平均长度为:  $n/b = \alpha$ 
  - $U_n \approx (\alpha + 1)/2$   $\alpha \geq 1$
  - $U_n \leq \alpha$   $\alpha < 1$

## 与线性探查比较

平均性能  $S_n$  :

- 假定 : 各元素是按升序插入的。
  - 当插入第  $i$  个元素时, 其所在链表的平均长度是:
    - $(i-1)/b$ , 第  $i$  个元素插入到该链表的尾部
    - 搜索 该元素:  $(i-1)/b + 1$
    - 假定  $n$  个标识符中每一个被搜索的概率都相同
- $$S_n = \frac{1}{n} \sum_{i=1}^n ((i-1)/b + 1) = 1 + (n-1)/2b$$
- $$\sim 1 + \alpha/2$$
- 使用链表时的平均性能要优于使用线性探查。

10.6 应用—文本压缩

- LZW (Lempel、Ziv、Welch) 压缩
- 压缩方法:
  - 基于原始数据, 创建一个字典, 字典中所存放的是文本中的字符串与其编码的映射。然后, 用字典中的编码来代替原始数据中的相应字符串。
- 例:要压缩的文本文件为 :
  - aaabbbbaabaaba
- LZW压缩字典:

code	0	1	2	3	4	5	6	7
key	a	b	aa	aab	bb	bbb	bbba	aaba

10.6 应用—文本压缩

- LZW规则 :
  - 开始, 为该文本文件中所有可能出现的字母分配一个代码, 构成初始字典。
  - LZW压缩器不断地在输入文件的未编码部分中寻找在字典中出现的**最长的前缀p**, 输出前缀p相应的代码; 所谓p是最长的前缀, 是指若在S中存在下一个字符c, 则需要为pc分配一个代码,并插入字典。

LZW方法示例

- 用LZW方法来压缩字符串“aaabbbbaabaaba”。

code	0	1	2	3	4	5	6	7
key	a	b	aa	aab	bb	bbb	bbba	aaba

压缩串 = 0214537

LZW方法示例

- 用LZW方法来压缩字符串“aaabbbbaabaaba”。

代码  
0 1  
关键字  
a b

aaabbbbaabaaba

压缩串=unll  
a) 初始布局

0	1	2	3	4
a	b	aa	aab	bb

aaabbbbaabaaba

压缩串=021  
d) aaab被压缩

0	1	2	3	4	5	6	7
a	b	aa	aab	bb	bbb	bbba	aaba

aaabbbbaabaaba

压缩串=021453  
g) aaabbbbaabaaba被压缩

代码  
0 1 2  
关键字  
a b aa

aaabbbbaabaaba

压缩串=0  
b) a被压缩

0	1	2	3	4	5
a	b	aa	aab	bb	bbb

aaabbbbaabaaba

压缩串=0214  
e) aaabbb被压缩

0	1	2	3	4	5	6	7
a	b	aa	aab	bb	bbb	bbba	aaba

aaabbbbaabaaba

压缩串=021453  
g) aaabbbbaabaaba被压缩

代码  
0 1 2 3  
关键字  
a b aa aab

aaabbbbaabaaba

压缩串=02  
c) aaa被压缩

0	1	2	3	4	5	6
a	b	aa	aab	bb	bbb	bbba

aaabbbbaabaaba

压缩串=02145  
f) aaabbbbaabaaba被压缩

压缩串 = 0214537

10.6 应用—文本压缩

- LZW压缩实现:
- 设当前前缀p为空, 读取下一个字符c;  
循环 (当前字符串pc不为空)
  - { if (当前字符串pc在字典中)
  - 当前前缀p=当前字符串pc;
  - else { 将当前字符串pc插入到字典中;
  - 输出当前前缀p的编码;
  - p=c;
  - }
  - 读取下一个字符c;
  - }
  - 输出最后一个当前前缀p的编码;

LZW解压缩方法

- 0 2 1 4 5 3 7 --> aaabbbbaabaaba
- fc(p): text(p)的第一个字符
- 1) 首先, 初始化字典, 在其中插入 (0, a)和(1, b)。压缩文件的第一个代码为0, 则应用a代替。
- 2) 下一个代码2不在字典中。因前一个代码为0, 且text(0)=a, fc(0)=a, 则text(2)=text(0)fc(0)=aa。因此用aa代替2, 并把(2, aa)插入字典中。
- 3) 下一个代码1在字典中, 直接由b来替换且把(3, text(2)fc(1))=(3, aab)插入字典中。
- 4) 下一代码4不在字典中。其前面的代码为1, 则text(4)=text(1)fc(1)=bb。把(4, bb)插入字典中, 且在解压文件中输出bb。
- 5) 下一个代码5不在字典中, 其前面的代码为5, 则text(5)=text(4)fc(4)=bbb, 把(5, bbb)被插入字典中, 同时把bbb输出到解压文件中。
- 6) 再下一代码为3在字典中, text(3)=aab则把aab输出, 并将(6, text(5)fc(3))=(6, bbba)插入字典。
- 7) 当遇到7时不在字典中, 把(7, text(3)fc(3))=(7, aaba)插入字典中并输出aaba。把分配给单一字母的代码插入字典中。

## LZW解压缩方法

1. 把分配给单一字母的代码插入字典中。
2. 输入第一个代码，用相应的文本（第一个代码一定对应于一个单一的字母）代替。
3. 设当前输入代码为p，q为p前面的代码。  
两种情况：1) 在字典中；2) 不在字典中。
  - 1) 当p在字典中时，
    - 找到与p相关的文本text(p)并输出。
    - text(q)fc(p)插入字典
  - 2) 当p不在字典中，此时只会是一种情况，
    - text(p)=text(q)fc(q)，(p, text(p)) 插入字典
    - 代码串:qp 对应文本串:text(q)text(q)fc(q)。因此p被解码为text(p)并输出。

## LZW解压缩实现：

1. 把分配给单一字母的代码插入字典中。
2. 输入第一个代码q，输出相应的文本（第一个代码一定对应于一个单一的字母）。
3. 循环：
  - ①输入下一个代码p；
  - ②If (p在字典中)
    - {输出代码p对应的文本串text(p)；
    - 将text(q)fc(p)及代码插入到字典中}
  - Else{【此时只会是一种情况：  
text(p)=text(q)fc(q)；  
即代码串qp对应文本串text(q)text(q)fc(q)】  
输出代码p对应的文本串text(q)fc(q)；  
将text(q)fc(q)及代码插入到字典中}
  - ③q=p；

### ■ 作业 P258 17