



高级数据结构

湖南大学信息科学与工程学院



提纲

2.1 堆

2.2 散列表



提纲

2.1 堆

2.2 散列表



提纲

2.3.1 堆的概述

2.3.2 堆的调整

2.3.3 维护堆的性质

2.3.4 建堆

2.3.5 堆排序算法

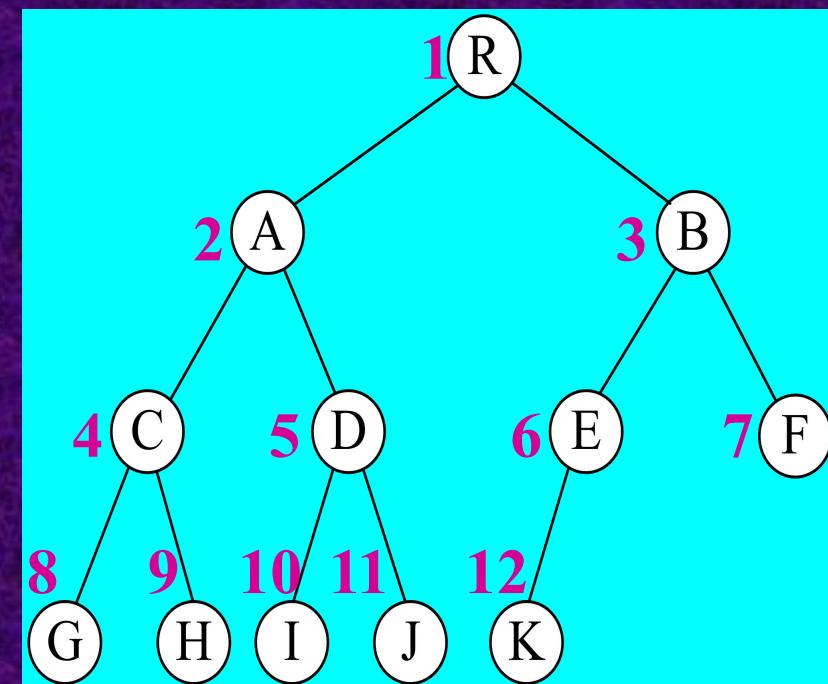
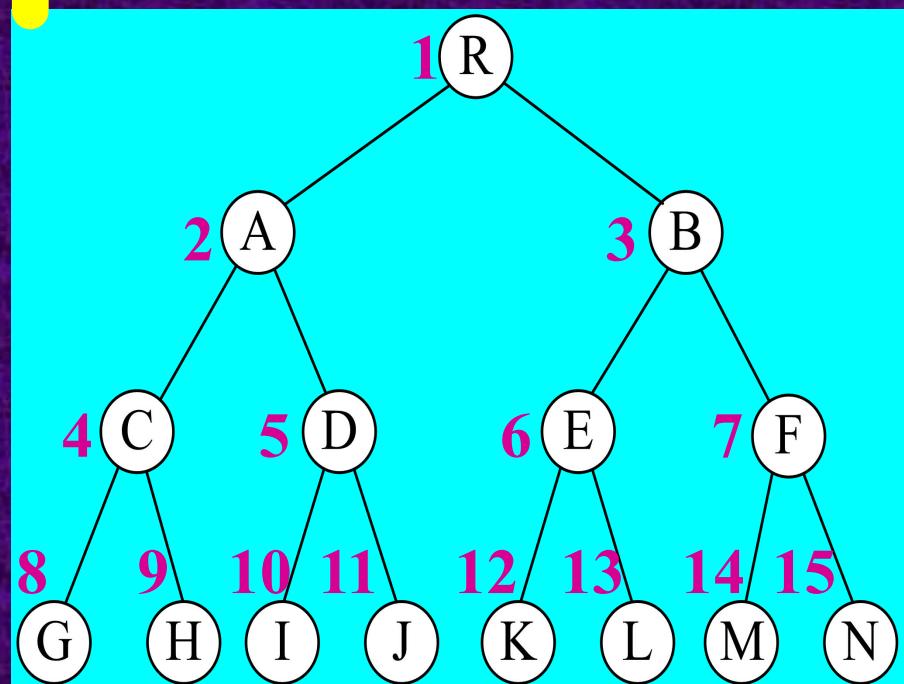
2.3.6 优先队列

■ 满二叉树

二叉树的所有分支结点都有左子树和右子树，并且所有叶子结点都在二叉树的最下一层；

■ 完全二叉树

深度为h的完全二叉树：前h-1层为满二叉树，最后一层的结点必须从左向右连续出现。





堆的定义

n 个元素的序列 $\{k_1, k_2, \dots, k_n\}$ 当且仅当满足如下关系时，称之为堆（heap）。若满足条件 $k_i \leq k_{i/2}$ 且则称大根堆（或最大堆）；若满足条件 $k_i \geq k_{i/2}$ 则称小根堆（或最小堆）。

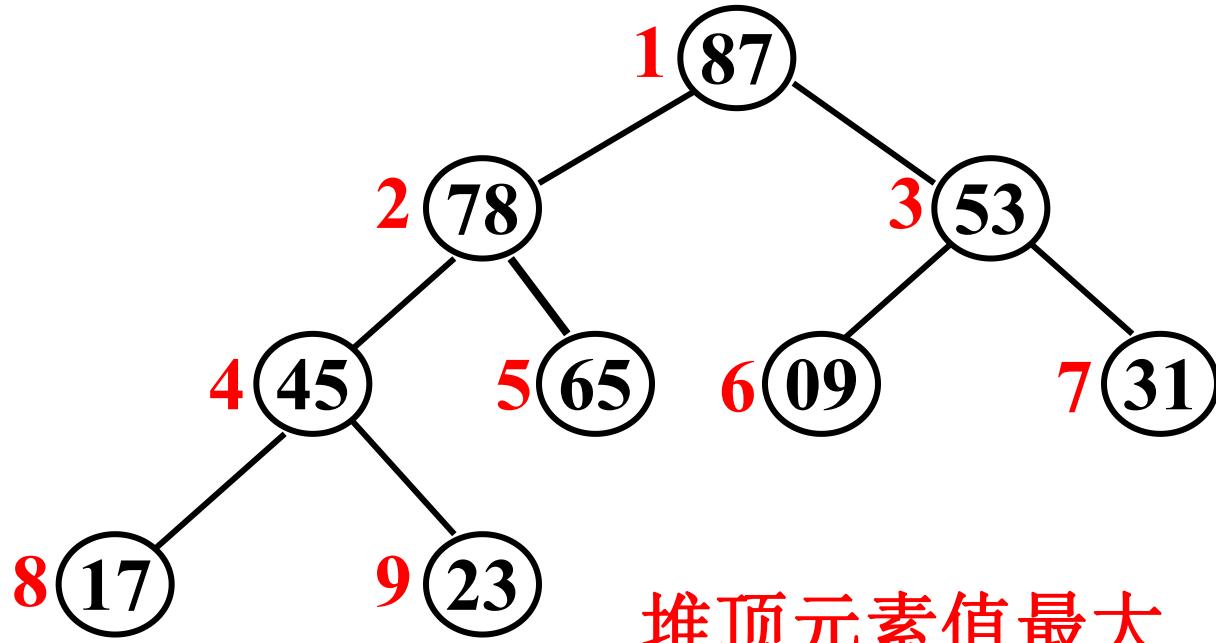
若将和此序列对应的一维数组(即以一维数组作此序列的存储结构)看成是一个完全二叉树，则堆实质上是满足如下性质的完全二叉树：树中任一非叶结点的关键字均不大于(或不小于)其左右孩子(若存在)结点的关键字。



2.3.1 堆的概述

最大堆

1	2	3	4	5	6	7	8	9
87	78	53	45	65	09	31	17	23

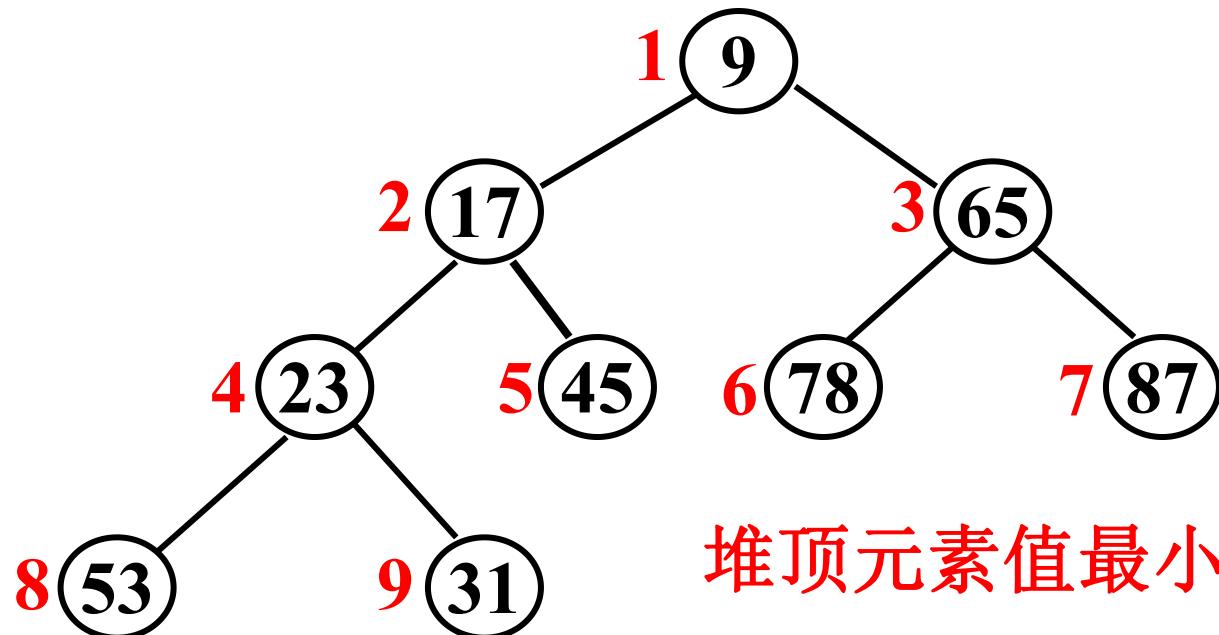


2.3.1 堆的概述



最小堆

1	2	3	4	5	6	7	8	9
9	17	65	23	45	78	87	53	31



堆顶元素值最小

将堆用顺序存储结构来存储，则堆对应一组序列



2.3.1 堆的概述

1、优先级队列的使用场景

- 1)、定时任务轮训问题
- 2)、合并有序小文件

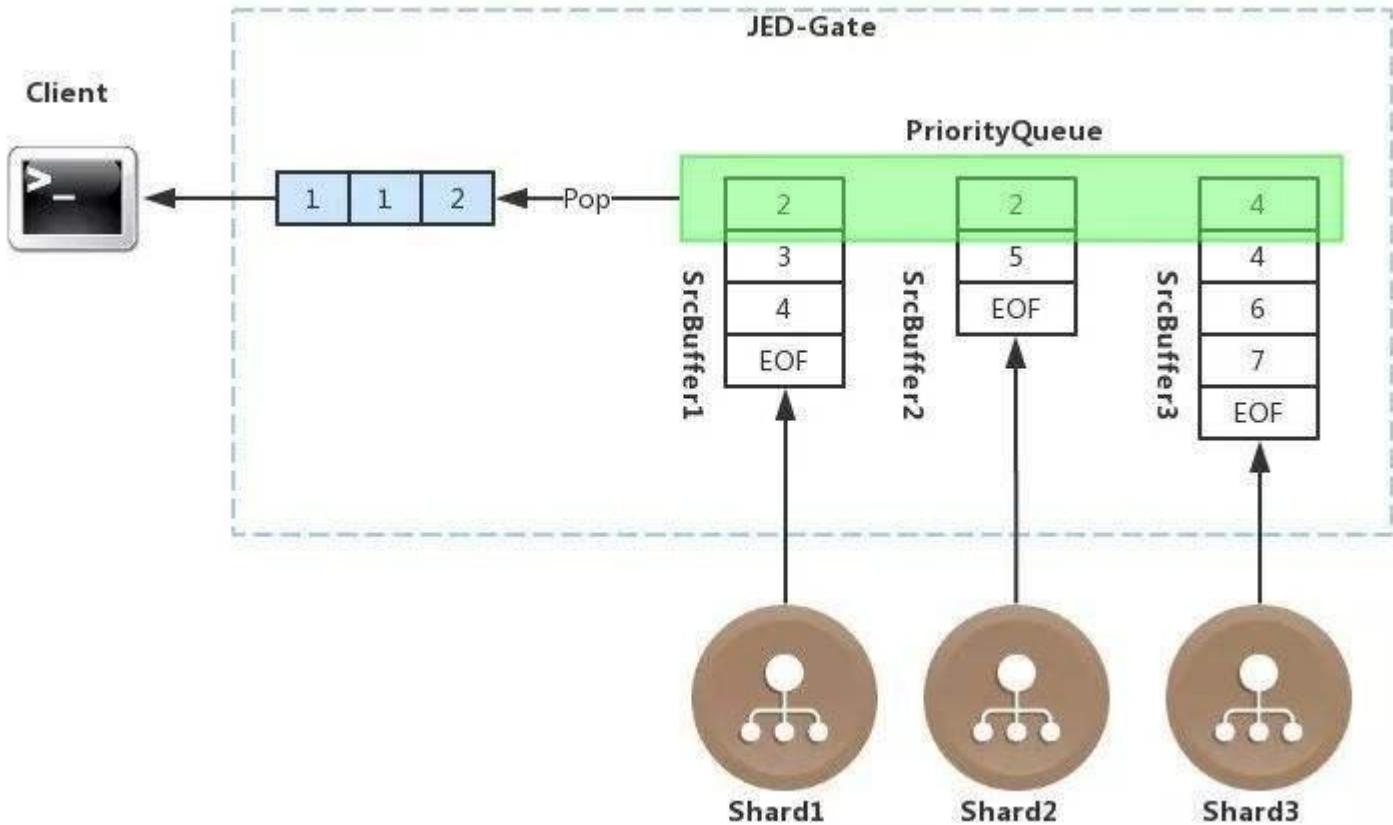
2、求Top K值问题

3、求中位数、百分位数

4、大数据量日志统计搜索排行榜



2.3.1 优先队列应用



优先级队列有两种：最大优先级队列和最小优先级队列

2.3.1 优先队列应用 - 最大优先级队列操作实现



优先队列是一种用来维护由一组元素构成的集合 S 的数据结构，其中每一个元素都有一个关键字(key)，关键字赋予了一个元素的优先级，故名为优先队列。

2.3.1 优先队列应用 - 最大优先级队列操作实现



优先队列的操作：

最大：

- 1) $\text{Insert}(S, x)$: 插入 x 到集合 S 中；
- 2) $\text{Maximum}(S)$: 返回 S 中具有最大关键字的元素；
- 3) $\text{Extract_Max}(S)$: 去掉并返回集合 S 中具有最大关键字的元素；
- 4) $\text{Increase_Key}(S, x, key)$: 将元素 x 的关键字增加到 key 。

最小

- 1) $\text{Insert}(S, x)$: 插入 x 到集合 S 中；
- 2) $\text{Minimum}(S)$: 返回 S 中具有最小关键字的元素；
- 3) $\text{Extract_Min}(S)$: 去掉并返回集合 S 中具有最小关键字的元素；
- 4) $\text{Increase_Key}(S, x, key)$: 将元素 x 的关键字增加到 key 。

2.3.1 优先队列应用 - 最大优先级队列操作实现



堆实现优先队列：

- (1) `HEAP_MAXIMUM` 返回最大堆第一个元素的值 **(堆的调整)**
- (2) `HEAP_EXTRACT_MAX` 实现 `EXTRACT_MAX` 操作，删除最大堆中第一个元素，然后调整堆。 **(堆的调整)**
- (3) `MAX_HEAP_INSERT` 实现 `INSERT` 操作，向最大堆中插入新的关键字。 **(堆的调整)**
- (4) `HEAP_INCREASE_KEY` 实现 `INCREASE_KEY`，通过下标来标识要增加的元素的优先级 `key`，增加元素后需要调整堆，从该节点的父节点开始自顶向上调整。 **(维护堆的性质)**

2.3.1 堆的概述



完全二叉树的存储方式（逻辑表达）

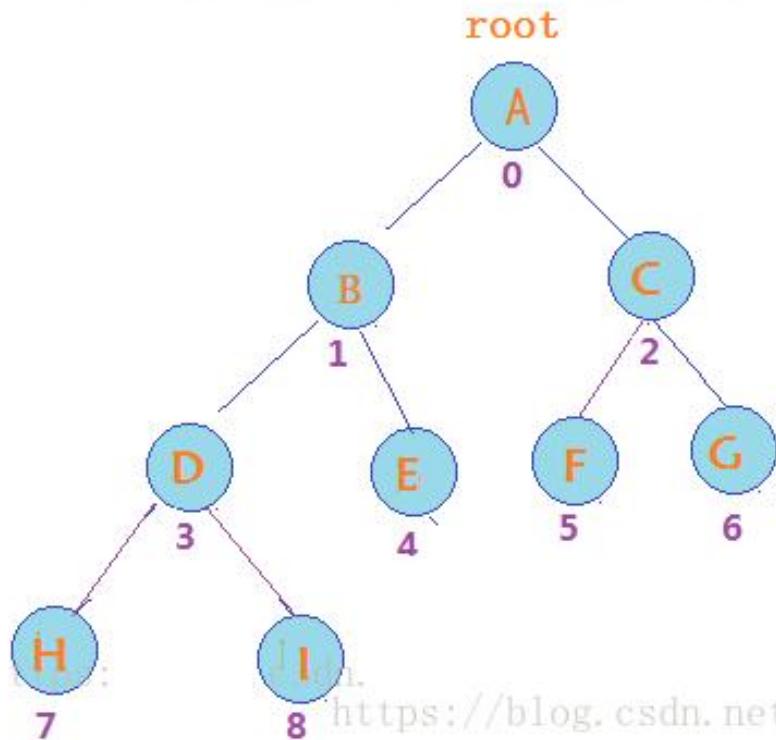
- ✓ 顺序存储（数组）（物理实现）
- ✓ 链式存储（链表）（物理实现）



2.3.1 堆的概述

完全二叉树的存储方式：顺序存储

A	B	C	D	E	F	G	H	I
0	1	2	3	4	5	6	7	8

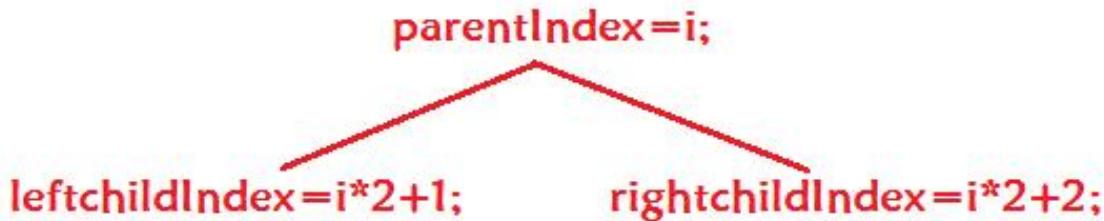




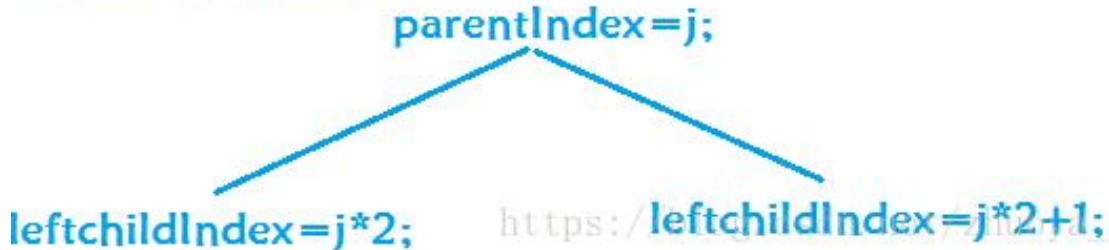
2.3.1 堆的概述

完全二叉树的存储方式：顺序存储

从数组的0号下标开始存储时
逻辑意义上根结点和左、
右孩子的索引关系：



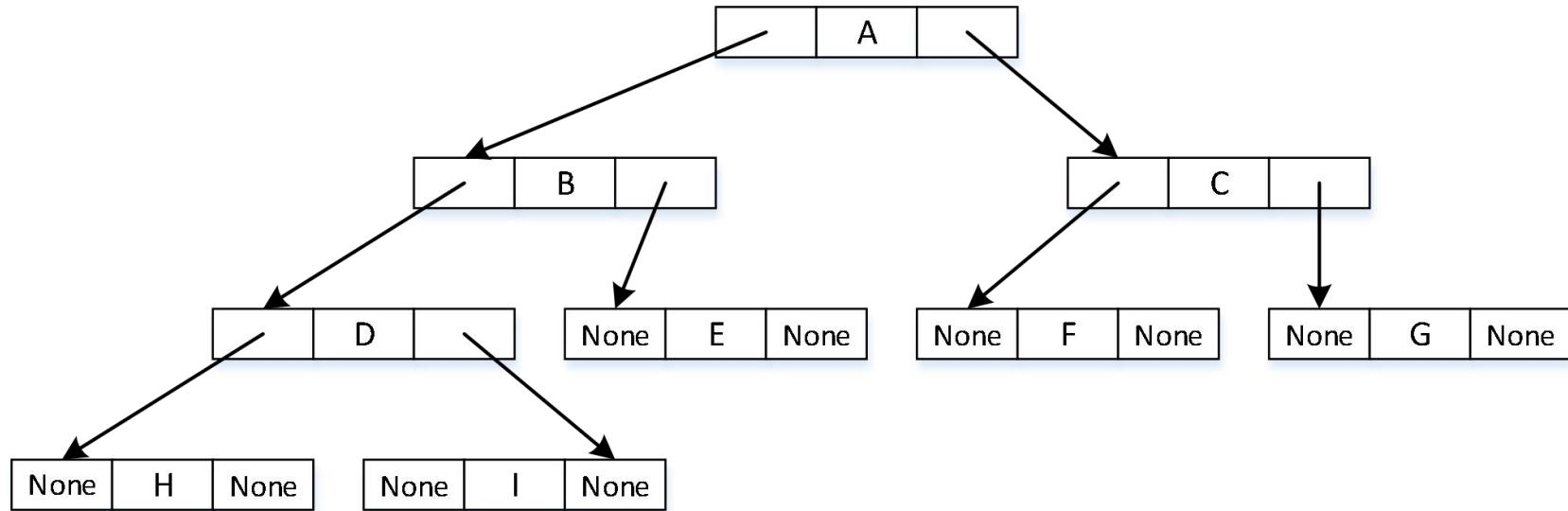
从数组的1号下标开始存储时
逻辑意义上根结点和左、
右孩子的索引关系：



2.3.1 堆的概述



完全二叉树的存储方式：链式存储





提纲

2.3.1 堆的概述

2.3.2 堆的调整

2.3.3 维护堆的性质

2.3.4 建堆

2.3.5 堆排序算法

2.3.6 优先队列

2.3.2 堆的调整



如何在输出堆顶元素后，调整剩余元素为一个新的堆？

解决方法：

输出堆顶元素之后，以堆中最后一个元素替代之；然后将根结点值与左、右子树的根结点值进行比较，并与其中小者进行交换；重复上述操作，直至叶子结点，将得到新的堆，称这个从堆顶至叶子的调整过程为“筛选”

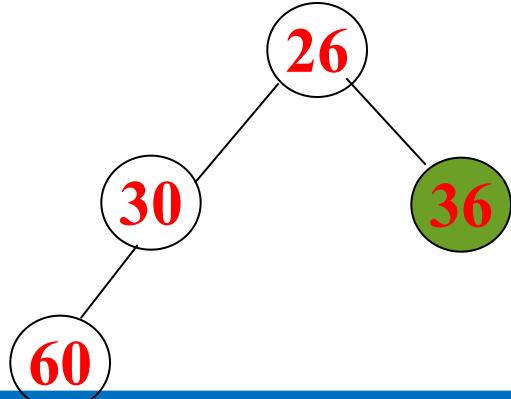
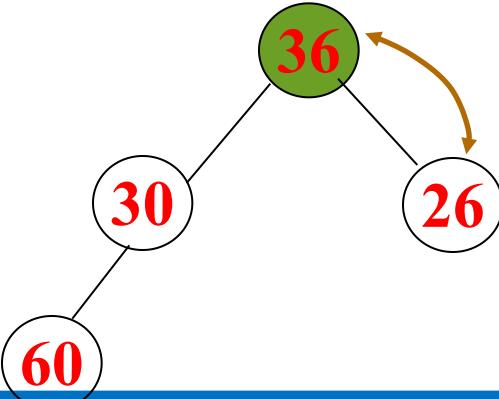
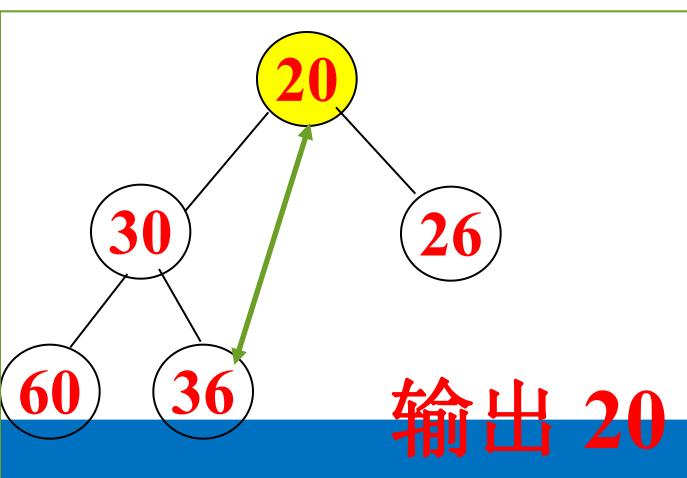
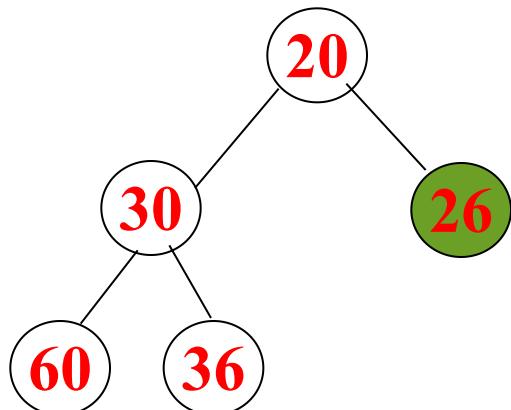
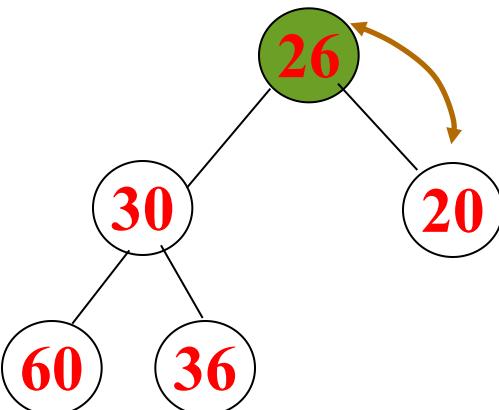
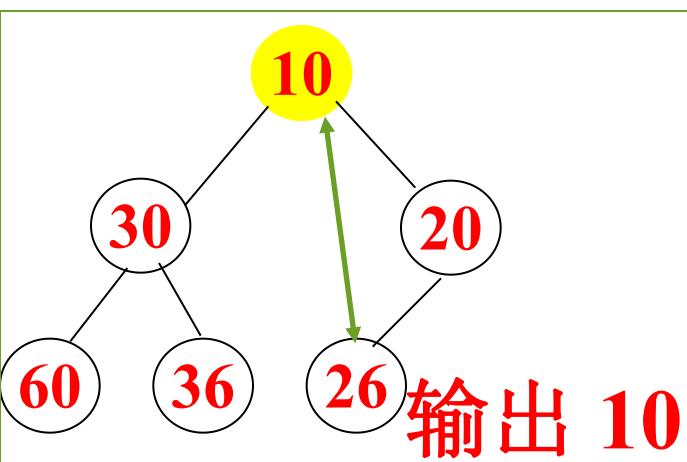
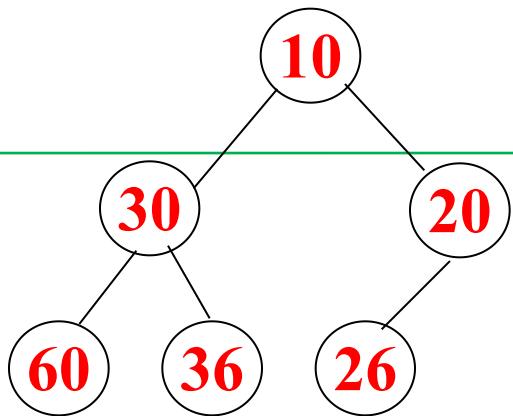


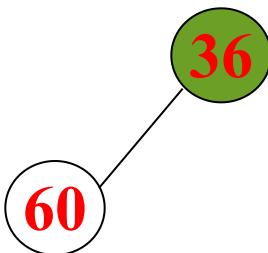
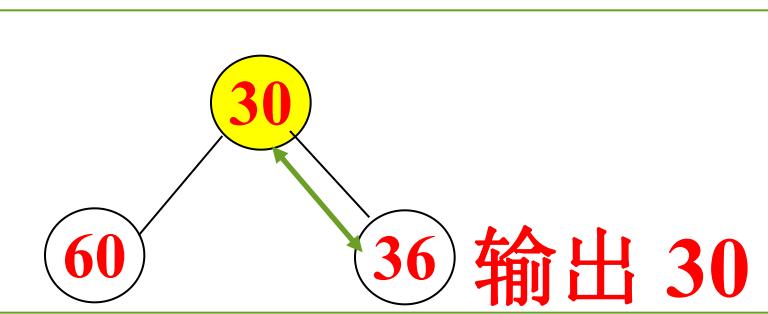
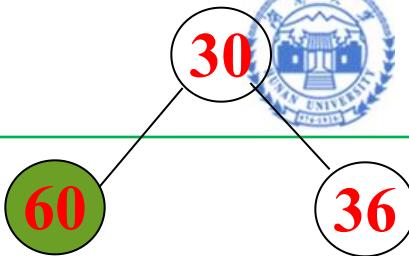
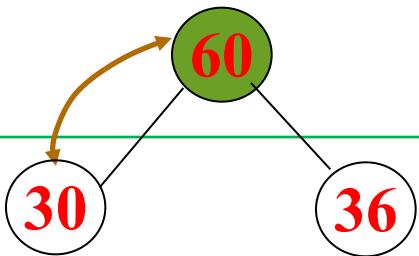
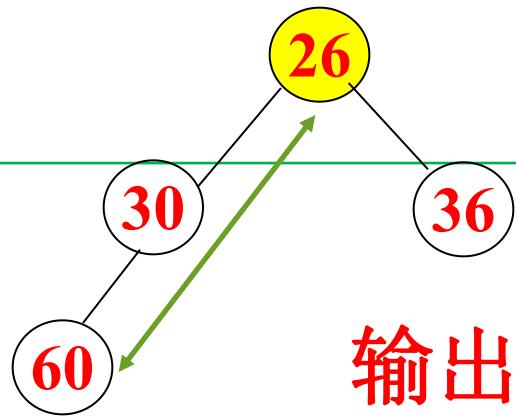
2.3.2 堆的调整

动画演示



最小堆输出序列





递增输出序列为: 10 20 26 30 36 60

2.3.2 堆的调整



如何在堆中插入一个元素？

解决方法：

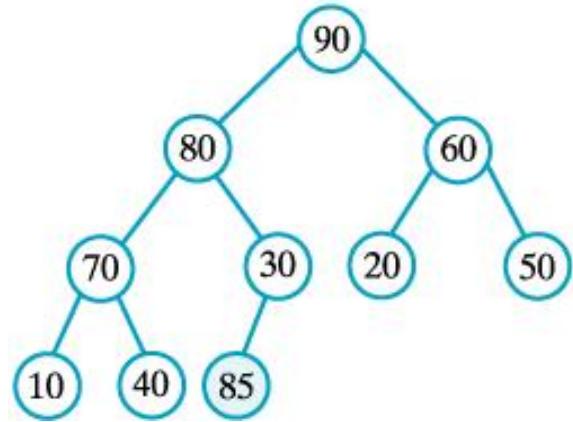
在堆的最后增加一个元素，然后新添加的元素不能比它的父元素大，如果比它的父元素大，就涉及到上浮的动作

2.3.2 堆的调整

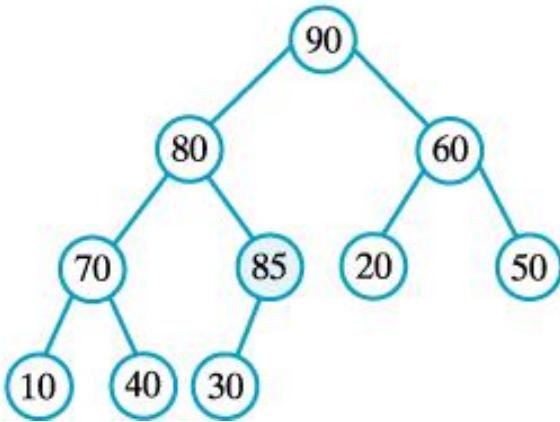


最大堆中加入元素 85

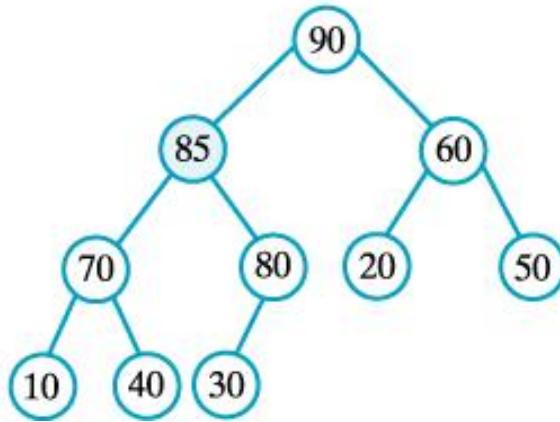
(a)



(b)



(c)





2.3.2 堆的调整

最大堆中加入元素 85

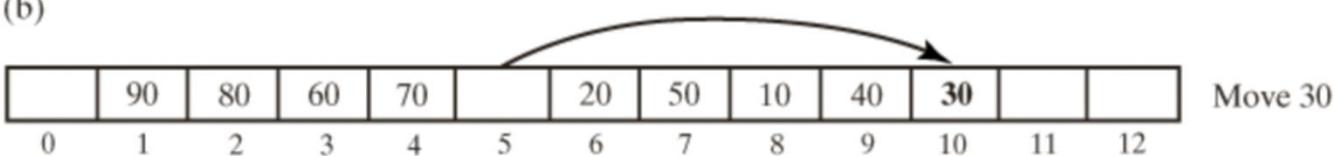
(a)

	90	80	60	70	30	20	50	10	40			
0	1	2	3	4	5	6	7	8	9	10	11	12

$(10/2)$

$85 > 30$

(b)



Move 30

(c)

	90	80	60	70		20	50	10	40	30		
0	1	2	3	4	5	6	7	8	9	10	11	12

$(5/2)$

$85 > 80$



2.3.2 堆的调整

最大堆中加入元素 85

(d)

	90		60	70	80	20	50	10	40	30		
0	1	2	3	4	5	6	7	8	9	10	11	12

Move 80

(e)

	90		60	70	80	20	50	10	40	30		
0	1	2	3	4	5	6	7	8	9	10	11	12

$85 < 90$

(f)

		90	85	60	70	80	20	50	10	40	30		
0	1	2	3	4	5	6	7	8	9	10	11	12	

Insert 85



2.3.2 堆的调整

加入元素时向上调整
删除元素时向下调整



提纲

2.3.1 堆的概述

2.3.2 堆的调整

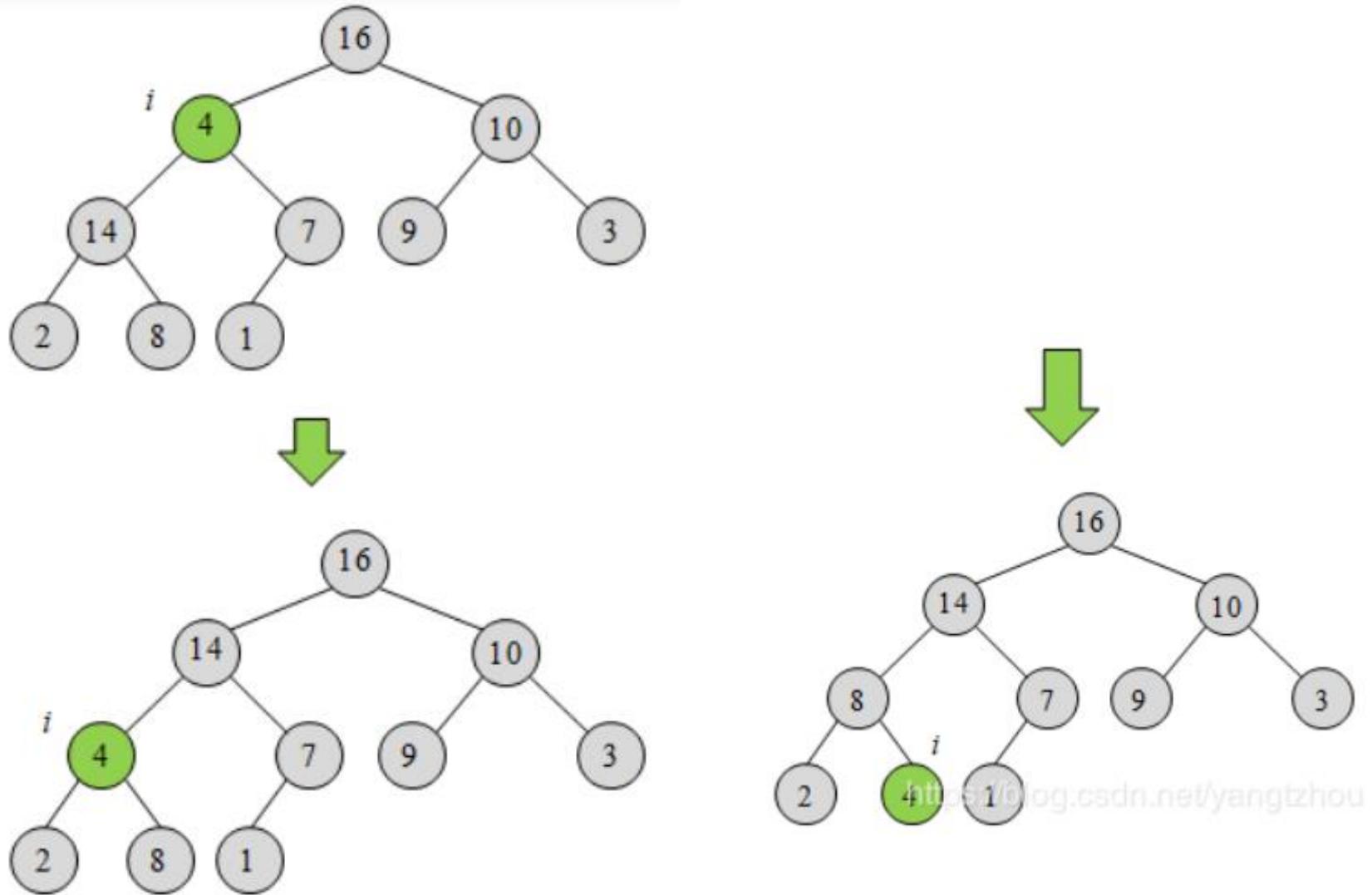
2.3.3 维护堆的性质

2.3.4 建堆

2.3.5 堆排序算法

2.3.6 优先队列

2.3.3 维护堆的性质



tips / blog.csdn.net/yangzhou



2.3.3 维护堆的性质

代码 6.2-1：维护最大堆性质

// 参数 A : 一个最大堆的数组

// 参数 i : 堆中的一个元素

MAX-HEAPIFY(A, i)

$l = \text{LEFT}(i)$

$r = \text{RIGHT}(i)$

if $l \leq A.\text{heap_size}$ and $A[l] > A[i]$

$largest = l$

else

$largest = i$

if $r \leq A.\text{heap_size}$ and $A[r] > A[largest]$

$largest = r$

if $largest \neq i$

exchange $A[i]$ with $A[largest]$

MAX-HEAPIFY($A, largest$)



提纲

2.3.1 堆的概述

2.3.2 堆的调整

2.3.3 维护堆的性质

2.3.4 建堆

2.3.5 堆排序算法

2.3.6 优先队列



2.3.4 建堆

建堆方法：

- ✓ 自底向上
- ✓ 插入建堆



2.3.4 建堆

建堆方法：自底向上（迭代实现）

代码 6.3-1：构建最大堆

// 参数 A : 一个用来构建最大堆的数组

BUILD-MAX-HEAP(A)

$A.heap_size = A.length$

for $i = \lfloor A.length / 2 \rfloor$ **down to** 1

MAX-HEAPIFY(A, i)

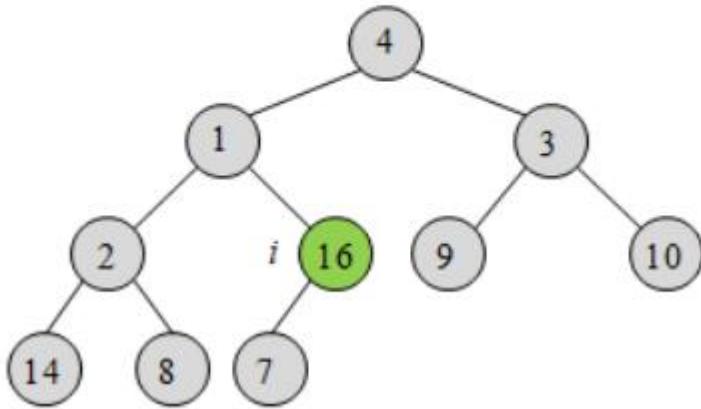


2.3.4 建堆

建堆方法1：自底向上（迭代实现）

原始数组 A

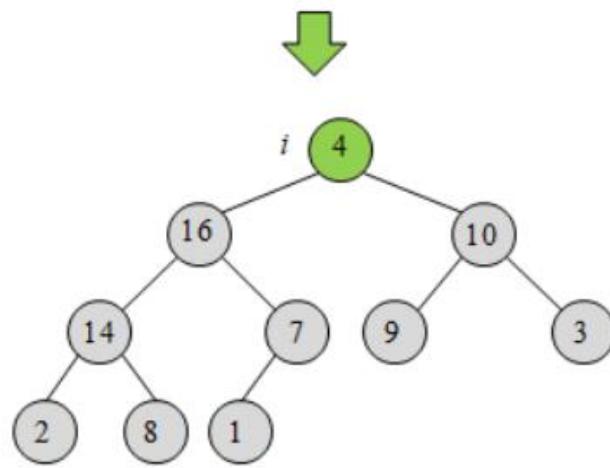
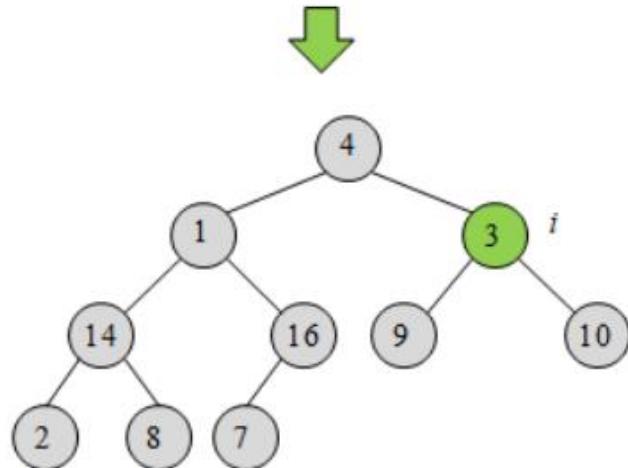
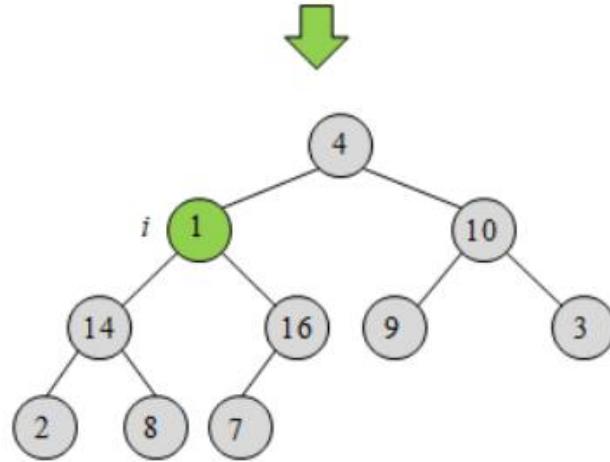
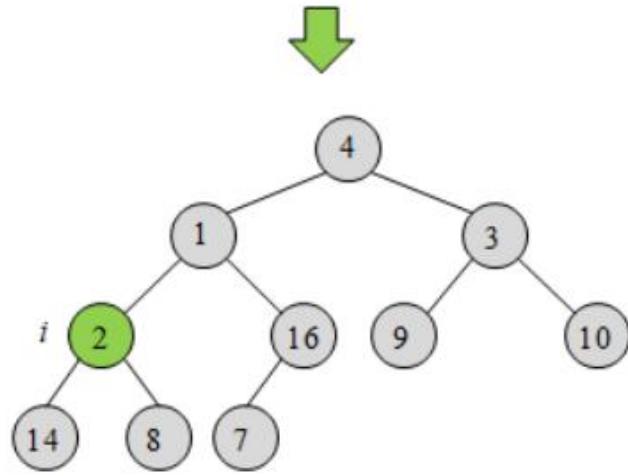
4	1	3	2	16	9	10	14	8	7
---	---	---	---	----	---	----	----	---	---



2.3.4 建堆



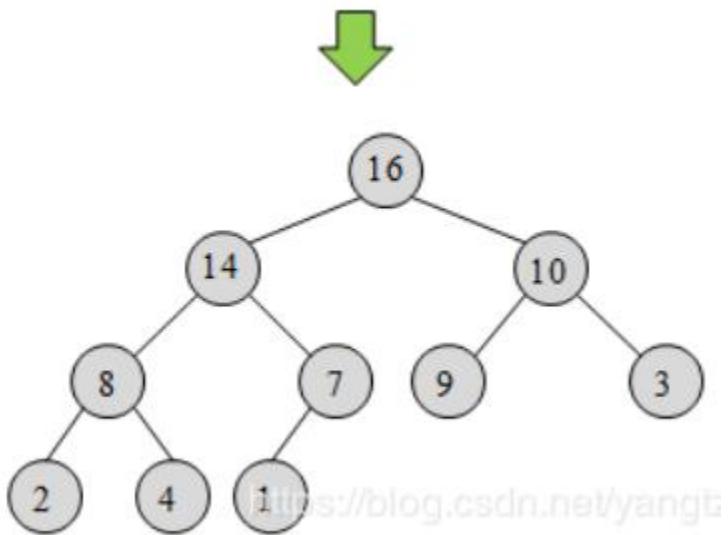
建堆方法1：自底向上（迭代实现）





2.3.4 建堆

建堆方法1：自底向上（迭代实现）



2.3.4 建堆



建堆方法1：自底向上（迭代实现）

```
def max_heap(to_adjust_list, heap_size, index):  
    "调整列表中的元素以保证以index为根的堆是一个最大堆"
```

将当前结点与其左右子节点比较，将较大的结点与当前结点交换，然后递归地调整子树

```
    left_child = 2 * index + 1  
    right_child = left_child + 1  
    if left_child < heap_size and to_adjust_list[left_child] >  
        to_adjust_list[index]:  
        largest = left_child  
    else:  
        largest = index  
    if right_child < heap_size and to_adjust_list[right_child] >  
        to_adjust_list[largest]:  
        largest = right_child  
    if largest != index:  
        to_adjust_list[index], to_adjust_list[largest] = \  
            to_adjust_list[largest], to_adjust_list[index]  
        max_heap(to_adjust_list, heap_size, largest)
```

2.3.4 建堆



建堆方法1：自底向上（迭代实现）

```
def build_max_heap(to_build_list):
    """建立一个堆"""
    # 自底向上建堆
    for i in range(len(to_build_list)//2 - 1, -1, -1):
        max_heap(to_build_list, len(to_build_list), i)

if __name__ == '__main__':
    to_sort_list = [4, 1, 3, 2, 16, 9, 10, 14, 8, 7]
    build_max_heap(to_sort_list)
    print(to_sort_list)
```

2.3.4 建堆



建堆方法2：插入建堆

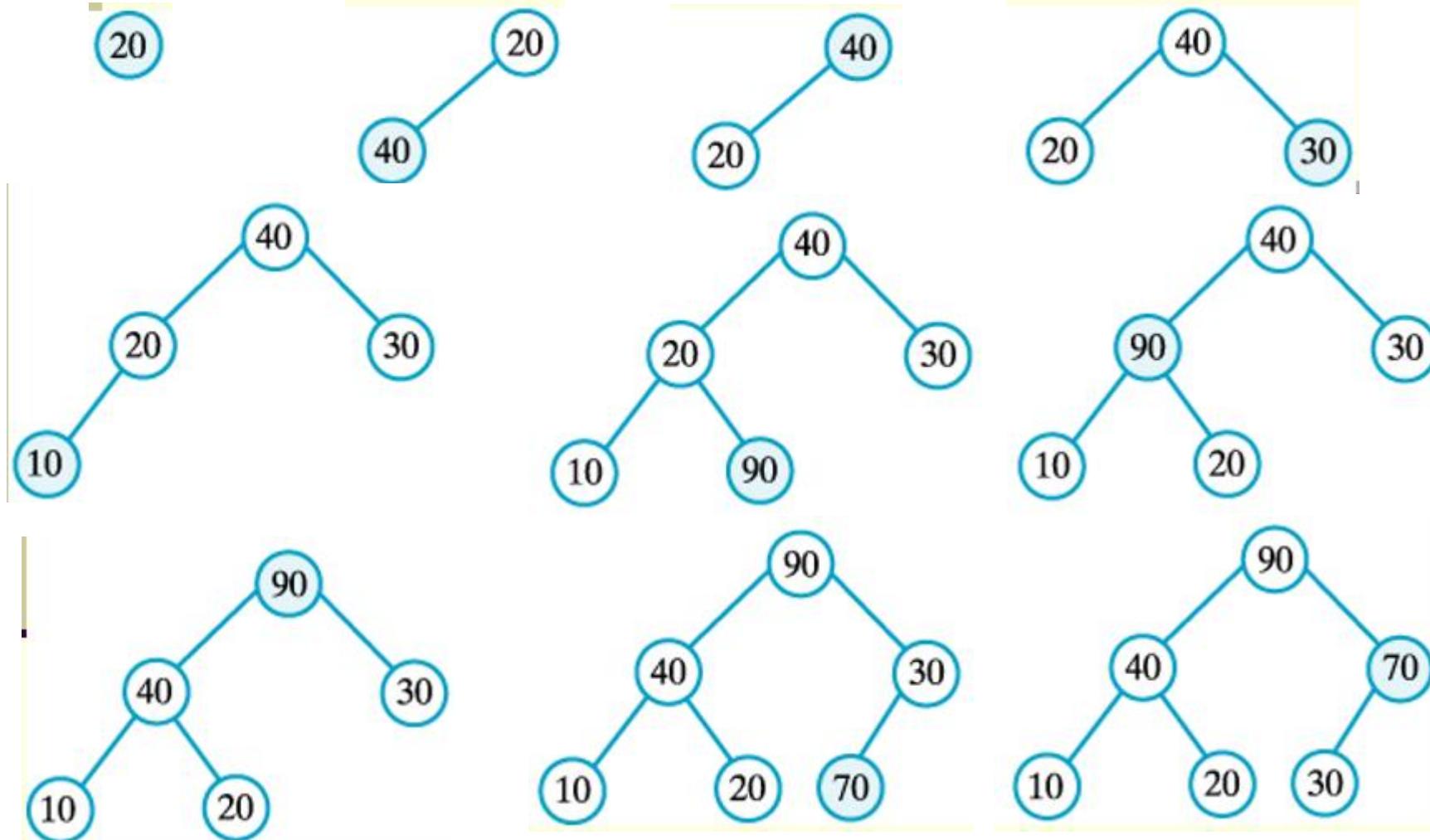
把一个数组看成两部分，左边是堆，右边是还没加入堆的元素，步骤如下：

- 1、数组里的第一个元素自然地是一个堆
- 2、然后从第二个元素开始，一个个地加入左边的堆，当然，每加入一个元素就破坏了左边元素堆的性质，得重新把它调整为堆

2.3.4 建堆



建堆方法2：插入建堆





提纲

2.3.1 堆的概述

2.3.2 堆的调整

2.3.3 维护堆的性质

2.3.4 建堆

2.3.5 堆排序算法

2.3.6 优先队列



2.3.5 堆排序算法

若在输出堆顶的最小值（最大值）后，使得剩余 $n-1$ 个元素的序列重又建成一个堆，则得到 n 个元素的次小值（次大值）……如此反复，便能得到一个有序序列，这个过程称之为**堆排序算法**。

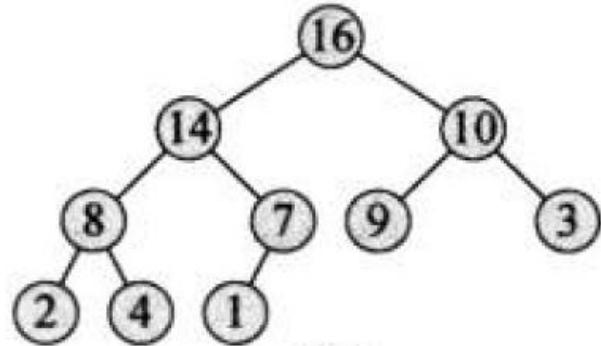


2.3.5 堆排序算法

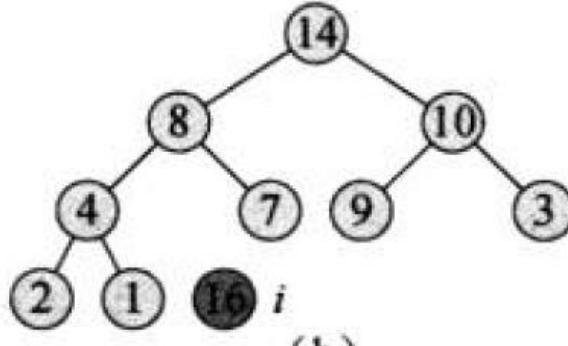
HEAPSORT(A)

- 1 BUILD-MAX-HEAP(A)
- 2 for $i = A.length$ down to 2
- 3 exchange $A[1]$ with $A[i]$
- 4 $A.heap-size = A.heap-size - 1$
- 5 MAX-HEAPIFY($A, 1$)

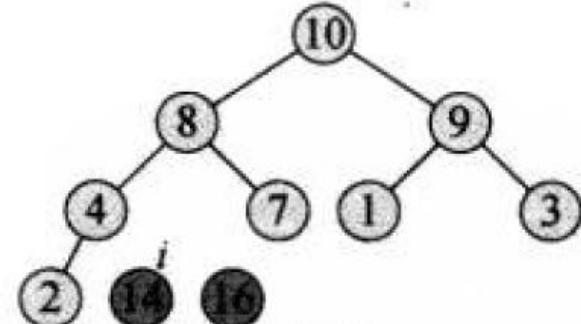
2.3.5 堆排序算法



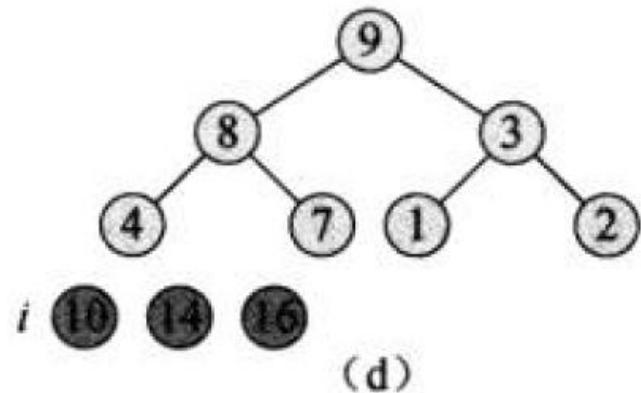
(a)



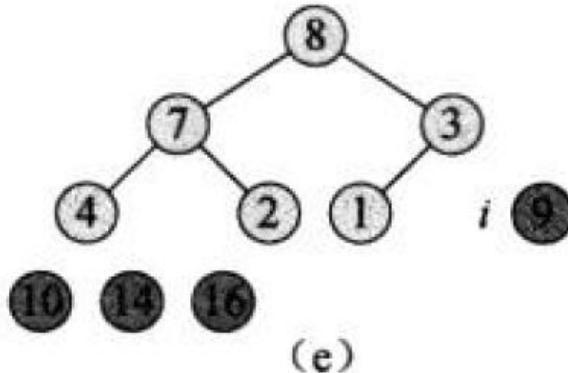
(b)



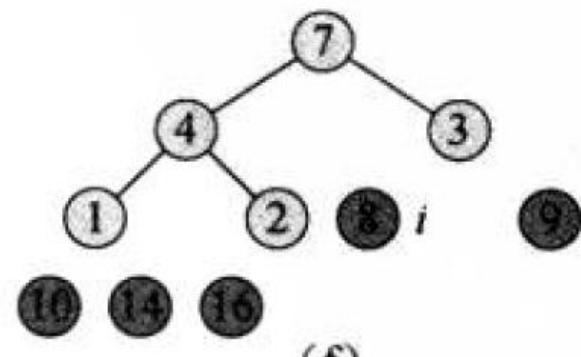
(c)



(d)



(e)



(f)

2.3.5 堆排序算法



堆排序的时间复杂度

堆排序 = 堆调整 + 堆排序

堆排序时间复杂度 = 堆调整时间复杂度 + 堆排序时间复杂度

堆调整的时间复杂度为 $O(n)$

排序重建堆的时间复杂度为 $n \lg(n)$ (递归)

总的时间复杂度为 $O(n+n \lg n) = O(n \lg n)$



提纲

2.3.1 堆的概述

2.3.2 堆的调整

2.3.3 维护堆的性质

2.3.4 建堆

2.3.5 堆排序算法

2.3.6 优先队列

2.3. 6 优先队列-最大优先级队列操作实现



优先队列的操作：

最大：

- 1) $\text{Insert}(S, x)$: 插入 x 到集合 S 中；
- 2) $\text{Maximum}(S)$: 返回 S 中具有最大关键字的元素；
- 3) $\text{Extract_Max}(S)$: 去掉并返回集合 S 中具有最大关键字的元素；
- 4) $\text{Increase_Key}(S, x, key)$: 将元素 x 的关键字增加到 key 。

最小

- 1) $\text{Insert}(S, x)$: 插入 x 到集合 S 中；
- 2) $\text{Minimum}(S)$: 返回 S 中具有最小关键字的元素；
- 3) $\text{Extract_Min}(S)$: 去掉并返回集合 S 中具有最小关键字的元素；
- 4) $\text{Increase_Key}(S, x, key)$: 将元素 x 的关键字增加到 key 。

2.3. 6 优先队列-最大优先级队列操作实现



堆实现优先队列：

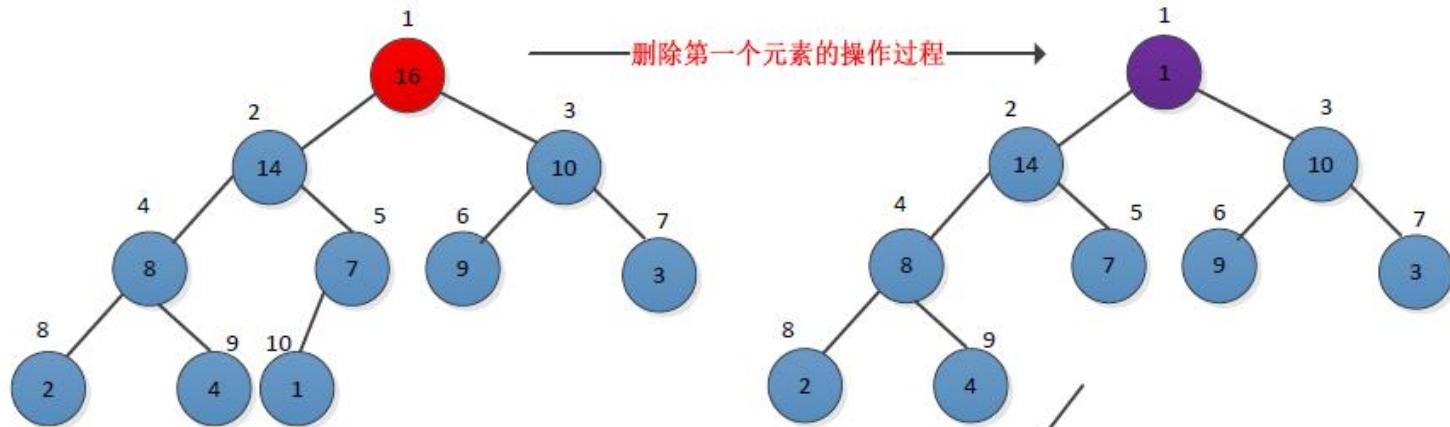
- (1) `HEAP_MAXIMUM`返回最大堆第一个元素的值
(堆的调整)
- (2) `HEAP_EXTRACT_MAX`实现`EXTRACT_MAX`操作，删除最大堆中第一个元素，然后调整堆。
(堆的调整)
- (3) `MAX_HEAP_INSERT`实现`INSERT`操作，向最大堆中插入新的关键字。
(堆的调整)
- (4) `HEAP_INCREASE_KEY`实现`INCREASE_KEY`，通过下标来标识要增加的元素的优先级`key`，增加元素后需要调整堆，从该节点的父节点开始自顶向上调整。
(维护堆的性质)

2.3. 6 优先队列-返回第一个元素



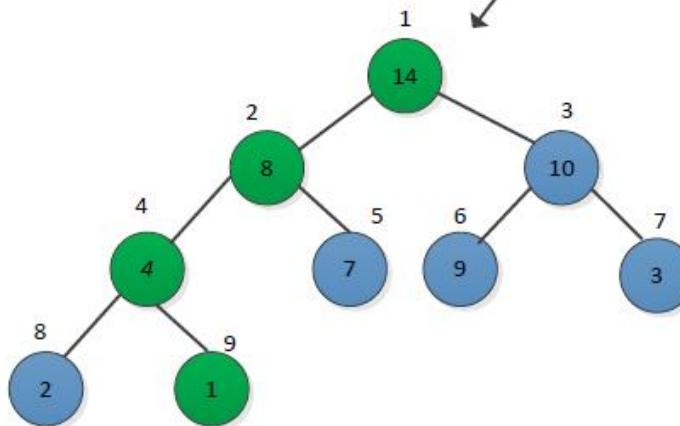
```
MAX-HEAP-MAXIMUM( $A$ )
    if  $A.heap\_size < 1$ 
        error “heap underflow”
    else
        return  $A[1]$ 
```

2.3. 6 优先队列-删除操作

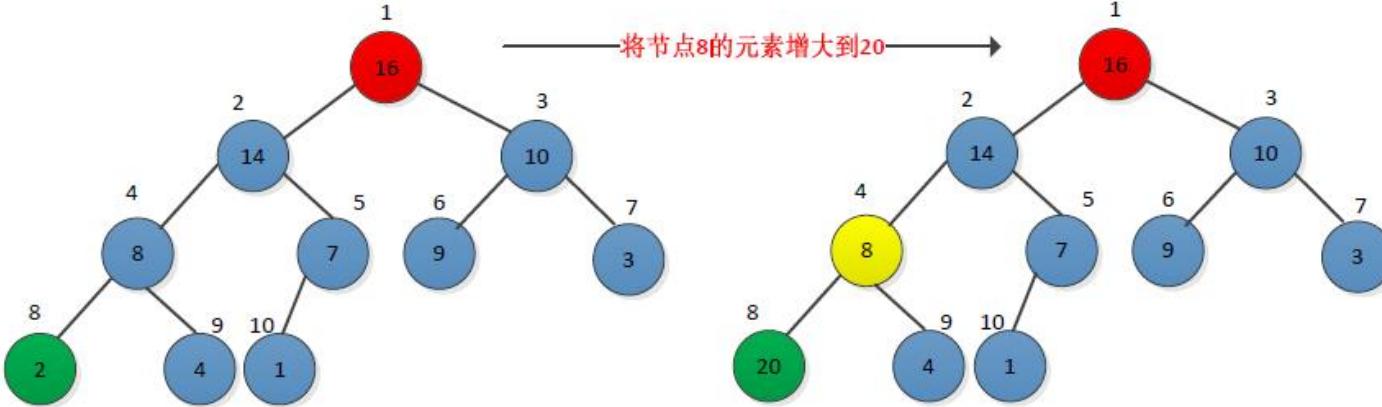


http://blog.csdn.net/chenyun_2010

从第1个元素调整堆

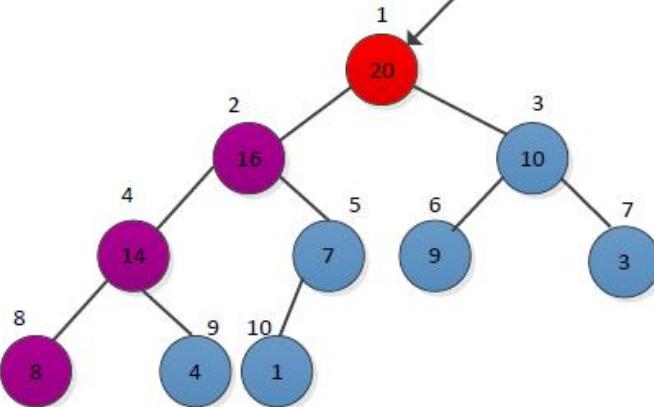


2.3. 6 优先队列-元素值增加操作



http://blog.csdn.net/chenjun_2010

从节点8的父节点开始调整堆



2.3. 6 优先队列-插入元素操作



MAX-HEAP-INSERT(A , key)

if $A.heap_size = A.length$

error “heap is full”

$A.heap_size = A.heap_size + 1$

$A[A.heap_size] = -\infty$

MAX-HEAP-INCREASE-KEY(A , $A.heap_size$, key)



提纲

2.1 堆

2.2 散列表



提纲

2.1.1 概述

2.1.2 直接寻址

2.1.3 散列表

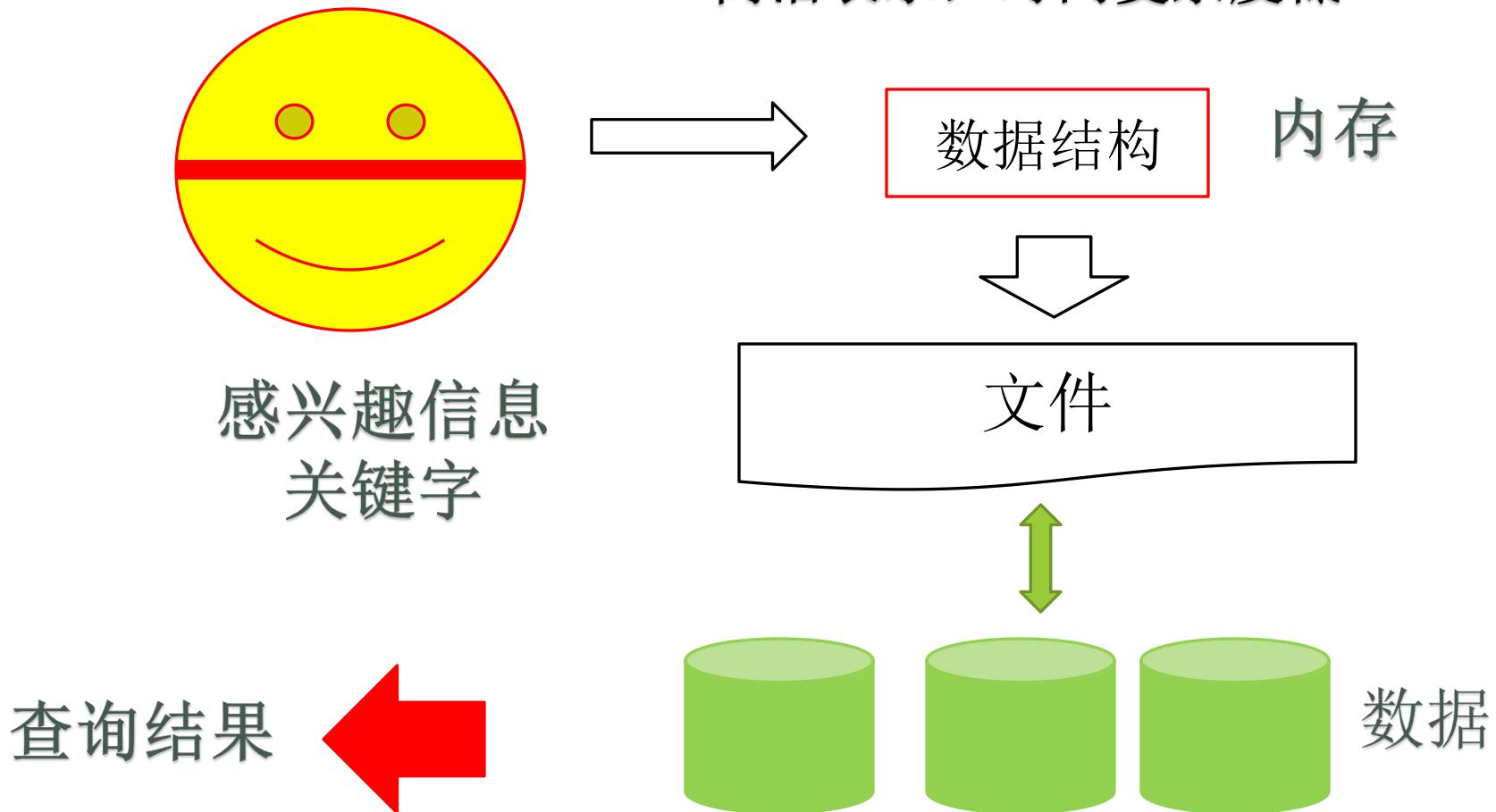
2.1.4 散列函数

2.1.5 开放寻址法

2.1.6 *Cuckoo Hash* 布谷鸟哈希

2.1.7 哈希的扩展应用

2.1.1 概述





提纲

2.1.1 概述

2.1.2 直接寻址

2.1.3 散列表

2.1.4 散列函数

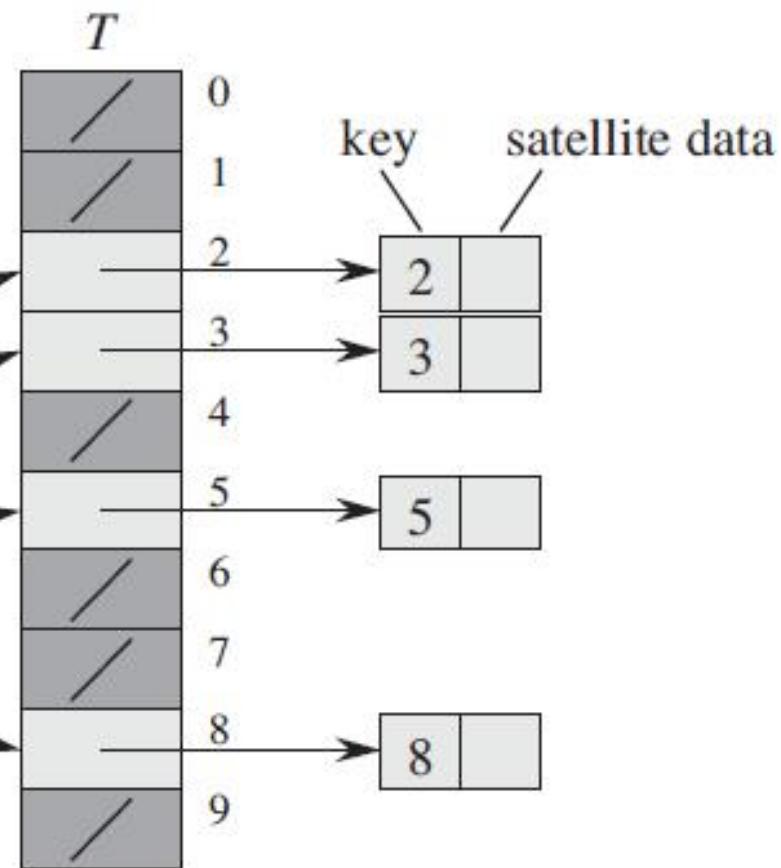
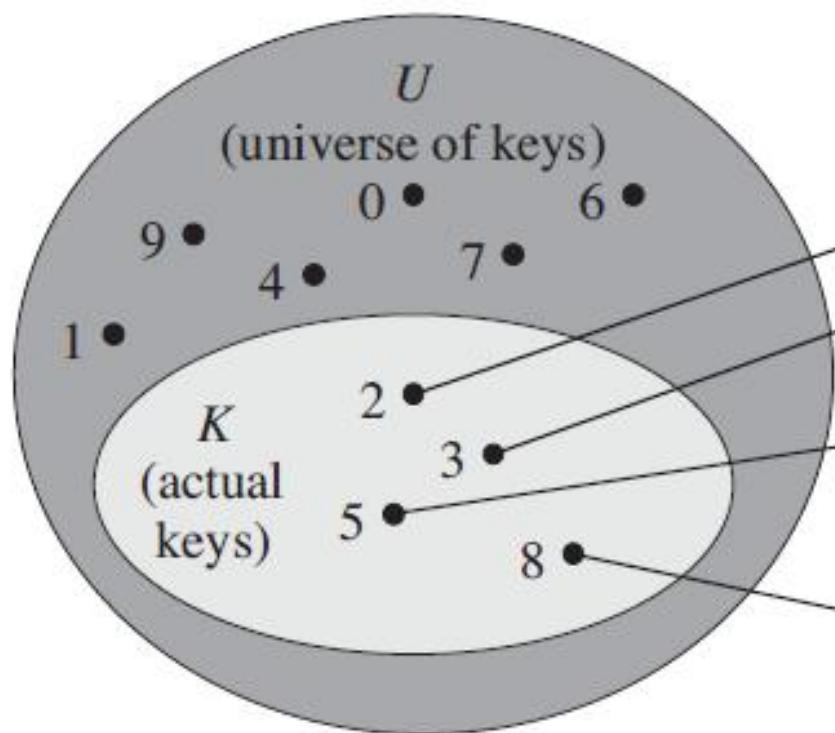
2.1.5 开放寻址法

2.1.6 *Cuckoo Hash* 布谷鸟哈希

2.1.7 哈希的扩展应用



2.1.2 直接寻址表





2.1.2 直接寻址表

DIRECT-ADDRESS-SEARCH(T, k)

1 **return** $T[k]$

DIRECT-ADDRESS-INSERT(T, x)

1 $T[x.key] = x$

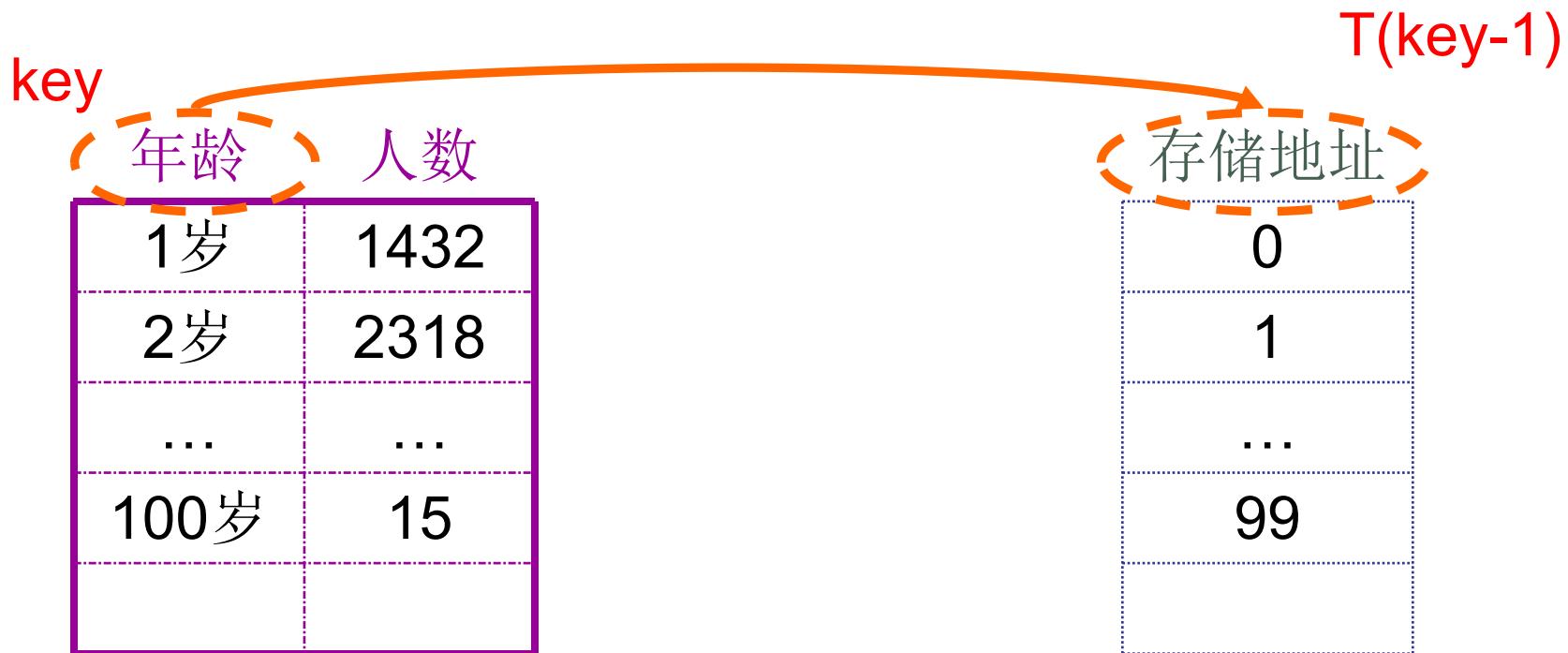
DIRECT-ADDRESS-DELETE(T, x)

1 $T[x.key] = \text{NIL}$



2.1.2 直接寻址表

举例



2.1.2 直接寻址表



直接寻址的缺点

对于全域较大，但是元素却十分稀疏的情况，使用这种存储方式将浪费大量的存储空间。



提纲

2.1.1 概述

2.1.2 直接寻址

2.1.3 散列表

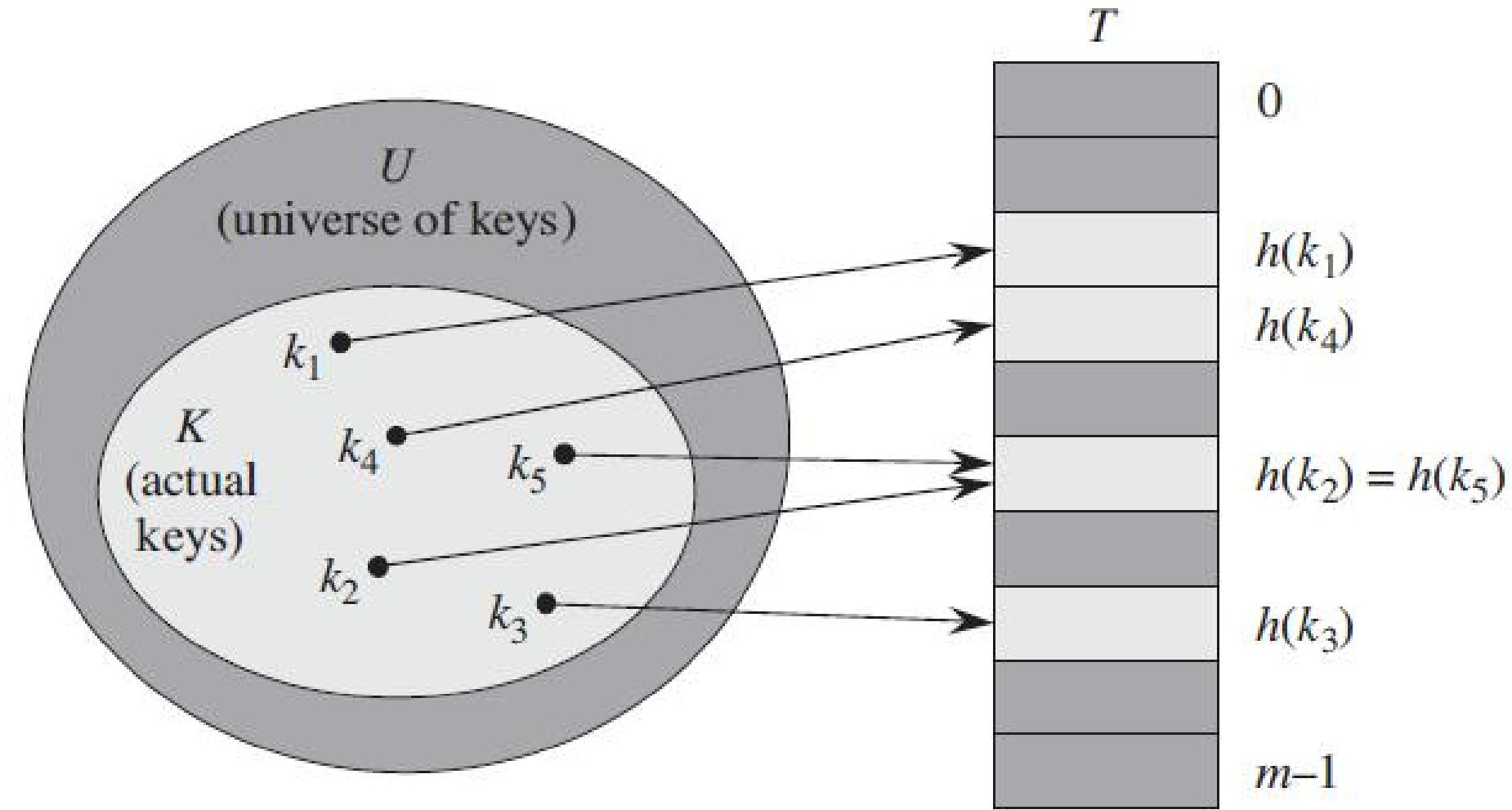
2.1.4 散列函数

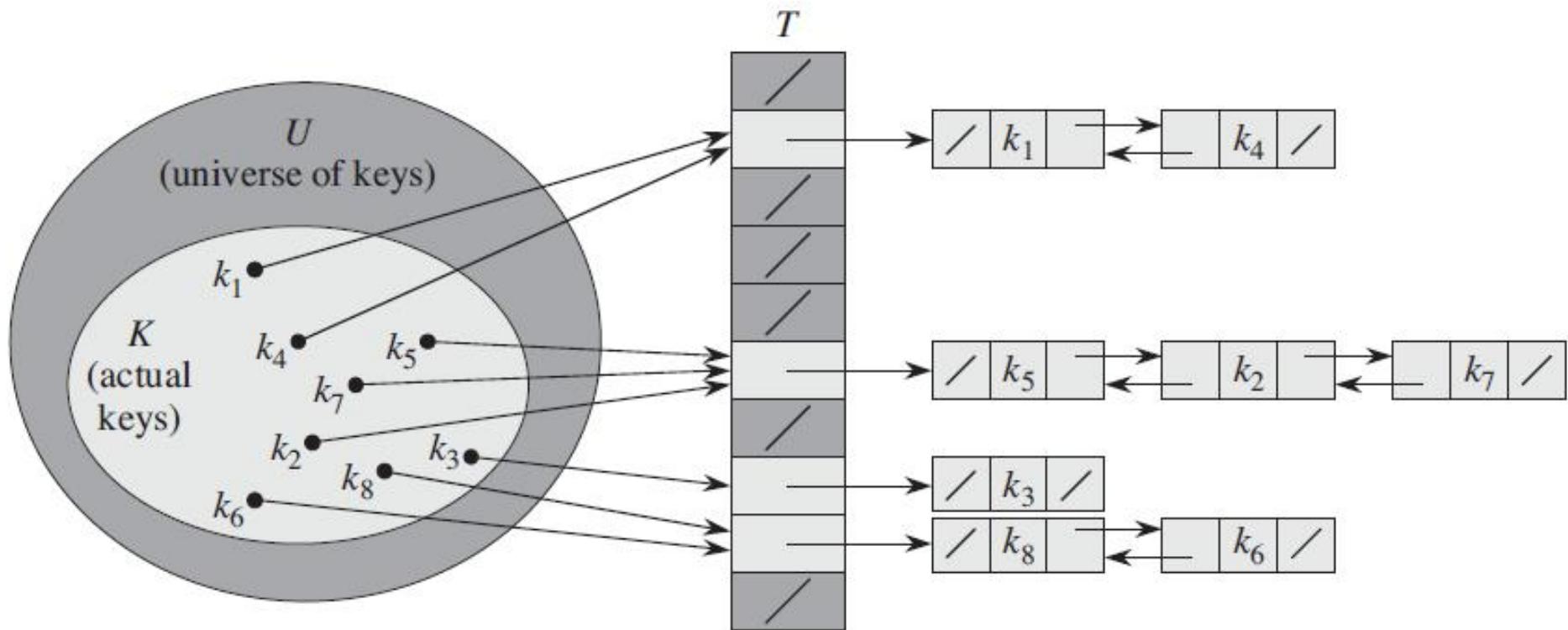
2.1.5 开放寻址法

2.1.6 *Cuckoo Hash* 布谷鸟哈希

2.1.7 哈希的扩展应用

2.1.3 散列表







CHAINED-HASH-INSERT(T, x)

- 1 insert x at the head of list $T[h(x.key)]$

CHAINED-HASH-SEARCH(T, k)

- 1 search for an element with key k in list $T[h(k)]$

CHAINED-HASH-DELETE(T, x)

- 1 delete x from the list $T[h(x.key)]$



链接法散列的分析（算法导论）

①一次不成功查找平均时间分析

在查找不成功的情况下，我们需要遍历链表 $T[j]$ 的每一个元素，而链表 $T[j]$ 的长度是 α ，因此需要时间 $O(\alpha)$ ，加上索引到 $T(j)$ 的时间 $O(1)$ ，总时间为 $\Theta(1 + \alpha)$ 。

②一次成功查找平均时间分析

在查找成功的情况下，我们无法准确知道遍历到链表 $T[j]$ 的何处停止，因此我们只能讨论平均情况。平均时间都为 $\Theta(1+\alpha)$



法)

链接法散列的分析（数据结构）

表 9.2 不同处理冲突的平均查找长度

处理冲突的方法	平均查找长度	
	查找成功时	查找不成功时
线性探测法	$S_{nl} \approx \frac{1}{2} \left(1 + \frac{1}{1-\alpha} \right)$	$U_{nl} \approx \frac{1}{2} \left(1 + \frac{1}{(1-\alpha)^2} \right)$
二次探测法与双哈希法	$S_{nr} \approx -\frac{1}{\alpha} \ln(1-\alpha)$	$U_{nr} \approx \frac{1}{1-\alpha}$
链地址法	$S_{nc} \approx 1 + \frac{\alpha}{2}$	$U_{nc} \approx \alpha + e^{-\alpha}$

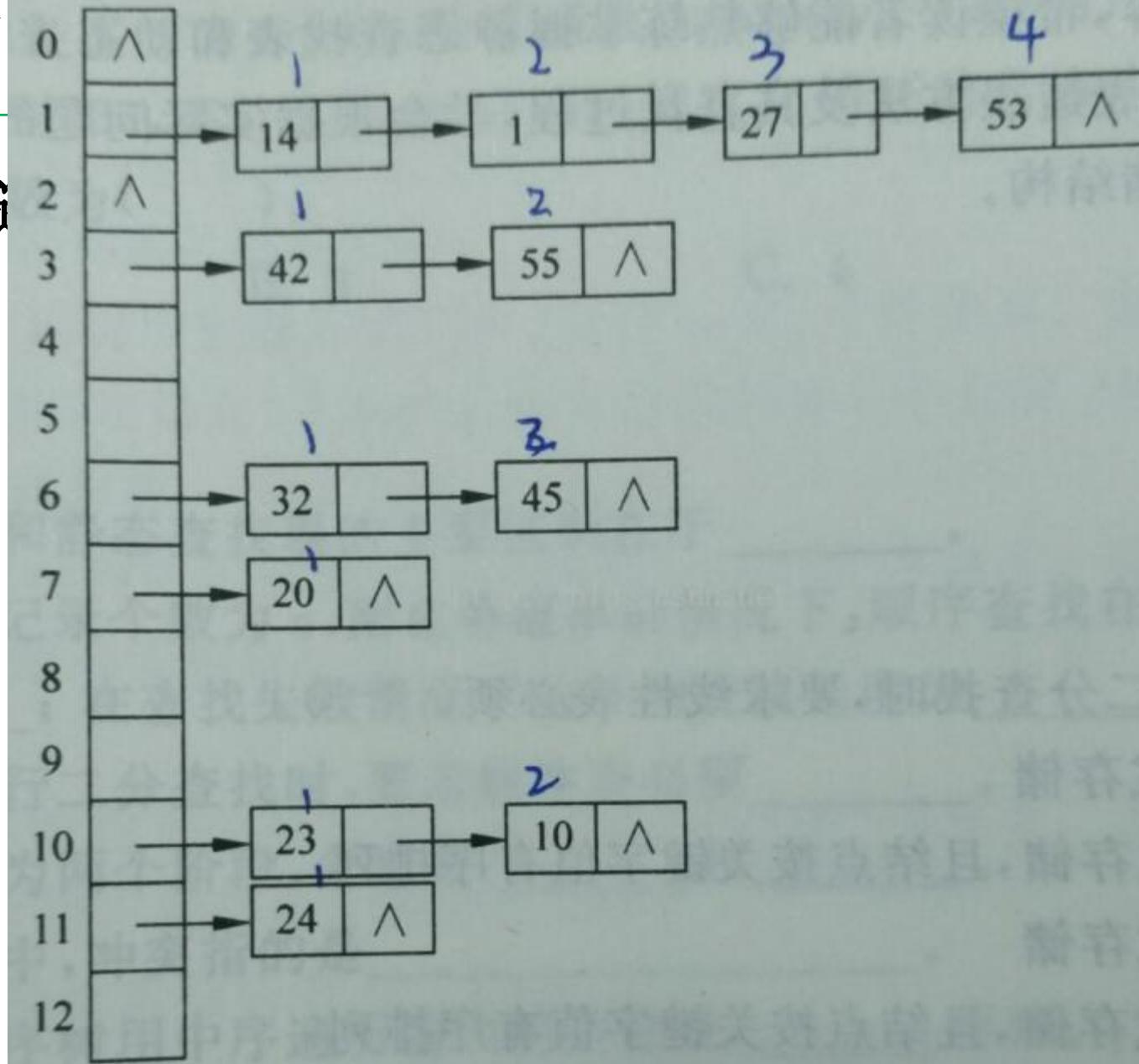


图 9.24 用拉链法解决冲突时的哈希表



提纲

2.1.1 概述

2.1.2 直接寻址

2.1.3 散列表

2.1.4 散列函数

2.1.5 开放寻址法

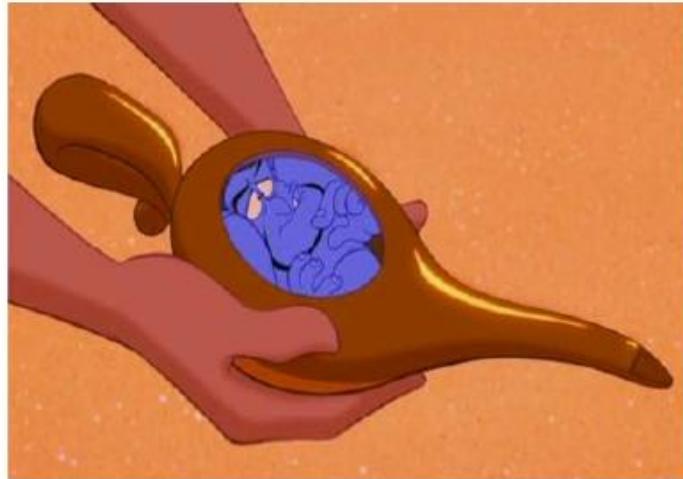
2.1.6 Cuckoo Hash 布谷鸟哈希

2.1.7 哈希的扩展应用

2.1.4 散列函数

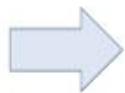


The ‘magic’ of hash functions



**PHENOMENAL
COSMIC
POWERS!!**

惊人的、广大无边的、强势（拉丁神）



itty bitty living space

微人模式、小的居住空间
(拉丁神灯)

2.1.4 散列函数



构造散列函数的基本方法

▶ 基本要求

设关键字集 K 中有 n 个关键字，哈希表长为 m ，即哈希表地址集为 $[0, m-1]$ ，则哈希函数 H 应满足：

1. 对任意 $k_i \in K$, $i=1, 2, \dots, n$, 有 $0 \leq H(k_i) \leq m-1$;
2. 对任意 $k_i \in K$, $H(k_i)$ 取 $[0, m-1]$ 中任一值的概率相等。

2.1.4.1 散列函数:除法散列法



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

A	B	C
十进制 (k)	二进制 (k)	$h(k)$ (其中 $m=8$)
20	10100	4
98	1100010	2
204	11001100	4
71	1000111	7
67	1000011	3
1234	10011010010	2

2.1.4.1 除法散列法的不足



Problems with division method

对于规则性的键值集合

- Regularity
 - Suppose keys are $x, 2x, 3x, 4x, \dots$ 假如 x 与 m 有个公约数 d
 - Suppose x and chosen m have common divisor d
 - Then only use $1/d$ fraction of table
 - x series cycles back, leaving $d-1$ out of d entries blank
 - E.g., m power of 2 and all keys are even, only use half
- So make m a prime number 所以 m 得是个质数
 - But finding a prime number is hard 但是找个质数很难
 - And now you have to divide (slow)

那么我们将只能有效利用到表格 $1/d$ 部分空间

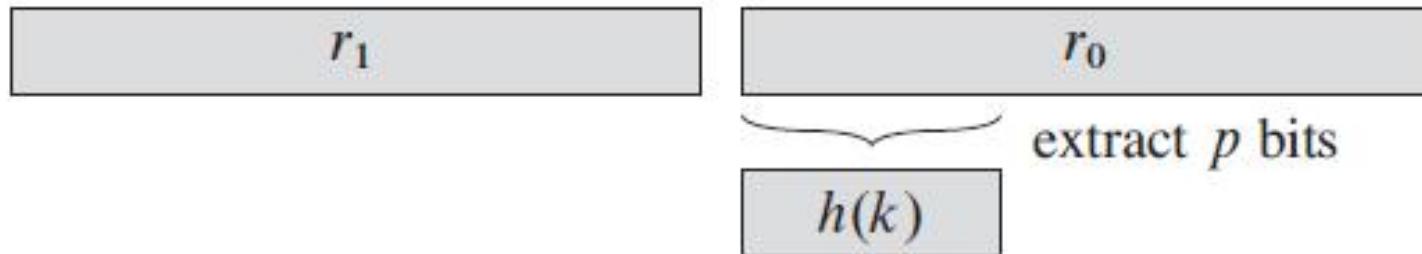
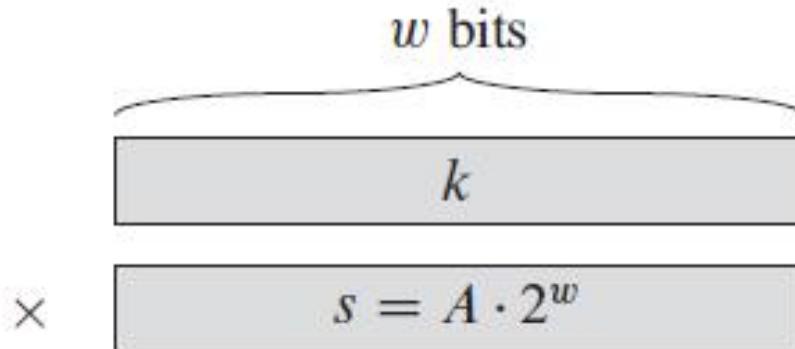
并且你要去相除（比乘法，位移都慢很多）

2.1.4.2 散列函数:乘法散列法



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

$m=2^p$



$$A \approx (\sqrt{5} - 1) / 2 = 0.618033988$$

2.1.4.3 散列函数: 散列值

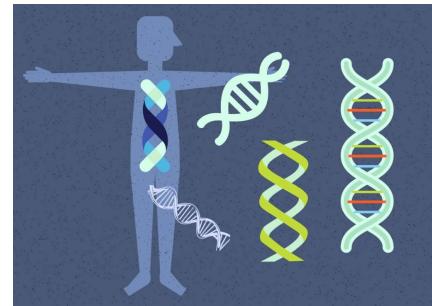


Key不是数字，是字符串呢？

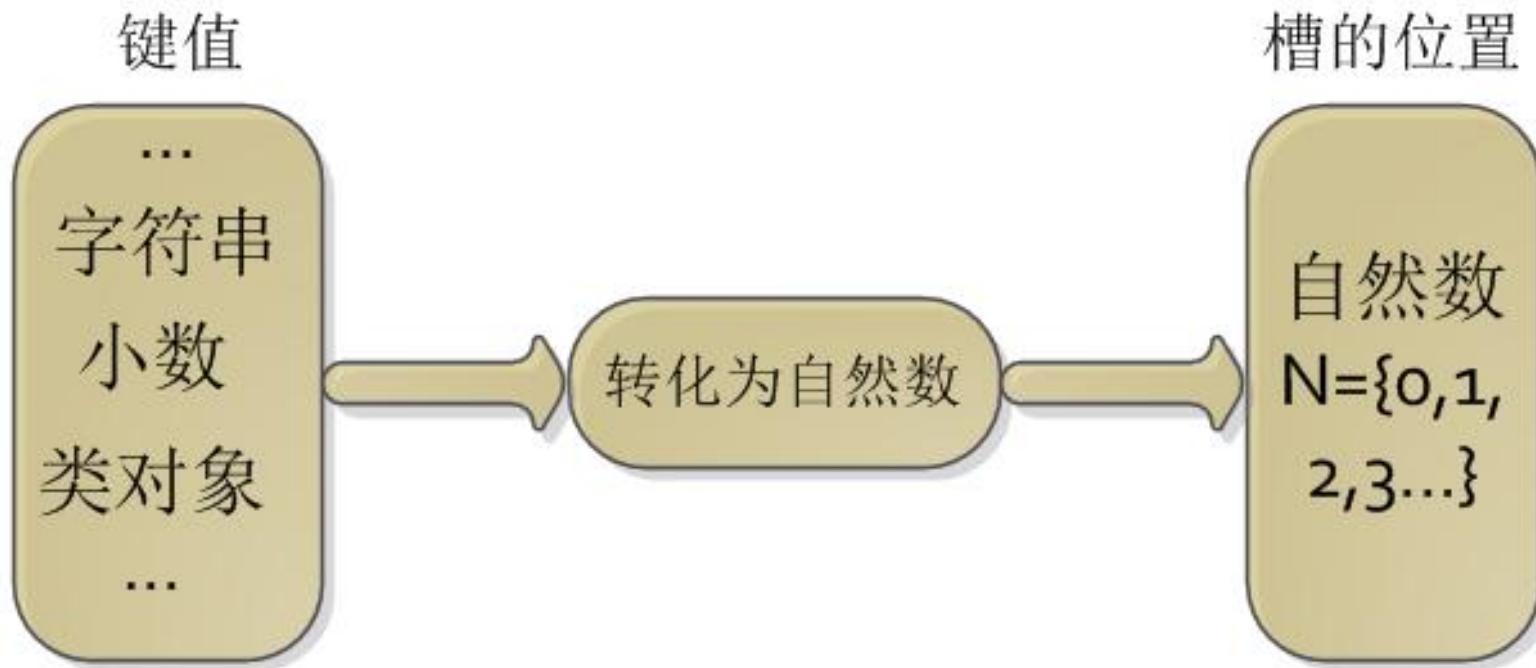
Key = “abcdef”



Hash值



2.1.4.3 散列函数: 散列值



Java: *hashcode*
C++: ASCII码相加

Python:
Key=“abcdef”
A=hash(Key)



提纲

2.1.1 概述

2.1.2 直接寻址

2.1.3 散列表

2.1.4 散列函数

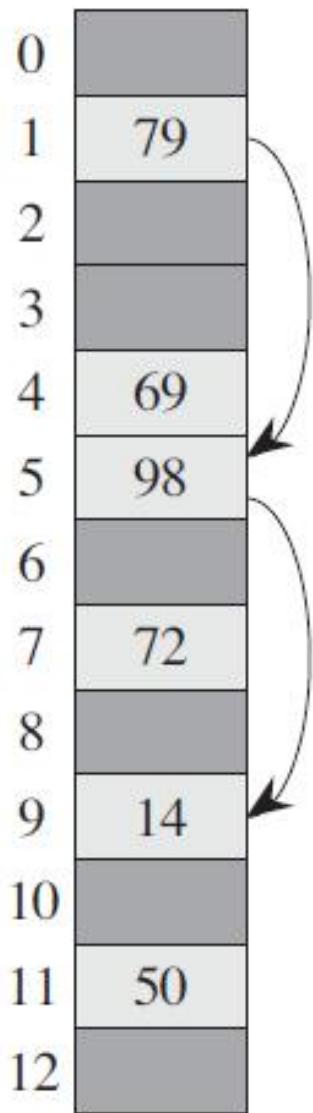
2.1.5 开放寻址法

2.1.6 Cuckoo Hash 布谷鸟哈希

2.1.7 哈希的扩展应用



2.1.5 开放寻址法





2.1.5 开放寻址法

HASH-INSERT(T, k)

```
1    $i = 0$ 
2   repeat
3        $j = h(k, i)$ 
4       if  $T[j] == \text{NIL}$ 
5            $T[j] = k$ 
6           return  $j$ 
7       else  $i = i + 1$ 
8   until  $i = m$ 
9   error “hash table overflow”
```

2.1.5 开放寻址法



HASH-SEARCH(T, k)

```
1   i = 0
2   repeat
3       j = h( $k, i$ )
4       if  $T[j] == k$ 
5           return  $j$ 
6       i =  $i + 1$ 
7   until  $T[j] == \text{NIL}$  or  $i == m$ 
8   return NIL
```

2.1.5 开放寻址法



HASH-DELETE(T, k)

$i=0$

repeat

$j = h(k, i)$

if ($T[j] == k$)

$T[j] = \text{DELETED}$

return

$i = i + 1$

until $T[j] == \text{NIL}$ or $i == m$

return

2.1.5 开放寻址法



HASH-INSERT(T, k)

$i = 0$

repeat

$j = h(k, i)$

if ($T[j] == NIL$ or $T[j] == DELETED$)

$T[j] = k$

return j

else $i = i + 1$

until $i == m$

error "hash table overflow"

2.1.5 开放寻址法



三种常用技术来计算开放寻址法中的探查序列

- 1.线性探查
- 2.二次探查
- 3.双重散列

2.1.5 开放寻址法 - 线性探查



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

关键词 (key)	47	7	29	11	9	84	54	20	30
散列地址 h(key)	3	7	7	0	9	7	10	9	8
冲突次数	0	0	1	0	0	3	1	3	6

关键词 (key)	47	7	29	11	9	84	54	20	30
散列地址 h(key)	3	7	7	0	9	7	10	9	8
冲突次数	0	0	1	0	0	3	1	3	6

操作 地址	0	1	2	3	4	5	6	7	8	9	10	11	12	说明
插入47				47										无冲突
插入7				47				7						无冲突
插入29				47				7	29					$d_1 = 1$
插入11	11			47				7	29					无冲突
插入9	11			47				7	29	9				无冲突
插入84	11			47				7	29	9	84			$d_3 = 3$
插入54	11			47				7	29	9	84	54		$d_1 = 1$
插入20	11			47				7	29	9	84	54	20	$d_3 = 3$
插入30	11	30		47				7	29	9	84	54	20	$d_6 = 6$

2.1.5 开放寻址法 - 二次探查



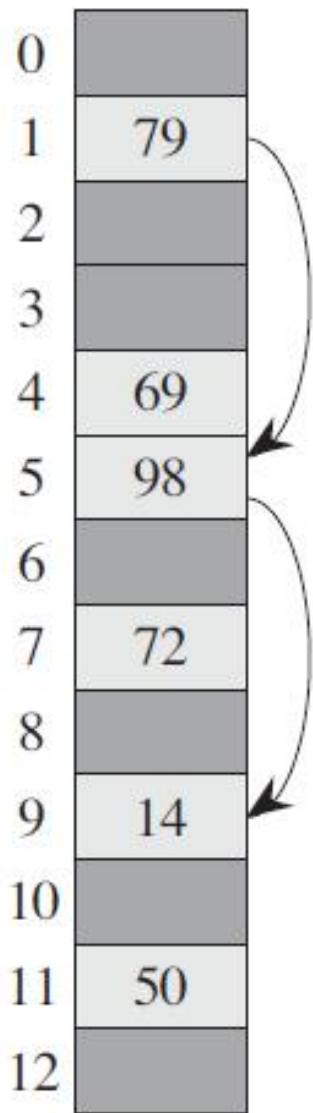
$$h(k, i) = (h'(k) + c_1 i + c_2 i^2) \bmod m$$

2.1.5 开放寻址法-双重散列



$$h(k, i) = (h1(k) + ih2(k)) \bmod m$$

2.1.5 开放寻址法





提纲

2.1.1 概述

2.1.2 直接寻址

2.1.3 散列表

2.1.4 散列函数

2.1.5 开放寻址法

2.1.6 *Cuckoo Hash* 布谷鸟哈希

2.1.7 哈希的扩展应用



定义：一种解决 *hash* 冲突的方法，其目的是使用简单的 *hash* 函数来提高 *hash table* 的利用率，同时保证 $O(1)$ 的查询时间。基本思想是使用 2 个 *hash* 函数来处理碰撞，从而每个 *key* 都对应到 2 个位置。



操作：

- 1) 对 key 值 $hash$, 生成两个 $hash$ key 值, $hashk_1$ 和 $hashk_2$, 如果对应的两个位置上有一个为空, 那么直接把 key 插入即可。
- 2) 否则, 任选一个位置, 把 key 值插入, 把已经在那个位置的 key 值踢出来。
- 3) 被踢出来的 key 值, 需要重新插入, 直到没有 key 被踢出为止。



衍生背景：

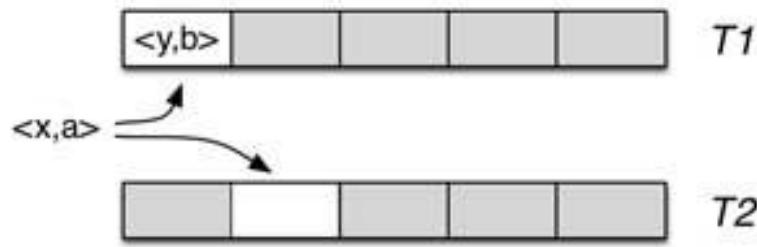
*Cuckoo*中文名叫布谷鸟，这种鸟有一种即狡猾又贪婪的习性，它不肯自己筑巢，而是把蛋下到别的鸟巢里，而且它的幼鸟又会比别的鸟早出生，布谷幼鸟天生有一种残忍的动作，幼鸟会拼命把未出生的其它鸟蛋挤出窝巢，今后以便独享“养父 母”的食物。

借助生物学上这一典故，*cuckoo hashing*处理碰撞的方法，就是把原来占用位置的这个元素踢走，不过被踢出去的元素还要比鸟蛋幸运，因为它还有一个备用位置可以安置，如果备用位置上 还有人，再把它踢走，如此往复。直到被踢的次数达到一个上限，才确认哈希表已满，并执行*rehash*操作。



Insertion when one of the two buckets is empty

Step 1: Both buckets for $\langle x, a \rangle$ are tested, the one in T2 is empty.

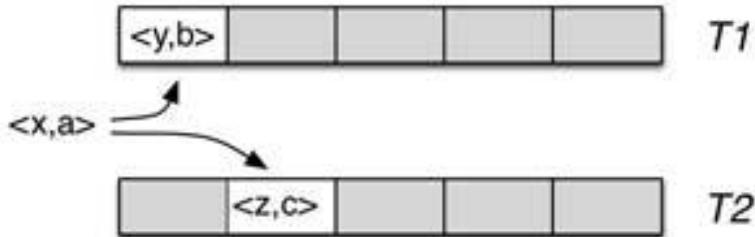


Step 2: $\langle x, a \rangle$ is stored in the empty bucket in T2.



Insertion when the two buckets already contain entries

Step 1: Here $\langle y, b \rangle$ will be withdrawn from T1 so that $\langle x, a \rangle$ can be stored.



Step 2: After $\langle x, a \rangle$ has been stored in T1, $\langle y, b \rangle$ needs to be moved to T2. The bucket in T2 may already contain an entry, if so this entry will need to be moved.





其他：

1) *Cuckoo hash* 有两种变形。一种通过增加哈希函数进一步提高空间利用率；另一种是增加哈希表，每个哈希函数对应一个哈希表，每次选择多个张表中空余位置进行放置。三个哈希表可以达到80% 的空间利用率。

2) *Cuckoo hash* 的过程可能因为反复踢出无限循环下去，这时候就需要进行一次循环踢出的限制，超过限制则认为需要添加新的哈希函数。



提纲

2.1.1 概述

2.1.2 直接寻址

2.1.3 散列表

2.1.4 散列函数

2.1.5 开放寻址法

2.1.6 Cuckoo Hash 布谷鸟哈希

2.1.7 哈希的扩展应用

2.1.7 哈希的扩展应用

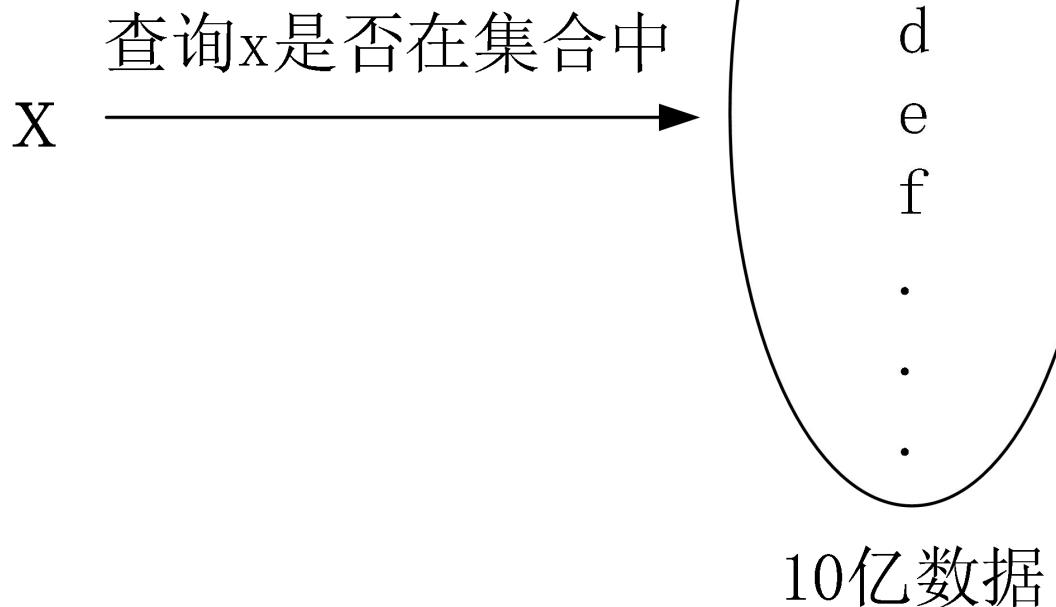


- 布鲁姆过滤器 (*BloomFilter*)
- 最小哈希
- 区块链

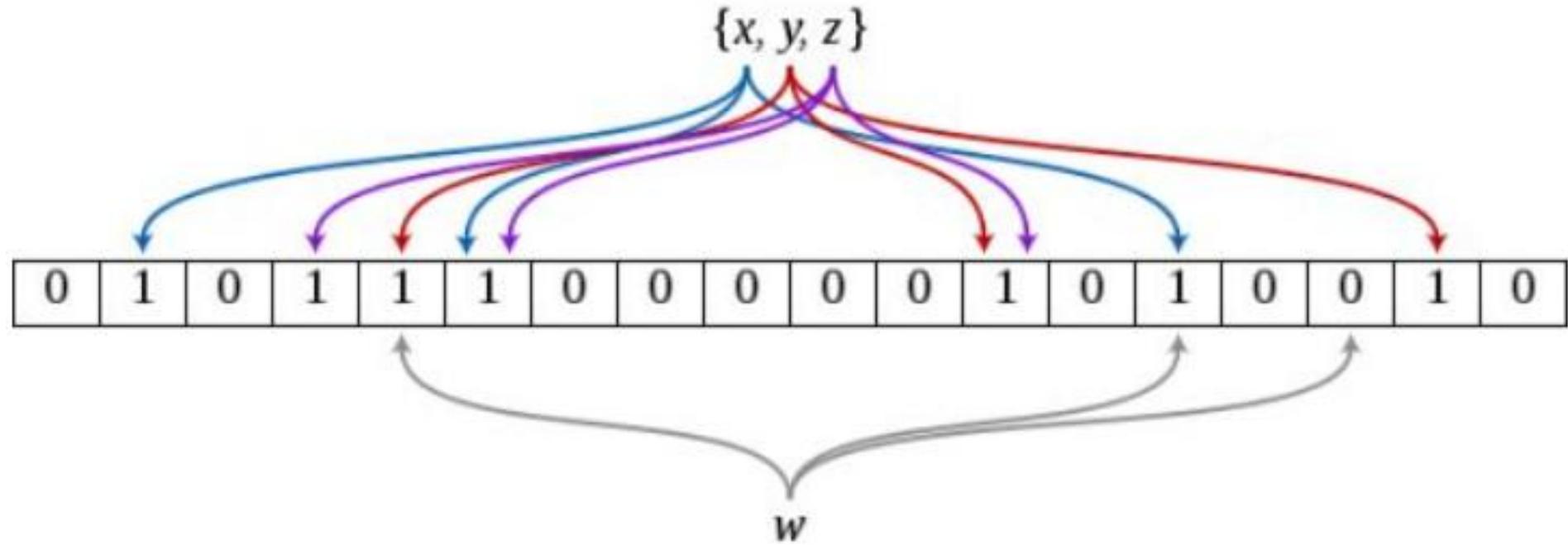
布隆过滤器 (bloomfilter)



问题1



(bloomfilter)



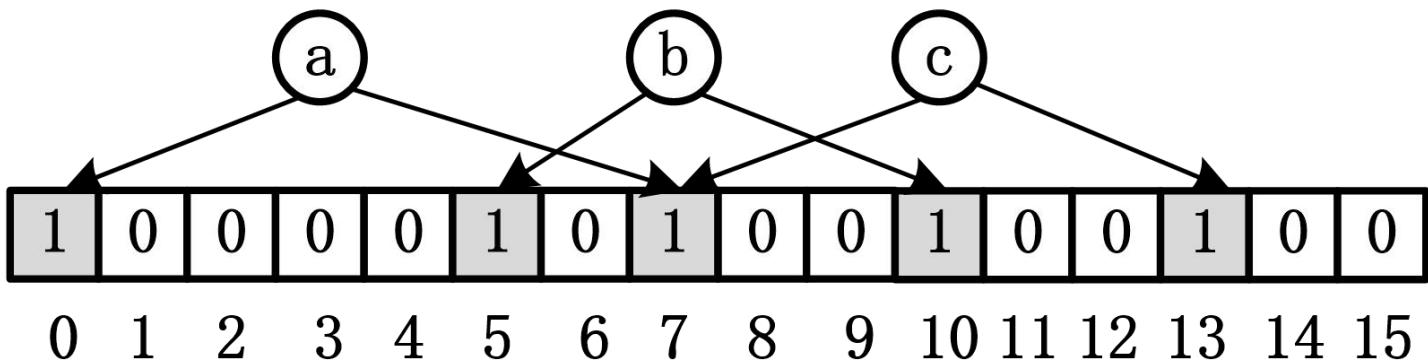


filter)

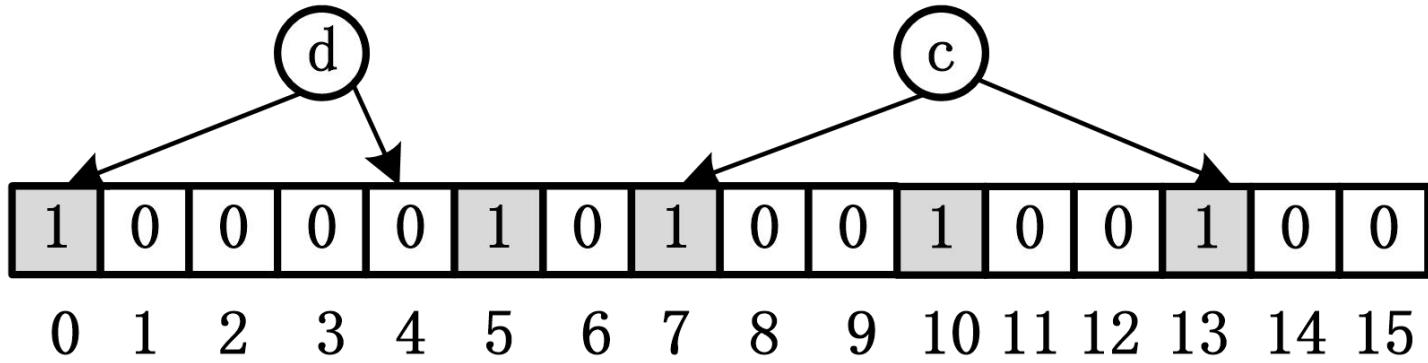
初始化

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

插入a, b, c



查询d, c





布鲁姆过滤器算法关键指标

- ▶ 空间复杂度
- ▶ 时间复杂度
- ▶ 误判率（假阳性）

将不属于集合的元素误判断成属于集合中的这种假阳性误判称为一次误判。发生误判的概率称为误判率。



误判率分析（假阳性）

(1) 理论分析，假设k个哈希函数、m位bitset，n个元素之后

$$f_{\text{BF}}(m, k, n) = (1 - e^{-k \cdot n / m})^k$$

2.1.7 布鲁姆过滤器



误判率分析（假阳性）

(3) 最优分析

$$k = \lceil \ln 2(m/n) \rceil$$

$$k_{\min} = (\ln 2) \left(\frac{m}{n} \right)$$

2.1.7 最小哈希



MinHash是一种局部敏感哈希技术
可以用来快速估算两个集合的相似度



2.1.7 Jaccard相似度

Jaccard相似度是用来计算集合相似性，也就是距离的一种度量标准

假如有集合A、B，那么 $J(A, B) = |A \cap B| / |A \cup B|$



2.1.7 最小哈希

$h(x)$: 把 x 映射成一个整数的哈希函数

$h_{\min}(S)$: 集合 S 中的元素经过 $h(x)$ 哈希后，具有最小哈希值的元素

那么对集合 A 、 B , $h_{\min}(A) = h_{\min}(B)$ 成立的条件是 $A \cup B$ 中具有最小哈希值的元素也在 $A \cap B$ 中

2.1.7 最小哈希性质



假设 $h(x)$ 是一个良好的哈希函数，它具有很好的均匀性，能够把不同元素映射成不同的整数。

所以有， $\Pr[h_{\min}(A) = h_{\min}(B)] = J(A, B)$ ，即集合A和B的相似度为集合A、B经过hash后最小哈希值相等的概率。

2.1.7 区块链



假设 $h(x)$ 是一个良好的哈希函数，它具有很好的均匀性，能够把不同元素映射成不同的整数。

所以有， $\Pr[h_{\min}(A) = h_{\min}(B)] = J(A, B)$ ，即集合A和B的相似度为集合A、B经过hash后最小哈希值相等的概率。