

CSE 326: Data Structures

Hash Tables

Brian Curless
Spring 2008

1

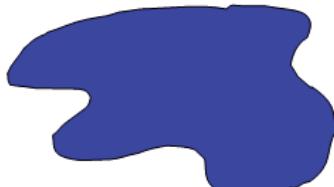
Announcements (5/7/08)

- Project 2B due today.
- Homework #4 due on Friday, beginning of class
- Project #3 assigned on Friday
 - Partner signups by 3pm Friday
- Section: project warm-up, midterms returned, ...
- Reading for this lecture: Chapter 5.

2

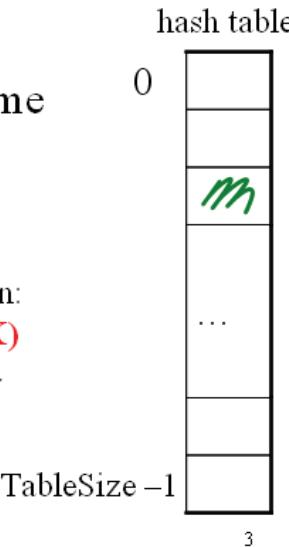
Hash Tables

- Find, insert, delete: constant time on average!
- A **hash table** is an array of some fixed size.
- General idea:



key space (e.g., integers, strings)

hash function:
 $\text{index} = h(K)$



Hash Tables

Key space of size M , but we only want to store subset of size N , where $N \ll M$.

- Keys are identifiers in programs. Compiler keeps track of them in a symbol table.
- Keys are student names. We want to look up student records quickly by name.
- Keys are chess configurations in a chess playing program.
- Keys are URLs in a database of web pages.

4

Simple Integer Hash Functions

- key space = integers
- TableSize = 10
- $h(K) = K \% 10$
- **Insert:** 7, 18, 41, 94

0	
1	41
2	
3	
4	94
5	
6	
7	7
8	18
9	

5

Simple Integer Hash Functions

- key space = integers
- TableSize = 6
- $h(K) = K \% 6$
- **Insert:** 7, 18, 41, 34

0	18
1	7
2	
3	
4	34
5	41

6

Aside: Properties of Mod

To keep hashed values within the size of the table,
we will generally do:

$$h(K) = \text{function}(K) \% \text{TableSize}$$

(In the previous examples, $\text{function}(K) = K$.)

$$\begin{array}{r} 24 \rightarrow 4 \\ +57 \rightarrow 7 \\ \hline 1 \end{array}$$

It's worth noting a couple properties of the mod
function:

- $(a + b) \% c = [(a \% c) + (b \% c)] \% c$
- $a \% c = b \% c \rightarrow (a - b) \% c = 0$
- $(a \cdot b) \% c = [(a \% c) \cdot (b \% c)] \% c$

$$\begin{array}{r} 24 \quad 4 \\ 57 \quad 7 \\ \hline 8 \end{array}$$

7

Some String Hash Functions

$$s_i \in [0 \dots 127]$$

key space = strings

$$K = s_0 s_1 s_2 \dots s_{m-1} \text{ (where } s_i \text{ are characters)}$$

$$1. \quad h(K) = s_0 \% \text{TableSize}$$

$$2. \quad h(K) = \left(\sum_{i=0}^{m-1} s_i \right) \% \text{TableSize}$$

$$3. \quad h(K) = \left(\sum_{i=0}^{m-1} s_i \cdot 128^i \right) \% \text{TableSize}$$

$$= s_0 + s_1 \cdot 128 + s_2 \cdot 128^2 + \dots$$

letter aliasing
anagrams spot, post, stop

8

Hash Function Desiderata

What are some desirable properties for a hash function?

EVENLY distributes values throughout the table.

→ MINIMIZES collisions

FAST to compute

9

A Fancier Hash Function

Some experimental results indicate that modular hash functions with prime tables sizes are not ideal.

Instead, we can work on designing a really good hash function:

```
jenkinsOneAtATimeHash(String key, int keyLength) {  
    hash = 0;  
    for (i = 0; i < key_len; i++) {  
        hash += key[i];  
        hash += (hash << 10);  
        hash ^= (hash >> 6);  
    }  
    hash += (hash << 3);  
    hash ^= (hash >> 11);  
    hash += (hash << 15);  
  
    return hash % TableSize;  
}
```

11

Designing Hash Functions

We've seen a few possibilities. The simplest is **modular hashing**:

$$h(K) = K \% P$$

where P is usually just the TableSize.

P is often chosen to be prime:

- Reduces likelihood of collisions due to patterns in data
- Is useful for guarantees on certain hashing strategies (as we'll see)

But what would be a more convenient value of P? 2^k

10

Collision Resolution

Collision: when two keys map to the same location in the hash table.

How can we cope with collisions?

open addressing / Look at next spot in table, if empty → use
if full, → go to next

Separate chaining / List at each table entry
Overflow table?

12

Separate Chaining

0	10
1	
2	22
3	
4	
5	
6	
7	107
8	
9	

$$\text{TableSize} = 10$$

$$h(K) = K \% 10$$

Separate chaining:

All keys that map to the same hash value are kept in a list (or “bucket”).

Insert:
10
22
107
12
42

13

Analysis of Separate Chaining

The **load factor**, λ , of a hash table is

$$\lambda = \frac{N}{\text{TableSize}} \leftarrow \text{no. of elements}$$

Separate chaining: $\lambda = \text{average } \# \text{ of elems per bucket}$

Average cost of:

- Unsuccessful find? λ
- Successful find? $1 + \frac{\lambda}{2}$
- Insert? 1

14

Alternative: Use Empty Space in the Table

0	8
1	109
2	10
3	
4	
5	
6	
7	
8	38
9	19

$$h(k) = k \% 10$$

Insert:
38
19
8
109
10

Try $h(K)$.

If full, try $h(K)+1$.

If full, try $h(K)+2$.

If full, try $h(K)+3$.

Etc...

15

Open Addressing

The approach on the previous slide is an example of **open addressing**:

After a collision, try “next” spot. If there’s another collision, try another, etc.

Finding the next available spot is called **probing**:

0th probe = $h(k) \% \text{TableSize}$

1st probe = $(h(k) + f(1)) \% \text{TableSize}$

2nd probe = $(h(k) + f(2)) \% \text{TableSize}$

...

ith probe = $(h(k) + f(i)) \% \text{TableSize}$

f(i) is the probing function. We’ll look at a few...

16

Terminology Alert!

- **Separate chaining** is sometimes called **open hashing**.
- **Open addressing** is sometimes called **closed hashing**.

17

Open Addressing Example, Revisited

0	8
1	109
2	10
3	
4	
5	
6	
7	
8	38
9	19

Insert:
38
19
8
109
10

Try $h(K)$
If full, try $h(K)+1$.
If full, try $h(K)+2$.
If full, try $h(K)+3$.
Etc...

What is $f(i)$? $f(i) = i$

18

Linear Probing

$$f(i) = i$$

- Probe sequence:

0th probe = $h(K) \% \text{TableSize}$

1th probe = $(h(K) + 1) \% \text{TableSize}$

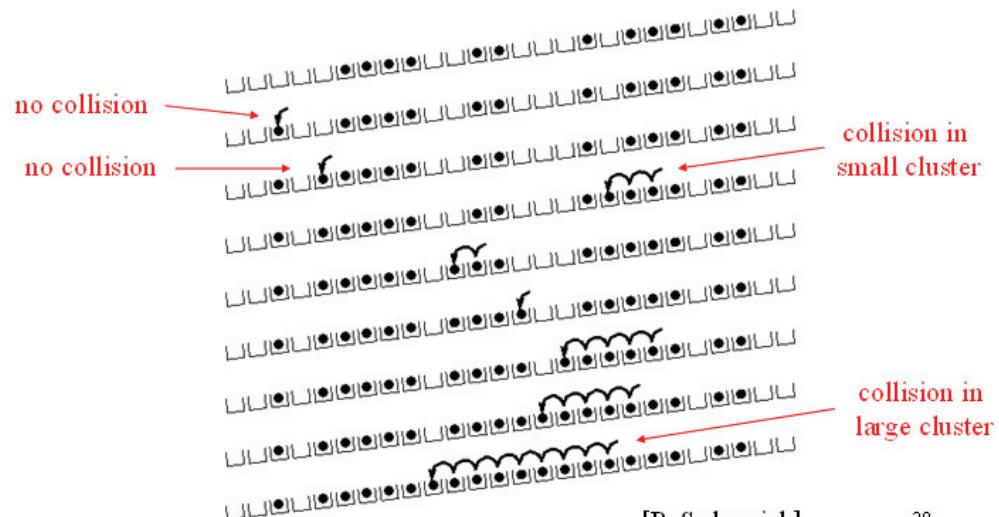
2th probe = $(h(K) + 2) \% \text{TableSize}$

...

ith probe = $(h(K) + i) \% \text{TableSize}$

19

Linear Probing – Clustering



[R. Sedgewick]

20

Analysis of Linear Probing

- For any $\lambda < 1$, linear probing *will* find an empty slot
- Expected # of probes (for large table sizes)
 - unsuccessful search: $\frac{1}{2} \left(1 + \frac{1}{(1-\lambda)^2} \right)$
 - successful search: $\frac{1}{2} \left(1 + \frac{1}{(1-\lambda)} \right)$
- Linear probing suffers from *primary clustering*
- Performance quickly degrades for $\lambda > 1/2$

21

Less likely
to encounter
Primary
Clustering

Quadratic Probing

- $$f(i) = i^2$$
- Probe sequence:

$$\begin{aligned} 0^{\text{th}} \text{ probe} &= h(K) \% \text{TableSize} \\ 1^{\text{th}} \text{ probe} &= (h(K) + 1) \% \text{TableSize} \\ 2^{\text{th}} \text{ probe} &= (h(K) + 4) \% \text{TableSize} \\ 3^{\text{th}} \text{ probe} &= (h(K) + 9) \% \text{TableSize} \\ \dots \\ i^{\text{th}} \text{ probe} &= (h(K) + i^2) \% \text{TableSize} \end{aligned}$$

22

Quadratic Probing Example

0	49
1	
2	58
3	79
4	
5	
6	
7	
8	18
9	89

Insert:
89
18
49
58
58
79

23

Another Quadratic Probing Example

0	48
1	
2	5
3	55
4	
5	40
6	76

TableSize = 7
 $h(K) = K \% 7$

$$\begin{aligned} \text{insert}(76) \quad 76 \% 7 &= 6 \\ \text{insert}(40) \quad 40 \% 7 &= 5 \\ \text{insert}(48) \quad 48 \% 7 &= 6 \\ \text{insert}(5) \quad 5 \% 7 &= 5 \\ \text{insert}(55) \quad 55 \% 7 &= 6 \\ \text{insert}(47) \quad 47 \% 7 &= 5 \end{aligned}$$

24

$$5 + \underbrace{0, 1, 4, 2, 2,}_{4, 1} 0, 1, 4, 9, 16, 25, 36$$

$$5, 6, 9, 7$$

$$5, 6, 2, 0$$

Quadratic Probing: Success guarantee for $\lambda < \frac{1}{2}$

Assertion #1: If $T = \text{TableSize}$ is **prime** and $\lambda < \frac{1}{2}$, then quadratic probing will find an empty slot in $T/2$ probes or fewer.

Assertion #2: If the following holds

$$\rightarrow (h(K) + i^2) \% T \neq (h(K) + j^2) \% T$$

For prime T and all $0 \leq i, j \leq T/2$ and $i \neq j$,

then assertion #1 is true.

Assertion #3: If assertion #2 true, then so is assertion #1.

25

Quadratic Probing: Properties

- For *any* $\lambda < \frac{1}{2}$, quadratic probing will find an empty slot; for bigger λ , quadratic probing *may* find a slot.
- Quadratic probing does not suffer from *primary clustering*: keys hashing to the same *area* are not bad.
- But what about keys that hash to the same *spot*?
Secondary Clustering!

27

Quadratic Probing: Success guarantee for $\lambda < \frac{1}{2}$

We can prove assertion #2 by contradiction.

Suppose that for some $i \neq j$, $0 \leq i, j \leq T/2$, prime T :

$$(h(K) + i^2) \% T = (h(K) + j^2) \% T$$

$$(h(K) + i^2 - h(K) - j^2) \% T = 0$$

$$(i^2 - j^2) \% T = 0$$

$$[(i+j)(i-j)] \% T = 0$$

$$i+j = 0 \% N$$

$$i+j = T \% N$$

$$i-j = 0 \% N$$

$$i-j = T \% N$$

26

Double Hashing

Idea: given two different (good) hash functions $h(K)$ and $g(K)$, it is unlikely for two keys to collide with both of them.

So... let's try probing with a second hash function:

$$f(i) = i * g(K)$$

- Probe sequence:

$$0^{\text{th}} \text{ probe} = h(K) \% \text{TableSize}$$

$$1^{\text{th}} \text{ probe} = (h(K) + g(K)) \% \text{TableSize}$$

$$2^{\text{th}} \text{ probe} = (h(K) + 2*g(K)) \% \text{TableSize}$$

$$3^{\text{th}} \text{ probe} = (h(K) + 3*g(K)) \% \text{TableSize}$$

...

$$i^{\text{th}} \text{ probe} = (h(K) + i*g(K)) \% \text{TableSize}$$

28

Double Hashing Example

0
1
2
3
4
5
6

TableSize = 7
 $h(K) = K \% 7$
 $g(K) = 5 - (K \% 5)$

Insert(76) $76 \% 7 = 6$ and $5 - 76 \% 5 = 3$
 Insert(93) $93 \% 7 = 2$ and $5 - 93 \% 5 = 1$
 Insert(40) $40 \% 7 = 5$ and $5 - 40 \% 5 = 0$
 Insert(47) $47 \% 7 = 5$ and $5 - 47 \% 5 = 3$
 Insert(10) $10 \% 7 = 3$ and $5 - 10 \% 5 = 2$
 Insert(55) $55 \% 7 = 6$ and $5 - 55 \% 5 = 5$

29

Another Example of Double Hashing

0
1
2
3
4
5
6
7
8
9

Hash Functions:
 $T = \text{TableSize} = 10$
 $h(K) = K \% T$
 $g(K) = 1 + (K/T) \% (T-1)$

Insert these values into the hash table in this order. Resolve any collisions with double hashing:

13
 28
 33 $h=3, g=4$
 147 $h=7, g=6$ $\Sigma = 13 + 3 + 28 + 147 = 196$
 43 $h=3, g=5$

30

Analysis of Double Hashing

- Double hashing is safe for $\lambda < 1$ for at least one case:
 - $h(k) = k \% p$
 - $g(k) = q - (k \% q)$
 - $2 < q < p$, and p, q are primes
- Expected # of probes (for large table sizes)
 - unsuccessful search: $\frac{1}{1-\lambda}$
 - successful search: $\frac{1}{\lambda} \log_e \left(\frac{1}{1-\lambda} \right)$

Never let $g(k)=0$

31

Deletion in Separate Chaining

How do we delete an element with separate chaining?

find
 remove from list

32

Rehashing

Deletion in Open Addressing

Can we do something similar for open addressing?

- Delete → mark "occupied"
- Find
- Insert

	$h(k) = k \% 7$
0	
1	
2	16
3	23 → occupied → 30
4	59
5	
6	76

33

Idea: When the table gets too full, create a bigger table (usually 2x as large) and hash all the items from the original table into the new table.

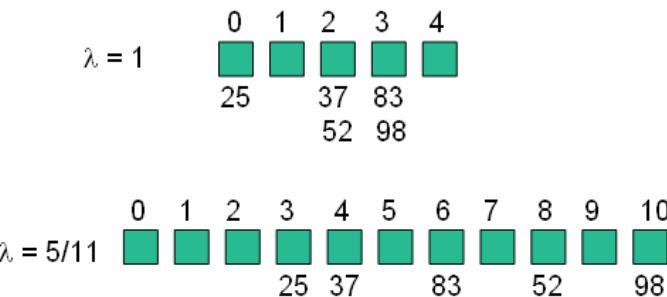
- When to rehash?
 - Separate chaining: full ($\lambda = 1$)
 - Open addressing: half full ($\lambda = 0.5$)
 - When an insertion fails
 - Some other threshold
- Cost of a single rehashing?

34

Rehashing Example

- Separate chaining example:

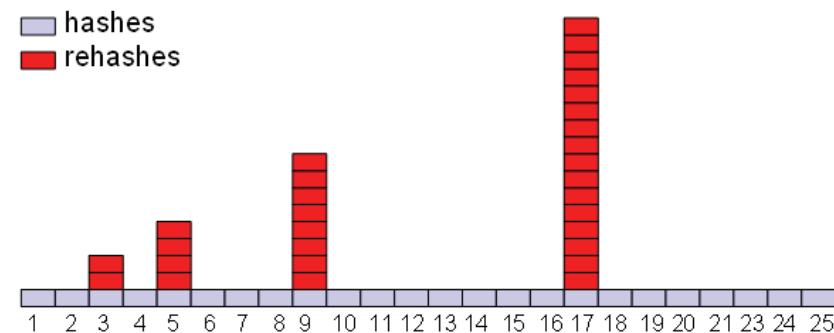
$$h_1(x) = x \% 5 \text{ rehashes to } h_2(x) = x \% 11.$$



35

Rehashing Picture

- Starting with table of size 2, double when load factor > 1.



36

Amortized Analysis of Rehashing

- Cost of inserting n keys is $< 3n$
- $2^k + 1 \leq n \leq 2^{k+1}$
 - Hashes = n
 - Rehashes = $2 + 2^2 + \dots + 2^k = 2^{k+1} - 2$
 - Total = $n + 2^{k+1} - 2 < 3n$
- Example
 - $n = 33$, Total = $33 + 64 - 2 = 95 < 99$

37

Hashing Summary

- Hashing is one of the most important data structures.
- Hashing has many applications where operations are limited to find, insert, and delete.
 - But what is the cost of doing, e.g., findMin?
- Can use:
 - Separate chaining (easiest)
 - Open hashing (memory conservation, no linked list management)
 - Java uses separate chaining
- Rehashing has good amortized complexity.
- Also has a big data version to minimize disk accesses: extendible hashing. (See textbook.)

38