

# Hashing

CSE 326  
Data Structures  
Lecture 14

## Readings and References

- Reading
  - Chapter 5

2/28/05

Lecture 14 - Hashing

2

## Hashing

- Hashing is a family of data structures used to efficiently support insert, delete, find.
- It cannot be used efficiently for other operations where the order of data is important. No list-all, range queries, successor, predecessor.

2/28/05

Lecture 14 - Hashing

3

## General Idea

- 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.

2/28/05

Lecture 14 - Hashing

4

## Simple Hash Table

T	
0	
1	
2	John Smith
3	
4	Judy Jones
5	
6	Martha Lee
7	Jerry Lee
8	
9	

Hash function:

$h : U \rightarrow \{0, 1, \dots, \text{HSize} - 1\}$

U is the universe of keys

$h(\text{"name"})$  is the hash value of "name"

$h(\text{Judy Jones}) = 4$

$h(\text{Jerry Lee}) = 7$

$\text{Find}(\text{"name"}) = T[h(\text{"name"})]$

2/28/05

Lecture 14 - Hashing

5

## Hashing Properties

- Load Factor  $= \lambda = \frac{N}{\text{HSize}}$ 
  - Hash tables may have unused entries  $\lambda < 1$
- Good quality hash function distribute data as evenly as possible over the keys.
- Collisions:  $h(\text{inserted key}) = h(\text{existing key})$ .
  - Open hashing - linked lists
  - Closed hashing - find a new place to put inserted key

2/28/05

Lecture 14 - Hashing

6

## Good Hash Functions

- Integers: Division method
  - Choose Hsize to be a prime
  - $h(n) = n \bmod \text{Hsize}$
  - Example. Hsize = 23,  $h(50) = 4$ ,  $h(1257) = 15$
- Character Strings
  - $x = a_0a_1a_2\dots a_m$  is a character string. Define  $\text{int}(x) = a_0 + a_1 128 + a_2 128^2 + \dots + a_m 128^{m-1}$
  - $h(x) = \text{int}(x) \bmod \text{Hsize}$
  - Compute  $h(x)$  using Horner's Rule
    - $h := 0$
    - for  $i = m$  to 0 by -1 do  $h := (a_i + 128h) \bmod \text{Hsize}$
    - return  $h$

2/28/05

Lecture 14 - Hashing

7

## A Bad Hash Function

- Keys able1, able2, able3, able4
  - Hsize = 128
  - $\text{int}(\text{able}x) \bmod 128 = \text{int}(a) = 97$
  - Thus,  $h(\text{able}x) = h(\text{able}y)$  for all  $x$  and  $y$
- Why use primes for hash table sizes?
  - Primes have no nontrivial divisors
  - Numbers relatively prime to 128 will also work for character strings

2/28/05

Lecture 14 - Hashing

8

## Multiplication Method

- Hash function defined by HSize and a floating point number A.
  - Integer case
  - $h(k) = \lfloor \text{HSize} * (k * A \bmod 1) \rfloor$
  - Example: HSize = 10, A = .485
  - $h(50) = \lfloor 10 * (50 * .485 \bmod 1) \rfloor$
  - $= \lfloor 10 * (24.25 \bmod 1) \rfloor$
  - $= \lfloor 10 * .25 \rfloor$
  - $= 2$
  - + HSize need not be prime
  - More computation than division method
- Another alternative – Universal Hashing

2/28/05

Lecture 14 - Hashing

9

## What about Collisions?

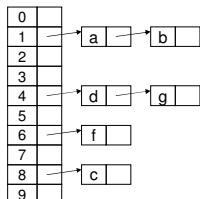
- Open Hashing - Collisions overflow into linked lists.
  - Load factors  $> 1$  are possible
- Closed Hashing - if a collision occurs find another place in the hash table for the entry.
  - Load factor must be  $\leq 1$

2/28/05

Lecture 14 - Hashing

10

## Open Hashing (Chaining)



- $h(a) = h(b)$  and  $h(d) = h(g)$
- Chains may be ordered or unordered. Little advantage to ordering.

2/28/05

Lecture 14 - Hashing

11

## Open Hashing Properties

- Load factor  $= \lambda$ 
  - Unsuccessful searches cost  $\lambda$  comparisons on average
  - Successful searches cost  $1 + \lambda/2$  comparisons on average
- Comparisons can be expensive so choosing  $\lambda$  between  $1/2$  and  $1$  is wise.

2/28/05

Lecture 14 - Hashing

12

## Closed Hashing (Open Addressing)

- No chaining, every key fits in the hash table.
- Probe sequence
  - $h(k)$
  - $(h(k) + f(1)) \bmod HSize$
  - $(h(k) + f(2)) \bmod HSize, \dots$
- Insertion: Find the first probe with an empty slot.
- Find: Find the first probe that equals the query or is empty. Stop at HSize probe, in any case.
- Deletion: lazy deletion is needed. That is, mark locations as deleted, if a deleted key resides there.

2/28/05

Lecture 14 - Hashing

13

## Linear Probing

- $f(i) = i$
- Probe sequence
  - $h(k)$
  - $(h(k) + 1) \bmod HSize$
  - $(h(k) + 2) \bmod HSize, \dots$
- Insertion (assuming  $\lambda < 1$ )
 

```
h := h(k)
while T(h) not empty do
    h := (h + 1) mod HSize;
insert k in T(h)
```

2/28/05

Lecture 14 - Hashing

14

## Linear Probing Example

	76	93	40	47	10	55
0				47	47	47
1						55
2		93	93	93	93	93
3					10	10
4						
5			40	40	40	40
6	76	76	76	76	76	76
Probes	1	1	1	3	1	3

2/28/05

Lecture 14 - Hashing

15

## Performance of Linear Probing

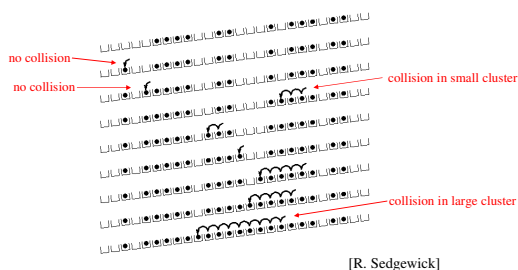
- If there is an available slot linear probing will find it.
- For large hash tables the expected number of probes on insertion is:
 
$$\frac{1}{2} \left( 1 + \frac{1}{(1-\lambda)^2} \right)$$
- The expected number of probes on successful searches is:
 
$$\frac{1}{2} \left( 1 + \frac{1}{1-\lambda} \right)$$
- Linear probing suffers from primary clustering.
- Not a good idea to use linear probing with  $\lambda > 1/2$ .
- Lazy deletion needed.

2/28/05

Lecture 14 - Hashing

16

## Linear Probing – Clustering



2/28/05

Lecture 14 - Hashing

17

## Quadratic Probing

- $f(i) = i^2$
- Probe sequence
  - $h(k)$
  - $(h(k) + 1) \bmod HSize$
  - $(h(k) + 4) \bmod HSize$
  - $(h(k) + 9) \bmod HSize, \dots$
- Insertion (assuming  $\lambda < 1/2$ )
 

```
h := h(k);
i := 0;
while T(h) not empty do {
    h := (h + 2*i + 1) mod HSize;
    i := i + 1;
}
insert k in T(h)
```

Note:  $(i+1)^2 - i^2 = 2i + 1$

2/28/05

Lecture 14 - Hashing

18

## Quadratic Probing Works for $\lambda < 1/2$

- If HSize is prime then  $(h(x) + i^2) \bmod \text{HSize} \neq (h(x) + j^2) \bmod \text{HSize}$  for  $i \neq j$  and  $0 \leq i, j < \text{HSize}/2$ .
- Proof**

$$(h(x) + i^2) \bmod \text{HSize} = (h(x) + j^2) \bmod \text{HSize}$$

$$(h(x) + i^2) - (h(x) + j^2) \bmod \text{HSize} = 0$$

$$(i^2 - j^2) \bmod \text{HSize} = 0$$

$$(i-j)(i+j) \bmod \text{HSize} = 0$$

$$\Rightarrow \Leftarrow \text{HSize does not divide } (i-j) \text{ or } (i+j)$$

2/28/05

Lecture 14 - Hashing

19

## Quadratic Probing may Fail if $\lambda \geq 1/2$

51	
0	
1	
2	16
3	45
4	59
5	
6	76

$51 \bmod 7 = 2; i = 0$   
 $(2 + 1) \bmod 7 = 3; i = 1$   
 $(3 + 3) \bmod 7 = 6; i = 2$   
 $(6 + 5) \bmod 7 = 4; i = 3$   
 $(4 + 7) \bmod 7 = 4; i = 4$   
 $(4 + 9) \bmod 7 = 6; i = 5$   
 $(6 + 11) \bmod 7 = 3; i = 6$   
 $(3 + 13) \bmod 7 = 2; i = 7$   
 ...

2/28/05

Lecture 14 - Hashing

20

## Performance of Quadratic Probing

- Although quadratic probing can fail for  $\lambda \geq 1/2$ , it is not likely to do so. We can use load factors greater than  $1/2$ , but load factors close to 1 should be avoided.
- Quadratic hashing does not suffer from primary clustering, but has only minor secondary clustering.
- With load factors near  $1/2$  the expected number of probes per successful search is about 1.5.
- Lazy deletion must be used.

2/28/05

Lecture 14 - Hashing

21

## Double Hashing

- $f(i) = i \cdot g(k)$  where  $g$  is a second hash function
- Probe sequence
  - $h(k)$
  - $(h(k) + g(k)) \bmod \text{HSize}$
  - $(h(k) + 2g(k)) \bmod \text{HSize}$
  - $(h(k) + 3g(k)) \bmod \text{HSize}, \dots$
- In choosing  $g$  care must be taken so that it never evaluates to 0.
- A good choice for  $g$  is to choose a prime  $R < \text{HSize}$  and let  $g(k) = R - (k \bmod R)$ .

2/28/05

Lecture 14 - Hashing

22

## Double Hashing Example

$h(k) = k \bmod 7$  and  $g(k) = 5 - (k \bmod 5)$

76	93	40	47	10	55
0	0	0	0	0	0
1	1	1	1	1	1
2	2	2	2	2	2
3	3	3	3	3	3
4	4	4	4	4	4
5	5	5	5	5	5
6	6	6	6	6	6
76	93	40	47	10	55

Probes 1

2/28/05

Lecture 14 - Hashing

23

## Double Hashing is Safe for $\lambda < 1$

- Let  $h(k) = k \bmod p$  and  $g(k) = q - (k \bmod q)$  where  $2 < q < p$  and  $p$  and  $q$  are primes. The probe sequence  $h(k) + ig(k) \bmod p$  probes every entry of the hash table.
  - Let  $0 \leq m < p$ ,  $h = h(k)$ , and  $g = g(k)$ . We show that  $h + ig \bmod p = m$  for some  $i$ .  $0 < g < p$ , so  $g$  and  $p$  are relatively prime. By extended Euclid's algorithm there are  $s$  and  $t$  such that  $sg + tp = 1$ . Choose  $i = (m-h)s \bmod p$
  - $(h + ig) \bmod p =$
  - $(h + (m-h)sg) \bmod p =$
  - $(h + (m-h)sg + (m-h)tp) \bmod p =$
  - $(h + (m-h)(sg + tp)) \bmod p =$
  - $(h + (m-h)) \bmod p = m \bmod p = m$

2/28/05

Lecture 14 - Hashing

24

## Deletion in Hashing

- Open hashing (chaining) – no problem
- Closed hashing – must do lazy deletion. Deleted keys are marked as deleted.
  - Find: done normally
  - Insert: treat marked slot as an empty slot and fill it.

$h(k) = k \bmod 7$   
Linear probing

Find 59

0		0		Insert 30
1		1		
2	16	2	16	
3	23	3	30	
4	59	4	59	
5		5		
6	76	6	76	

2/28/05

Lecture 14 - Hashing

25

## Rehashing

- Build a bigger hash table of approximately twice the size when  $\lambda$  exceeds a particular value
  - Go through old hash table, ignoring items marked deleted
  - Recompute hash value for each non-deleted key and put the item in new position in new table
  - Cannot just copy data from old table because the bigger table has a new hash function
- Running time is  $O(N)$  but happens very infrequently
  - Not good for real-time safety critical applications

2/28/05

Lecture 14 - Hashing

26

## Rehashing Example

- Open hashing –  $h_1(x) = x \bmod 5$  rehashes to  $h_2(x) = x \bmod 11$ .

$\lambda = 1$

0	1	2	3	4
25		37	83	
		52	98	

$\lambda = 5/11$

0	1	2	3	4	5	6	7	8	9	10
			25	37		83		52		98

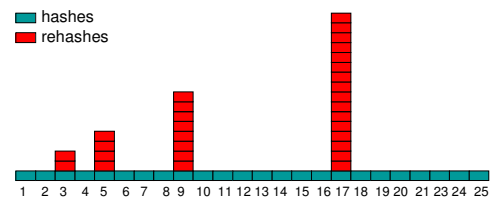
2/28/05

Lecture 14 - Hashing

27

## Rehashing Picture

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



2/28/05

Lecture 14 - Hashing

28

## 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 =  $33 + 64 - 2 = 95 < 99$

2/28/05

Lecture 14 - Hashing

29

## Case Study

- Spelling Dictionary - 30,000 words
- Goals
  - Fast spell checking
  - Minimal storage
- Possible solutions
  - Sorted array and binary search
  - Open hashing (chaining)
  - Closed hashing with linear probing
- Notes
  - Almost all searches are successful
  - 30,000 word average 8 bytes per word, 240,000 bytes
  - Pointers are 4 bytes

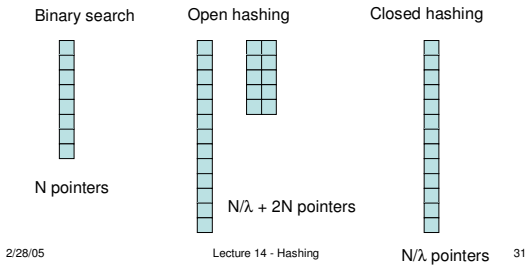
2/28/05

Lecture 14 - Hashing

30

## Storage

- Assume words are stored as strings and entries in the arrays are pointers to the strings.



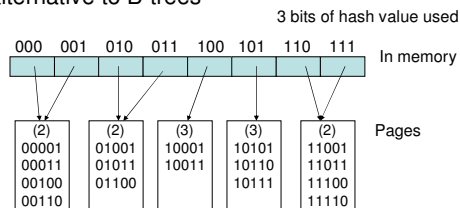
## Analysis

- Binary Search**
  - Storage =  $N$  pointers + words = 360,000 bytes
  - Time =  $\log_2 N \leq 15$  probes in worst case
- Open hashing**
  - Storage =  $2N + N/\lambda$  pointers + words
  - $\lambda = 1$  implies 600,000 bytes
  - Time =  $1 + \lambda/2$  probes per access
  - $\lambda = 1$  implies 1.5 probes per access
- Closed hashing**
  - Storage =  $N/\lambda$  pointers + words
  - $\lambda = 1/2$  implies 480,000 bytes
  - Time =  $(1/2)(1 + 1/(1-\lambda))$  probes
  - $\lambda = 1/2$  implies 1.5 probes per access

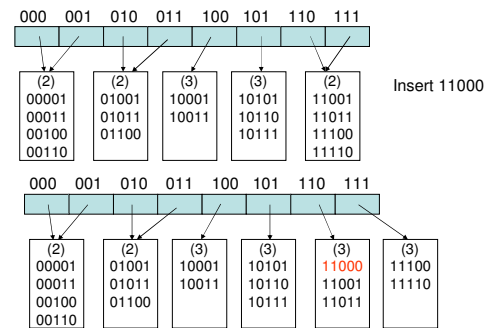
2/28/05 Lecture 14 - Hashing 32

## Extendible Hashing

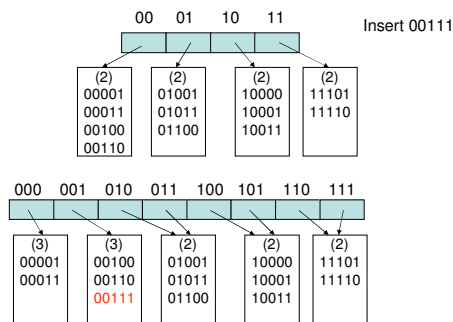
- Extendible hashing is a technique for storing large data sets that do not fit in memory.
- An alternative to B-trees



## Splitting



## Rehashing



## Analysis of Extendible Hashing

- On deletion neighbors can be merged.
- If table uses  $k$  bits but all pages use  $k-1$  bits then rehashing to a smaller table can be done. Not normally an issue with large databases.
- Rehashing does not touch pages.
- Splitting and merging touch only two pages.

2/28/05 Lecture 14 - Hashing 36

## Hashing Summary

- Hashing is one of the most important data structures.
- Hashing has many applications where operations are limited to find, insert, and delete.
- Dynamic hash tables have good amortized complexity.
- Extendible hashing is useful in databases.