

CSE 326: Data Structures Lecture #12 Hashing II

Henry Kautz
Winter 2002

Load Factor in Linear Probing

- For any $\lambda < 1$, linear probing will find an empty slot
- Search cost (for large table sizes)

- successful search: $\frac{1}{2} \left(1 + \frac{1}{(1-\lambda)} \right)$

- unsuccessful search: $\frac{1}{2} \left(1 + \frac{1}{(1-\lambda)^2} \right)$

- Performance quickly degrades for $\lambda > 1/2$

Linear Probing – Expected # of Probes

Load factor	failure	success
.1	1.11	1.06
.2	1.28	1.13
.3	1.52	1.21
.4	1.89	1.33
.5	2.5	1.50
.6	3.6	1.75
.7	6.0	2.17
.8	13.0	3.0
.9	50.5	5.5

Open Addressing II: Quadratic Probing

- Main Idea: Spread out the search for an empty slot – Increment by i^2 instead of i

- $h_i(X) = (\text{Hash}(X) + i^2) \% \text{TableSize}$

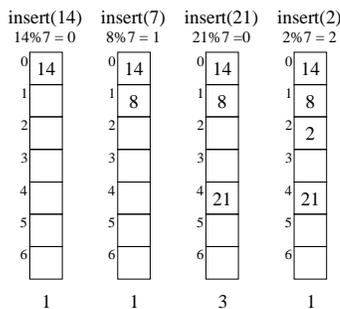
$h_0(X) = \text{Hash}(X) \% \text{TableSize}$

$h_1(X) = \text{Hash}(X) + 1 \% \text{TableSize}$

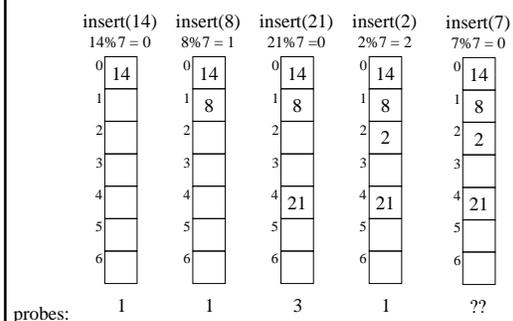
$h_2(X) = \text{Hash}(X) + 4 \% \text{TableSize}$

$h_3(X) = \text{Hash}(X) + 9 \% \text{TableSize}$

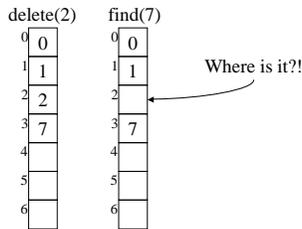
Quadratic Probing Example



Problem With Quadratic Probing

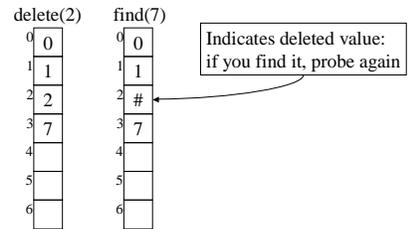


Deletion in Open Addressing



What should we do instead?

Lazy Deletion



But *now* what is the problem?



The Squished Pigeon Principle

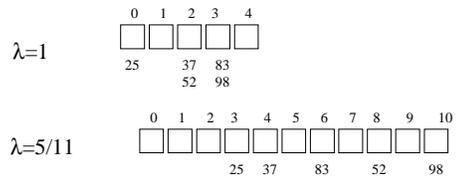
- An insert using open addressing *cannot* work with a load factor of 1 or more.
 - Quadratic probing can *fail* if $\lambda > \frac{1}{2}$
 - Linear probing and double hashing *slow* if $\lambda > \frac{1}{2}$
 - Lazy deletion never frees space
- Separate chaining becomes slow once $\lambda > 1$
 - Eventually becomes a linear search of long chains
- How can we relieve the pressure on the pigeons?

REHASH!

Rehashing Example

Separate chaining

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



Stretchy Stack Amortized Analysis

- Consider sequence of n operations
push(3); push(19); push(2); ...
- What is the max number of stretches? $\log n$
- What is the total time?
 - let's say a regular push takes time a , and stretching an array contain k elements takes time bk .

$$an + b(1 + 2 + 4 + 8 + \dots + n) = an + b \sum_{i=0}^{\log n} 2^i$$

$$= an + b(2n - 1)$$

- Amortized time = $(an + b(2n - 1))/n = O(1)$



Rehashing Amortized Analysis

- Consider sequence of n operations
insert(3); insert(19); insert(2); ...
- What is the max number of rehashes? $\log n$
- What is the total time?
 - let's say a regular hash takes time a , and rehashing an array contain k elements takes time bk .

$$an + b(1 + 2 + 4 + 8 + \dots + n) = an + b \sum_{i=0}^{\log n} 2^i$$

$$= an + b(2n - 1)$$

- Amortized time = $(an + b(2n - 1))/n = O(1)$



Rehashing without Stretching

- Suppose input is a mix of inserts and deletes
 - Never more than TableSize/2 active keys
 - Rehash when $\lambda=1$ (half the table must be deletions)
- Worst-case sequence:
 - T/2 inserts, T/2 deletes, T/2 inserts, Rehash, T/2 deletes, T/2 inserts, Rehash, ...
- Rehashing at most doubles the amount of work – still $O(1)$

Case Study

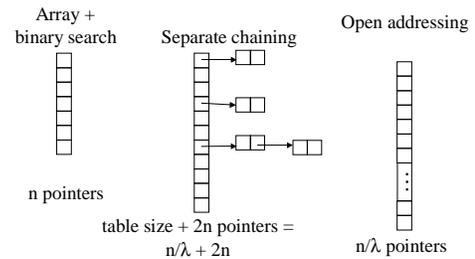
- Spelling dictionary
 - 30,000 words
 - static
 - arbitrary(ish) preprocessing time
 - Goals
 - fast spell checking
 - minimal storage
 - Practical notes
 - almost all searches are successful
 - words average about 8 characters in length
 - 30,000 words at 8 bytes/word is 1/4 MB
 - pointers are 4 bytes
 - there are *many* regularities in the structure of English words
- Why?

Solutions

- Solutions
 - sorted array + binary search
 - separate chaining
 - open addressing + linear probing

Storage

- Assume words are strings and entries are pointers to strings



Analysis

- Binary search
 - storage: n pointers + words = 360KB
 - time: $\log_2 n \leq 15$ probes per access, worst case
- Separate chaining
 - storage: $2n + n/\lambda$ pointers + words ($\lambda = 1 \Rightarrow 600\text{KB}$)
 - time: $1 + \lambda/2$ probes per access on average ($\lambda = 1 \Rightarrow 1.5$)
- Open addressing
 - storage: n/λ pointers + words ($\lambda = 0.5 \Rightarrow 480\text{KB}$)
 - time: $\frac{1}{2} \left(1 + \frac{1}{1-\lambda} \right)$ probes per access on average ($\lambda = 0.5 \Rightarrow 1.5$)

Which one should we use?

A Random Hash...

- Universal hashing
 - Given a particular input, pick a hash function parameterized by some random number
 - Useful in proving average case results – instead of randomizing over inputs, randomize over choice of hash function
- Minimal perfect hash function: one that hashes a given set of n keys into a table of size n with no collisions
 - Always exist
 - Might have to search large space of parameterized hash functions to find
 - Application: compilers
- One way hash functions
 - Used in cryptography
 - Hard (intractable) to invert: given just the hash value, recover the key

Coming Up

- Wednesday: Nick leads the class
- Try all the homework problems BEFORE Thursday, so you can ask questions in section!
- Friday: Midterm