

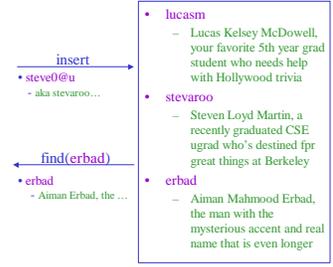
CSE 326: Data Structures

It's an open-and-closed hash!

Luke McDowell
Summer Quarter 2003

Reminder: Dictionary ADT

- Dictionary ADT:
 - Maps *values* to user-specified *keys*
 - Or: a set of (key, value) pairs
 - Keys may be any (homogeneous) type
 - Values may be any (homogeneous) type
- Operations:
 - Insert (key, value)
 - Find (key)
 - Remove (key)



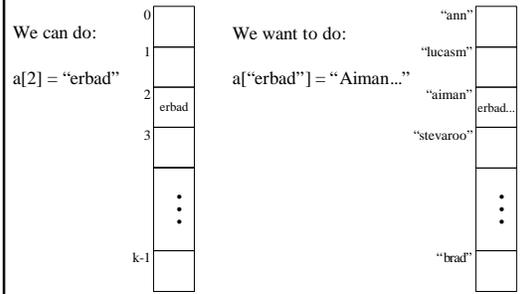
The Dictionary ADT is sometimes called the "Map ADT"

Implementations So Far

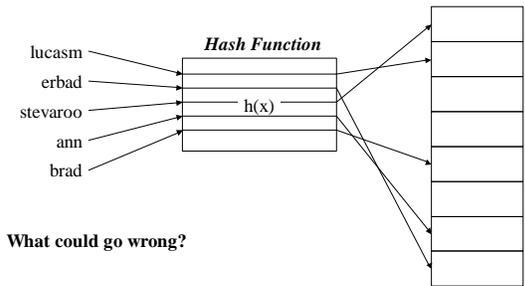
	insert	find	delete
• Unsorted list	O(1)	O(n)	O(n)
• Sorted list	O(n)	O(log n)?	O(n)
• Trees	O(log n)	O(log n)	O(log n)

How about O(1) insert/find/delete?

Hash Table Goal



Hash Table Approach



Hash Table Code First Pass

```

Value find(Key k) {
    int index = hash(k) % tableSize;
    return Table[index];
}
    
```

Key Questions:

1. What should the hash function be?
2. How should we resolve collisions?
3. What should the table size be?

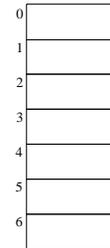
A Good Hash Function...

- ...is easy (fast) to compute ($O(1)$ and practically fast).
- ...distributes the data evenly ($\text{hash}(a) \% \text{size} \neq \text{hash}(b) \% \text{size}$).
- ...uses the whole hash table (for all $0 \leq k < \text{size}$, there's an i such that $\text{hash}(i) \% \text{size} = k$).

Good Hash Function for Integers

- Choose
 - tableSize is prime
 - $\text{hash}(n) = n$
- Example:
 - tableSize = 7

insert(4)
insert(17)
find(12)
insert(9)
delete(17)



Easy/boring stuff we're going to skip

- Why does the table size have to be prime?
- Picking a good hash function for strings

Read Section 5.2 of the text!

Collisions

- Pigeonhole principle says we can't avoid all collisions
 - try to hash without collision m keys into n slots with $m > n$
 - e.g., try to put 7 pigeons into 5 holes

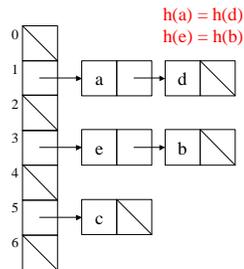


- What do we do when two keys hash to the same entry?
 - Separate chaining: put little dictionaries in each entry
 - Open addressing: pick a next entry to try
- Frequency depends on load factor

load factor $= \frac{\# \text{ of entries in table}}{\text{tableSize}}$

Separate Chaining

- Put a mini-Dictionary at each entry
 - Usually a linked list
 - Why not a search tree?
- Properties
 - ? can be greater than 1
 - performance degrades with length of chains



Load Factor in Separate Chaining

- Search cost
 - unsuccessful search: ?
 - successful search: ?
- Desired load factor: ?

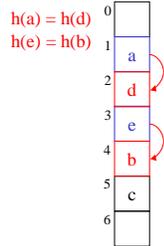
Open Addressing

What if we only allow one Key at each entry?

- two objects that hash to the same spot can't both go there
- first one there gets the spot
- next one must **probe** for another spot

- Properties

- ? ? 1
- performance degrades with difficulty of finding right spot



Salary-Boosting Obfuscation

“Open Hashing” equals “Separate Chaining”
 “Closed Hashing” equals “Open Addressing”

Probing

- Probing how to:
 - First probe - given a key k , hash to $h(k)$
 - Second probe - if $h(k)$ is occupied, try $h(k) + f(1)$
 - Third probe - if $h(k) + f(1)$ is occupied, try $h(k) + f(2)$
 - And so forth
- Probing properties
 - we force $f(0) = 0$
 - the i^{th} probe is to $(h(k) + f(i)) \bmod \text{size}$
- When does the probe fail?
- Does that mean the table is full?

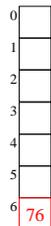
Linear Probing

$$f(i) = i$$

- Probe sequence is
 - $h(k) \bmod \text{size}$
 - $h(k) + 1 \bmod \text{size}$
 - $h(k) + 2 \bmod \text{size}$
 - ...

Linear Probing Example

insert(76) insert(93) insert(40) insert(47) insert(10) insert(55)
 $76\%7 = 6$ $93\%7 = 2$ $40\%7 = 5$ $47\%7 = 5$ $10\%7 = 3$ $55\%7 = 6$



Problem?

Load Factor in Linear Probing

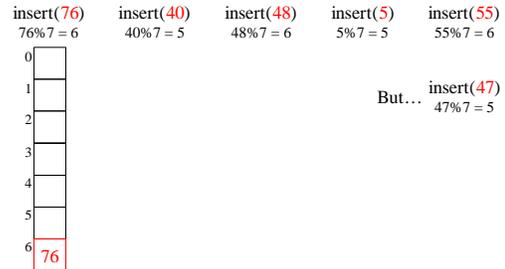
- For any $\alpha < 1$, linear probing will find an empty slot
- Search cost (for large table sizes)
 - successful search: $\frac{1}{2} \frac{1-\alpha^2}{1-\alpha}$
 - unsuccessful search: $\frac{1}{2} \frac{1+\alpha}{1-\alpha}$
- Linear probing suffers from *primary clustering*
- Performance quickly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

Quadratic Probing

$$f(i) = i^2$$

- Probe sequence is
 - $h(k) \bmod \text{size}$
 - $(h(k) + 1) \bmod \text{size}$
 - $(h(k) + 4) \bmod \text{size}$
 - $(h(k) + 9) \bmod \text{size}$
 - ...

Quadratic Probing Example



Quadratic Probing Succeeds (for $\alpha < 1/2$)

- If size is prime and $\alpha < 1/2$, then quadratic probing will find an empty slot in size/2 probes or fewer.
 - show for all $0 \leq i, j < \text{size}/2$ and $i \neq j$

$$(h(x) + i^2) \bmod \text{size} \neq (h(x) + j^2) \bmod \text{size}$$
 - by contradiction: suppose that for some $i \neq j$:

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

$$i^2 \bmod \text{size} = j^2 \bmod \text{size}$$

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

$$[(i + j)(i - j)] \bmod \text{size} = 0$$
 - but how can $i + j = 0$ or $i + j = \text{size}$ when $0 \leq i, j < \text{size}/2$?
 - same for $i - j \bmod \text{size} = 0$

Load Factor in Quadratic Probing

- For any $\alpha < 1/2$, quadratic probing will find an empty slot; for greater α , quadratic probing *may* find a slot
- Quadratic probing does not suffer from *primary* clustering
- But what about keys that hash to the same spot?

Double Hashing

$$f(i) = i \cdot \text{hash}_2(x)$$

- Probe sequence is
 - $h_1(k) \bmod \text{size}$
 - $(h_1(k) + 1 \cdot \text{hash}_2(x)) \bmod \text{size}$
 - $(h_1(k) + 2 \cdot \text{hash}_2(x)) \bmod \text{size}$
 - ...
- Goal?

A Good Double Hash Function...

- ...is quick to evaluate.
- ...differs from the original hash function.
- ...never evaluates to 0 (mod size).

One good choice is to choose a prime $R < \text{size}$ and:

$$\text{hash}_2(x) = R - (x \bmod R)$$

Extendible Hashing

Hashing technique for huge data sets

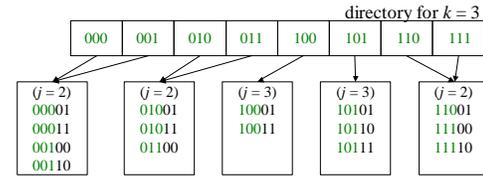
- Optimizes to reduce disk accesses
- Each hash bucket fits on one disk block
- Better than B-Trees if order is not important – *why?*

Table contains:

- Buckets, each fitting in one disk block, with the data
- A directory is used to hash to the correct bucket

Extendible Hash Table

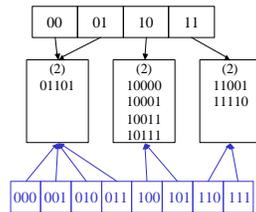
- Directory entry: *key prefix* (first k bits) and a pointer to the bucket with all keys starting with its prefix
- Each bucket contains keys matching on first j k bits, plus the value associated with each key



insert(11011)?
insert(11011)?

Splitting the Directory

1. insert(10010)
But, no room to insert and *no adoption!*
2. Solution: **Expand directory**
3. Then, it's just a normal split.



How to ensure this uncommon?

If Extendible Hashing Doesn't Cut It

Option 1: Store only pointers/references to the items:
(key, value) pairs are in disk

Option 2: Improve Hash + Reshash

The One-Slide Hash

Hash function: maps keys to integers

Collision resolution

- Separate Chaining
 - Expand beyond hashtable via secondary Dictionaries
 - Allows $? > 1$
- Open Addressing
 - Expand within hashtable
 - Secondary probing: {linear, quadratic, double hash}
 - ? ? 1 (by definition!)
 - ? ? ½ (by preference!)

Choosing a Hash Function

- Make sure table size is prime
- Careful choice for strings
- **"Perfect hashing"**
 - If keys known in advance, tune hash function for them!

Reshashing

- Tunes up hashtable when ? crosses the line

Extendible hashing

- For disk-based data
- Combine with B-tree directory if needed

Implementations So Far

	insert	find	delete
• Unsorted list	$O(1)$	$O(n)$	$O(n)$
• Sorted list	$O(n)$	$O(\log n)?$	$O(n)$
• Trees	$O(\log n)$	$O(\log n)$	$O(\log n)$
• Hash Table	$O(1)$	$O(1)$	$O(1)$

Is there anything a hash table can't do fast?