MIT OpenCourseWare <u>http://ocw.mit.edu</u>

6.006 Introduction to Algorithms Spring 2008

For information about citing these materials or our Terms of Use, visit: <u>http://ocw.mit.edu/terms</u>.

6.006 Recitation

Build 2008.14

Coming up next...

- Open addressing
- Karp-Rabin
 - coming back from the dead to hunt us

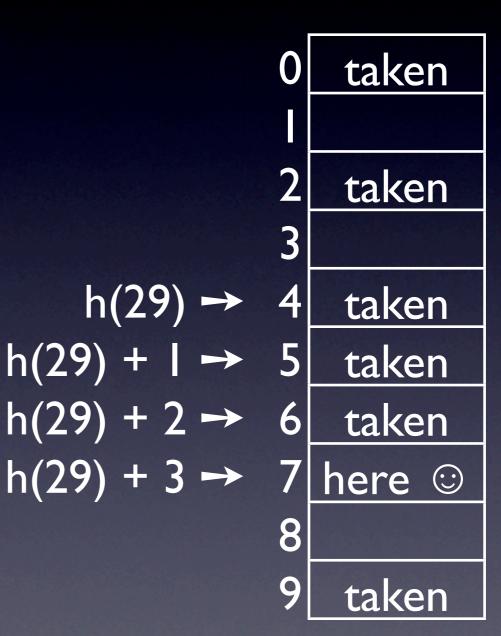
Open Addressing

• Goal: use nothing but the table

- Hoping for less code, better caching
- Hashing \Rightarrow we must handle collisions
 - Solution: try another location

Easy Collision handling

- h(x) = standard hash function
- if T[h(x)] is taken
 - try T[h(x)+1]
 - then T[h(x) + 2]
 - then T[h(x) + 3]
- just like parking a car



H(29, 0)

- h(k) grows up to H(k, i) where i is the attempt number
- first try T[H(k, 0)]

	0	taken
		taken
	2	taken
	3	taken
	4	taken
	5	taken
	6	taken
	7	taken
	8	taken
\rightarrow	9	taken

- h(k) grows up to H(k, i) where i is the attempt number
- first try T[H(k, 0)]
 - then T[H(k, I)]

taken 0 H(29, I) taken taken 3 taken taken 4 5 taken taken 6 taken 8 taken H(29, 0) → taken

- h(k) grows up to H(k, i) where i is the attempt number
- first try T[H(k, 0)]
 - then T[H(k, I)]
 - then T[H(k, 2)]

taken 0 H(29, I) → taken taken 3 taken H(29, 2)taken 5 taken taken 6 taken 8 taken H(29, 0) → taken

- h(k) grows up to H(k, i) where i is the attempt number
- first try T[H(k, 0)]
 - then T[H(k, I)]
 - then T[H(k, 2)]
- stop after trying all

 $H(29,3) \rightarrow 0$ taken H(29, I) → I taken $H(29,4) \rightarrow 2$ taken $H(29,9) \rightarrow 3$ taken $H(29,2) \rightarrow 4$ taken H(29, 5) → 5 taken H(29, 6) → 6 taken $H(29,7) \rightarrow 7$ taken H(29, 8) → 8 taken H(29,0) → 9 taken

Collision Handling: Abstracting it Up $H(29,3) \rightarrow 0$ taken

- H(k) =
 <H(k, 0), H(k, 1), H(k, 2) ... >
- Linear probing, h(29) =
 4, H_{linear}(29) = ?

<4, 5, 6, 7, 8, 9, 0, 1, 2, 3>

• General properties?

 $H(29,3) \rightarrow 0$ H(29, I) → I $H(29,4) \rightarrow 2$ $H(29,9) \rightarrow 3$ $H(29,2) \rightarrow 4$ $H(29,5) \rightarrow 5$ H(29, 6) → 6 $H(29,7) \rightarrow 7$ H(29, 8) → 8 $H(29,0) \rightarrow 9$

taken taken taken taken taken taken taken taken taken

- Any collision handling strategy comes to:
 - for key k, probe H(k,0), then H(k,1) etc.
- No point in trying the same place twice
- Probes should cover the whole table (otherwise we raise 'table full' prematurely)
- Conclusion: H(k, 0), H(k, 1) ... H(k, m-1) are a permutation of {1, 2, 3 ... m}

Linear Probing and Permutations

h(29) = 4; H(29) =

<4, 5, 6, 7, 8, 9, 0, 1, 2, 3>

• $h(k) = h_0 (mod m); H(k) =$

 $<h_0 \mod m$, $(h_0 + I) \mod I$ m, $(h_0 + 2) \mod m$, ... $(h_0 + m - 1) \mod m >$

m permutations (max m!)

taken taken 3 h(29) → 4 taken h(29) + I → 5 taken h(29) + 2 → 6 taken $h(29) + 3 \rightarrow 7$ here 😳 8 9

taken

Ideal Collision Handling

- Simple Hashing (collision by chaining)
 - Ideal hashing function: uniformly distributes keys across hash values
- Open Addressing
 - Ideal hashing function: uniformly distributes keys across permutations
 - a.k.a. uniform hashing

Uniform Hashing: Achievable?

- Simple mapping between permutations of m and numbers I ... m!
- Convert key to big number, then use permutation number (bignum mod m!)
- ... right?

n	k mod 6	Permutation
	0	<1, 2, 3>
]	<1,3,2>
	2	<2, , 3>
	3	<2, 3, >
	4	<3, I, 2>
	5	<3.2. >

Uniform Hashing: Achievable?

 Number of digits in m! 	k mod 6	Permutation
 O(log(m!)) 	0	<1,2,3>
 O(m*log(m)) 	l	<1,3,2>
	2	<2, I, 3>
 Working mod m! is slow 	3	<2, 3, I>
 check your Python cost model 	4	<3, I, 2>
	5	<3, 2, >

Working Compromise

• Why does linear probing suck?

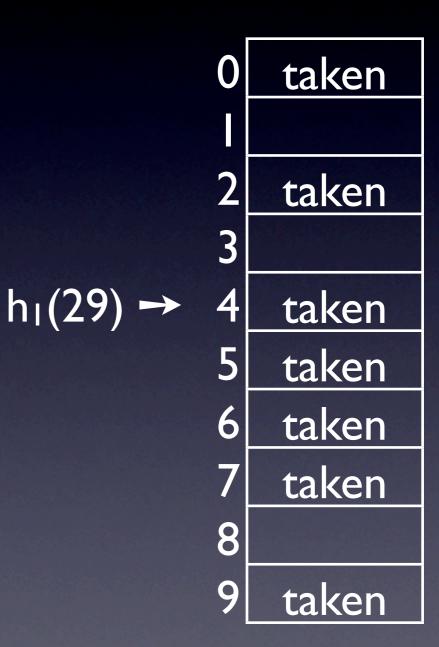
- We jump in the table once, then walk
- Improvement
 - Keep jumping after the initial jump
 - Jumping distance: 2nd hash function
 - Name: double hashing

 h₁(k) and h₂(k) are hashing functions



 h₁(k) and h₂(k) are hashing functions

• $H(k, 0) = h_1(k)$



 $h_1(k)$ and $h_2(k)$ are 0 taken hashing functions 2 taken $H(k,0) = h_1(k)$ 3 h₁(29) → $H(k, I) = h_1(k) + h_2(k)$ 4 taken 5 taken taken 6 $h_1(29)+h_2(29) \rightarrow$ taken 8 9 taken

•	$h_1(k)$ and $h_2(k)$ are	h₁(29)+2·h₂(29) →	0	taken
	hashing functions		I	
•	$H(k, 0) = h_1(k)$		2	taken
			3	
•	$H(k, I) = h_1(k) + h_2(k)$	h₁(29) →	4	taken
			5	taken
			6	taken
		$h_1(29)+h_2(29) \rightarrow$	7	taken
			8	
			9	taken

•	h ₁ (k) and h ₂ (k) are hashing functions	$h_1(29)+2 \cdot h_2(29) \rightarrow$	0 	taken
•	$H(k, 0) = h_1(k)$		2	taken
		$h_1(29) + 3 \cdot h_2(29) \rightarrow$	3	here 😳
•	$H(k, I) = h_1(k) + h_2(k)$	h₁(29) →	4	taken
			5	taken
			6	taken
		h₁(29)+h₂(29) →	7	taken
			8	
			9	taken

- $h_1(k)$ and $h_2(k)$ are $h_1(29)+2 \cdot h_2(29) \rightarrow$ 0 taken hashing functions taken $H(k, 0) = h_1(k)$ $h_1(29) + 3 \cdot h_2(29) \rightarrow$ 3 here 🙂 h₁(29) → $H(k, I) = h_1(k) + h_2(k)$ 4 taken 5 taken • $H(k, i) = h_1(k) + i \cdot h_2(k)$ 6 taken $h_1(29) + h_2(29) \rightarrow$ taken mod m 8 9 taken
 - you knew that, right?

Double Hashing Trap

 gcd(h₂(k), m) must be I 	$h_1(29)+2 \cdot h_2(29) \rightarrow$	0	taken
• colution (coov to got)			
 solution I (easy to get) 		2	taken
• m prime, $h_2(k) = k$	$h_1(29) + 3 \cdot h_2(29) \rightarrow$	3	here 🙂
mod q (where q < m)	h₁(29) →	4	taken
		5	taken
 solution 2 (faster, better) 		6	taken
• m = 2 ^r (table can grow)	h₁(29)+h₂(29) →	7	taken
• $\Pi = Z$ (table call grow)		8	
b.(k) is odd (not over)		9	takon

taken

• h₂(k) is odd (not even)

Open Addressing: Deleting Keys

- Suppose we want to delete k_d stored at 7
- Can't simply wipe the entry, because key 29 wouldn't be found anymore
 - rember H(29) =
 <4, 7, 0, 3 ...>

 $h_1(29) + 3 \cdot h_2(29) \rightarrow h_1(29) \rightarrow$

 $h_1(29) + 2 \cdot h_2(29) \rightarrow$

 $h_1(29) + h_2(29) \rightarrow$

)	taken
2	taken
3	here 😊
4	taken
5	taken
6	taken
7	kd
3	
9	taken

Open Addressing: Deleting Keys

- Entry meaning 'deleted'
- Handling 'deleted'
 - Search: Keep looking
 - Insert: Stop, replace 'deleted' with the new key/value

 $h_1(29) + 3 \cdot h_2(29) \rightarrow h_1(29) \rightarrow$

 $h_1(29) + 2 \cdot h_2(29) \rightarrow$

 $h_1(29) + h_2(29) \rightarrow$

0	taken
2	taken
3	here 😊
4	taken
5	taken
6	taken
7	deleted
8	
9	taken

- Design: implementing a collection in Python
 - ______(self, key)
 - return key item or raise KeyError(key)
 - _____(self, key, item)
 - insert / replace (key, item)
 - __delitem__(self, key)

- Closures: not for n00bs
- def compute_modulo is local to the mod_m call
- the function created by def compute_modulo is returned like any object
- the object remembers the context around the def (the value of m)

```
1 def mod_m(m):
2     def compute_modulo(n):
3         return (n % m)
4         return compute_modulo
5
6 >>> m5 = mod_m(5)
7 >>> m3 = mod_m(3)
8 >>> m5(9)
9 4
10 >>> m3(9)
11 0
```

```
1 def linear_probing(m = 1009):
      def hf(key, attempt):
 2
          return (hash(key) + attempt) % m
 3
 4
      return hf
 5
   def double_hashing(hf2, m = 1009):
 6
      def hf(key, attempt):
 7
          return (hash(key) + attempt * hf2(key)) % m
 8
 9
      return hf
10
   class DeletedEntry:
11
12
      pass
13
   class OpenAddressingTable:
14
15
      def __init__(self, hash_function, m = 1009):
          self.entries = [None for i in range(m)]
16
          self.hash = hash_function
17
18
          self.deleted_entry = DeletedEntry()
```

```
14 class OpenAddressingTable:
15
      def __init__(self, hash_function, m = 1009):
16
          self.entries = [None for i in range(m)]
          self.hash = hash_function
17
18
          self.deleted_entry = DeletedEntry()
19
20
      def get_entry(self, key):
21
          for attempt in xrange(len(self.entries)):
22
              h = self.hash(key, attempt)
23
              if self.entries[h] is None:
24
                  return None
25
              if self.entries[h] is not self.deleted_entry and \setminus
26
                 self.entries[h][0] == key:
27
                   return self.entries[h]
28
29
      def __getitem__(self, key):
30
          entry = self.get_entry(key)
31
          if entry is None:
32
              raise KeyError(key)
33
          return entry[1]
34
35
      def __contains__(self, key):
36
          return self.get_entry(key) is not None
```

```
14 class OpenAddressingTable:
      def __init__(self, hash_function, m = 1009):
15
          self.entries = [None for i in range(m)]
16
          self.hash = hash function
17
          self.deleted_entry = DeletedEntry()
18
19
37
      def __setitem__(self, key, value):
38
          if value is None: raise 'Cannot set value to None'
39
          del self[key]
          for attempt in xrange(len(self.entries)):
40
              h = self.hash(key, attempt)
41
              if self.entries[h] is None or \
42
                 self.entries[h] is self.deleted_entry:
43
                  self.entries[h] = (key, value)
44
45
                  return
46
          raise 'Table full'
```

```
14 class OpenAddressingTable:
15
      def __init__(self, hash_function, m = 1009):
16
          self.entries = [None for i in range(m)]
17
          self.hash = hash_function
18
          self.deleted_entry = DeletedEntry()
19
47
      def __delitem__(self, key):
48
          for attempt in xrange(len(self.entries)):
49
              h = self.hash(key, attempt)
50
              if self.entries[h] is None:
51
                   return
52
              if self.entries[h] is not self.deleted_entry and \setminus
53
                 self.entries[h][0] == key:
54
                   self.entries[h] = self.deleted_entry
55
                   return
56
          return
```

Getting Rolling Hashes Right

Modular Arithmetic

• Foundation:

- (a + b) mod m = ((a mod m) + (b mod m)) mod m
- From that, it follows that:
 - (a · b) mod m = ((a mod m) · (b mod m)) mod m
 - induction: multiplication is repeated +

Modular Gotcha

- Never give mod a negative number
 - want q = $(a b) \mod m$, but a b < 0
 - $q \mod m = (a (b \mod m)) \mod m$
 - but (b mod m) is < m
 - so (a + m (b mod m)) > 0
 - $q = (a + m (b \mod m)) \mod m$

Modular Arithmetic-Fu

- Multiplicative inverses: assume p is prime
- For every a and p, there is a⁻¹ so that:
 - $(a * a^{-1}) \mod p = 1$
 - example: p = 23, $a = 8 \Rightarrow a^{-1} = 3$
 - check: 8 * 23 = 24, 24 mod 23 = 1
- Multiplying by a⁻¹ is like dividing by a

Modular Arithmetic-Fu

- How do we compute a⁻¹?
- Fermat's Little Theorem:
 - p prime $\Rightarrow a^{a-1} \mod p = 1$
- Huh?
 - $a^{a-1} \mod p = a * a^{a-2} \mod p = 1$
 - so (for p) $a^{-1} \mod p = a^{a-2} \mod p$

Back to Rolling Hashes

- Data Structure (just like hash table)
 - start with empty list
 - append(val): appends val at the end of list
 - skip(): removes the first list element
 - hash(): computes a hash of the list