Week 3: Balls and Bins	Jan 29, 2018
	LS 5 330
Last Time:	
• QuickSort	
• Average-Case Analysis Stable Marriage	
• The Coupon Collector's Problem.	
Today: Hashing	
• Balls-in-Bins	
• Hashing with Chaining	
• Linear Probing	
Cuckoo Hashing	
Puzzle of the Day	
-> 2 envelopes	
-> one has 2× money	
-> Pick an envelope at random	
Should suitch?	
let A= money in envelope	
E[other envelope] = 1/2 2A + 1/2 A	$f = \frac{5}{0} A$
Alucos sutch!	
7	
What is Wrong!	

## Balls-and-Bins











### Hashing with Chaining

The simplest form of hashing is to resolve collisions with chaining. Each key x is hashed using hash function f to a table  $A[1, \ldots, m]$ . Each slot in the table stores a linked list containing all the keys hashed to that slot. To insert an item, simply compute f(x) and add the key to the front of the linked list in slot A[f(x)]. This takes O(1) time (assuming you can compute the hash function in O(1) time). To search for an item, compute f(x) and then search through the entire linked list at A[f(x)]. The cost of a search depends on the size of the linked lists.

Using the balls in bins analysis, you should be able to show that:

- If the hash table has size m and there are n items in the table, then the expected cost of an insert is O(n/m).
- If n = m, the expect cost of an insert is O(1).
- Using the analysis above, you can show that if m = n, then with high probability, every linked list has length  $O(\log n / \log \log n)$ , and so each search operation terminates in time  $O(\log n / \log \log n)$  with high probability.
- You should be able to show that if m = n, then the following statement is *false*: with high probability, each linked list has length O(1).
- In fact, using the analysis above, you can show that in expectation, many search operations will be cost  $\Omega(\log n / \log \log n)$ , and you can calculate the probability that at least one searc operation has cost  $\Omega(\log n / \log \log n)$ .



### •









-> S in at most 5 level 1-2 nodes

Claim: ≥1 of first 4 level l-2 nodes for S are crouded

Why?  
1) All Keys in first 4 nodes hash to  
first 4 nodes (cannot hash earlier:  
A[S\_1-] is empty)  
2) 
$$\geq 3 \cdot 2^{p-2}$$
+1 Keys hash to first 4 nodes  
 $\Rightarrow \geq 1$  in node 1  
 $\Rightarrow$  nodes 2,3,4 are full  $\Rightarrow$  covered by cluster  
 $\Rightarrow \geq 3 \cdot 2^{p-2}$ +1 in nodes 1,2,3,4  $\Rightarrow \geq 3 \cdot 2^{p-2}$   
 $\Rightarrow \geq 3 \cdot 2^{p-2}$ +1 in nodes 1,2,3,4  $\Rightarrow \geq 3 \cdot 2^{p-2}$   
hash to first 4 nodes  
Contradiction  $\Rightarrow \geq 1$  first 4 nodes crouded

# Long streak => crouded nodes



Claim: Prlq is in cluster of size [2, 2"] ≤ 12 Po-2 For each UEN, Pr(In crowded) ≤ Po-2 Union bound ... E cost q < < K. Pr(q in cluster of size K) ∠ Ž 2 · Pr(q in cluster of size [2,2"]

< 3 + 2 2 · 12 P1-2

If Poz is small, then E[cost] is small!

Balls - and - Bins  
D balls  
b bins 
$$(b = \frac{m}{2^{p-1}} \leftarrow one \ bin \ per \ l-2 \ node)$$

Assume:

Fix one bin Z.  $bin crouded 2^{2-2}$  $Z = \frac{r}{b} = M$ E Markov: Pr(Z>3M) ≤ 1/3 Not good:  $\sum_{q=2}^{1} 2^{q+1} \cdot \frac{1}{3} = G(m)$ Chebychev:  $Pr(Z > 3M) \leq Pr(|Z-M| > 2M)$  $\leq E (z-m)^{2}$ 

Define 
$$X_{j} = 1$$
 if bill  $j$  in bin  
 $= 0$  other use  
Pr  $\{X_{j}\} = E[X_{j}] = \frac{1}{b}$   
 $Z = \sum_{j=1}^{n} X_{j}$   
Define  $Y_{j} = X_{j} - \frac{1}{b}$   
 $E Y_{j} = \sum_{j=1}^{n} X_{j} - \frac{1}{b} = Z - M$   
 $\Rightarrow E[(Z - M^{L})] = E[(EY_{j})^{L}]$   
 $E[[Z - M^{L}]] = E[(EY_{j})^{L}]$   
 $E[[Z - M^{L}]] = E[Y_{i}, Y_{j}]$   
 $i) i \neq j$ :  $Y_{i}, Y_{j}$  independent  
 $E[Y_{i}, Y_{j}] = E[Y_{i}] E[Y_{j}] = 0 \cdot 0 = 0$   
 $(E[Y_{i}] = E[X_{i}] - \frac{1}{b} = \frac{1}{b} - \frac{1}{b} = 0$   
 $(E[Y_{i}] = E[Y_{i}]) = \frac{1}{b} (1 - \frac{1}{b})^{L} + (1 - \frac{1}{b}) [-\frac{1}{b}]^{L}$   
 $= \frac{1}{b} - \frac{1}{b} + \frac{1}{b} + \frac{1}{b} - \frac{1}{b}$ 

Conclusion:  $E[(\xi, Y_i)^2] \leq \xi E[Y_i^2] \leq \xi = M$  $Pr(|z-M| > 2M) \leq \frac{M}{4M^2} \leq \frac{1}{4M}$  $\frac{1}{4M} = \frac{1}{4} \cdot \frac{1}{2^{\ell-2}} \cdot \frac{1}{2} = \frac{1}{4} \cdot \frac{1}{2^{\ell-2}}$  $= \sum_{n=2}^{1} 2^{n+1} \cdot \frac{1}{42^{n-2}} = \sum_{p=2}^{1} \frac{1}{4} \cdot 8 \approx 2\log(m)$ Better... On to the fourth moment! Why not 3??









$$E\left[\cos t \ 2\right] \leq 3 + 12 \cdot \sum_{p=2}^{t} 2^{n+1} P_{p-2}$$

$$\leq 3 + \frac{12}{9} \sum_{p=2}^{t} 2^{n+1} \frac{1}{2^{1+n}}$$

$$\leq 3 + \frac{12}{9} \sum_{p=2}^{t} 2^{n+1} \frac{1}{2^{1+n}}$$

$$\leq 3 + \frac{12}{9} \cdot \frac{32}{2} \leq 21 = O(1)$$

$$E\left[\cos t \text{ of insert/access}\right] = O(1)$$

$$More \ care = smaller \ constant$$

$$More \ care \Rightarrow M = n(1 + 2)$$

$$In \ practice: \ Very \ fast!$$

$$Cache \ performance \ easily$$

$$More \ important \ than \ constant$$

$$f \ actors.$$

$$More \ important \ than \ constant$$

$$f \ actors.$$

$$More \ linear \ Probing$$

$$[Don't \ believe \ und \ they \ tell \ pour...)$$

#### Cuckoo Hashing

This week, we considered Cuckoo Hashing. The basic idea behind Cuckoo Hashing is that we use to tables A and B of size m = 2n to store n items. For table A, we use hash function f, and for table B, we use hash function g. To insert an item, we proceed as follows, executing insert(x, A, f):

```
insert(item x, Table T, hashfunction h)
    slot = h(x)
    z = T[slot]
    T[slot] = x
    if (z != null)
        if (T == A)
            insert(z, B, g)
        else if (T == B)
            insert(z, A, f)
```

That is, we put the new item in the proper slot in table A. If the slot is already occupied by some item z, then we kick z out of table A and insert it (recursively) into table B. Eventually, if this proceeds for more than 2m steps, then we know there is a problem and we abort.

Searching for an item is simple, and requires at most 2 queries:

search(x)

```
if (A[f(x)] == x) then return A[f(x)];
if (B[g(x)] == x) then return B[g(x)];
else return NOT_IN_TABLE;
```

Item x can be in at most two locations: table A or table B, so searching is very fast: O(1) time.

An insertion, on the other hand, can be more expensive. If  $m \ge 2n$ , we can show that most of the time, each insertion completes in  $O(\log n)$  with high probability. But sometimes, on occasion, you may have a problem. For example, imagine you have three keys x, y, z where f(x) = f(y) = f(z) and g(x) = g(y) = g(z). In this case, the insertion will loop forever—or until it aborts. More generally, we have a problem when the Cuckoo Graph has two cycles<sup>1</sup>:



In either of these cases, we will not be able to insert all the items. When this happens, we choose new hash functions and try again to insert all the items. We can show that for a given set of n elements, there is a probability of O(1/n) of failure. (Below, we will suggesting proving the easier bound, i.e., that there is a probability 1/2 of a failure.) Thus, with probability 1/n we will have to re-hash (and re-build) the tables, which has cost O(n). Thus the expected cost due to rebuilding is O(1). Overall, then, we can show that insertions have expected amortized cost O(1).

<sup>&</sup>lt;sup>1</sup>Figure taken from: http://infoweekly.blogspot.sg/2010/02/cuckoo-hashing.html

### Analysis

To analyze the performance of Cuckoo Hashing, we consider the Cuckoo Graph which is a bipartite graph defined as follows: there are m left node representing array A and m right nodes representing array B; for every key x, add one edge to the graph connecting node f(x) on the left side to node g(x) on the right side.

Notice that when analyzing the Cuckoo Graph, we are not actually *executing* the insertion algorithm. We are simply looking at the edges that would be created by the keys. By looking at the Cuckoo Graph for n keys, we can determine both whether or not it is possible to insert them all (without an infinite loop) and the expect cost of each insertion.

As an exercise, prove the following facts regarding Cuckoo Hashing. Assume throughout that m = 4n.

- Given a Cuckoo Graph containing n edges, show that any path of length p in the Cuckoo Graph that is generated by an insert operation that terminates has a simple subpath (with no repeated nodes) of length at least p/3.
- Given a Cuckoo Graph containing *n* edges, show that for all nodes *i* and *j*, for all  $\ell$ , the probability that the graph contains a path of length exactly  $\ell$  from *i* to *j* is at most  $\frac{1}{4\ell} \frac{1}{m}$ .
- Prove that the expected cost of an insert operation is O(1) (as long as the Cuckoo Hash Table contains at most n items).
- Given a Cuckoo Graph containing n edges, prove that the probability it contains a cycle is at most 1/2.
- Whenever an insert does not terminate (i.e., ends in an infinite loop), you need to create a new Cuckoo Hash Table, choose new hash functions, and reinsert all n items. (If your insert operation continues for 2m steps, you can assume you have entered an infinite loop!). Assume that you continue this process, rehashing as necessary, until you successfully insert all n items. Prove that the expected total cost of this process is O(n) (including the cost of the repeated rehashing).