

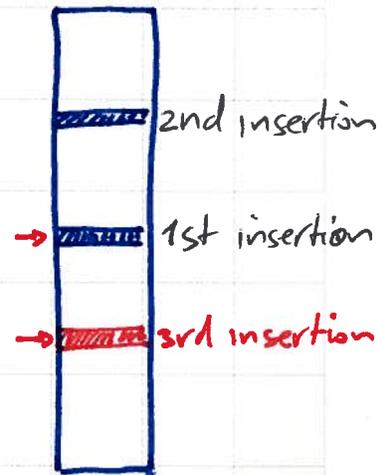
## 6.006 Lecture 7: Hashing 3

- Rolling hashes
- Open Addressing
- Uniform hashing + analysis
- Advanced topics: universal, perfect, & cryptographic hashing

□ CLRS 11.4, 11.3.3, 11.5

### Open Addressing

- No linked lists
- Collision? Store elsewhere in hash table
- More collisions? Probe more; may need to probe  $m-1$  times to find an empty slot



- Hash function ~~of key k~~ <sup>of key k</sup> is now a sequence of probes  $\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$   
sequence must be a permutation of  $0, 1, 2, \dots, m-1$   
a key maps to a permutation
- Clearly, load factor  $\alpha \leq 1$ .



Example of a rolling hash

Alphabet     <sup>0</sup> "A"    <sup>1</sup> "C"    <sup>2</sup> "G"    <sup>3</sup> "T"

$h = \text{value of string} \% p$  for prime  $p$

T = C T A T T A C G T  
1 3 0 3 3 0 1 2 3  $4 \leftarrow \text{base } 4$   
18 4 0 9 1<sub>10</sub>  
4 5 3 mod  $p = 1009$

$\downarrow$  delete C      $\leftarrow$  add G  
T A T T A C G T G

value of leading C =  $1 \cdot 4^8 = 1 \cdot 960$   $\leftarrow \text{mod } p, \text{ precompute!}$

so dropping leading C  $\Rightarrow h = 453 - 690$   
 $= 502 \text{ mod } p$

shifting and adding a G:

$$h = \underset{\substack{\uparrow \\ \text{shift}}}{4} \cdot 502 + \underset{\substack{\uparrow \\ G}}{2} = 1001 \text{ mod } p$$

subtract  $\times 1$   
multiply  $\times 1$   
add  $\times 1$  } cost of moving  $h$  to next window

## Insert (k, v):

for i in range(m):

if  $T[h(k, i)] == \text{None}$ :  $T[h(k, i)] = (k, v)$

return

raise Exception('full')

## Search(k):

for i in range(m)

if  $T[h(k, i)] == \text{None}$ : return None # not in table

if  $T[h(k, i)][0] == k$ : return  $T[h(k, i)]$

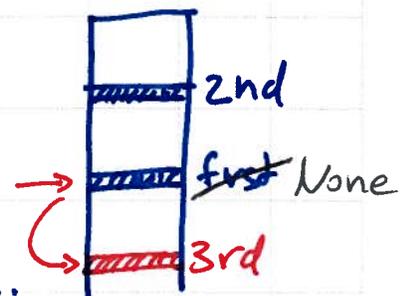
return None

## Delete(k):

# tricky! setting  $T[h(k, i)] = \text{None}$  may cause

search to fail (e.g. deleting the first element inserted in the example causes the third not to be found)

- Find key
- Replace by 'Deleted'
- Skip ~~over~~ over 'Deleted' in search but use 'Deleted' slots in insert.

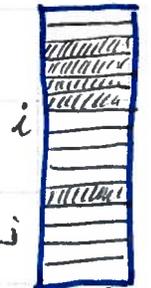


## How to construct $h(k,i)$

- What do we want?  
for chaining, we want simple uniform hashing:  
each key is equally likely to hash to any slot.
- For open addressing, we want uniform hashing:  
each key is equally likely to hash to any of the  $m!$  probe sequences (permutations of  $0, 1, \dots, m-1$ ).
- Harder to achieve, but double hashing works well.

## Linear probing

- Start with an ordinary hash function  $h'(k)$
- $h(k,i) = (h'(k) + i) \bmod m$
- Start at  $h'(k)$  and scan sequentially
- Not good: only  $m$  possible sequences, leads to clustering



which slot is more likely to get filled next?  
 $i$ , leads to clustering

## Double hashing

- $h(k,i) = (h_1(k) + i \cdot h_2(k)) \bmod m$
- to ensure  $h(k,*)$  hits all slots, make  $h_2(k)$  and  $m$  relatively prime. Ex:  $m = 2^r$ ,  $h_2(k)$  odd

## Open addressing

vs

## Chaining

cost explodes as  $\alpha$  approaches 1

No memory allocation (except to resize), cache efficient

Hard to find a really good  $h$

easier to implement in hardware

cost rises gently with  $\alpha$

allocates memory as chains grow (constant overhead)

Easy to find a good  $h$

## Advanced topics in hashing

### Universal hashing (back to chaining)

For any  $h$  there are collisions; if we are unlucky all the keys may hash to 1 slot.

Solution: Don't use a fixed  $h$ ; choose it at random

Universal hashing: random selection of  $h$  should guarantee  $\Pr[k_1, k_2 \text{ collide}] \leq \frac{1}{m}$

$h(k) = ((ak + b) \bmod p) \bmod m$  (last lecture) works.  
↑  
random

### Perfect hashing: $\Theta(1)$ worst-case search (fixed set of keys)

• Primary table stores pointers to secondary tables + their hash fn

• Make secondaries large enough so prob. of any collision  $\leq \frac{1}{2}$

• Try secondary hashes until no collisions at all

Space is still  $\Theta(n)$  (expected)

