

CS180 - Hashing (part 3)

Note Title

5/2/2011

Announcements

- Checkpoint due today
- Program due Friday
- Last HW out today
due Monday (not graded)

Data Storage - Dictionary : insert find remove

Ex. key →

Locker #	Name
26	Dan
355	Kevin
101	Tracy
53	Nitish
201	David
⋮	⋮

← data

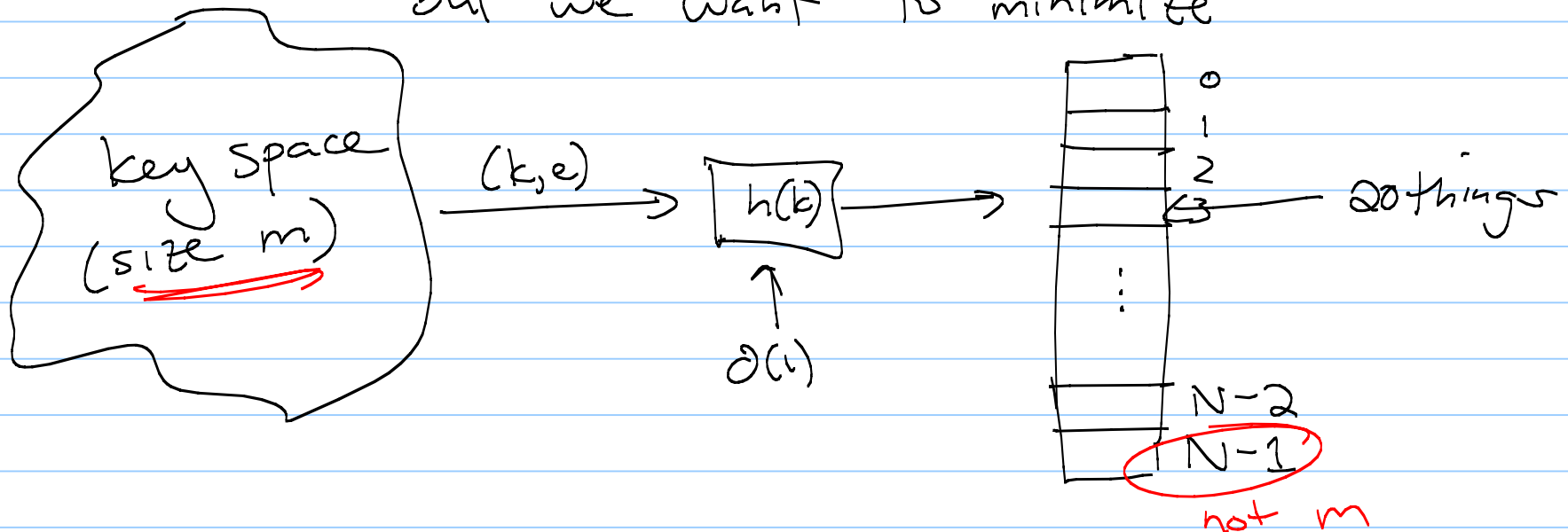
We want to be able to retrieve a name quickly when given a locker number.

(Let $n = \#$ of people, &
 $m = \#$ of lockers)

$$n \approx m$$

Good hash functions:

- Are fast goal: $O(1)$
- Don't have collisions - when $k_1 \neq k_2$ but $h(k_1) = h(k_2)$
these are unavoidable, but we want to minimize



Step 1: Turn key into an integer
polynomial hashing \Rightarrow integer (32-bits)

Step 2: Compression map

Take integer and make it $< N$.

mod (or %)

MAD

Collisions

Can we ever totally avoid collisions?

No

m is bigger than n or N

Step 3: Handle collisions
(gracefully & quickly)

So how can we handle collisions?

[Hint: Do we have any data structures that can store more than 1 element?]

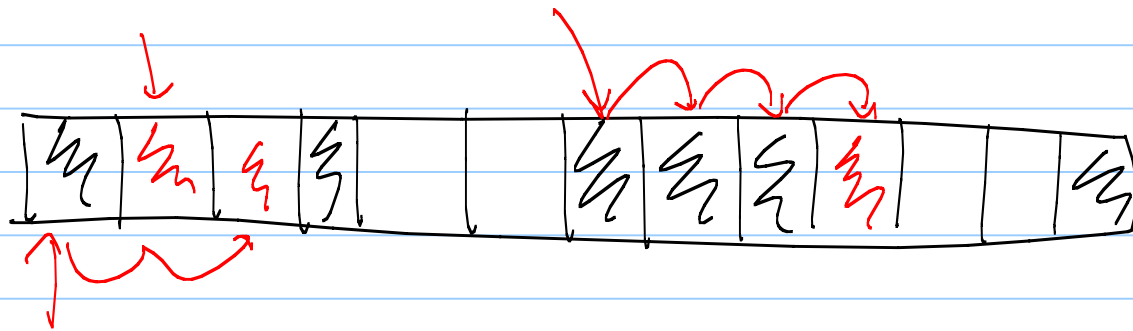
inserted a list in each
array space

lists = separate chaining
 $O(n)$

AVL tree
Vector

Linear Probing

Instead of lists, if we hash to a full spot, just keep checking next spot (as long as the next spot is not empty).

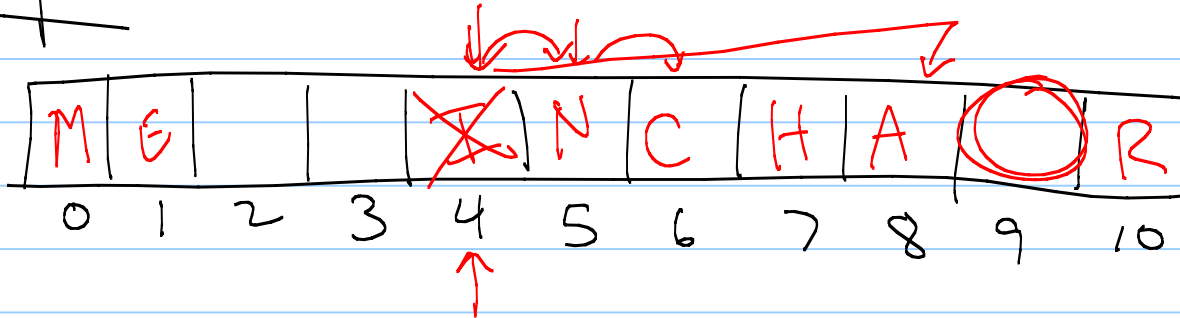


Example

$$h(k) = k \bmod 11$$

$$h(k) = (k + i) \bmod 11$$

$i = 0, 1, 2, \dots$



Insert:

remove →

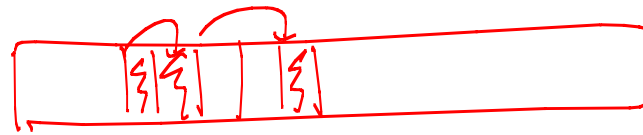
- (12, E) $12 \bmod 11 = 1$
 - (21, R) $21 \bmod 11 = 10$
 - (37, I) $37 \bmod 11 = 4$
 - (26, N) $26 \bmod 11 = 4$
 - (16, C) $16 \bmod 11 = 5$
 - (5, H) $5 \bmod 11 = 5$
 - (15, A) $15 \bmod 11 = 4$
 - (10, M)
- find (48, -)

Running Time for Linear Probing

Insert: $O(n)$ (not $O(N)$)
 ↑ ↑
 # of elements in array

Remove: remove bit (don't actually remove)
 "dirty bits"
 $O(n)$ if enough of array is dirty,
 re hash everything

Find: hash to value & keep checking
 next spot while not found &
 $O(n)$ next not empty



Quadratic Probing

Linear probing checks $A[h(k) + j \text{ mod } N]$ if $A[h(k) \text{ mod } N]$ is full.
 j , or $j=1, 2, 3, 4, \dots$

To avoid clusters, try \leftarrow primary clustering

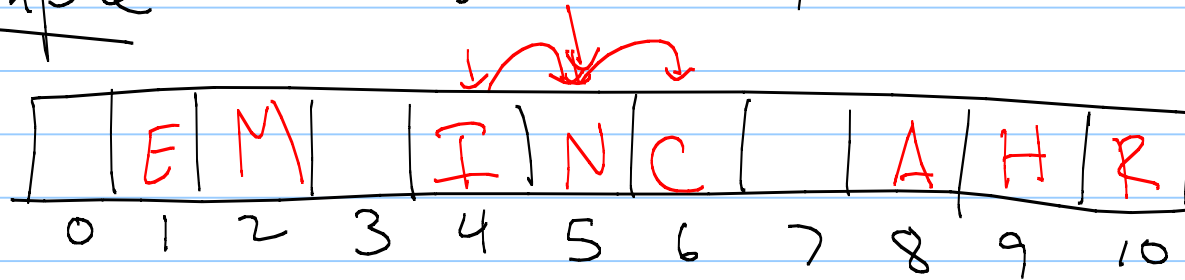
$$A[h(k) + j^2 \text{ mod } N]$$

where $j=0, 1, 2, 3, 4, \dots$

$$\begin{aligned} h(k) \\ h(k) + 1^2 &= h(k) + 1 \\ h(k) + 2^2 &= h(k) + 4 \\ h(k) + 3^2 &= h(k) + 9 \\ &\vdots \end{aligned}$$

Example

$$h(k) = k \bmod 11$$



Insert:

- (12, E)
- (21, R)
- (37, I)
- (26, N)
- (16, C)
- (5, H)
- (15, A)
- (4, M)

$$\begin{array}{l} 12 \bmod 11 = 1 \\ 21 \bmod 11 = 10 \\ 37 \bmod 11 = 4 \\ 26 \bmod 11 = 4 \\ 16 \bmod 11 = 5 \\ 5 \bmod 11 = 5 \\ 15 \bmod 11 = 4 \\ 4 \bmod 11 = 4 \end{array}$$

Issues with Quadratic Probing:

- Can still cause secondary clustering
- N really must be prime for this to work
- Even with N prime, starts to fail when array gets half full

(Runtimes are essentially the same)

Can fail entirely!

Secondary Hashing

- Try $A[h(k)]$

- If full, try $A[h(k) + f(j) \bmod N]$
for $j = 1, 2, 3, \dots$

where

$$f(j) = j \cdot h'(k)$$

with h' a different
hash function

Best if we want to avoid clusters.

Load Factors

Separate chaining actually works as well as most others in practice, although it does use more space.

Most of these methods only work well if $\frac{n}{N} < 0.5$.

(Even chaining starts to fail if $\frac{n}{N} > 0.9$)

usually $0.25 < \frac{n}{N} < 0.5$

Rehashing

Because we need $\frac{n}{N} < 0.5$, most hash code checks if the array has become more than half full.

If so, it stops & recomputes everything for a larger N , usually at least twice as big.

(Still not too bad in an amortized sense - think vectors.)