

Admin: → PSet 2 due tonight!

→ PSet 3 out

→ Looking at problems 3-1 to 3-4
before the quiz is a good review

→ No Part B this time

→ Quiz 1 next week!

See handout □
on Stellar ○

→ Review recitations on Wed

→ NO LECTURE on Thur

→ "Normal" recitations on Fri

→ Practice on past quizzes

courses.csail.mit.edu/6.006

→ Cookies today! ?

Today: → Table resizing

→ Amortized analysis

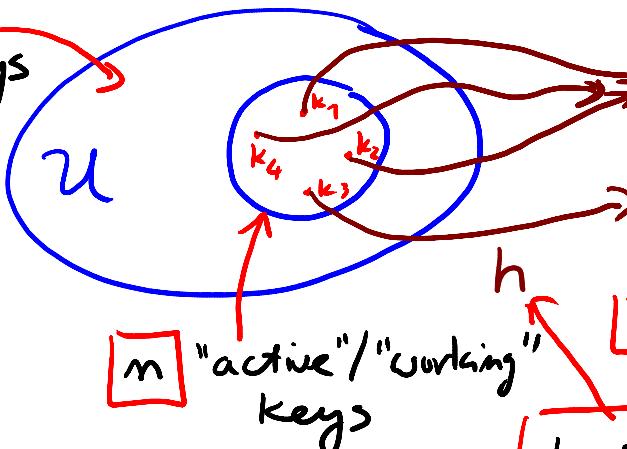
→ String matching & Karp-Rabin alg.

→ Rolling hash

Reading: CLRS Ch 17

Last time: Hashing with chaining:

all possible keys



table

colliding items

k_1

k_2

k_3

k_4

expected
chain length

n "active"/"working" keys

hash function

$$\alpha = \frac{n}{m}$$

load factor

→ Expected cost of insert/search/delete:

$$\Theta(1 + \alpha)$$

$O(1)$ if $m = \Omega(n)$

Important: This assumes that:

→ h satisfies Simple Uniform Hashing Assumption (or Universal Hashing)

→ h takes $O(1)$ time to compute

"Good" hash functions: → Division method $h(k) = k \bmod m$

Can be broken by certain sets of keys

→ Multiplication method:

$$128 = 2^7$$

$$h(k) = [(a_j \cdot k) \bmod 2^7] \gg (7-r)$$

→ Universal hashing: random $m = 2^r$

Provably good for any set of keys

$$h(k) = [(a \cdot k + b) \bmod p] \bmod m$$

random

large prime

How large should the table be?

- Want $m = \Theta(n)$ at all times
- But: we don't know at first what n will be!
- Our guess on m too small \Rightarrow slow (lots of collisions)
- ————— m too large \Rightarrow wasteful on space

Key idea: Start with table being small (e.g., $m=1$)
→ grow (& shrink) as necessary

How to grow/shrink a hash table?

Rehashing: → choose the new size m'
& new hash function h' (need a new hash fun. as size changed)

→ Build new hash table from scratch
insert each item from the old table
into the new one (& discard old one when done)

⇒ Time needed:

$$\underline{\Theta(n+m) = \Theta(n)} \quad \text{if } m = \Theta(n)$$

(This is NOT $O(1)$ time!)

How fast to grow the table?

Assume m reaches n (after an insert)

$\rightarrow m+1$?

\Rightarrow could need to rehash every step

\Rightarrow n inserts cost

$$\Theta(1+2+3+\dots+n) = \Theta(n^2) \gg n!$$

(Table doubling)

$\rightarrow m \times 2$? (Note: still $m = \Theta(n)$)

\Rightarrow rehash at each 2^i -th insertion

\Rightarrow n inserts cost

(really the next power of 2)

$$\Theta(1+2+4+8+\dots+n) = \Theta(n)$$

(Similar in spirit to build_heap [L4])

Result: A few inserts cost $\Omega(n)$ time
but $O(1)$ "on average"

Broader concept: Amortized analysis

- \rightarrow Common technique in DSs. Think: Like paying rent:
- \rightarrow Operation has amortized cost $T(n)$ if $\forall k \geq 1$, k operations cost $\leq k \cdot T(n)$

\Rightarrow " $T(n)$ amortized" \approx " $T(n)$ on average"; but averaged over all ops.

Note: Some of these ops might still take long time (but most don't!)

Back to hashing:

→ Can maintain $m = \Theta(n)$ under insertions
in $\Theta(1)$ amortized time

⇒ $\alpha = O(1) \Rightarrow$ search in $O(1)$ expected time
(under SUHA/univ. hashing
+ $O(1)$ time eval. of h)

Delete operations?

→ Already $O(1)$ expected time as is

→ But: m can get big wrt $n \rightarrow$ wasteful!
e.g.: $n \times$ insert, $n \times$ delete

→ Solution: When n decreases to $\frac{m}{4}$ (Any const. > 2 works)
shrink to half the size

⇒ $O(1)$ amortized cost to maintain
 $m = \Theta(n)$ for both inserts & deletions

But: Analysis harder → see CLRS 17.4

Exercise:
What goes wrong if we use $\frac{m}{2}$?

Sidenote : can make everything $O(1)$ time per op.
in an non-amortized sense, but a bit tricky

Python: The same approach gives us resizable
lists (arrays) in Python

⇒ `list.append` & `list.pop` in $O(1)$ amortized time

String matching problem:

Given two strings $S[1 \dots s]$ & $T[1 \dots t]$,
does S occur as a substring of T ?

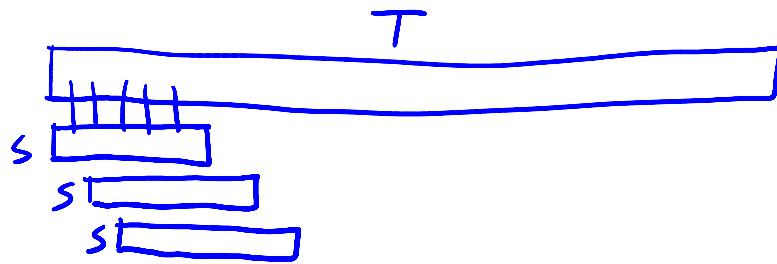
In other words: Does there exist a shift k ($0 \leq k \leq t-s$)
s.t. $\boxed{T[k+1 \dots k+s] = S[1 \dots s]}$? (s-character
match)

Think: $S \leftarrow$ pattern $T \leftarrow$ text

E.g. $S = 'G.006'$ & $T = \text{your entire INBOX}$

Naive algorithm:

Slide a length- s "window" through T ,
checking all ($t-s$) shifts for
a potential match with S



Running time?

$$O((t-s) \cdot s) = O(s \cdot t)$$

of possible shifts

Checking for a match each time

Potentially quadratic.



Can we do better?

Key bottleneck: Need to check for a match each time

Karp-Rabin algorithm:

Idea: Quickly "pre-screen" potential matches by comparing hashes, i.e.,

$$\text{check if } h(T[k+1 \dots k+s]) \stackrel{?}{=} h(S[1 \dots s])$$

Note: If hashes are different \Rightarrow no match

but if hashes are equal \Rightarrow there still might be NO match

"false positives"
due to collisions

\Rightarrow If hashes match, still need to verify for match

\Rightarrow But: probability p of "false positive" will be small for a suitable choice of h

linear time!

Running time?
 $(*)$

$$O(t(1 + p \cdot s) + s) = O(t) \text{ if } p \leq \frac{1}{s} !$$

Computing & dealing
Comparing hashes with false positives checking "final" match

Unfortunately: Above analysis cheats/abuses the model

Indication something is "fishy": If we take h to be a "trivial" identity hash function $h(x) = x$

Then: \rightarrow prob. of false positive is 0

\rightarrow K-R alg. becomes the naïve alg.

\rightarrow BUT: (*) would still indicate $O(1)$ _{run time}
(instead of $O(s)$)

Root of the problem: Can't really assume that computing $h(T[k+1..s])$ takes $O(1)$ time

Note: Size of $T[k+1..s]$ is s \leftarrow not a const.

\Rightarrow Computing $h(T[k+1..s])$ should take $\Omega(s)$ time

Resulting "real" running time is

$$O\left(+ \left(\underbrace{s}_{\text{"real" time}} + p \cdot s\right) + s\right) = \boxed{O(ps)}$$

"real" time
needed to compute
& compare hashes

No improvement!

What can we do now?

→ Use a hash function that can leverage the fact that every two consecutive sliding windows differ only by two characters

Rolling hash: Provides a quick way to update hash value when we slide the window by one position:

Given $h(T[k+1 \dots k+s])$ compute

$h(T[k+2 \dots k+s+1])$ in $O(1)$ time

How to do that?

Neat idea: Algebraization

→ Encode the window as a s-digit number
(base B, for $B >$ alphabet size)

$$x = d_1 d_2 \dots d_s \rightarrow \sum_{i=1}^s d_i B^{s-i} \quad h(x) \quad (\text{mod } q)$$

→ To keep sizes of #s reasonable use $\text{mod } q$

→ As long as q is some random prime $\approx s$

(one needs to prove that) \Rightarrow prob. p of false positives $\leq \frac{1}{s}$

6.006

3/3/16

p.10

Now:

$$y = h(x) = \sum_{i=1}^s d_i B^{s-i}$$

$$x = d_1 \dots d_s$$

→ To remove first character

$$x = 0d_2 \dots d_s \quad y \leftarrow y - d_1 B^{s-1} \pmod{q}$$

→ To shift everything left

$$x = d_2 \dots d_s 0 \quad y \leftarrow y \cdot B \pmod{q}$$

→ To specify new low-order digit l (after shift)

$$x = d_2 \dots d_s d_{s+1} \quad y \leftarrow y + l \pmod{q}$$

⇒ Updating a hash after a shift $O(1)$ time

Resulting running time of Karp-Rabin alg.:

$$O(s + (t-s)(1 + p \cdot s)) = O(t)$$

Compute
the initial hash
(& verify the
final match)

$(t-s)$ shifts

update
the rolling
hash

deal with
false positives