

6.006 Lecture 6: Hashing II

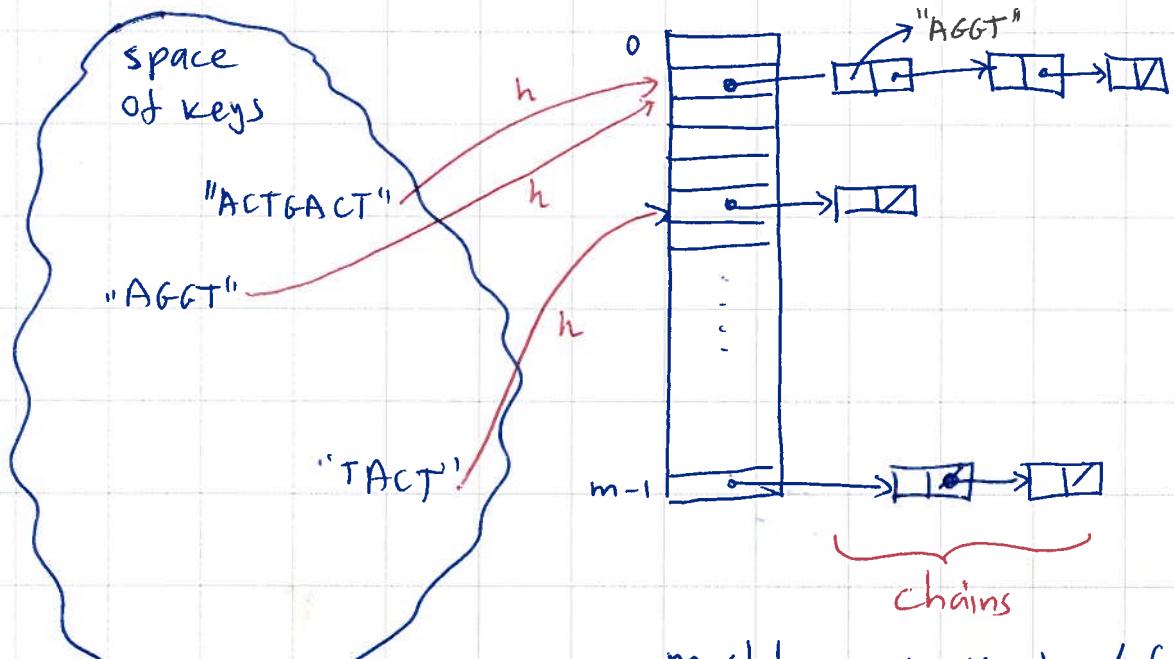
- Hash functions

11.3

- Resizing & amortized analysis (CLRS ch 17)

- Rolling hash

Recap of lecture 5:



m slots, n keys, load factor $\alpha = \frac{n}{m}$

Computing $h(x)$ (question ②)

Lots of ways. Here's a good one
for integer x 's (large integers)

Let p be a prime $p > m$

pick $a \quad 0 < a < p$

pick $b \quad 0 < b < p$

let $h(x) = ((ax + b) \bmod p) \bmod m$
only if $p > m$

e.g. $m = 1,000,000$

$p = 1,000,003$

$a = 314,159$

$b = 271,828$

} can reuse p, a, b for smaller m 's

If keys are not integers, convert them first
to integers

$x = "ATTGCTAC"$

treat as a base-4 integer

$x = "Boston"$

treat as a base-52 integer
or base-26
or base 128...

Textbook describes other methods.

Amortized analysis

Some insertions are really expensive (the ones that trigger doubling).

So instead of analyzing the worst case cost per operation, we ~~completely~~ analyze the total cost $T(n)$ for n operations, divided by n : amortized analysis.

A bit like analyzing the average case but without making any probabilistic assumptions.

Suppose we start with $m=5$

$$T(1) = 1 \xrightarrow{\text{insertion}}$$
$$T(2) = 1 + 1 \xrightarrow{\text{new insertion}} = 2$$

$$T(3) = 3$$
$$T(4) = 3 + 1 \xrightarrow{\substack{\text{new insertion} \\ \text{copying}}} + 4 = 8$$

$$T(5) = 9$$

$$T(6) = 10$$

$$T(7) = 11$$

$$T(8) = 11 + 1 + 8 = 20$$

$$\text{if } n=2^k : T(n) = n + \left(n + \frac{n}{2} + \frac{n}{4} + \dots + 8 + 4\right) \leq n + 2n = 3n$$

$$\text{so } T(n)/n \leq 3 : \text{worst case amortized cost per insertion is 3}$$

Resizing a hash table (and arrays in general)

We want $m = \Theta(n)$ at all times

m too small: large α , slow searching (long chains)

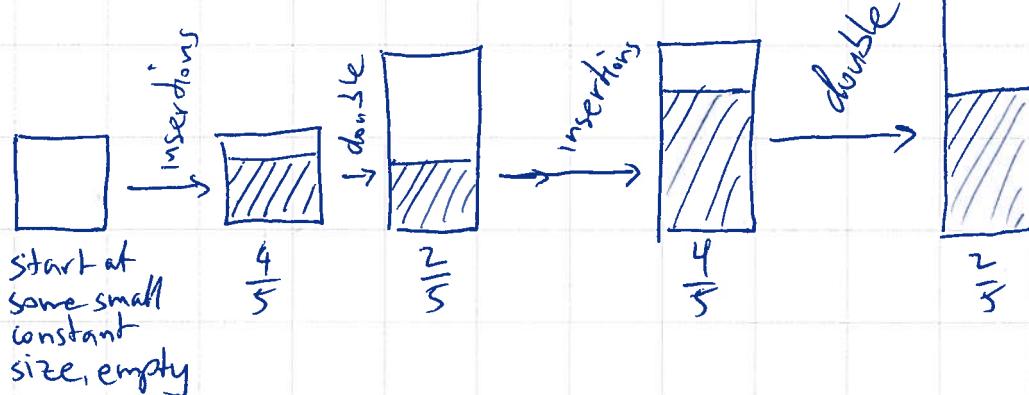
m too small: expensive to enumerate, initialize
wastes memory

How to set m if we do not know how large n
will be? (or how small)

Adjust m as n changes:

if $\alpha \geq 4/5$: double table size

now $\alpha = 2/5$



If no deletions, $\frac{2}{5} \leq \alpha < \frac{4}{5}$ or $n, m < \text{const}$

Another view of amortized analysis: every insertion deposits 3 "units of work" in a savings account, but ~~immediately~~ uses one unit to pay for the insertion.

When the account is full enough to copy entire table, do it (does not work exactly for our example starting from $m=5$, but can be made to work).

Deletions

If $\alpha \leq \frac{1}{5}$: halve table size

Amortized cost per operation is still ≤ 3

Example: deleting so n decreases from $\frac{2}{5}m$ to $\frac{1}{5}m$ leaves $2(\frac{1}{5}m)$ units in the account; this ~~pays~~ pays for copying the remaining $\frac{1}{5}m$ elements (and ~~leaves~~ ^{still} $\frac{1}{5}m$ in the account)

Why not halve the table when $\alpha < \frac{2}{5}$?

Rabin-Karp String matching & Rolling hashes

Given pattern $P[1..m]$ } lists of characters
text $T[1..n]$ }
does P occur in T ?

E.g. find "ATG" in "AATCGC..."

Idea:

- Compute $h(P)$
- for each length- m window of T , $T[i, \dots, i+m-1]$

A A T C G C
 \underbrace{ } \underbrace{ } \underbrace{ }
 \underbrace{ } \underbrace{ } \underbrace{ }
 ...

compute $h(T[i, \dots, i+m-1])$ and compare to $h(P)$

if $=$: check to see if really a match

if \neq : move on to next i

- Use a hash function h s.t. can compute $h(T[i, \dots, i+m-1])$ from $h(T[i-1, \dots, i+m-1])$ easily
 \Rightarrow rolling hash

Example of a rolling hash

Alphabet $\begin{matrix} 0 & 1 & 2 & 3 \\ \text{A} & \text{C} & \text{G} & \text{T} \end{matrix}$

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

$T = C T A T T A C G T$

$1 \ 3 \ 0 \ 3 \ 3 \ 0 \ 1 \ 2 \ 3 \ 4 \leftarrow \text{base 4}$

$1 \ 8 \ 4 \ 0 \ 9 \ 1_{10}$

$$453 \bmod p = 1009$$

\downarrow delete C

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

value of leading C = $1 \cdot 4^8 = 1 \cdot 65536 \bmod p$, precompute!

$$\begin{aligned} \text{so dropping leading C} \Rightarrow h &= 453 - 65536 \\ &= 502 \bmod p \end{aligned}$$

shifting and adding a G:

$$h = 4 \cdot 502 + 2 = 1001 \bmod p$$

$\begin{matrix} \text{shift} & \text{A} \\ & \text{G} \end{matrix}$

subtract +1

multiply +1

add +1

cost of moving \downarrow h to next window