

## 6.006 Lecture 5: Hashing I

- Naming data, fast docdist, direct addressing
- Hash functions - general idea
- 3 questions in the design of hash tables
- Chaining
- Analysis, simple uniform hashing analysis
- ~~An example of good hash functions~~
- Problem sets: how are students doing?
- Lab assistant: hours, give us feedback

## Forms of Naming

$L = []$

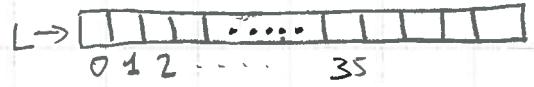
$L[0] = \text{True}$

:

if  $L[35]: \dots$

$\Theta(1)$

$\Theta(1)$



$\text{address}(L[35]) =$

$\text{address}(L) + 35 \cdot \text{width of elt}$

Naming elements of a list (array) using integers is very efficient, if the range of integers is not too large.

The ability to use arbitrary names and still be able to  $\text{insert(key,value)}$ ,  $\text{delete(key)}$ ,  $\text{search(key)}$  is very useful: in  $\Theta(1)$  time

`def count-frequency(word-list):`

$D = \{\}$  Python dictionary  $\equiv$  hash table

Search  $\rightarrow$  for word in word-list:

    if word in D:  $D[\text{word}] += 1$   $\leftarrow$  search, insert  
    else:  $D[\text{word}] = 1$   $\leftarrow$  insert

`def inner-product(D1, D2):`

$\text{sum} = 0$

for word in D1:

    if word in D2:

$\text{sum} += D1[\text{word}] * D2[\text{word}]$

$\Theta(n)$  if insert, search take  $\Theta(1) \Rightarrow$  optimal

## Hash Functions

Keys (names) that are natural for the application are not always small integers:

- set of length-20 strings over alphabet A,T,G,C
- credit-card #'s
- English words

PS2  $\uparrow$

Idea: define a "random-looking" function  $h$  from set  $U$  of keys to the set  $\{0, 1, \dots, m-1\}$  of indices of table  $T$ , store  $(x, v)$  in  $T[h(x)]$

Such a function is built into Python:  $\text{hash}(x)$  gives a 32-bit integer, and  $\text{hash}(x) \% m$  gives a number in  $\{0, 1, \dots, m-1\}$  ( $x$  must be immutable; more in recitation)

## Three questions

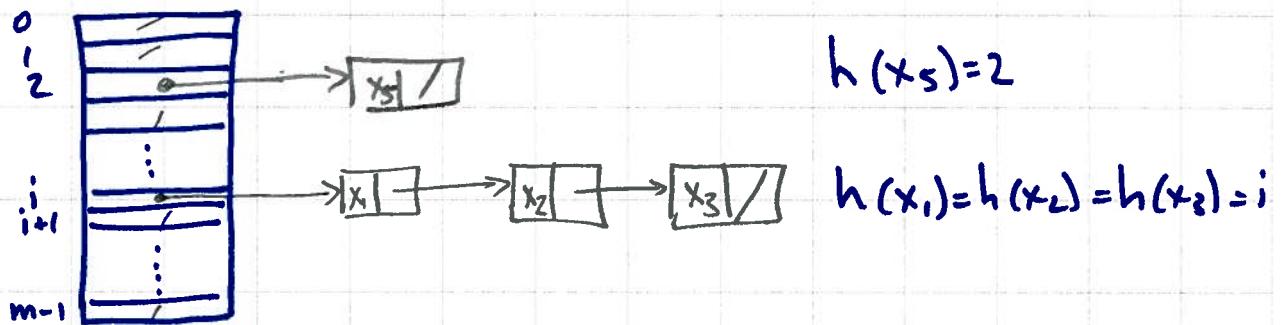
- ① How to choose  $m$ ?
- ② How to compute  $h(x)$ ?
- ③ How to insert both  $x$  and  $y$  if  $h(x) = h(y)$ ? (collision)

One approach for dealing with collisions: chaining

(another approach, open addressing, next week)

Hash table  $T[0..m-1]$

$T[i]$  is a list of elements that have  $h(x)=i$



$T[i]$  can be a Python list (array that Python grows dynamically) or a linked list.

## Analysis of Hashing with Chaining

n items in a table of size m

worst-case: all hash to position i

$T[i]$  has length n

search, delete takes  $\Theta(n)$  time

(insert is still  $\Theta(1)$ )

Let's demand less: expected time under  
an assumption

### Simple Uniform Hashing Assumption

Each key is equally likely to hash  
to any slot, independent of where  
other keys are hashed to

Load Factor  $\alpha = n/m = \text{average keys/slot}$

- time to do search/delete:  $\Theta(1+\alpha)$

↑ search list  $T[i]$   
compute  $i = h(x)$

-  $m = \lceil L(n) \rceil \Rightarrow \alpha = O(1) \Rightarrow \Theta(1+\alpha) = \Theta(1)$

- time to enumerate:  $\Theta(m+n)$

↑ ↑ enumerate all  
enumerate  $T$  the lists  
in  $T$

-  $m = O(n) \Rightarrow \Theta(m+n) = \Theta(n)$

- We want both  $m = \lceil L(n) \rceil$  and  $m = O(n)$

so we aim for  $m = \Theta(n)$ , ~~say~~

say  $\frac{n}{4} \leq m \leq 4n$

- This almost addresses question (1), but next time we'll see how to resize dynamically

## 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  $a$   $0 < a < p$

pick  $b$   $0 < b < p$

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.