

Lecture 24: hashing

- Map interface

- Hash tables

- Hash codes

- Chaining

Map (also a dictionary, lookup table)

store values associated w/ keys.

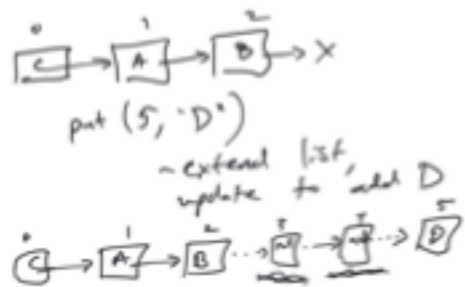
- put a new value with a given key
- check to see if my value assoc. w/ a given key
- look up value for a given key

Eg keys = words values = definitions.

Possible implementations:

- linked list

key: position in list (integer)
 value: whatever is stored there.



- heap?

look up element w/ given index?



doesn't help, looking still $O(n)$ have to search everywhere.

- array: if keys are small integers
 put (i, v) $a[i] = v$
 get (i) return $a[i]$.

extra work to check if $a[i]$ actually exists,

$\} O(1)$

$\} O(1)$

- Binary search tree

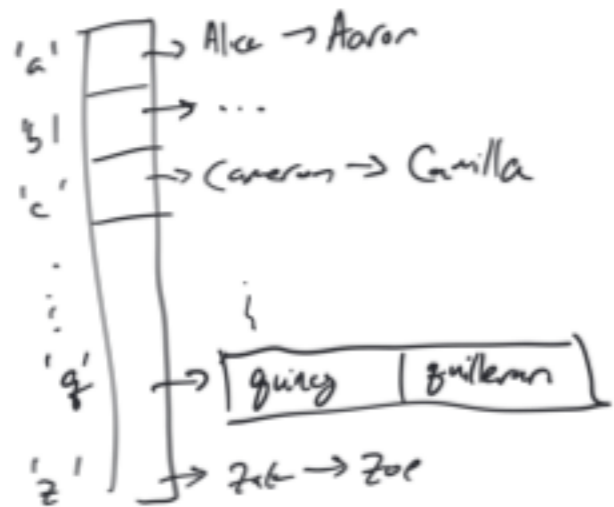
good for sparse, ordered keys

put: enter into tree $O(\log n)$

lookup: follow path to desired key, $O(\log n)$ (hard!)
assuming balanced tree.

Big simplifying assumption: assume data is perfectly "spread out"

ex: assume data are names, assume as many names starting with "A" as with "B"



n entries in table

N buckets

each bucket has 2 things, $\alpha = 2$

$n = 52$
 $N = 26$

(load factor: how full is my table?)

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

to look up "John", what do I do?

- go to John's bucket (bucket J)

$O(1)$ time

- search through J's, comparing each to "John"

By "magic" assumption: each bucket has about $\frac{n}{N}$ elements,

$O(\frac{n}{N}) = O(d)$ time.

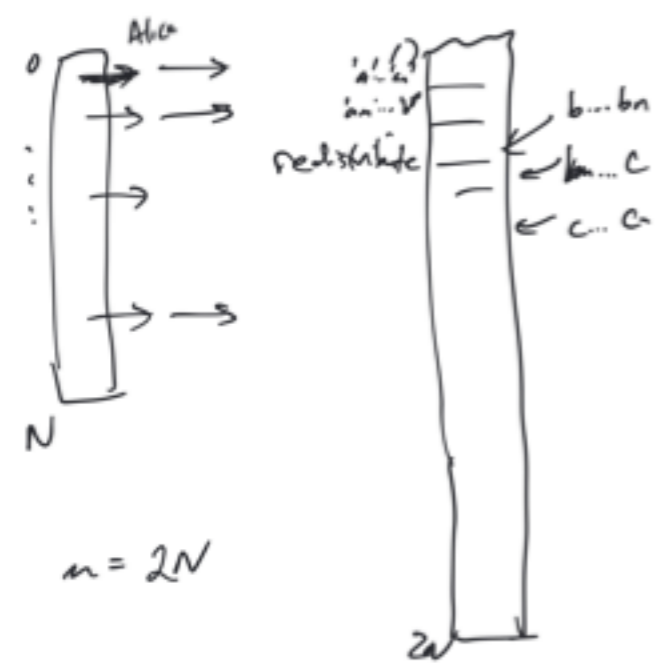
idea: let's make α constant

by increasing N as n increases.

Maintain $\alpha \leq 2$.

When adding an element,

if $n > 2N$
i.e. if $\alpha > 2$,
we'll double N .



$n = 2N$

- make a new array of buckets (size $2N$)
- copy everything from old array to new. (in correct bucket)

$O(n)$ for each elt of old table:
compute bucket in new tbl, $O(1)$
add elt to new bucket $O(1)$ } n times

$O(1)$ per element in table.

takes $O(n)$ time only after doing $2N$ insertions.

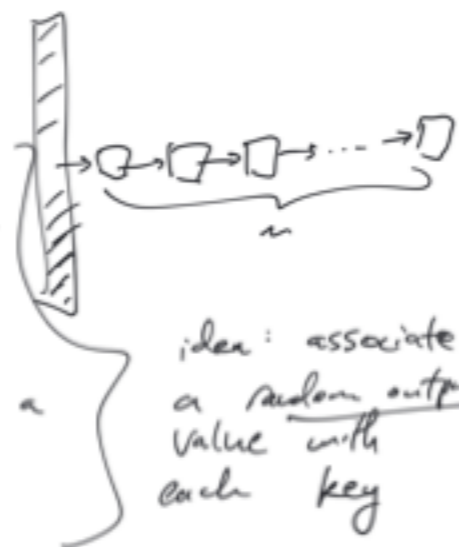
Can do n insertions in $O(n)$ time, even though in worst case, I need $O(n)$ for one insertion (this only happens once per n insertions)

Amortized Constant time

In practice, data is not uniformly distributed.
 Worst-case: all my entries are in same bucket, lookup takes $O(n)$ for every entry (not $O(1)$).

clustering: lots of data in the same bucket.

generalization: need to associate a "bucket" with each key in a general way.



idea: associate a random output value with each key. hash code.

Example: a "random" function assigns a number to each key

Alice \rightarrow 17

Karen \rightarrow 32

Alfred \rightarrow 6

⋮

Zak \rightarrow 0

Zoe \rightarrow 1000

↑
keys

↑
hash codes

- Can easily compute hash of any key.

- if two keys are equal, hashes should be equal.

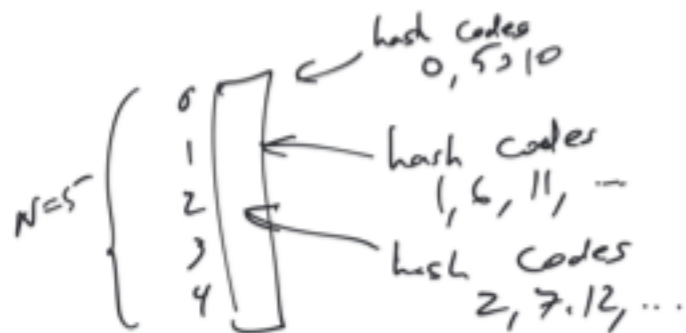
- should be "spread out": if two keys are different, they are unlikely to hash to same bucket (i.e. to have same hash code)

Hash table: array of N buckets,
each containing a list of entries.

- each entry in bucket b has a key whose hash code is $b \pmod{N}$

if we wrap around
i.e. remainder of
hashcode / N is
 b .

- No special structure
within buckets
(unsorted linked list)



Amortized
constant
time $O(1)$.

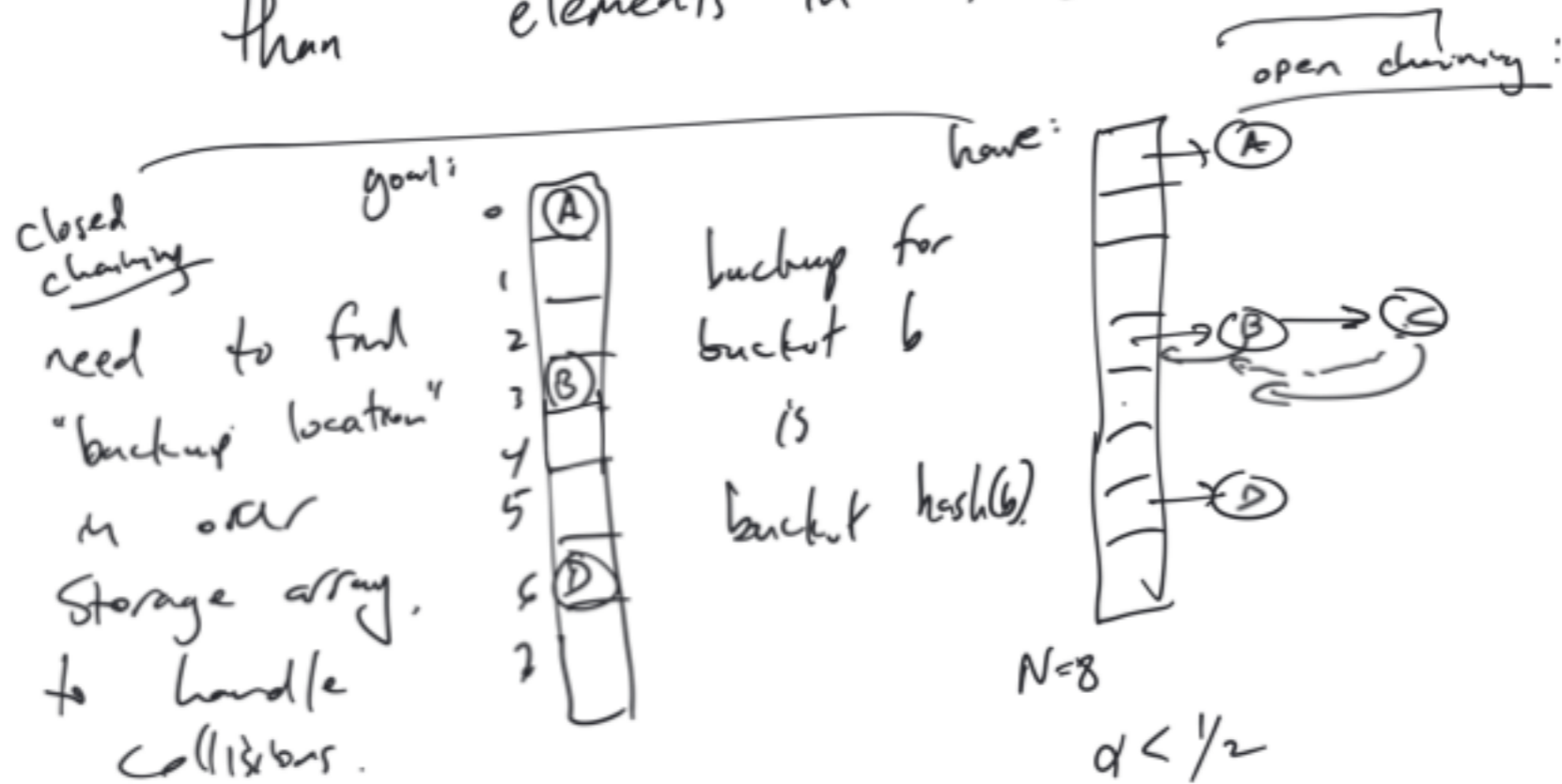
- to lookup key k : $O(1)$
 - hash k , take the remainder to find bucket b . if hash f^m is good.
 - linear search in bucket for key k . (average running time)

to insert value v with key k
• check if need to double size (i.e. if $n \geq 2N$)

- compute bucket by hashing k , taking remainder
- insert (k, v) in bucket b .

Closed chaining

if load factor is < 1 , there's more space in the array of buckets than elements in table



strategy for finding bucket location:

- if a slot is full, go to next slot.
- hash the bucket id to find "bucket" bucket