

Hash tables - how dictionaries are implemented

Basic idea

myMap["^{name}dave"] = ^{age}37

myMap["^{name}beccc"] = ^{age(?)}19



array
(fixed in size)

key → # which

tells you where to put it in the array

hash function:

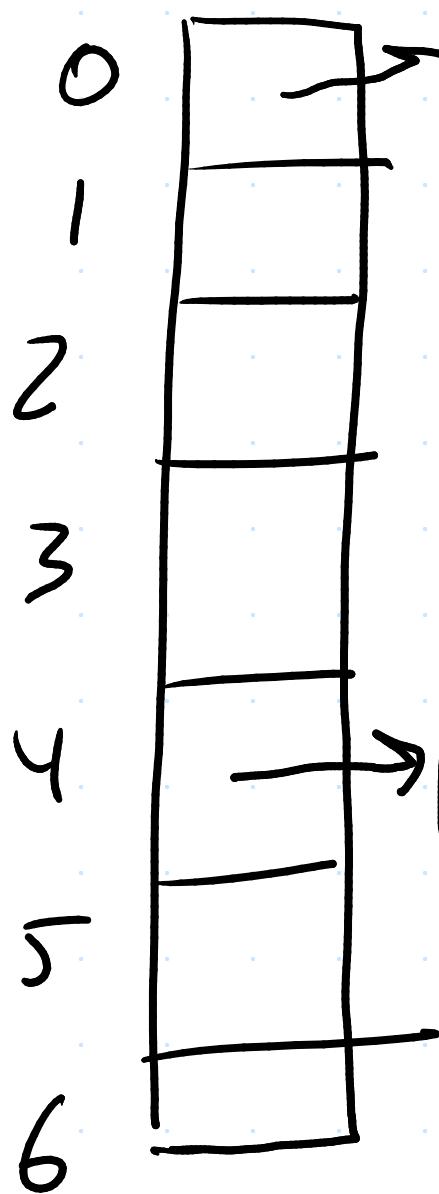
converts key to a number

e.g. on a string, maybe we add up the numeric codes for all letters in the word

(assume $a=1, b=2, \dots$)

$$\text{dave} \rightarrow 4+1+22+5 = \underline{\underline{32}}$$

$$\text{becca} \rightarrow 2+5+3+3+1 = 14$$



$$\text{my Map["dave"]} = 37$$

$$h("dave") = 32$$

hash code

Convert hash code
to a hash index,
which is in array
range

fixed size
array

$$h("dave") \% \text{ size of array}$$

$$32 \% 7 = 4$$

$$h("becca") \% \text{ size of array}$$

$$14 \% 7 = 0$$

How old is Dave?

println(myMap["dave"])

Recalculate:

h("dave") % size

$$32 \% 7 = 4$$

Look at spot 4 (instant)

-is Dave there? Yes!

Get the object, print
the age.

println(myMap["allyb"])

h("allyb") % size

$$1 + 12 + 12 + 25 + 2 \% \text{size}$$

$$52 \% 7 = 3$$

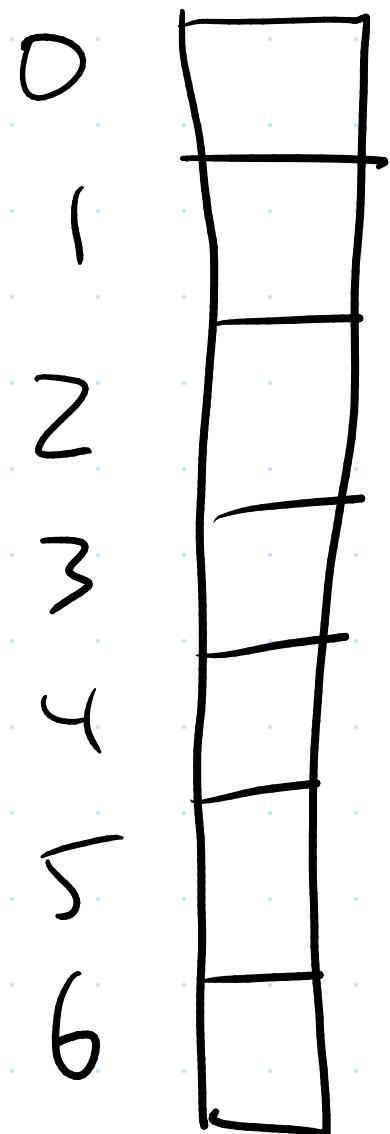
Look at loc 3, no allyb,
so print null (or something)

Lots to fix!

Collisions: when two keys land at same index.

myMap["dave"] "dave" \rightarrow 4

myMap["aj"] "aj" $1 + 10 = 11 \% 7 \rightarrow 4$

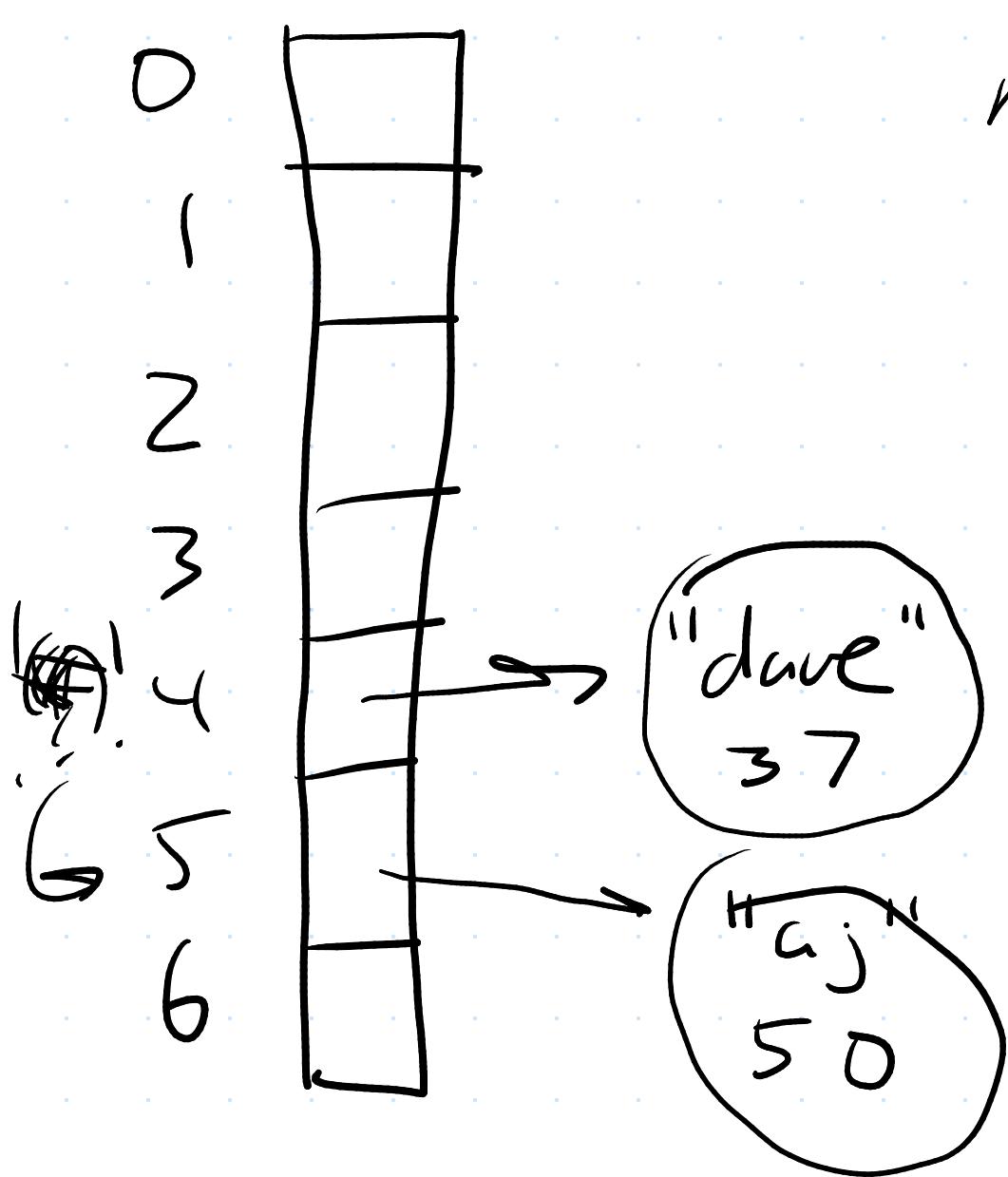


Two
Index
brood approaches
still just have
an array of
objects.

If a collision
jump to another
spot in the array,
and use that
(probing)

* ouch

Linear probing - if you collide,
go to next spot



myMap["dave"] = 37

myMap["aj"] = 50

Supposed to be
4, but
Collision
so go to
next one.

When you do a lookup, if you find something else than what you're looking for, try the next one, repeat, until you hit a blank spot.

To look up "aj"

↳ 4

Ugh, that's "dave"

So try next spot,
repeat as necessary until find
it, or hit a blank.

(; if hit bottom, circle around
to top)

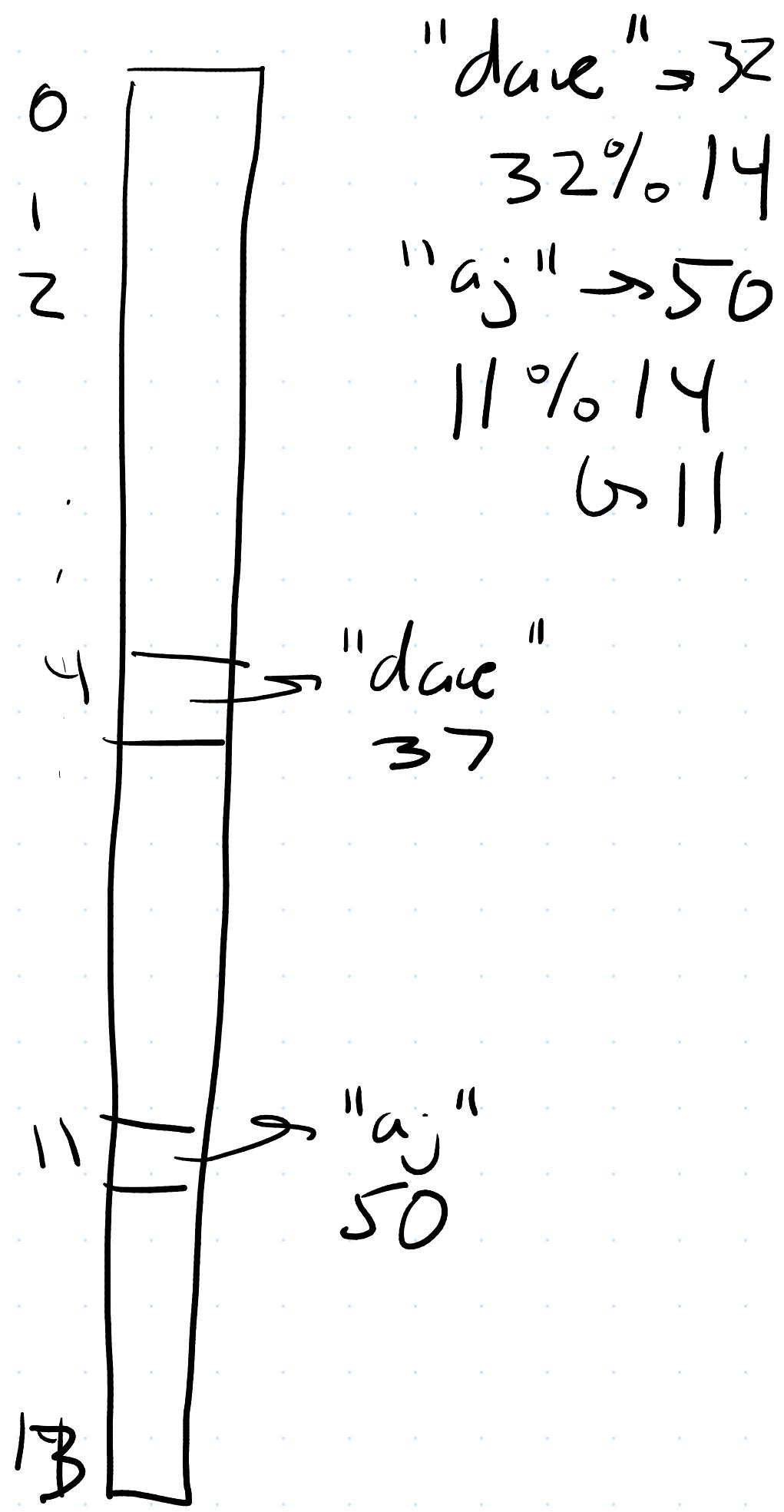
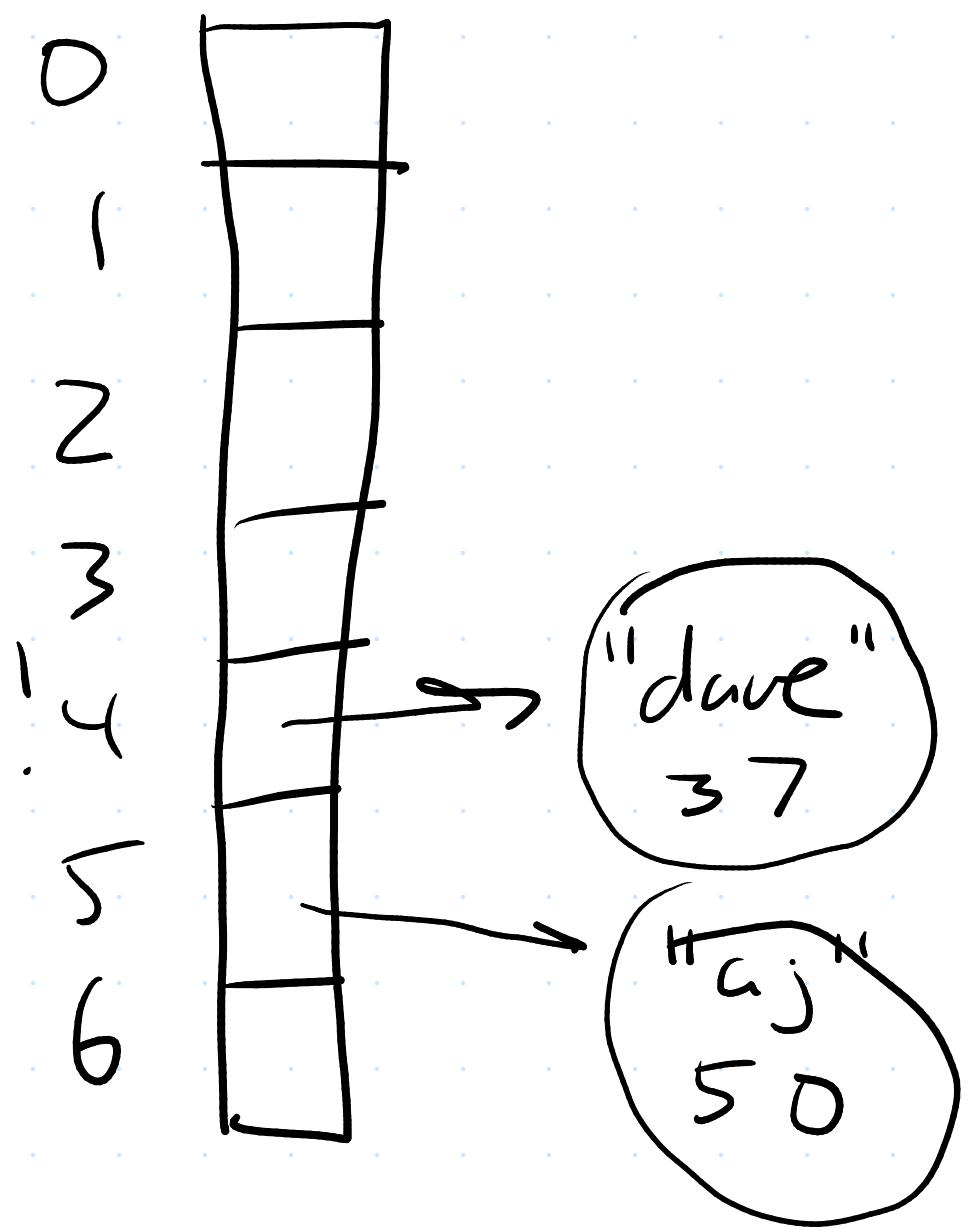
Over time, clusters develop,
and slow things down.

How do you avoid this?

- best solution is for table
to be large relative to the
number of entries.

Space vs time tradeoff

What you can also do is rebuild your hash table from scratch, if it gets filled than you like



too full?
→ build bigger
and rehash
(slow, but
hopefully not
very often)

There are other probing techniques that offer some improvements

Linear probing (if collision, go to next) results in a problem called primary clustering

→ any new entry that lands in that cluster adds to that cluster

Quadratic Probing: if collision,

try hash index ←

- hash index + 1 (linear probing)
 - hash index + 4
 - hash index + 9
 - hash index + 16
 - etc
- } wrap around
as necessary

This doesn't fix fundamental problems of collisions

Quadratic probing still results
in secondary clustering

↳ any entries that start at
the same hash index will
add to the same cluster.

Double hashing (probing)

- two hash functions constant

Sequence : hash index 1
 ↓
 hash index 1 + a hash index 2
 hash index 1 + 2a hash index 2
 hash index 1 + 3a hash index 2

These all help avoid clustering,
but it still happens when
of entries is high.
size of table

New groups (3 people)

Manager - keep everyone engaged

reader - reads each question loud

Scribe - writes answers for group

closed hashing

open addressing

"Approach I"

everything is an array



Approach II

chaining

(closed
addressing)

open hashing

