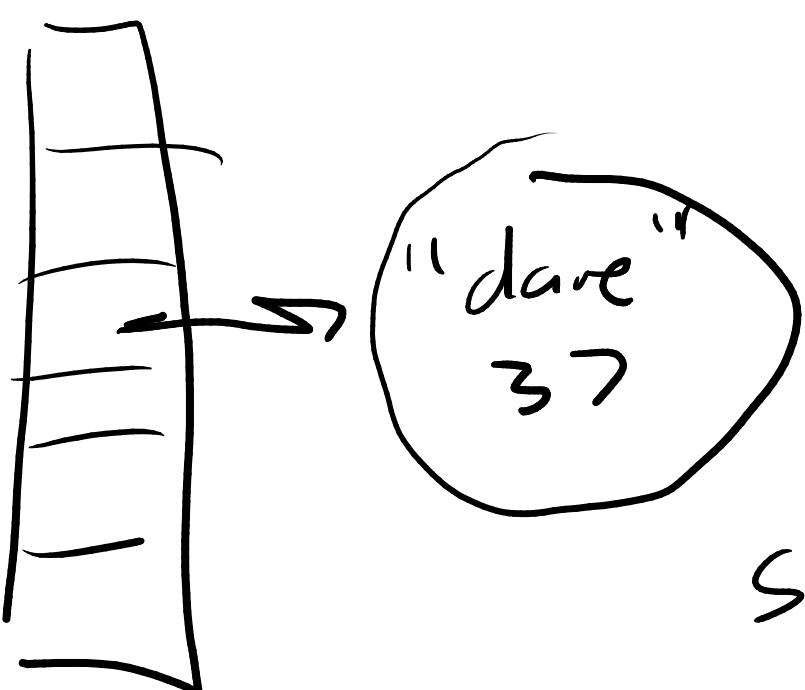


MVN clean

Hash tables

- hash functions (good ones)
 - hashing w/ chaining
 - Performance
 - assignment

↳ convert to a number
(hash function)



Want hash
function to
spread out your
keys (avoid
collisions)

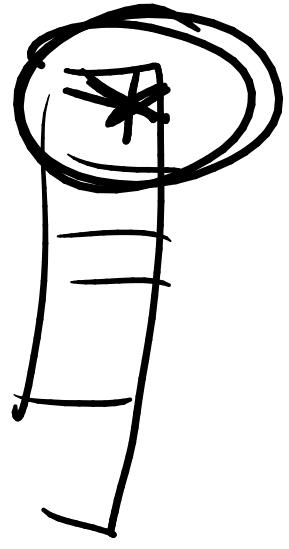
$h(\text{key}) \rightarrow \underline{\text{hash value}}$

$h(key) \rightarrow \underline{\text{hash value}}$ hash
 \sim hash value % size of array $\Rightarrow \underline{\text{index}}$

Worst hash function ever:

$$h(\text{key}) = 0$$

everything collides at 0.



Good ones:

Keys are integers:

$$h(\text{key}) = \text{key}$$

Most applications: this is fine
(but in some particular case,
could be otherwise).

Student ids would work fine

(but what if all student ids
ended in 5?)

In Kotlin, for built-in hash tables,
this is what it does.

Keys are strings:

Last week we used:

- take numeric value of first char
or
- add up numeric values of all chars
 ↳ better, because uses all data,
 but rearrangements of the
 same word collide

cat, tac give same hash
 Kotlin/Java code

if your strings $c_0 c_1 c_2 c_3 c_4 \dots$
 each char in string

hash value =

$$c_0 + c_1 * 3^1 + c_2 * 3^2 +$$

$$c_3 * 3^3 + c_4 * 3^4 \\ + \dots$$

(Horner's method)

Avoids collision for rearrangements
Why 31? Not a lot of great explanations

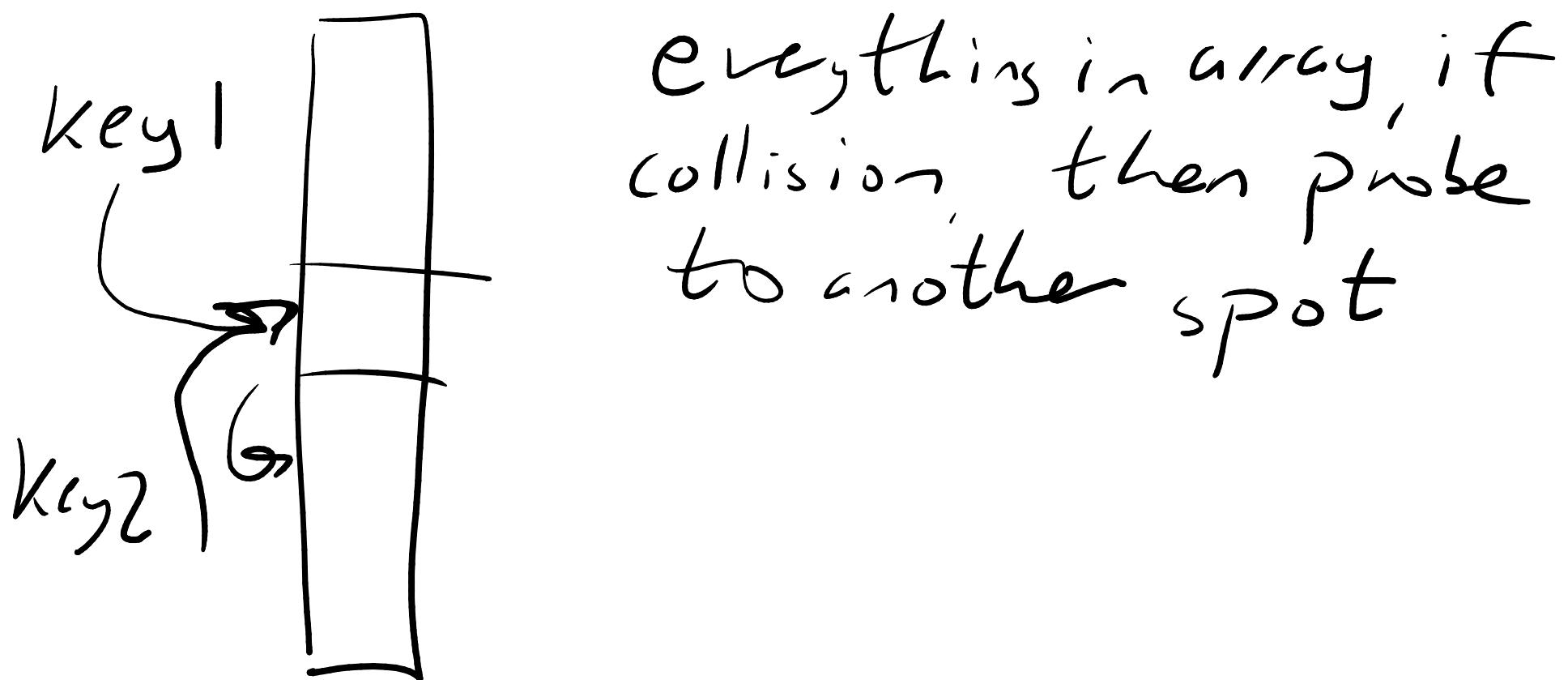
- 31 is prime, usually a helpful thing
- might be efficient to do arithmetic.

key is Double ... lots of crazy ideas

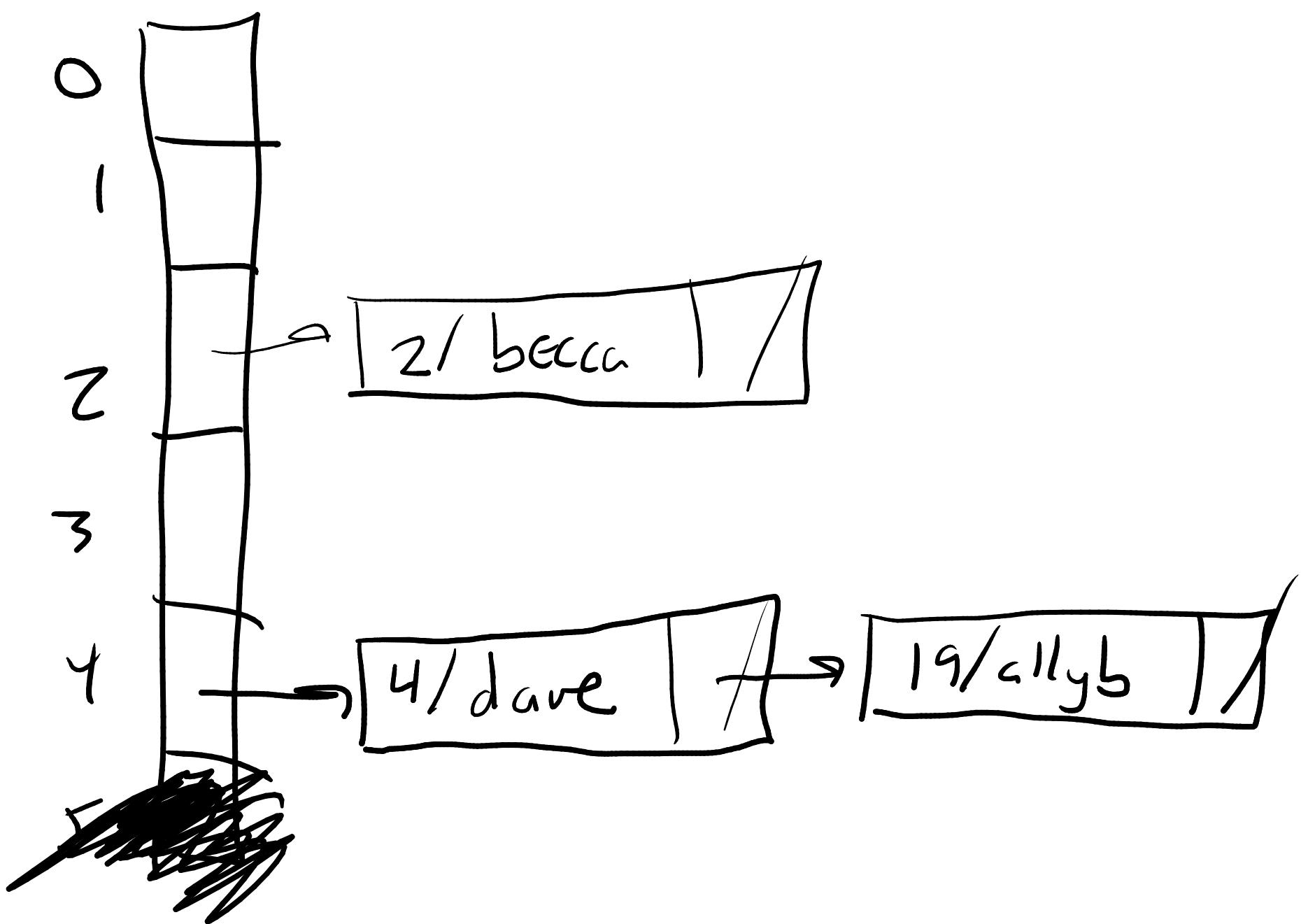
Chaining - alternative approach
+ the array

Last week, we did "hashing in an array"

- open addressing
-



Chaining: each spot in array is a linked list of entries



insert myMap[4] = "dave"

assume h(key) = key

myMap[2] = "becca"

myMap[19] = "ally b"

↳ goes to hash index 4

If a collision, just add on to linked list. No probing!

(in practice, add to beginning of list → faster)

To lookup:

print(myMap[19])

hash index: 4

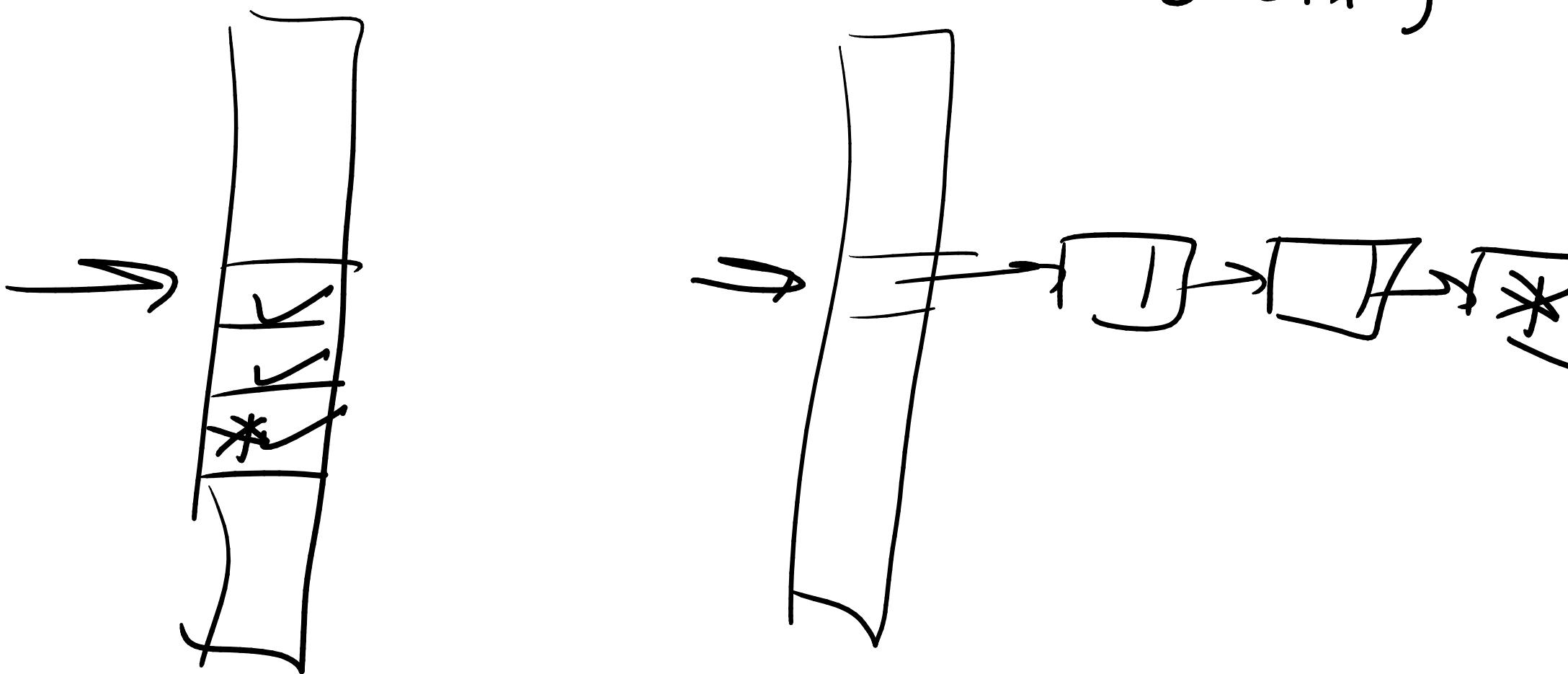
→ search linked list

Pros and cons of chaining vs all array:

Pros/cons:

- Memory tradeoffs: chaining takes slightly more memory (need addl pts for linked list)
- In chaining, worst case searches all collisions on a single hash index, but in an array, multiple hash index might mix in a cluster even worse
- chaining insertion guarantees $O(1)$ (at beginning of list), in an array does not
- chaining conceptually easier (?) (no probing algorithms)
- chaining can't fill up (but will slow down)

Searching in chaining is $O(\text{length of chain})$



Search in array is $O(\text{size of cluster})$

Performance

Worst case is (search, insert, etc)

$O(n)$ $n = \# \text{ of entries}$

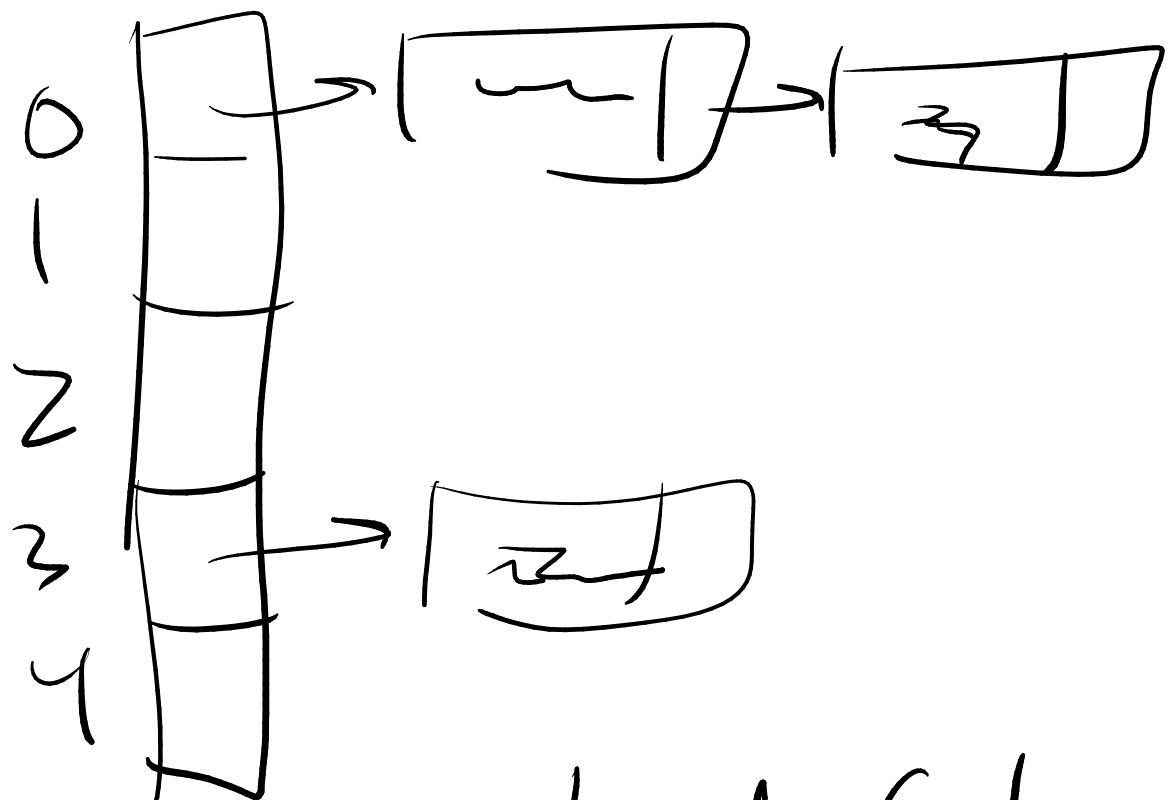
But, that should never happen
if you pick a reasonable
hash function

In general, if your hash table
is big enough, and your hash
function is good, it's
more like $O(1)$

More precise:

load factor = how full hash
table is

$$= \frac{\# \text{ of entries}}{\text{size of array}} = \lambda$$



$$\text{load factor} = \frac{\# \text{ of entries}}{\text{size of array}}$$

All in an array (open addressing)

\Rightarrow maxes out at 1

Chaining: \Rightarrow is unlimited

Performance for chaining:

Assume hash function works well -
evenly spreads things out.

Typical chain length = # entries = $\frac{\# \text{ entries}}{\text{size of array}}$

Typical lookup time: $O(\bar{\lambda})$

If you can keep $\bar{\lambda}$ small, it's fast!

How do you keep $\bar{\lambda}$ small?

$\bar{\lambda} = \frac{\# \text{ of entries}}{\text{size of array}} \rightarrow$ smaller
 \rightarrow bigger

Open addressing: theory is hard

Avg lookups for an unsuccessful search

$$= \frac{1}{2} \left(1 + \frac{1}{(1-\gamma)^2} \right)$$

(not obvious)

(approx for small γ)

As γ grows, the # of lookups grows

Cuckoo Hashing

bird - variation designed to make lookups guaranteed to be $O(1)$, no matter what at the cost of potentially slow insertions

Two hash tables each of which uses a different hash function.
(drawn horizontally)

Example, 2 tables of size $\underline{5}$

T1: $\underline{\quad}$ $\underline{\quad}$ $\underline{\quad}$ $\underline{\quad}$ $\underline{\quad}$

T2: $\underline{\quad}$ $\underline{\quad}$ $\underline{\quad}$ $\underline{\quad}$ $\underline{\quad}$

For T1: $h_1(\text{key}) = \text{key}$

T2: $h_2(\text{key}) = (\text{key} + 1)^2$

(not realistic hash functions, it's just for an easy example)

To insert, you start as usual on the first table

- Insert a 3. $h_1(3) = 3$
 $\text{index} = 3 \% \underline{5} = 3$

(only writing keys, assume values are there).

T1: — — — — 3 —

T2: — — — = —

Insert an 11

$$h_1(11) \% 5 = 11 \% 5 = 1$$

T1: 0 11 3 ←

T2: — — — = —

Insert a 13

$$h_1(13) \% 5 = 13 \% 5 = 3$$

collision!

When you have a collision, you kick out the value there, and put in new one:

T1: 0 11 13 → Kicked out 3

T2: — — — — —

The 3 which got kicked out,
we now try to insert in second
table

$$h_2(3) \% 5 = (3+1)^2 \% 5 \\ = 16 \% 5 = 1$$

T1:	0	1	2	3	4
T2	0	3	2	3	4

To do a lookup

"Find key 3"

- Try first table. If don't succeed, try second table.

On insert, if you collide on second table, kick out again, and try to add first using first table's hash function.

Insert: 23

Try T1:

$$h_1(23) \% 5 = 3$$

T1:

$\begin{array}{c} \text{||} \\ \hline 0 & 1 & 2 & 3 & 4 \end{array}$

kicked
out
13

T2

$\begin{array}{c} 3 \\ \hline 0 & 1 & 2 & 3 & 4 \end{array}$

Try to insert 13 in table Z

$$\begin{aligned} h_2(13) &= (13+1)^2 \% 5 \\ &= 14^2 \% 5 \\ &= 196 \% 5 = 1 \end{aligned}$$

This will break a cycle, but cycles don't always happen.

If you end up in a cycle, then
you need to start all over with
new hash functions.