

THE PROBLEM.

Most of the data we care about come in logical pairs of a (key, value)

- (home Address, Person)
- (UserID, Facebook Profile)
- (Trading Book ID, Book Metadata)
- (course code, (210, marks))
- Cryptography (identification)
- Systems (Virtual Memory)
- Bioinformatics (DNA search)

How can we allow for efficient inserts and lookups given a "key".?

Q: How could we represent this?

LINKED LIST

TREES.

GRAPHS?

HEAPS.

lookup? $O(n)$	$O(\log n)$	$O(V+E)$	$O(\log n)$
insert? $O(1)$	$O(\log n)$	$O(1)$	$O(\log n)$

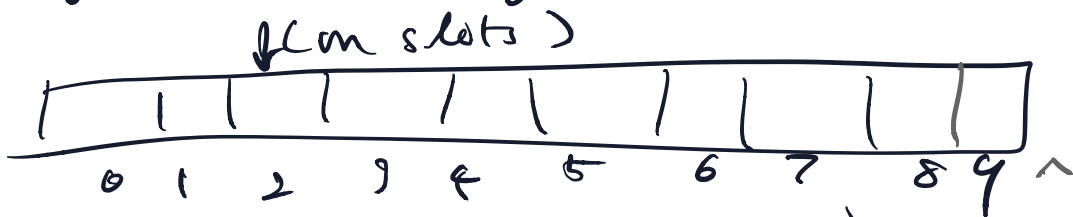
Goal?

LOOKUP: $\Theta(1)$ (amortized)

INSERT: $\Theta(1)$ (amortized)

Idea: Say we have an array of m slots.

①



deterministic? (key, value)

②

We have a function $h: \underline{K} \rightarrow \mathbb{Z}_m$ that is computed in $O(1)$ time.

③

For (key, value) to insert we go by rule.

$$A[h(\text{key})] = \text{value}.$$

④ For lookups, given a key where do I look?

$O(1)$

insert

$O(1)$

lookup.

PROBLEM!

$$|K| = |\Pi_i|$$

$$|K| > |\Pi_i|$$

"injecting"

Handling collisions : $(key \in \mathbb{Z}, char^* value) \cdot (k, v)$

e.g. STUDENT RECORDS

ZID

NAME

$$h(5162945)$$

$$h(4201337)$$

$$h(4)$$

"ADRIAN R. MARTINEZ"

"DORITOS X. MTNDEW"

"A"

h

$k(2)5$
slots



TAKEAWAY:

Set of keys does not have same cardinality as the range of our hash h .

Separate chaining

$$h_1(x) = x \bmod 7$$

where $m = 7$

\Rightarrow there must exist a $k_1, k_2 \in \mathbb{Z}$ s.t.

$$h(k_1) = h(k_2)$$

aka. collisions!

Q1. What is the worst case when inserting k (key, value) pairs into our hash table? Imagine our table has size $n/2$ and we measure search cost by NUMBER of key comparisons.

If we assume keys are evenly distributed what's average cost?

DOUBLE HASHING.

n = num keys stored
 m = num slots in table

$$\Omega\left(\frac{n}{m}\right)$$

load factor $\alpha = n/m$ = # keys per slot. lower bound of $\Omega(\alpha)$.

Idea: Have two hash functions h_1, h_2 and if we reach a collision add on $h_2(k)$.

Pseudocode:

```
base =  $h_1(k)$ , shift =  $h_2(k)$ 
while (no collision on  $A[\text{base}]$ ):
    base += shift.
insert  $(k, v)$  in  $A$  at base.
```

Q2. consider inserting the sequence

~~11~~ ~~16~~ ~~27~~ ~~35~~ 22 20 15 24 29 19 13

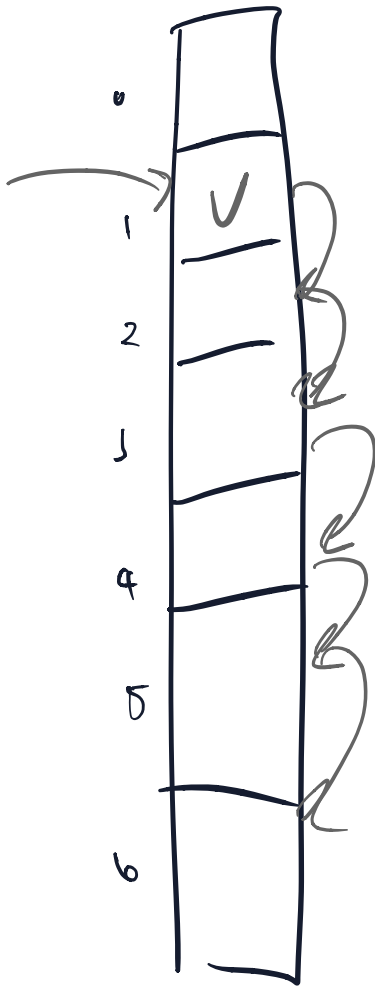
into an empty table with $h_1(x) = x \% 11$, $h_2(x) = (x \% 3) + 1$

why?

11	v''	35	v'''	v	16	27	..	v'	...	v'
0	1	2	3	4	5	6	7	8	9	10

LINEAR HASHING.

Idea: use $h(k)$ to define starting index, then +1 the index cell to as our point to insert our value.



$$h(x) = x \bmod 6$$

- ① insert(8, "hello")
- ② insert(4, "hey").
- ③ insert(17, "whatsup").

$$h_2(x) = 1$$

PROBLEMS

Why am I learning all three, when I only need one?

① Easy / Pragmatic

2. Using 2 complicated ↓ ↓ number
collisions

$h_1 \sim h_2$ ↑ distribution very. Cache