Lecture 7
6.046

<mark>Symbol - table problem (compilers)</mark>

Table S holding records

<mark>state set is only lookups, etc...</mark>



record

satellite data

<mark>Operations:</mark>
- insert (S,x): $S \leftarrow S \cup \{x\}$
- delete (S,x): $S \leftarrow S - \{x\}$
- search (S,k)

} dynamic set

return x ■ s.t.
key[x] = k  or
nil if no such x

<mark>Direct access table</mark>

Suppose keys are drawn from $u = \{0, 1, \dots, m-1\}$
Assume keys are distinct
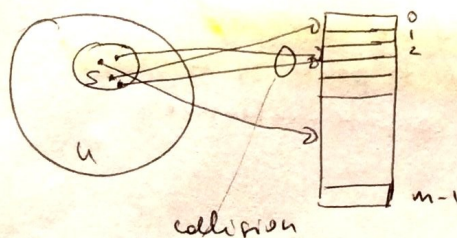Set up array $T[0 \dots m-1]$ to represent dyn. set S

$$T[k] = \begin{cases} x & \text{if } x \in S \text{ and } key[x] = k \\ nil & \text{otherwise} \end{cases}$$

Ops take $\Theta(1)$ time

e.g. 64 bit #s drawn, need a table of size $2^{64}$
↳ need hashing

<mark>Hashing</mark>

Hash function h maps keys "randomly" into slots of table T.



collision

when a record to be inserted maps to an already occupied slot, a <u>collision</u> occurs.

<mark>Resolving collisions by chaining</mark>
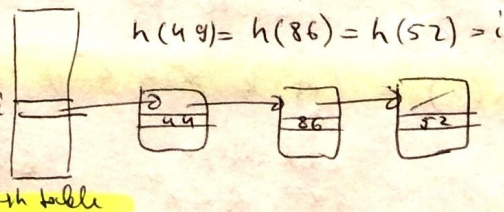idea: link records in same slot into list.

$h(49) = h(86) = h(52) = i$



hash table

<mark>Analysis</mark>

<u>Worst case</u>: every key hashes to the same slot. → long link list
access takes $\Theta(n)$ time if $|S| = n$

<mark>Average case</mark>
<mark>assumption of simple uniform hashing</mark>
- each key is equally likely to be hashed to any slot in T, independent where other keys are hashed

Def. the load factor of a hash table with n keys and m slots is $\alpha = n/m$ = ▮ ave # keys per slot.

Expected unsuccessful search time : $\Theta(1 + \alpha)$

Expected search time = $\Theta(1)$
if $\alpha = O(1)$, i.e. $n = O(m)$

↑ cost of hash access to slot  ↑ cost of lost search

if n grows
m has to grow

Represent a dynamic set with const ▮ time ops as long as:
- $n = O(m)$
- simple uniform hashing

Choosing a hash function :
- Should distribute keys uniformly into slots
- Regularity in key distribution should not affect uniformity

Division method

$h(k) = k \bmod m$  ← does not have to be eq. to size of table, can be around the size (len)
- don't pick m with small divisor d

Ex. d = 2 and all keys even
=> odd slots never used

Ex  $m = 2^r$ => hash does not depend on all bits of k

$k = 1011\ 0001\ 11\underline{011010}$  $r = 6$  $m = 2^6$

$h(k)$

h does not depend on other bits

Pick m = prime not too close to power of 2 or 10

However, the division method is compute-intensive.

Multiplication method :

$m = 2^r$, computer has w-bit words

$h(k) = (A \cdot k \bmod 2^w)\ \text{rsh}\ (w-r)$

↑ odd integer $2^{w-1} < A < 2^w$    ↑ right shifted by bitwise

- Don't pick A too close to $2^{w-1}$ or $2^w$
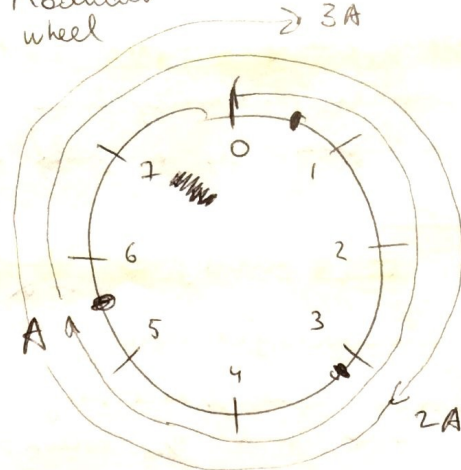- Fast method : mult. mod $2^w$ faster than division
  rsh is fast

Ex: $m = 8 = 2^3$, $w = 7$
← full word

$$1011001 = A$$
$$1101011 = k$$
$$\overline{1001010011\ 0011}$$

high-order ↑ two-word product
bits ignored
with mod $2^w$

$$011\ 0011$$
$\underbrace{\quad}_{h(k)}$ $\xrightarrow{\quad}$ rsh removes
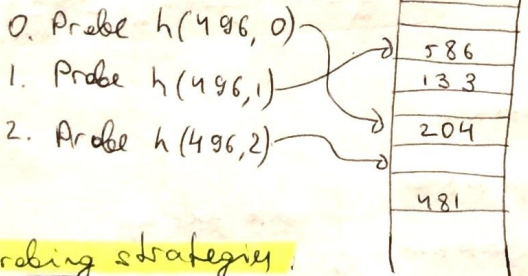4 bits

Modular wheel



$\to$ 3A

A ↑

2A

k spins the wheel...
↳ heuristic

---

## Resolving collisions by open addressing

- No storage for links
- Probe table systematically until an empty slot is found ( hash function after hash function)
- $h: U \times \{0, 1, \dots, m-1\} \to \{0, 1, \dots, m-1\}$

  universe ↑ probe # ← index of a
  of keys hash function

- Probe seq. should be permutation ← given a key, the sequence of slots hit, no repeating
- Table may fill up, $n \leq m$
- Deletion is difficult

Ex: **Insert** $k = 496$

0. Probe $h(496, 0)$
1. Probe $h(496, 1)$
2. Probe $h(496, 2)$

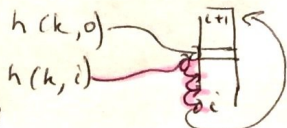| |
|---|
| |
| 586 |
| 133 |
| |
| 204 |
| |
| 481 |

**Search** - same probe sequence
- successful - find record
- unsuccessful - find nil

**My comment:**
- need to save the ID of the hash function together with the $k$

## Probing strategies

- Linear - $h(k, i) = (h(k, 0) + i) \bmod m$
  "primary clustering" - long runs of filled slots

  $h(k, 0)$
  $h(k, i)$

- Double hashing - $h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m$
  excellent method, usually pick $m = 2^r$ and $h_2(k)$ odd

Analysis of open addressing  ← stronger than simple uniform hashing

==Assumption of uniform hashing==:

  each key is equally likely to have any one of the
  $m!$ perms as its probe seq., indep. of other keys

## Theorem

$$E[\# \text{ probes}] \leq \frac{1}{1-\alpha} \quad \text{if } \alpha < 1 \quad (\text{i.e. } n < m)$$

Pf. (unsuccessful search):

  1 probe always necessary
    with prob. $n/m$ collision $\Rightarrow$ 2. probe is necessary

  2. probe prob. of collision $\frac{n-1}{m-1} \Rightarrow$ 3. probe is necessary

  3. probe prob. of collision $\frac{n-2}{m-2}$ $\cdots$ —

  Note: $\frac{n-i}{m-i} < \frac{n}{m} = \alpha$ for $i = 1, 2, \ldots, n-1$, assumption
  $n < m$ (necessary for open addressing)

$$E[\# \text{ probes}] = 1 + \frac{n}{m}\left(1 + \left(\frac{n-1}{m-1}\right)\left(1 + \frac{n-2}{m-2}\left(\cdots \left(1 + \frac{1}{m-n}\right)\right)\cdots\right)\right)$$

$$\leq 1 + \alpha\left(1 + \alpha\left(1 + \alpha\left(\cdots 1 + \alpha\right)\cdots\right)\right)$$

$$\leq 1 + \alpha + \alpha^2 + \alpha^3 \cdots$$

$$= \sum_{i=0}^{\infty} \alpha^i = \frac{1}{1-\alpha}$$

$1 + \frac{n}{m} + \frac{n}{m}\left(\frac{n-1}{m-1}\right) + \frac{n}{m}\left(\frac{n-1}{n-1}\right)\left(\frac{n-2}{m-2}\right)\cdots$

$\alpha < 1$ const $\Rightarrow O(1)$ probes in expectation

  – Table 50% full $\Rightarrow$ 2 probes in expectation
  – Table 90% full $\Rightarrow$ 10 probes in expectation

  | cost goes up
  ↓ need to keep $\alpha$ low