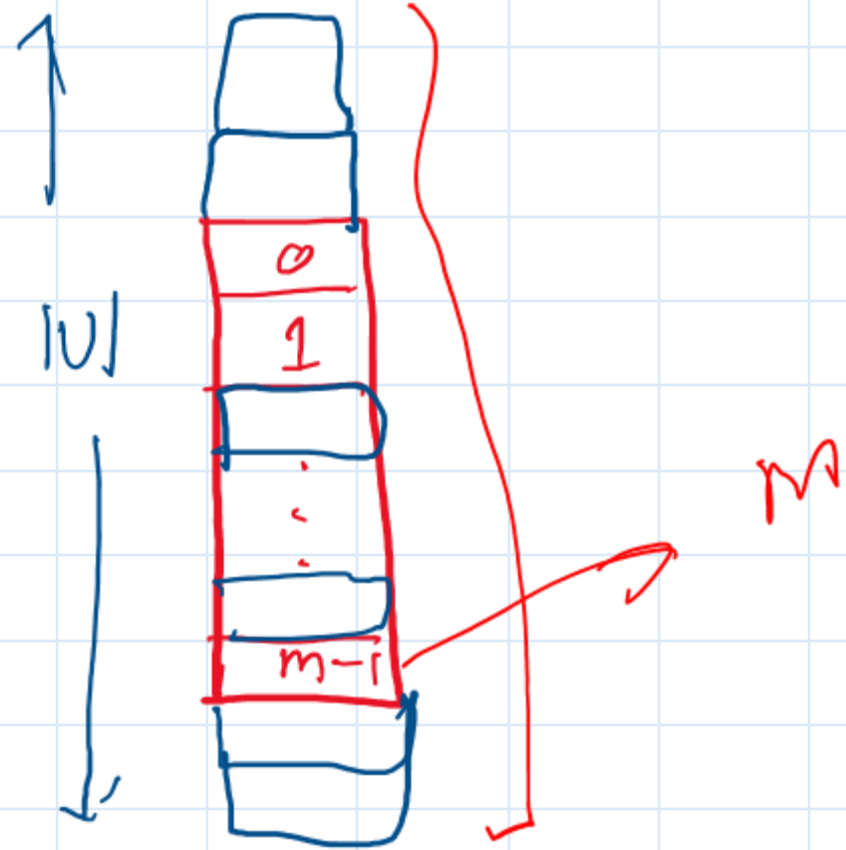
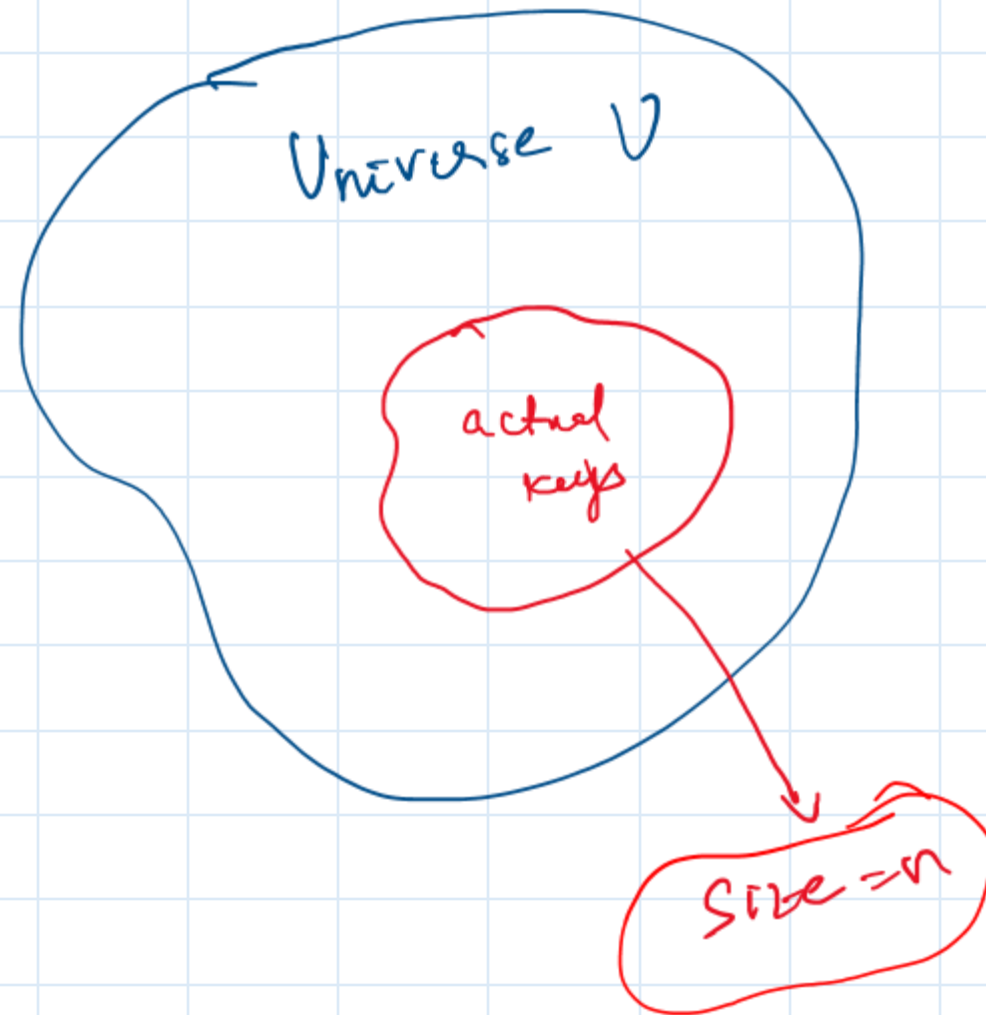


COL106 - Data Structures and Algorithms

Hashtables (contd), other related ADTs

Hash Tables



Hash functions = Hash code + Compression map
key \rightarrow integer integer $\rightarrow [0, m-1]$

Hash Code:

- Integer cast (interpret bits of key as integer)
- Component sum (partition bits of key into components & add up)
- Memory address (interpret memory address as integer)

- Polynomial accumulation (partition key to components a_0, \dots, a_{n-1} of fixed length & evaluate polynomial)

$$p(z) = a_0 + a_1 z + \dots + a_{n-1} z^{n-1}$$

Horner's rule:

$$p(z) = a_0 + z(a_1 + z(a_2 + \dots + a_{n-1}))$$
$$p_0(z) \leftarrow a_{n-1}, \quad p_i(z) \leftarrow a_{n-i-1} + z p_{i-1}(z)$$

$z=23$ gives
 ≤ 6 collisions
on a set of
50,000 English words

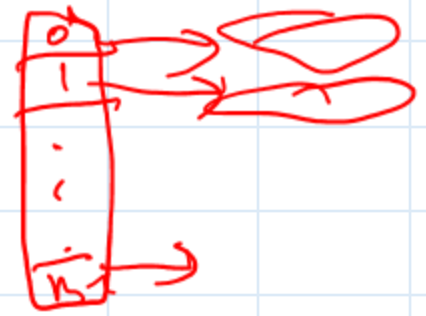
Compression Maps

①

$$h(k) = k \pmod{m} \quad (\text{simple division})$$

choosing $m \neq b^e$ (bad!)

Typically m is chosen to be prime - gives a uniform distribution.
not close to a power of 2.



Compression Maps

① $h(k) = k \pmod{m}$ (simple division)

choosing $m \neq b^e$ (bad!)

Typically m is chosen to be prime - gives a uniform distribution.

not close
to a power of 2.

② $h(k) = \lfloor m (kA \pmod{1}) \rfloor$ where $0 < A < 1$

Eg: $A = \frac{\sqrt{5} - 1}{2}$

(Fibonacci hashing)
(Knuth TAOCP vol. 2)

Compression maps

③

"MAD" (Multiply, Add & Divide)

$$h(k) = |ak + b| \bmod m$$

↓
not a multiple of m

(often used in
pseudo random number
generators)

Universal Hashing:

$$\mathcal{H} = \{h: U \rightarrow \{0, \dots, m-1\}\}$$

\mathcal{H} is universal if for any randomly chosen h from \mathcal{H} ,
& keys k_x, k_y

$$\Pr[h(k_x) = h(k_y)] \leq \frac{1}{m}$$

There are
efficient
constructions of
UHF

Handling collisions

$$h: U \rightarrow \{0, 1, \dots, m-1\}$$

Open addressing:

In chaining all elements were stored outside the hash table.

- Put all the elements in the table $\Rightarrow n \leq m$
- Have a systematic way to probe elements of the table.

$$h(x) = 1$$
$$h(y) = 1$$

Modify hash function: $h: U \times \{0, 1, \dots, m-1\} \rightarrow \{0, 1, \dots, m-1\}$

Probe sequence:

$$\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$$

specify the probe number

- h gives a sequence of slots examined for a given key.

Linear Probing

(uses less memory than chaining, but slower)
causes clustering

Example:

$$h(x) = x \bmod 13$$

Insert: 18, 41, 22, 44, 59, 32, 31, 72

$n \bmod 13$
 $n \bmod 13$

| | | | | | | | | | | | | |
|---|---|----|---|---|----|----|----|----|----|----|----|----|
| | | 41 | | | 18 | 44 | 59 | 32 | 22 | 31 | 72 | |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

Linear Probing

Example:

$$h(x) = x \bmod 13$$

Insert: 18, 41, 22, 44, 59, 32, 31, 73

| | | | | | | | | | | | | |
|---|---|----|---|---|----|----|----|----|----|----|----|----|
| | | 41 | | | 18 | 44 | 59 | 32 | 22 | 31 | 73 | |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

Red arrows point from the value 32 at index 8 to indices 9, 10, and 11, indicating the probing sequence.

Delete 32

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|----|----|----|
| | | | | | | | | | | | | |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

Linear Probing

$$ax + b \text{ mod } 10$$

Example:

$$h(x) = x \text{ mod } 13$$

$$n \leq m$$
$$h(x) = 12$$

Insert: 18, 41, 22, 44, 59, 32, 31, 73

| | | | | | | | | | | | | |
|---|---|----|---|---|----|----|----|---------------|----|----|----|----|
| | | 41 | | | 18 | 44 | 59 | 32 | 22 | 31 | 73 | |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

Delete 32

| | | | | | | | | | | | | |
|---|---|----|---|---|----|----|----|---|----|----|----|----|
| | | 41 | | | 18 | 44 | 59 | X | 22 | 31 | 73 | 74 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

↑ Try new
↑ Defunct
↑

lookup → ignore X
insert → replace X

Rehash if there are too many X's

Double Hashing

Quadratic probing:
 $(h(k) + i^2) \bmod m$

Use two functions: primary hash $h(k)$
secondary hash $d(k)$

- table size m
must be prime.

↓
cannot take 0 values.

handles collision
by placing item in
first available cell
in

$$(i + j d(k)) \bmod m$$
$$j \in \{0, 1, \dots, m-1\}$$

DoubleHashInsert(k)

if (table is full) error

probe = $h(k)$; offset = $d(k)$

while (table[probe] occupied)

probe = (probe + offset) mod m

table[probe] = k

Double Hashing Example

$k \geq V \times \underbrace{\{0, 1, \dots, m-1\}}_{\rightarrow \{0, \dots, m-1\}}$

$$m = 13$$

$$h(k) = k \bmod 13$$

$$d(k) = (7 - k) \bmod 7$$

$$h(k)$$

Insert: 18, 41, 22, 44, 59, 32, 81, 72

Hand-drawn array diagram with 13 slots. Slots 2, 5, 9, and 10 are highlighted with red boxes and contain the values 41, 18, 22, and 44 respectively. Slot 5 is also circled in red.

$\Rightarrow f_{\text{fix}} = d(f) = 5$

2

Analysis of Double Hashing

- let α (load factor) be < 1
- Assume every probe looks at a random location in the table
- $1 - \alpha$ fraction of table is empty.
- Expected # probes to find an empty location

$$= \frac{1}{1 - \alpha}$$

$$L \approx \frac{\alpha}{1 - \alpha}$$

$$\uparrow \frac{1}{1 - \alpha}$$

required for unsuccessful search.

Analysis of Double Hashing

- Average number of probes for a successful search
= Average number of probes required to insert all elements

- To insert an element we need to find an empty location.

$$\{1, \dots, m/2\}$$

$$\{\frac{m}{2}+1, \dots, \frac{m}{2}+\frac{m}{4}\}$$

$$\{\frac{m}{2}+\frac{m}{4}+1, \dots, \frac{m}{2}+\frac{m}{4}+\frac{m}{8}\}$$

next \vdots $m/2$

probes

$$\leq 2$$

$$\leq 4$$

$$\leq 8$$

$$\leq 2^i$$

Total # probes

$$\left. \begin{matrix} m \\ m \\ m \\ m \end{matrix} \right\}$$