

COL106 - Data Structures and Algorithms

Map
ADT:

HashTables

Implementing Map ADT

java.util. Map
java.util. Dictionary

- Arrays, LinkedList (inefficient)
- Unordered sequence:
 - $\text{get}(k)$, $\text{remove}(k)$ takes $O(n)$ time
 - $\text{put}(k, v)$ takes $O(1)$ time
 - Could be useful if there are not too many $\text{get}()$, $\text{remove}()$ ops required.
(e.g. log files in a computer)
- Ordered sequence (say array):
 - $\text{get}(k)$ takes $O(\log n)$ time, $\text{put}(k, v)$ & $\text{remove}(k)$ take $O(n)$ time
 - Good if all you need is search.

Hash tables

Direct addressing:

Array indexed by key: takes $O(1)$ time for all operations, but $O(r)$ space.
- e.g.: COL106 registry. $\hookrightarrow r$ is range of keys.

Hash Table:

- ↓
① Hash function
② Array (table)

- $O(1)$ expected time

- $O(n+m)$ space where m is table size.
 n is number of entries.

Store item (k, v) at
index $i = h(k)$

Example: Let keys be entry numbers of COL106 students.

↓
2022 PH 10140 \rightarrow 140

Hash function: Take last 3 digits: \rightarrow Collision

How to deal with alphanumeric keys?

Birthday Paradox

What is the probability that there are 2 people in this room who share a birthday?

340

365

Birthday Paradox

What is the probability that there are 2 people in this room who share a birthday?

$P(B)$ = Prob. that at least 2 people share a birthday

$$= 1 - \left(\frac{365}{365} \times \frac{364}{365} \times \dots \times \frac{365-360}{365} \right)$$

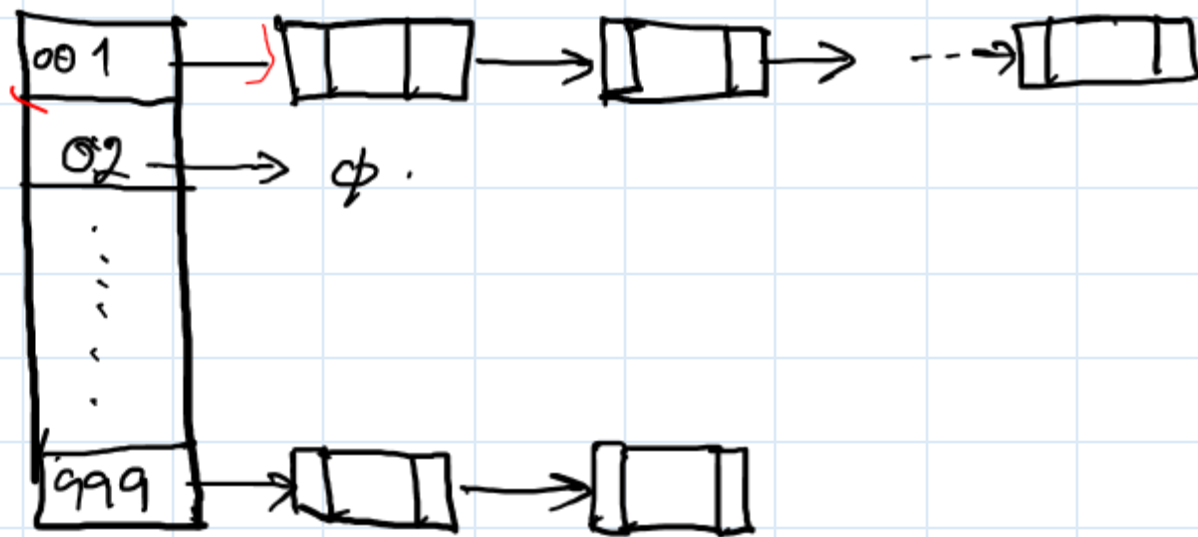
$$= 1 - 1.07 \times 10^{-118} \Rightarrow \approx 1$$

No matter what hash function you use, there will be collisions

Collision Resolution

- What if I have two keys hashing to same location?

- Chaining: Have an array of links indexed by keys, having list of items with same key.



To find/insert/delete an element, lookup position in table & search/insert/delete the element in the linked list of the hashed slot

$h: \mathbb{U} \rightarrow \{0, \dots, m-1\} \rightarrow m = \text{Size of hash table.}$

Analysis of Hashing

$$h(k) = [\text{rand}(k)] \rightarrow m$$

- Element with key k stored in slot $h(k)$.

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

$$\{0, 1, \dots, m-1\}$$

- Assume time taken to compute $h(k)$ is $\Theta(1)$.

- What is a good hash function?

- distributes keys "evenly" among slots, easy to compute, less space.
- ideally slot is picked uniformly at random & hash
- Simple uniform function (assumption)
- Load factor: $\alpha = \frac{n}{m} \rightarrow \begin{matrix} \# \text{ elements} \\ \text{size of table} \end{matrix}$

$$n = O(m)$$

↳ not actually.
we need to know
where a key got
mapped to!
↳ also not easy to store

Analysis of Hashing

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

- Simple Uniform hashing leads to a list of length $\alpha = \frac{1}{m}$ on average.

\Rightarrow Expected search time = $O(1 + \alpha)$ \rightarrow worst case unsuccessful search

- Efficiency relies on choice of hash function!

What about expected search time for a successful search?

$$h(k) = i \quad x[i]$$

$$O(1 + \frac{\alpha}{2})$$

r

Analysis of Hashing

- Simple Uniform hashing leads to a list of length $d = \frac{n}{m}$ on average.

\Rightarrow Expected search time = $O(1+d)$

\rightarrow worst case unsuccessful search

- Efficiency relies on choice of hash function!

What about expected search time for a successful search?

- You computed a hash & found a hit. How many elements did you have to look at?

Obs:

Expected length of list when i -th element was inserted = $\frac{i}{m}$

$O\left(1 + \frac{d}{2}\right)$

$a + bz + cz^2$ $c + az + bz^2$ Hash functions = Hash code + Compression map
 key \rightarrow integer integer $\rightarrow [0, m-1]$

Hash Code: 2021 CS 11011

- 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 of fixed length & evaluate polynomial)

a_0, \dots, a_{n-1}

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

numbers

Horner's rule: $p(z) = a_0 + z(a_1 + z(a_2 + \dots$

$$p_0(z) \leftarrow a_{n-1}, \quad p_i(z) \leftarrow a_{n-i-1} + z p_{i-1}(z)$$

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

Compression Maps

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

choosing $m \neq 2^e$ (bad!) ^{ge}

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)

④

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

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

Handling collisions

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

Open addressing:

In chaining all elements were stored outside the hash table.

Put all the elements in the table $\Rightarrow n \leq m$

$A[i]$

Have a systematic way to probe elements of the table.

Modify hash function:

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

specify the probe number

Probe sequence:

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

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