

CS21 WEEK 06

↳ Hash Tables

↳ Midterm: (90^{min} - 2 hrs) ← Expected Time

- In class, written exam

- Closed book

- Primarily Writing Functions

- Can memorize code or
memorize steps to then
write code

- Code in C++/Java for full cred.

or code in psuedo for partial cred.

- Won't have time to write comments

Not On Exam

- Insert Sort

- Merge Sort

- Counting Sort

↳ Def on Exam:

- Quicksort

- Priority Queue/Heaps

- Radix Bucket hybrid

- Hash Tables

- Review stuff

↳ Only exact var
multiplication method

↳ Linked Lists, Stacks, Queues

(Guessing ~7-10 q's)

HASH TABLES

DYNAMIC SET

- tracks membership that can change overtime

- Array PFA } Insert $O(1)$
search
Delete $O(n)$

- Linked List

insert

delete

Find/Search

Dictionary

ADT
(Abstract Data
Type)

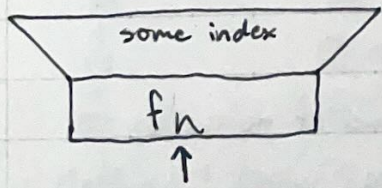
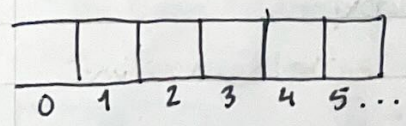
Hash Table

expected time

insert, remove, search $O(1)$

built on an array

(w.c. is $O(n)$)



Basic Idea: Hashing Function

- Given a value

- Returns an index (index into container - array)

Complications:

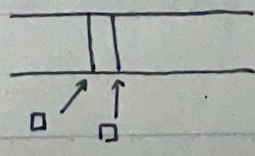
- Writing a good hashing function

- Collision

How resolve?

- Chaining

- Open addressing



Worst:

```
int hash(something){
    return 0;
    *A perfect hashing
    fn
    n data
}
```


* Required to support duplicates in our hash program *
 ↑
 causes collisions

Date

(Simplest?) technique For Collision Resolution

- Chaining

↳ Each 'spot' in an array is a list, not a single element.

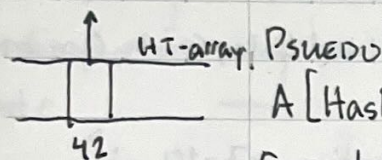
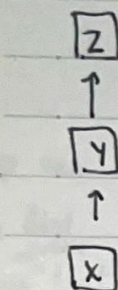
[x] [y] [z]

* Good HT Lists:

Vector - C++

Array List - Java

hashfn (~) → 42



PSEUDO


A[Hash(some-data)].insert(some-data)

Expected Time:

$O(1 + \alpha)$ α - load factor

:

$O(1)$

$\frac{n}{m} = \text{data}$ $m = \text{array}$ 
 ↑ size

RULE FOR SEARCH & DELETE → First Found Only

HASHING FUNCTIONS

- Def: Simple, Uniform Hashing

↳ All spots equally likely

↳ Simple concept, not simple implementation

- 3 Techniques: BAD, GOOD, BETTER

1) BAD, DON'T USE: Division Method

$h(k) = k \bmod m$ size of table

↳ Performance method is critical on value of m

↳ Usually a power of 2 is a bad idea

↳ Prime & not near a power of 2 (Eg. 701)

m = size of table

2) Good ← What to implement
Multiplication Method

$$\frac{1}{\phi} = \frac{\sqrt{5}-1}{2}$$

$$h(k) = \lfloor m (kC - \lfloor kC \rfloor) \rfloor$$

$$\phi + 1 = \phi^2$$

- Where $0 < C < 1 \leftarrow C$ is a const

- Optimal Choice of C depends on data ≈ 0.618034

$\hat{=}$ (inverse of golden ratio)

CODE

$$\text{hash}(k) = (\text{int})(m * \frac{(k * C - (\text{int})(k * C))}{[0.00 - 1.00]})$$

$k * C$

size of table

xxxxx.yyyyyy

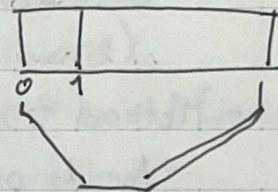
= xxxxx.00000

00000.yyyyyy

Value between

0-99%

key value of
data inserted in
hash table



3) Not Using (Protects You from Malicious attacks)

- Select C at runtime

Ideal Const: $\frac{1}{\phi} = 0.618034$
golden ratio

METHOD #2 FOR COLLISION RESOLUTION

- By Open Addressing (3 variations)

• On collision



A → hash func → 17

B → hash func → 17

17

collision

C →

- Method 1 Linear Probe

$$h'(k, i) = (h(k) + i) \% m$$

value attempt

- Easy to code, suffers from * primary clustering

- Bad performance

- Method 2 Quadratic Probe

$$h'(k, i) = (h(k) + c_1 * i + c_2 * (i * i)) \% m$$

$c_1, c_2 \neq 0$, they are constants

- Good choice of c_1, c_2 , and m is difficult :/

Add array simulation later

- Method 3 Double Hashing

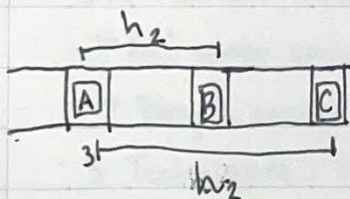
- Avoids primary & secondary clustering

$$h'(k, i) = (\underbrace{h_1(k)}_{\text{original hash func.}} + i \cdot \underbrace{h_2(k)}_{\text{second hash func., different c}}) \% m$$

original second hash

hash func. func., different c

$$i = 0$$



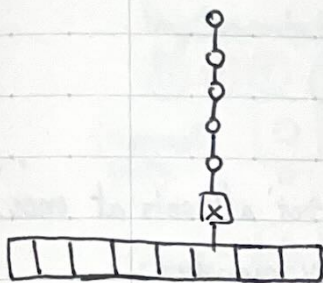
Outline/Pseudo

$\left[\begin{array}{l} \text{hash function (int)} \\ \text{(record)} \leftarrow \text{gets key value from record} \end{array} \right.$
 $\hookrightarrow \text{not required}$

record
 $\langle \text{int, string} \rangle$

insert (record)

$A[\text{hash}(\text{record})].\text{push_back}(\text{record})$



hash \leftarrow [x] rec

Search, Delete (record)

$\text{index} = \text{hash}(\text{record})$

loop through entries in $A[\text{index}]$

Time Complexity Matters

Performance Matters

\hookrightarrow Do a fast array delete