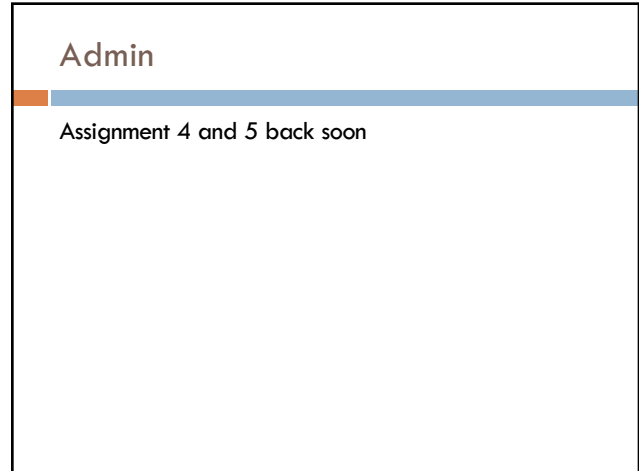# HASHTABLES

David Kauchak
CS 140 – Spring 2023
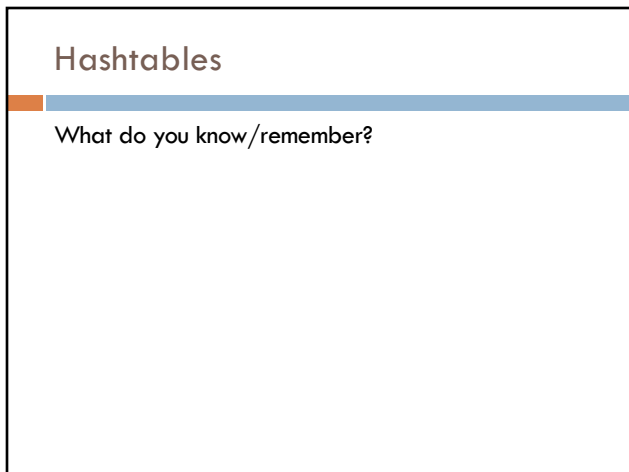
1

## Admin

Assignment 4 and 5 back soon
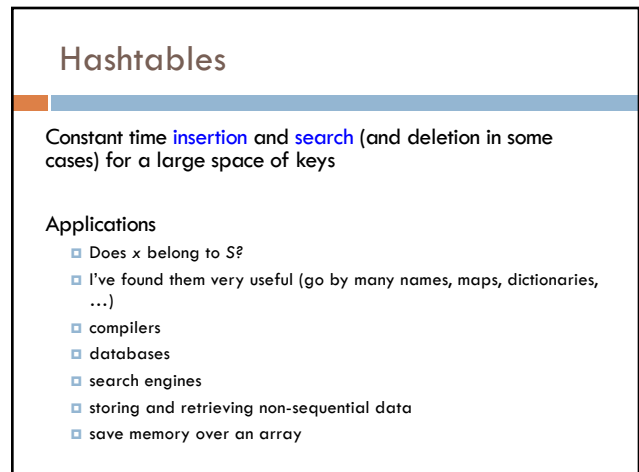
2

## Hashtables

What do you know/remember?

3

## Hashtables

Constant time insertion and search (and deletion in some cases) for a large space of keys

Applications
- Does *x* belong to *S*?
- I've found them very useful (go by many names, maps, dictionaries, …)
- compilers
- databases
- search engines
- storing and retrieving non-sequential data
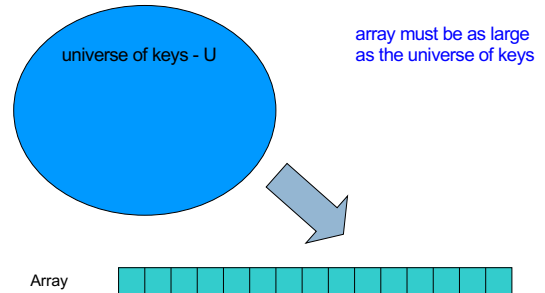- save memory over an array

4

## Hashtables

Constant time insertion and search (and deletion in some cases) for a large space of keys

For this class, we'll just think of them as a collection of keys

For many applications/implementations, there is a value associated with the key, i.e., key/value pair (though lookup is still exclusively based on the key)
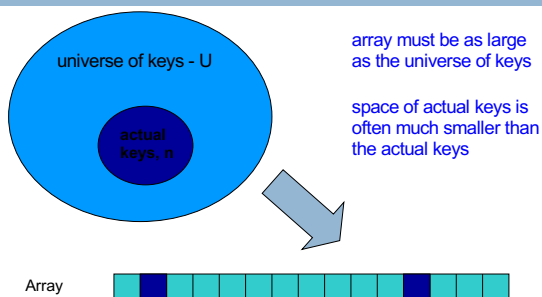
5

## Why not just arrays aka direct-address tables?

universe of keys - U

array must be as large as the universe of keys

Array

11

## Why not just arrays?

universe of keys - U

actual keys, n

array must be as large as the universe of keys

space of actual keys is often much smaller than the actual keys

Array

12

## Why not arrays?

Think of indexing all last names < 10 characters
- Census listing of all last names
  http://www.census.gov/genealogy/names/dist.all.last
  - 88,799 last names
- What is the size of our space of keys?
  - $26^{10}$ = a big number
- Not feasible!
- Even if it were, not space efficient

13

## The load of a table/hashtable

m = number of possible entries in the table

n = number of keys stored in the table

$\alpha = n/m$ is the **load factor** of the hashtable

**What is the load factor of the last example?**

- $\alpha = 88{,}799 / 26^{10}$ would be the load factor of last names using direct-addressing

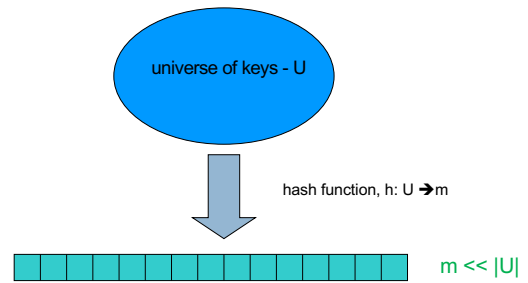The smaller $\alpha$, the more wasteful the table

The load also helps us talk about run time

14

## Hash function, *h*

A hash function is a function that maps the universe of keys to the slots in the hashtable



universe of keys - U

hash function, h: U ➔ m

m << |U|

15

## Hash function, *h*

A hash function is a function that maps the universe of keys to the slots in the hashtable



universe of keys - U

hash function, h: U ➔ m

m << |U|

16

## Hash function, *h*

A hash function is a function that maps the universe of keys to the slots in the hashtable



universe of keys - U

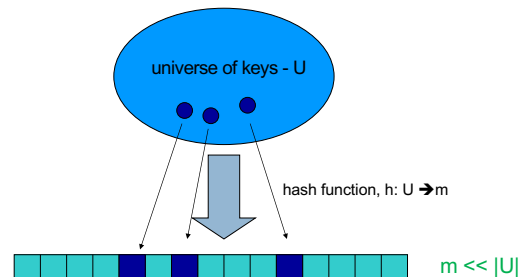hash function, h: U ➔ m

m << |U|

17

## Hash function, *h*

What can happen if  m < |U|?

universe of keys - U

hash function, h: U ➜ m

m << |U|

18

## Collisions

If m < |U|, then two keys can map to the same position in the hashtable (pidgeonhole principle)

universe of keys - U

hash function, h

m << |U|

19

## Collisions

A collision occurs when $h(x) = h(y)$, but $x \neq y$

A good hash function will minimize the number of collisions

Because the number of hashtable entries is less than the possible keys (i.e. m < |U|) collisions are inevitable!

Collision resolution techniques?

20

## Collision resolution by chaining

Hashtable consists of an array of linked lists

When a collision occurs, the element is added to linked list at that location

If two entries $x \neq y$ have the same hash value $h(x) = h(y)$, then $T(h(x))$ will contain a linked list with both values

21

## Insertion

ChainedInsert(x):
  entry = h(x)
  insert x at the head of T[entry]

ChainedHashInsert( 🟩 )

22

## Insertion

ChainedInsert(x):
  entry = h(x)
  insert x at the head of T[entry]

h( 🟩 )   hash function is a mapping from the key to some value < m

23

## Insertion

ChainedInsert(x):
  entry = h(x)
  insert x at the head of T[entry]

h( 🟩 )

24

## Deletion

ChainedDelete(x):
  entry = h(x)
  delete x at the list at T[entry]

What does that involve?

25

## Deletion

ChainedDelete(x):
  entry = h(x)
  delete x at the list at T[entry]

Search though the list!

26

## Deletion

ChainedDelete(x):
  entry = h(x)
  delete x at the list at T[entry]

ChainedHashDelete( )

27

## Deletion

ChainedDelete(x):
  entry = h(x)
  delete x at the list at T[entry]

h( )

28

## Deletion

ChainedDelete(x):
  entry = h(x)
  delete x at the list at T[entry]
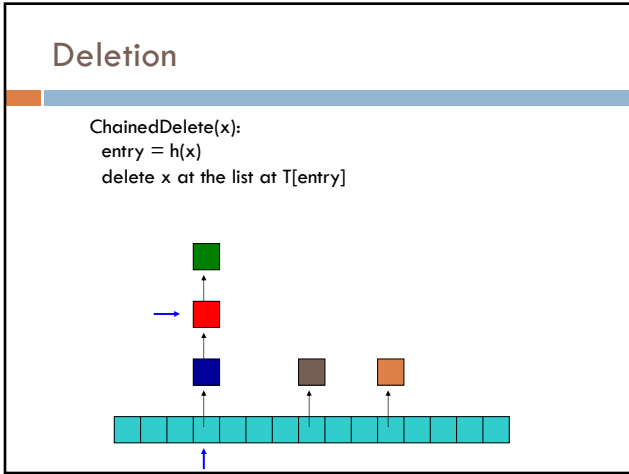
29

## Deletion

ChainedDelete(x):
  entry = h(x)
  delete x at the list at T[entry]

30

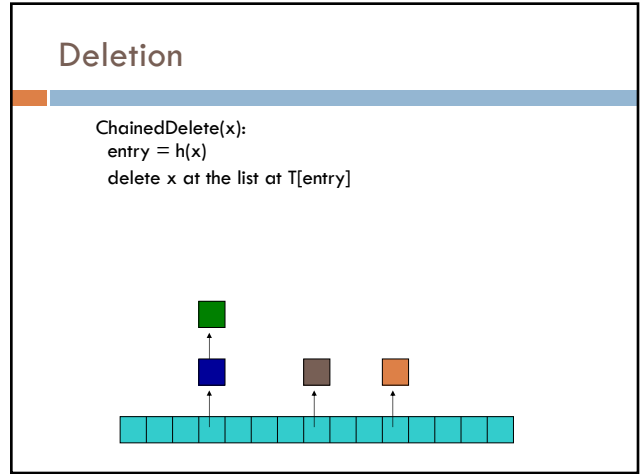## Deletion

ChainedDelete(x):
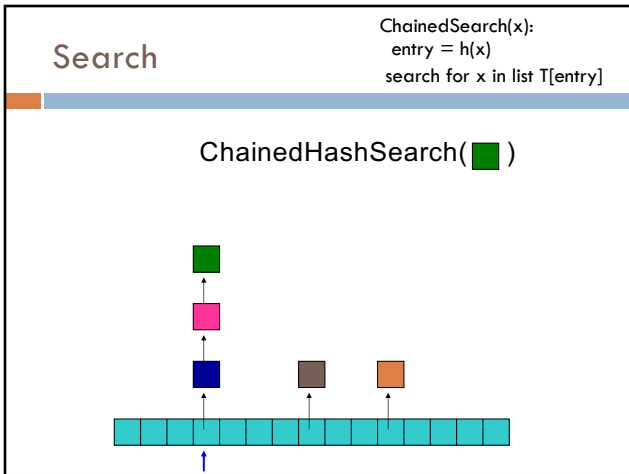  entry = h(x)
  delete x at the list at T[entry]

31

## Search

ChainedSearch(x):
  entry = h(x)
  search for x in list T[entry]

ChainedHashSearch( ▪ )

32

## Search

ChainedSearch(x):
  entry = h(x)
  search for x in list T[entry]

h( ▪ )

33

**Search**

ChainedSearch(x):
entry = h(x)
search for x in list T[entry]

ChainedHashSearch( ◼ )

34

**Search**

ChainedSearch(x):
entry = h(x)
search for x in list T[entry]

ChainedHashSearch( ◼ )

35

**Search**

ChainedSearch(x):
entry = h(x)
search for x in list T[entry]

ChainedHashSearch( ◼ )
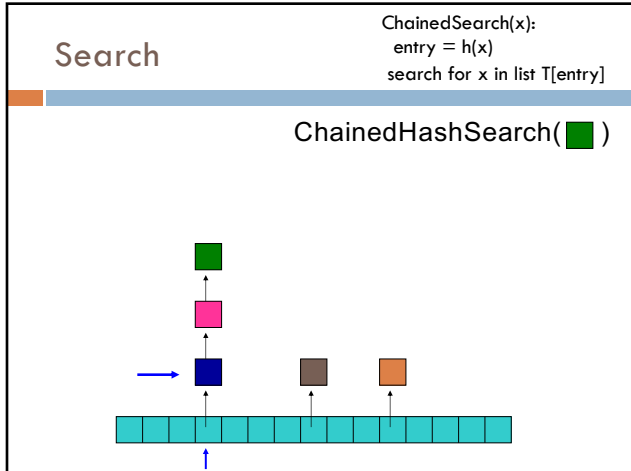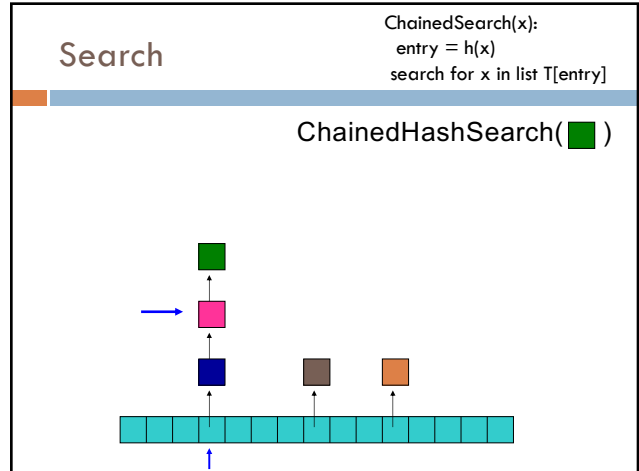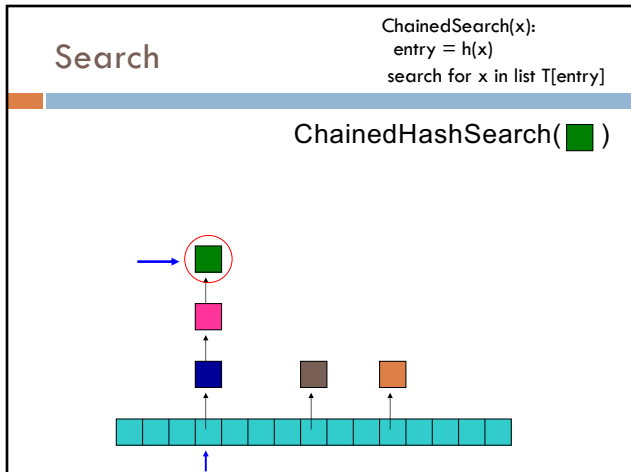
36

**Running time**

ChainedInsert(x):
entry = h(x)
insert x at the head of T[entry]

$\Theta(1)$

ChainedDelete(x):
entry = h(x)
delete x at the list at T[entry]

O(length of the chain)

ChainedSearch(x):
entry = h(x)
search for x in list T[entry]

O(length of the chain)

37

## Length of the chain

Worst case?

## Length of the chain

Worst case?
- All elements hash to the same location
- h(k) = 4
- **n**

## Length of the chain

Average case:

Depends on how well the hash function distributes the keys

What is the best we could hope for a hash function?
- simple uniform hashing: an element is equally likely to end up in any of the $m$ slots

Under simple uniform hashing what is the average length of a chain in the table?
- $n$ keys over $m$ slots = $n / m = \alpha$

## Average chain length

If you roll a fair $m$ sided die $n$ times, how many times are we likely to see a given value?

For example, 10 sided die:
- 1 time
  - 1/10
- 100 times
  - 100/10 = 10

## Search average running time

Two cases:

- Key is **not** in the table
  - must search all entries
  - $\Theta(1 + \alpha)$
- Key **is** in the table
  - on average search half of the entries
  - $O(1 + \alpha)$

## Hash functions

What makes a good hash function?

- Approximates the assumption of simple uniform hashing
- Deterministic – h(x) should always return the same value
- Low cost – if it is expensive to calculate the hash value (e.g. log n) then we don't gain anything by using a table

Challenge: we don't generally know the distribution of the keys

- Frequently data tend to be clustered (e.g. similar strings, run-times, SSNs). A good hash function should spread these out across the table

## Hash functions

What are some hash functions you've heard of before?

## Division method

$h(k) = k \bmod m$

| m | k | h(k) |
|---|----|------|
| 11 | 25 | |
| 11 | 1 | |
| 11 | 17 | |
| 13 | 133 | |
| 13 | 7 | |
| 13 | 25 | |

## Division method

$h(k) = k \bmod m$

| m | k | h(k) |
|---|---|------|
| 11 | 25 | 3 |
| 11 | 1 | 1 |
| 11 | 17 | 6 |
| 13 | 133 | 3 |
| 13 | 7 | 7 |
| 13 | 25 | 12 |

46

## Division method

**Don't** use a power of two.  Why?

| m | k | bin(k) | h(k) |
|---|---|--------|------|
| 8 | 25 | 11001 | |
| 8 | 1 | 00001 | |
| 8 | 17 | 10001 | |

47

## Division method

**Don't** use a power of two.  Why?

| m | k | bin(k) | h(k) |
|---|---|--------|------|
| 8 | 25 | 11001 | 1 |
| 8 | 1 | 00001 | 1 |
| 8 | 17 | 10001 | 1 |

if $h(k) = k \bmod 2^p$, the hash function is just the lower $p$ bits of the value

48

## Division method

Good rule of thumb for $m$ is a prime number not too close to a power of 2

Pros:
- quick to calculate
- easy to understand

Cons:
- keys close to each other will end up close in the hashtable

49

## Multiplication method

Multiply the key by a constant $0 < A < 1$ and extract the fractional part of $kA$, then scale by $m$ to get the index

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

extracts the fractional portion of $kA$

50

## Multiplication method

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

Common choice is for $m$ as a power of 2 and

$$A = (\sqrt{5} - 1)/2 = 0.6180339887$$

Why a power of 2?

Book has other heuristics

51

## Multiplication method

| m | k | A | kA | h(k) |
|---|----|-------|----|------|
| 8 | 15 | 0.618 | | |
| 8 | 23 | 0.618 | | |
| 8 | 100 | 0.618 | | |

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

52

## Multiplication method

| m | k | A | kA | h(k) |
|---|----|-------|--------|-------------------|
| 8 | 15 | 0.618 | 9.27 | floor(0.27*8) = 2 |
| 8 | 23 | 0.618 | 14.214 | floor(0.214*8) = 1 |
| 8 | 100 | 0.618 | 61.8 | floor(0.8*8) = 6 |

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

53

## Other hash functions

http://en.wikipedia.org/wiki/List_of_hash_functions

cyclic redundancy checks (i.e. disks, cds, dvds)

Checksums (i.e. networking, file transfers)

Cryptographic (i.e. MD5, SHA)

54

## Open addressing

Keeping around an array of linked lists can be inefficient and a hassle

Like to keep the hashtable as just an array of elements (no pointers)

How do we deal with collisions?
- compute another slot in the hashtable to examine

55

## Hash functions with open addressing

Hash function must define a **probe sequence** which is the list of slots to examine when searching or inserting

The hash function takes an additional parameter *i* which is the number of collisions that have already occurred

The probe sequence **must** be a permutation of every hashtable entry. Why?

{ h(k,0), h(k,1), h(k,2), …, h(k, m-1) } is a permutation of { 0, 1, 2, 3, …, m-1 }

56

## Hash functions with open addressing

Hash function must define a **probe sequence** which is the list of slots to examine when searching or inserting

The hash function takes an additional parameter *i* which is the number of collisions that have already occurred

The probe sequence **must** be a permutation of every hashtable entry. Why?

If not, we wouldn't explore all the possible location in the table!

57

## Probe sequence

h(k, 0)

58

## Probe sequence

h(k, 1)

59

## Probe sequence

h(k, 2)

60

## Probe sequence

h(k, 3)

61

## Probe sequence

h(k, …)

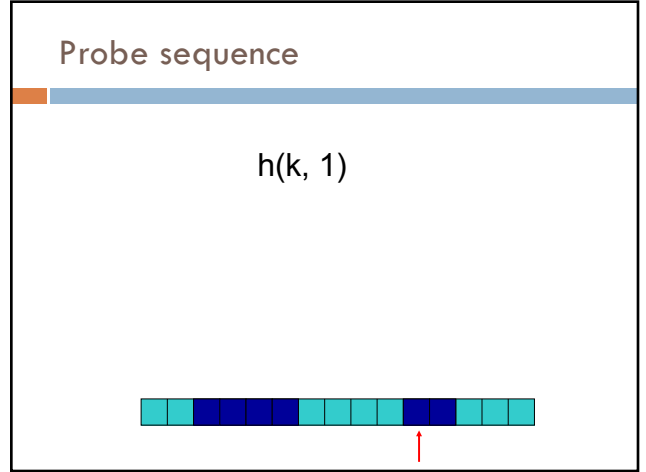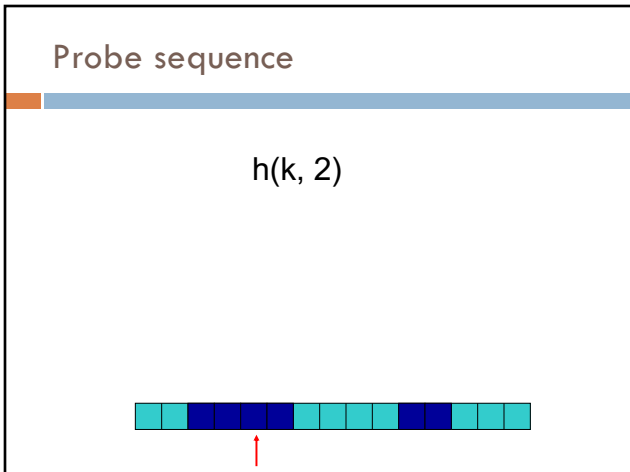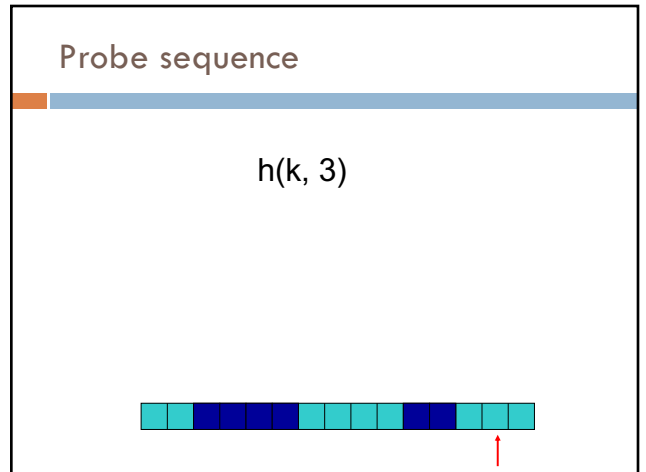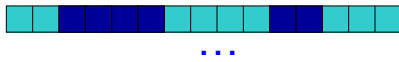must visit all locations

...

## Open addressing: Insert

HASH-INSERT$(T, k)$
1   $i \leftarrow 0$
2   $j \leftarrow h(k, i)$
3   **while** $i < m - 1$ and $T[j] \neq null$
4        $i \leftarrow i + 1$
5        $j \leftarrow h(k, i)$
6   **if** $T[j] = null$
7        **return** $j$
8   **else**
9        **error** "hash is full"

## Open addressing: Insert

HASH-INSERT$(T, k)$
1   $i \leftarrow 0$
2   $j \leftarrow h(k, i)$
3   **while** $i < m - 1$ and $T[j] \neq null$
4        $i \leftarrow i + 1$
5        $j \leftarrow h(k, i)$
6   **if** $T[j] = null$
7        **return** $j$
8   **else**
9        **error** "hash is full"

get the first hashtable
entry to look in

## Open addressing: Insert

HASH-INSERT$(T, k)$
1   $i \leftarrow 0$
2   $j \leftarrow h(k, i)$
3   **while** $i < m - 1$ and $T[j] \neq null$
4        $i \leftarrow i + 1$
5        $j \leftarrow h(k, i)$
6   **if** $T[j] = null$
7        **return** $j$
8   **else**
9        **error** "hash is full"

follow the probe
sequence until we find
an open entry

## Open addressing: Insert

HASH-INSERT$(T, k)$
1   $i \leftarrow 0$
2   $j \leftarrow h(k, i)$
3   **while** $i < m - 1$ and $T[j] \neq null$
4        $i \leftarrow i + 1$
5        $j \leftarrow h(k, i)$
6   **if** $T[j] = null$      return the open entry
7        **return** $j$
8   **else**
9        **error** "hash is full"

66

## Open addressing: Insert

HASH-INSERT$(T, k)$
1   $i \leftarrow 0$
2   $j \leftarrow h(k, i)$
3   **while** $i < m - 1$ and $T[j] \neq null$
4        $i \leftarrow i + 1$
5        $j \leftarrow h(k, i)$
6   **if** $T[j] = null$
7        **return** $j$
8   **else**
9        **error** "hash is full"      hashtable can fill up

67

## Open addressing: search

HASH-SEARCH$(T, k)$
1   $i \leftarrow 0$
2   $j \leftarrow h(k, i)$
3   **while** $i < m - 1$ and $T[j] \neq null$ and $T[j] \neq k$
4        $i \leftarrow i + 1$
5        $j \leftarrow h(k, i)$
6   **if** $T[j] = k$
7        **return** $j$
8   **else**
9        **return** $null$

68

## Open addressing: search

HASH-SEARCH$(T, k)$
1   $i \leftarrow 0$
2   $j \leftarrow h(k, i)$
3   **while** $i < m - 1$ and $T[j] \neq null$ and $T[j] \neq k$
4        $i \leftarrow i + 1$
5        $j \leftarrow h(k, i)$
6   **if** $T[j] = k$
7        **return** $j$
8   **else**
9        **return** $null$

HASH-INSERT$(T, k)$
1   $i \leftarrow 0$
2   $j \leftarrow h(k, i)$
3   **while** $i < m - 1$ and $T[j] \neq null$
4        $i \leftarrow i + 1$
5        $j \leftarrow h(k, i)$
6   **if** $T[j] = null$
7        **return** $j$
8   **else**
9        **error** "hash is full"

69

16

## Open addressing: search

HASH-SEARCH(T, k)
```
1  i ← 0
2  j ← h(k, i)
3  while i < m − 1 and T[j] ≠ null and T[j] ≠ k
4          i ← i + 1
5          j ← h(k, i)
6  if T[j] = k
7          return j
8  else
9          return null
```

HASH-INSERT(T, k)
```
1  i ← 0
2  j ← h(k, i)
3  while i < m − 1 and T[j] ≠ null
4          i ← i + 1
5          j ← h(k, i)
6  if T[j] = null
7          return j
8  else
9          error "hash is full"
```

"breaks" the probe sequence

70

## Open addressing: delete

Two options:

- mark node as "deleted" (rather than null)
  - modify search procedure to continue looking if a "deleted" node is seen
  - modify insert procedure to fill in "deleted" entries
  - increases search times

- if a lot of deleting will happen, use chaining

71

## Probing schemes

Linear probing – if a collision occurs, go to the next slot

- h(k,i) = (h(k) + i) mod m
- Does it meet our requirement that it visits every slot?
- for example, m = 7 and h(k) = 4

$$h(k,0) = 4$$
$$h(k,1) = 5$$
$$h(k,2) = 6$$
$$h(k,3) = 0$$
$$h(k,3) = 1$$

72

## Linear probing: search

h( ▮ , 0)



73

17

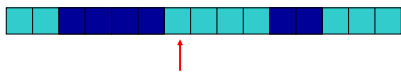Linear probing: search

h( ■, 1)

74

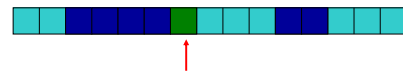Linear probing: search

h( ■, 2)

75

Linear probing: search

h( ■, 3)

76

Linear probing: search

h(   , 3)
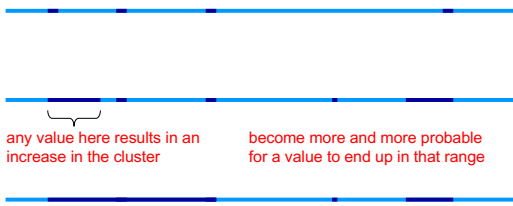
77

## Linear probing

Problem:

primary clustering – long runs of occupied slots tend to build up and these tend to grow



any value here results in an increase in the cluster

become more and more probable for a value to end up in that range

78

## Quadratic probing

$h(k,i) = (h(k) + c_1 i + c_2 i^2) \mod m$

Rather than a linear sequence, we probe based on a quadratic function

Problems:
- must pick constants and $m$ so that we have a proper probe sequence
- if $h(x) = h(y)$, then $h(x,i) = h(y,i)$ for all i
- secondary clustering

79

## Double hashing

Probe sequence is determined by a second hash function

$h(k,i) = (h_1(k) + i(h_2(k))) \mod m$

Problem:
- $h_2(k)$ must visit all possible positions in the table

80

## Running time of insert and search for open addressing

Depends on the hash function/probe sequence

Worst case?
- O(n) – probe sequence visits every full entry first before finding an empty

81

## Slide 82

### Running time of insert and search for open addressing

Average case?

We have to make at least one probe

82

## Slide 83

### Running time of insert and search for open addressing

Average case?

What is the probability that the first probe will **not** be successful (assume uniform hashing function)?

$$\alpha$$

83

## Slide 84

### Running time of insert and search for open addressing

Average case?

What is the probability that the first **two** probed slots will **not** be successful?

$\sim\alpha^2$

why
'~'?

84

## Slide 85

### Running time of insert and search for open addressing

Average case?

What is the probability that the first **two** probed slots will **not** be successful

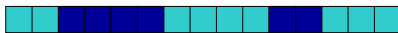Technically, second probe is: $\dfrac{n-1}{m-1}$     $\sim\alpha^2$

85

## Running time of insert and search for open addressing

Average case?

What is the probability that the first **three** probed slots will **not** be successful?

$$\sim \alpha^3$$



86

## Running time of insert and search for open addressing

Average case: expected number of probes

sum of the probability of making 1 probe, 2 probes, 3 probes, …

$$E[probes] = 1 + \alpha + \alpha^2 + \alpha^3 + \ldots$$
$$= \sum_{i=0}^{m} \alpha^i$$
$$< \sum_{i=0}^{\infty} \alpha^i$$
$$= \frac{1}{1-\alpha}$$

87

## Average number of probes

$$E[probes] = \frac{1}{1-\alpha}$$

| $\alpha$ | Average number of searches |
|------|------------------------|
| 0.1  | $1/(1-.1) = 1.11$ |
| 0.25 | $1/(1-.25) = 1.33$ |
| 0.5  | $1/(1-.5) = 2$ |
| 0.75 | $1/(1-.75) = 4$ |
| 0.9  | $1/(1-.9) = 10$ |
| 0.95 | $1/(1-.95) = 20$ |
| 0.99 | $1/(1-.99) = 100$ |

88

## How big should a hashtable be?

A good rule of thumb is the hashtable should be around half full

What happens when the hashtable gets full?

Copy: Create a new table and copy the values over
- results in one expensive insert
- simple to implement

Amortized copy: When a certain ratio is hit, grow the table, but copy the entries over a few at a time with every insert
- no single insert is expensive and can guarantee per insert performance
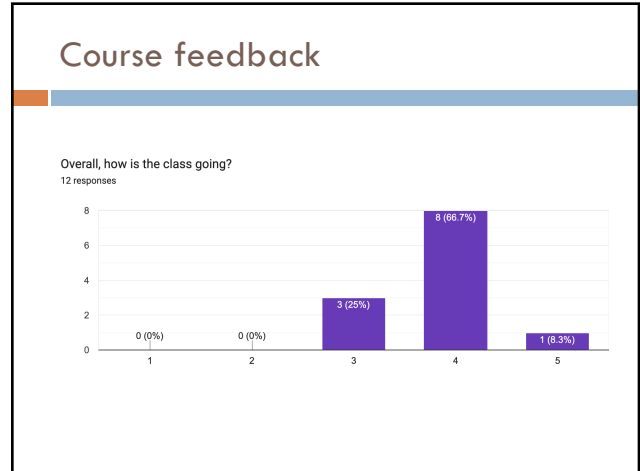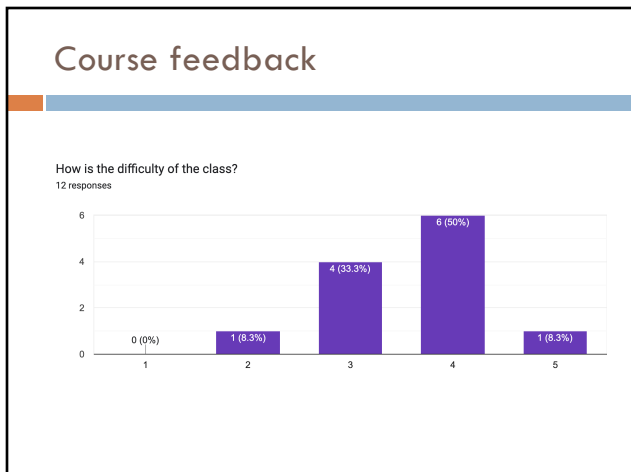- more complicated to implement

89

21

## Checkpoint 1

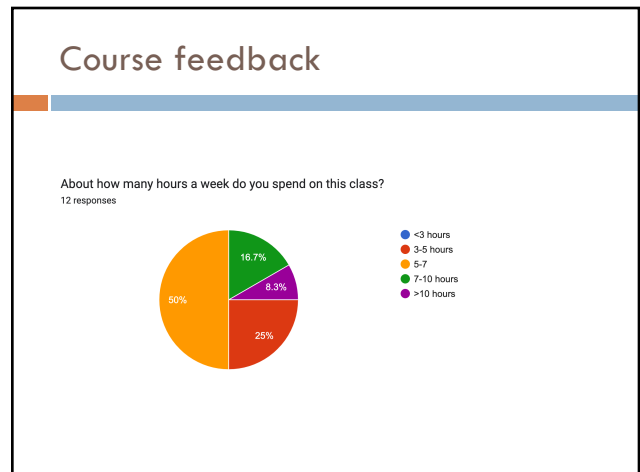Induction on trees

$$T(n) = T(\sqrt{n}) + c$$

90

## Course feedback

Overall, how is the class going?
12 responses



91

## Course feedback

How is the difficulty of the class?
12 responses



92

## Course feedback

About how many hours a week do you spend on this class?
12 responses



- <3 hours
- 3-5 hours
- 5-7
- 7-10 hours
- >10 hours

93

## Course feedback

I love proving things and looking at the Math behind the concepts from CS62.

the group assignments

Honestly I just really like the little comics at the start of every homework

94

## Course feedback

lectures are wayyy too fast, barely enough time to process things so it feels pointless to take notes; current course content is comprehensive and makes sense but it feels disorganized, like different content stitched together sort of so...

Having more examples, or going through the slides a bit slower

95

## Course feedback

The homeworks are a lot of work and the mentors are super helpful but someone's even they don't have the solutions and that wastes hours of our time. I think homeworks can have more straight forward problems that show we understand things rather than problems that we always have to scavenge the internet and bug mentors for understandings.

96

## Course feedback

During Class, could we have some more exercises along with the lecture contents?

97