# CPSC 221
# Basic Algorithms and Data Structures

# Hashing

Textbook References:
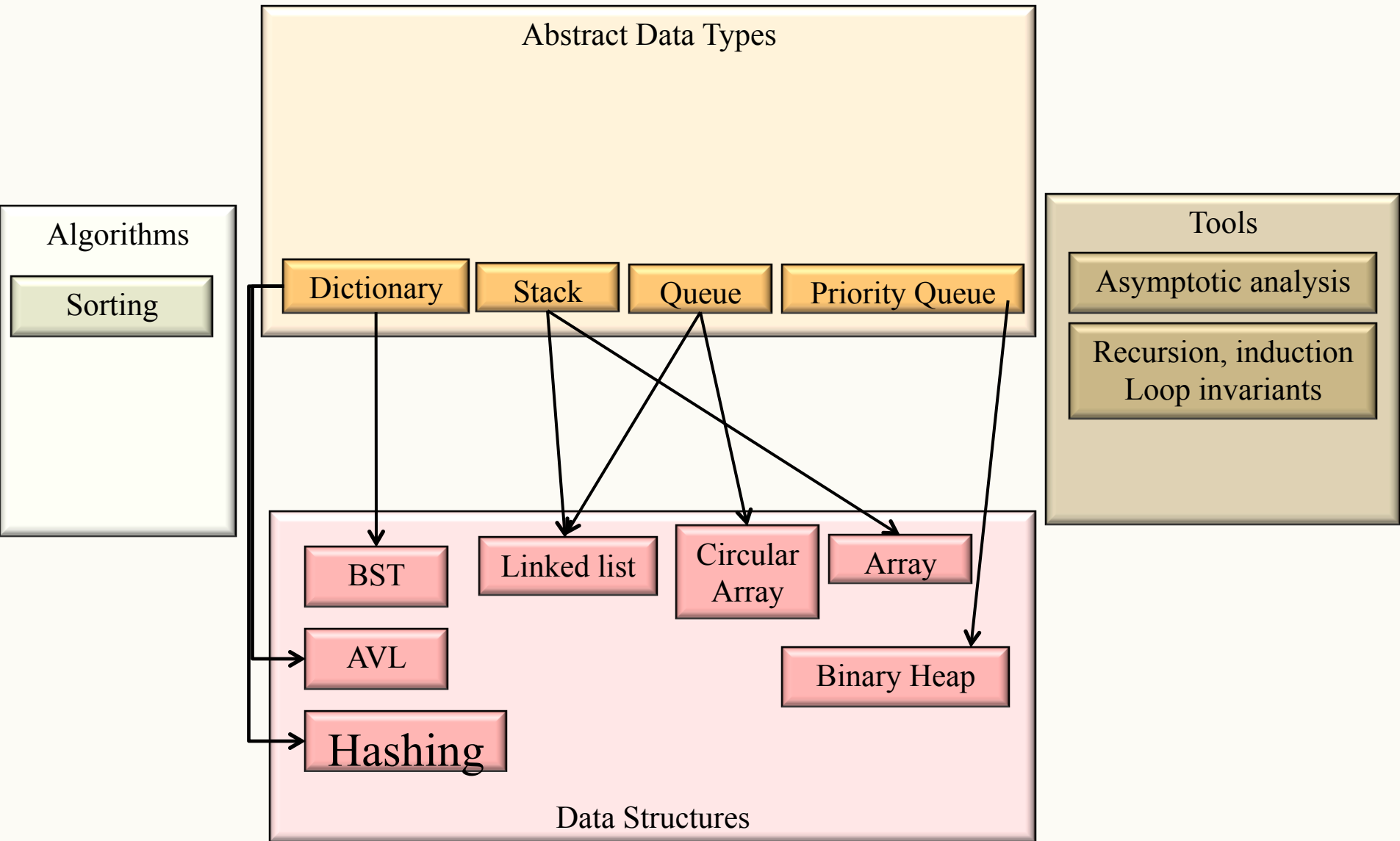Koffman: Chapter 9
EPP 3$^{rd}$ edition: 7.3
EPP 4$^{th}$ edition: 9.4

Hassan Khosravi
January – April  2015

# CPSC 221 Journey

Abstract Data Types

Algorithms

Sorting

Dictionary | Stack | Queue | Priority Queue

Tools

Asymptotic analysis

Recursion, induction
Loop invariants

BST

Linked list

Circular Array

Array

AVL

Binary Heap

Hashing

Data Structures

# Learning Goals

After this unit, you should be able to:

- Define various forms of the pigeonhole principle; recognize and solve the specific types of counting and hashing problems to which they apply.

- Provide examples of the types of problems that can benefit from a hash data structure.

- Compare and contrast open addressing and chaining.

- Evaluate collision resolution policies.

- Describe the conditions under which hashing can degenerate from $O(1)$ expected complexity to $O(n)$.

- Identify the types of search problems that do not benefit from hashing (e.g. range searching) and explain why.

- Manipulate data in hash structures both irrespective of implementation and also within a given implementation.

# CPSC 221 Administrative Notes

- Written Assignment 2 is due tomorrow(5pm)

- Lab 9 is posted, which is on Hashing
    - (Mar 20 – Mar 26)

- PeerWise  Call #4 is due Monday March 23.

# Reminder: Dictionary ADT

- Dictionary operations
  - create
  - destroy
  - insert
  - find
  - delete

insert
- brownies
  - tasty

find(wolf)
- wolf
  - the perfect mix of oomph and Scrabble value

- midterm
  - would be tastier with brownies
- prog-project
  - so painful… who invented templates?
- wolf
  - the perfect mix of oomph and Scrabble value

- Stores values associated with user-specified keys
  - values may be any (homogenous) type
  - keys may be any (homogenous) comparable type

# Implementations So Far

|  | insert | find | delete |
|---|---|---|---|
| • Unsorted list | O(1) | O(n) | O(n) |
| • Sorted Array | O(n) | O(log n) | O(n) |
| • AVL Trees | O(log n) | O(log n) | O(log n) |

Can we do better? O(1)?

# Example 1 (natural, numeric keys)

- In a small company of 100 employees, each employee is assigned an Emp_ID number in the range 0 – 99.

- To store the employee's records in an array, each employee's Emp_ID number acts as an index into the array where this employee's record will be stored as shown in figure

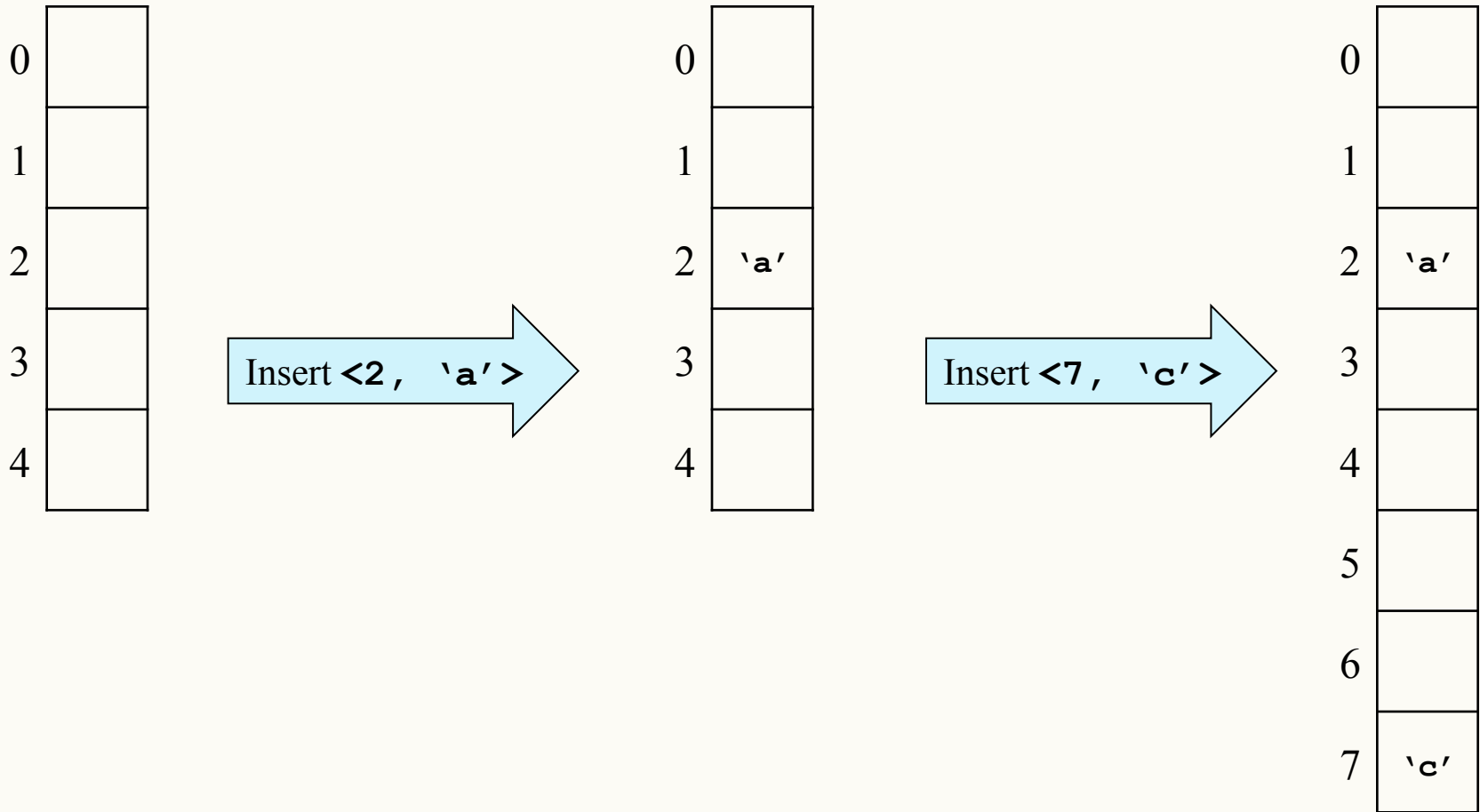| KEY | | ARRAY OF EMPLOYEE'S RECORD |
|---|---|---|
| Key 0 | [0] | Record of employee having Emp_ID 0 |
| Key 1 | [1] | Record of employee having Emp_ID 1 |
| ………………………….… | | ………………………………………….. |
| Key 99 | [99] | Record of employee having Emp_ID 99 |

# Follow-up example

- Let's assume that the same company uses a five digit Emp_ID number as the primary key. In this case, key values will range from 00000 to 99999. If we want to use the same technique as above, we will need an array of size 100,000, of which only 100 elements will be used.
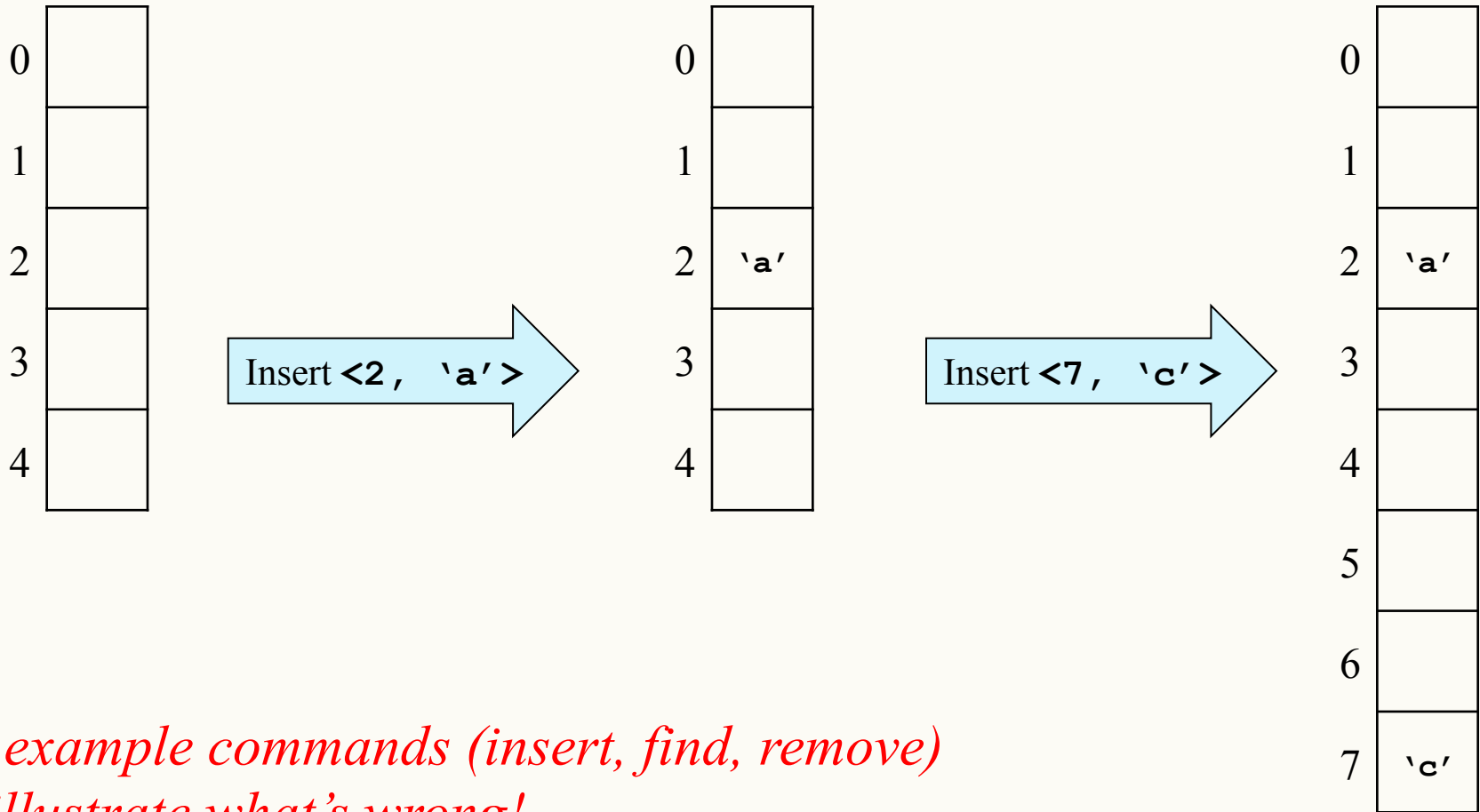
| KEY | ARRAY OF EMPLOYEE'S RECORD |
|---|---|
| Key 00000                    [0] | Record of employee having Emp_ID 00000 |
| …………………………… | ……………………………………………… |
| Key n                        [n] | Record of employee having Emp_ID n |
| …………………………… | …………………………………………….. |
| Key 99999              [99999] | Record of employee having Emp_ID 99999 |

- It is impractical to waste that much storage just to ensure that each employee's record is in a unique and predictable location.
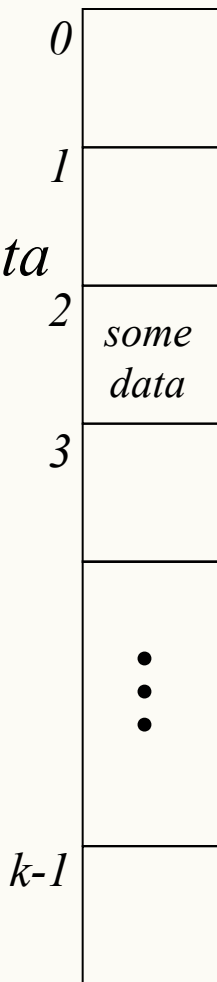
# First Pass: Resizable Vectors

0
1
2
3
4

Insert **<2, 'a'>**

0
1
2 **'a'**
3
4

Insert **<7, 'c'>**

0
1
2 **'a'**
3
4
5
6
7 **'c'**

# What's Wrong with Our First Pass?

```
0                  0                  0
                                      
1                  1                  1
                                      
2                  2  'a'             2  'a'
                                      
3    Insert <2, 'a'>   3   Insert <7, 'c'>   3
                                      
4                  4                  4
                                      
                                      5
                                      
                                      6
                                      
                                      7  'c'
```

*Give example commands (insert, find, remove)*
*that illustrate what's wrong!*

# Hash Table Goal

*We can do:*

*a[2] = some data*

| | |
|---|---|
| *0* | |
| *1* | |
| *2* | *some data* |
| *3* | |
| ⋮ | |
| *k-1* | |

*We want to do:*

*a["Steve"] = some data*

| | |
|---|---|
| *"Alan"* | |
| *"Hassan"* | |
| *"Steve"* | *some data* |
| *"Ed"* | |
| *"Will"* | |
| | ⋮ |
| *"Martin"* | |

*How will insert, find, and delete work?*

# Aside: How do arrays do that?

*We can do:*

*a[2] = some data*

```
0
1
2   some
    data
3

    ⋮

k-1
```

Q: If I know houses on a certain block in Vancouver are on 33-foot-wide lots, where is the 5th house?
A: It's from (5-1)*33 to 5*33 feet from the start of the block.
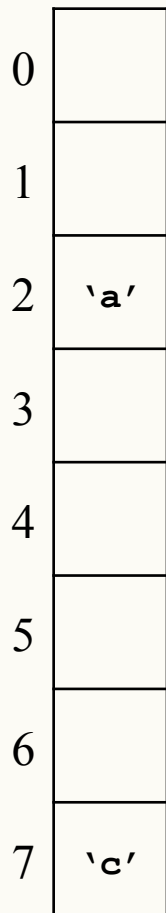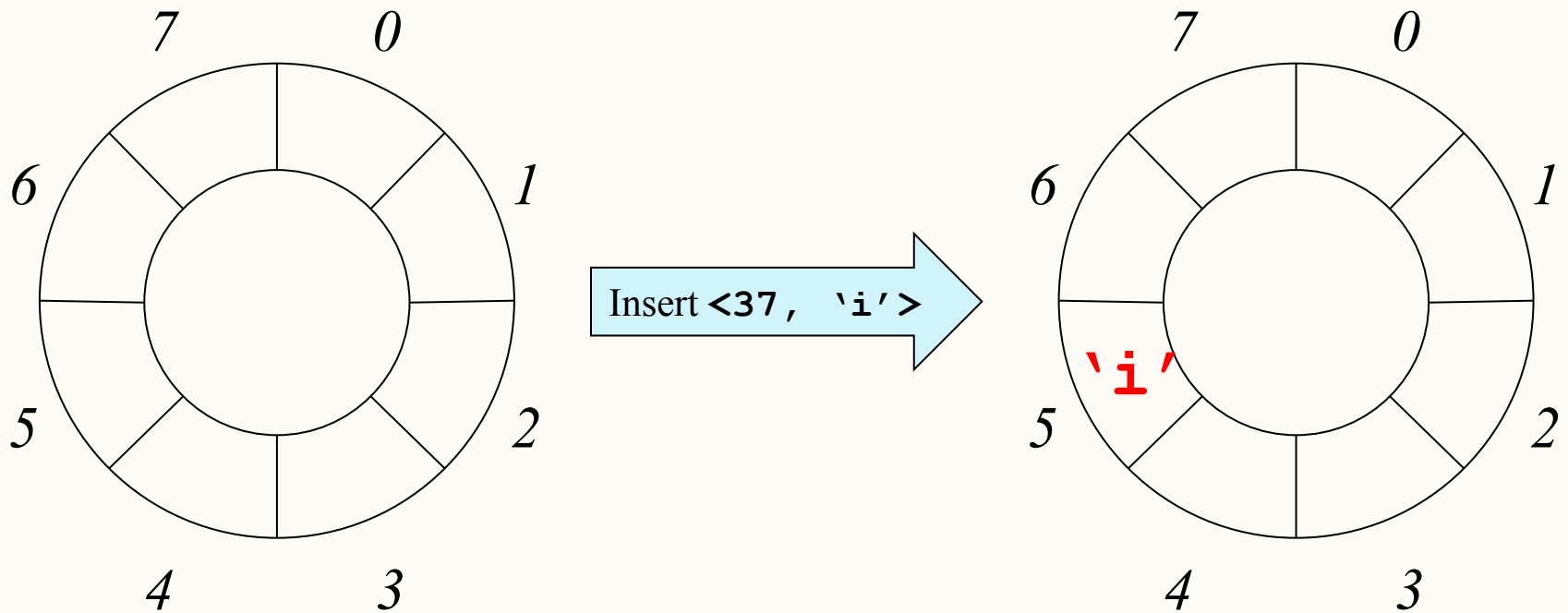
element_type a[SIZE];

Q: Where is a[i]?
A: start of a + i*sizeof(element_type)

Aside: This is why array elements have to be the same size, and why we start the indices from 0.

# What is the 25th Element?

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | `'a'` |
| 3 | |
| 4 | |
| 5 | |
| 6 | |
| 7 | `'c'` |

# What is the 25<sup>th</sup> Element Now?

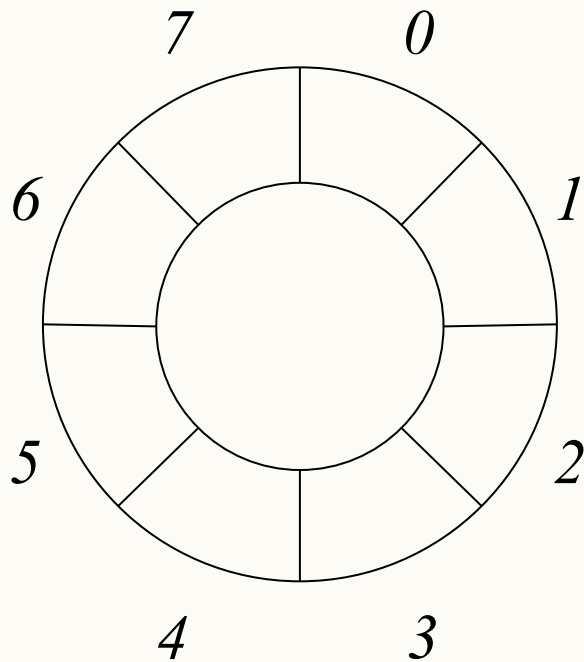|   |   |
|---|---|
| 0 |   |
| 1 |   |
| 2 | 'a' |
| 3 |   |
| 4 |   |
| 5 |   |
| 6 |   |
| 7 | 'c' |

considered as
a circular array

# Second Pass: Circular Array
# (For the Win?)



Insert **<37, 'i'>**

*Does this solve our memory usage problem?*

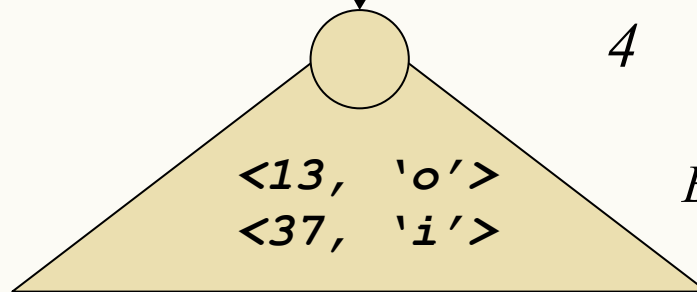# What's Wrong with our **Second** Pass?
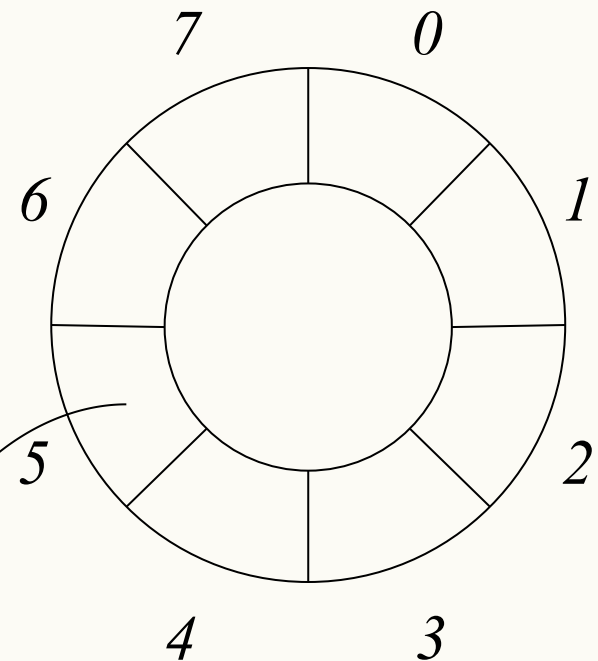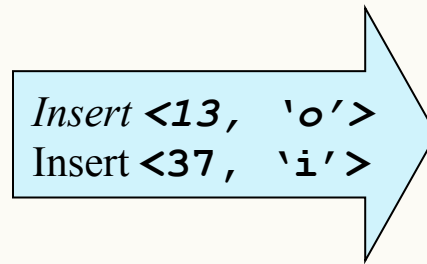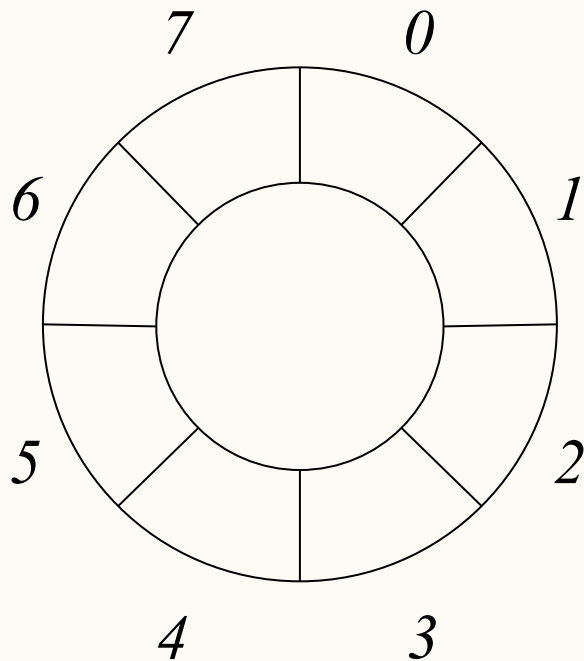
7　　　0
6　　　1
5　　　2
4　　　3

Let's insert　2 and 258?

258 % 8 = 2
258 % 16 = 2
258 % 32 = 2
258 % 64 = 2
258 % 128 = 2
258 % 256 = 2

Resize until they don't?

Solutions:
- Prime table sizes helps
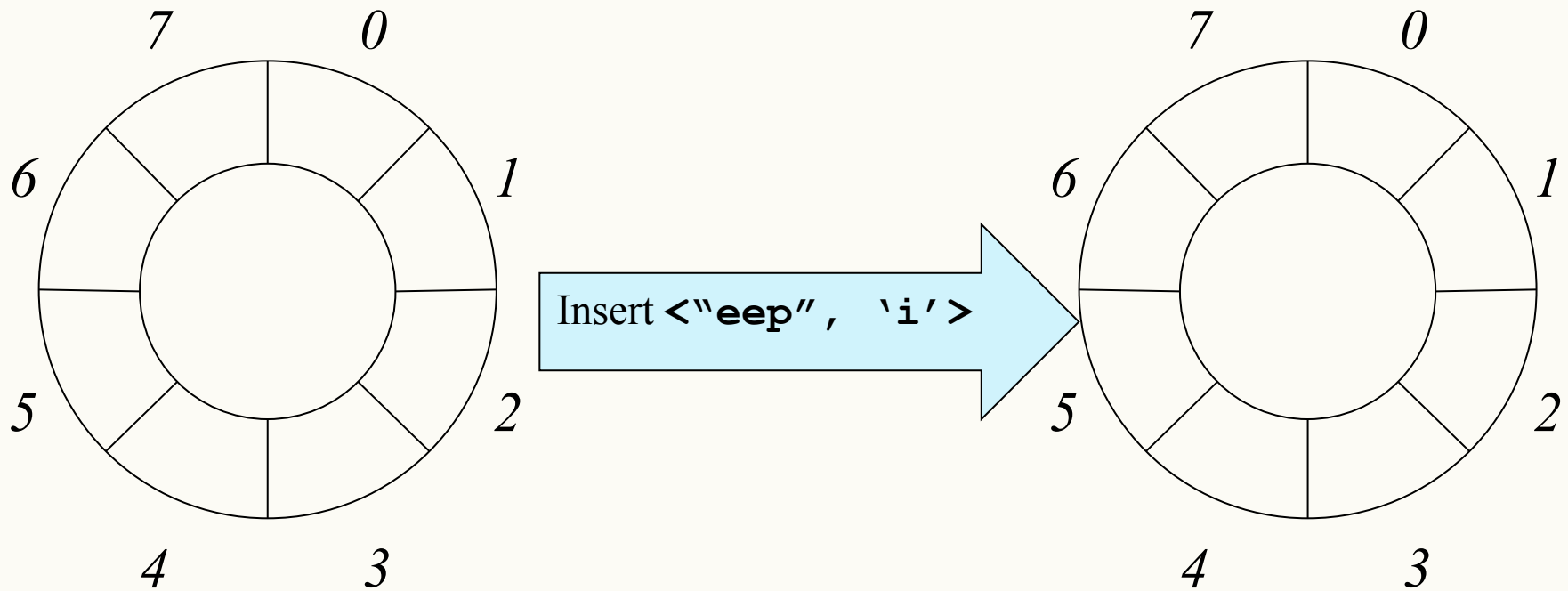- Some way to handle these collisions without resizing?

# Third Pass:
# Punt to Another Dictionary?



*Insert* `<13, 'o'>`
Insert `<37, 'i'>`

`<13, 'o'>`
`<37, 'i'>`

*BST, AVL, linked list, or other dictionary*

*When should we resize in this case?*

# How Do We Turn Strings into Numbers?

```
      7         0                          7         0
  6                 1                   6                 1

  5                 2                   5                 2
      4         3                          4         3
```
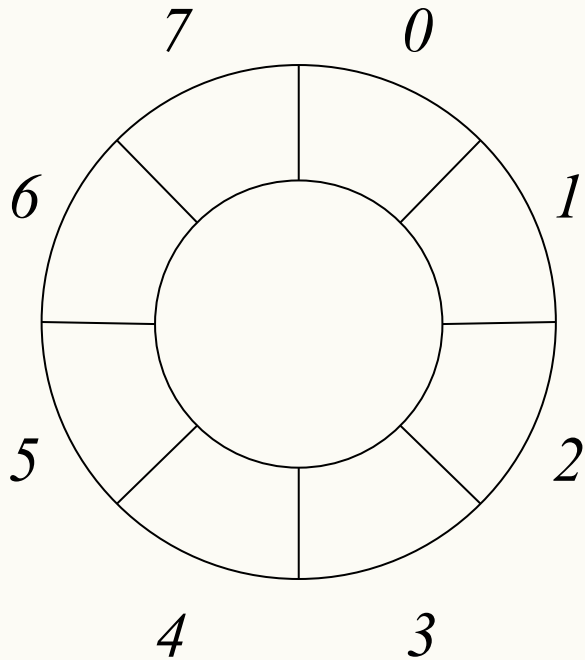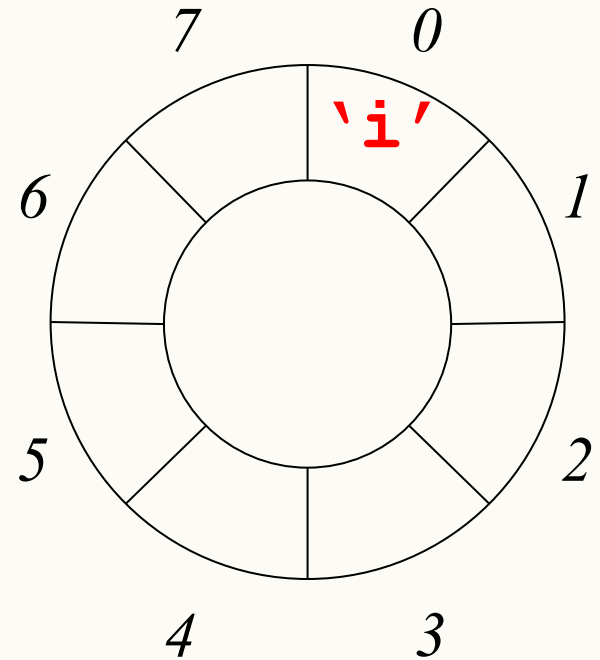
Insert **<"eep", 'i'>**

*What should we do?*

# Fourth Pass: Strings **ARE** Numbers

| e | e | p |
|---|---|---|
| 01100101 | 01100101 | 01110000 |

*= 6,645,104*

Insert **<"eep", 'i'>**

**'i'**

*6,645,104 % 8 = 0*

# Fourth Pass: Strings **ARE** Numbers

| e | e | p |
|---|---|---|
| 01100101 | 01100101 | 01110000 |

*= 6,645,104*

Insert **<"eep", 'i'>**

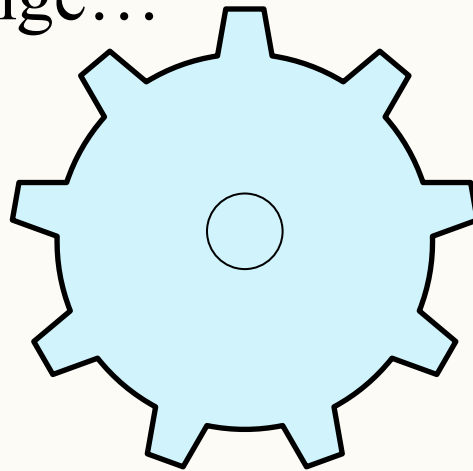Those numbers get REALLY big antidisestablishmentarianism. Just saying.

*6,645,104 % 8 = 0*

# Fifth Pass: Hashing!

- We only need perhaps a 64 (128?) bit number. There's no point in forming a **huge** number.

- We need a function to turn the strings into numbers, typically on a bounded range…

*antidisestablishmentarianism* → → *1,097,757,801*

*Maybe we can only use some parts of the string*

# Schlemiel, Schlemazel, Trouble for Our Hash Table?

- Let's try out:
  - "schlemiel" and "schlemazel"?
  - "microscopic" and "telescopic"?
  - "abcdefghijklmnopqrstuvwxyzyxwvutsrqponmlkjihgfedcba" and "abcdefghijklmnopqrstuvwxyzzyxwvutsrqponmlkjihgfedcba"

- Which bits of the string should we keep? Does our hash table care?

That's hashing! Take our data and turn it into a sorta-random number, ideally one that spreads out similar strings far apart!

# Third Pass, Take Two: Punt to Another Slot?



Insert **<13, 'o'>**
Insert **<37, 'i'>**

Slot 5 is full, but no "dictionaries in each slot" this time. Overflow to slot 6? When should we resize?

# Hash Table Approach

*Alan*

*Steve*

*Hassan*

*Will*

*Ed*

*f(x)*

*But… is there a problem in this pipe-dream?*

# Hash Table Dictionary Data Structure

- Hash function: maps keys to integers

  – result: can quickly find the right spot for a given entry

- Unordered and sparse table

  – result: cannot efficiently list all entries, *definitely* cannot efficiently list all entries in order or list entries between one value and another (a "range" query)

*Alan*

*Steve*

*Hassan*

*Will*

*Ed*

*f(x)*

# Hash Table Terminology

*hash function*

*Alan*

*Steve*

*Hassan*

*Will*

*Ed*

*f(x)*

*collision*

*keys*

$$\text{load factor } \lambda = \frac{\text{\# of entries in table}}{\text{tableSize}}$$

# Hash Table Code First Pass

```
Value & find(Key & key) {
  int index = hash(key) % tableSize;
  return Table[index];
}
```

- What should the hash function be?

- What should the table size be?

- How should we resolve collisions?

# A Good (Perfect?) Hash Function…

- is easy (fast) to compute
  - O(1) *and* fast in practice.

- distributes the data evenly
  - hash(a) % size ≠ hash(b) % size.

- uses the whole hash table. for all $0 \leq k <$ size, there's an i such that
  - hash(i) % size = k.

# Good Hash Function for Integers

- Choose
  - tableSize is
    - **prime** for good spread
    - **power of two** for fast calculations/convenient size
  - hash(n) = n (fast and good enough?)

| | |
|---|---|
| Insert 2 | 0: *14* |
| Insert 5 | 1: |
| Insert 10 | 2: *2* |
| Find 10 | 3: *10* |
| Insert 14 | 4: |
| Insert -1 | 5: *5* |
| | 6: *-1* |

# Good Hash Function for Strings?

Suppose we have a table capable of holding 5000 records, and whose keys consist of strings that are 6 characters long. We can apply numeric operations to the ASCII codes of the characters in the string in order to determine a hash index:

```
int hash(char * key){
  int  hashCode = 0;
  int  index    = 0;
  while (key[index] != '\0'){
    hashCode += (int)key[index];
    index++;
  }
  return  hashCode % 5000;
}
```

| C | A | M | E | C | O | /0 |
|---|---|---|---|---|---|----|

$C = 67$
$A = 65$
$M = 77$
$E = 69$
$C = 67$
$O = 79$
$\quad = 424$

| | | | CAMECO | | |
|---|---|---|---|---|---|

*0*                 *424*      *4999*

*Is this a good idea?*

# Good Hash Function for Strings?

- What is a significant problem with this approach?

  - Hash of any string with the same 6 letters is the same

  - ASCII values have a max of 255

    - 6*255 = 1530, which means [1531 – 4999] are wasted

- Alternative approach

- Let $s = s_1 s_2 s_3 s_4 \ldots s_5$: choose

  – $hash(s) = s_1 + s_2 128 + s_3 128^2 + s_4 128^3 + \ldots + s_n 128^n$

- Problems:

  – hash("really, really big") is really, really big!

  – hash("one thing") % 128 = hash("other thing") % 128

# Making the String Hash Easy to Compute

- Use Horner's Rule (Qin's Rule?)

$$a + bx + cx^2 = a + x(b + xc)$$

```
int hash(string s) {
  h = 0;
  for (i = s.length() - 1; i >= 0; i--) {
    h = (s_i + 31*h) % tableSize;
  }
  return h;
}
```

$$hash(help) = h+31(e+31(l+31*p))$$

*You would also need to %*

# Hash Function Summary

- Goals of a hash function
  - reproducible mapping from key to table entry
  - evenly distribute keys across the table
  - separate commonly occurring keys (neighbouring keys?)
  - complete quickly

# How to Design a Hash Function

- Know what your keys are *or*
  Study how your keys are distributed.

- Try to include all important information in a key in the construction of its hash.

- Try to make "neighbouring" keys hash to very different places.

- Balance complexity/runtime of the hash function against spread of keys (very application dependent).

# The Pigeonhole Principle (informal)

You can't put k+1 pigeons into k holes without putting two pigeons in the same hole.

**This place just isn't coo anymore.**



*Image by en:User:McKay, used under CC attr/share-alike.*

# Clicker question

Suppose we have 5 colours of Halloween candy, and that there's lots of candy in a bag. How many pieces of candy do we have to pull out of the bag if we want to be sure to get 2 of the same colour?

a. 2
b. 4
c. 6
d. 8
e. *None of these*

# Clicker question (answer)

Suppose we have 5 colours of Halloween candy, and that there's lots of candy in a bag. How many pieces of candy do we have to pull out of the bag if we want to be sure to get 2 of the same colour?

*a.* 2
*b.* 4
*c.* 6
*d.* 8
*e.* None of these

# The Pigeonhole Principle (formal)

Let X and Y be finite sets where $|X| > |Y|$.

If $f : X \rightarrow Y$, then $f(x_1) = f(x_2)$ for some $x_1, x_2 \in X$, where $x_1 \neq x_2$.

*Now that's coo!*

**f**

$x_1$

**X**

**Y**

$x_2$

$f(x_1) = f(x_2)$

# The Pigeonhole Principle (Example #2)

If there are 1000 pieces of each colour, how many do we need to pull to guarantee that we'll get 2 *black* pieces of candy (assuming that black is one of the 5 colours)?

a. 2

b. 6

c. 4002

d. 5001

e. None of these

# The Pigeonhole Principle (Example #2)

If there are 1000 pieces of each colour, how many do we need to pull to guarantee that we'll get 2 *black* pieces of candy (assuming that black is one of the 5 colours)?

a. 2

b. 6

c. 4002

d. 5001

e. None of these

This is not an appropriate problem for the pigeonhole principle! We don't know **which** hole has two pigeons!

# The Pigeonhole Principle (Example #3)

If 5 points are placed in a 6cm x 8cm rectangle, argue that there are two points that are not more than 5 cm apart.

*8cm*

*6cm*

Hint: How long is the diagonal?

# The Pigeonhole Principle (Example #4)

For $a, b \in \mathbf{Z}$, we write *a divides b* as $a|b$, meaning $\exists\, c \in \mathbf{Z}$ such that $b = ac$.

Consider $n + 1$ distinct positive integers, each $\leq 2n$. Show that one of them must divide on of the others.

For example, if $n = 4$, consider the following sets:

$\{1, 2, 3, 7, 8\}$   $\{2, 3, 4, 7, 8\}$   $\{2, 3, 5, 7, 8\}$

Hint: Any integer can be written as $2^k * q$ where k is an integer and q is odd. E.g., $129 = 2^0 * 129$; $60 = 2^2 * 15$.

# The Pigeonhole Principle (Example #4)

For $a, b \in \mathbf{Z}$, we write *a divides b* as $a|b$, meaning $\exists\, c \in \mathbf{Z}$ such that $b = ac$.

Consider $n + 1$ distinct positive integers, each $\leq 2n$. Show that one of them must divide on of the others.

Hint: Any integer can be written as $2^k * q$ where k is an integer and q is odd. E.g., $129 = 2^0 * 129$;  $60 = 2^2 * 15$.

- $x_1, x_2, x_3, \ldots x_{n+1}$  $\rightarrow$  $x_i = 2^{ki} * q_i$
  - Holes = Odd numbers $\leq 2n$  (n)
  - Pigeons = $q_i$ s (n+1)

- By the PHP, with n+1 numbers there exists $q_i = q_j$

$$\frac{2^{ki}\, q_i}{2^{kj}\, q_j}$$

- Therefore, one is some multiple of 2 times the other.
  - $x_i|x_j$ or $x_j|x_i$

# Example revisited

- In a small company of 100 employees, each employee is assigned an Emp_ID number in the range 00000 - 99999.
    - U  (number of potential keys)=100,000
    - m (number of keys) = 100
    - n (space allocated) =?
        - Hopefully not much bigger than m
        - Maybe 200 or 300

- By the Pigeonhole Principle(PHP) multiple potential keys are mapped to the same slot, which introduces the possibility of collisions.
- As m gets larger there is a higher probability of collision.

# Clicker question

- Consider n people with random birthdays (i.e., with each day of the ear equally likely). How large does n need to be before there is at least a 50% chance that two people have the same birthday.

A: 23

B: 57

C: 184

D: 367

E: None of the above

# Clicker question (Birthday Paradox)

- Consider n people with random birthdays. How large does n need to be before there is at least a 50% chance that two people have the same birthday.

A: 23 → 50%

B: 57 → 99%

C: 184

D: 367 → 100%

E: None of the above

- Corollary: Even if we randomly hash only $\sqrt{2m}$ keys into m slots, we get a collision with probability > **0.5**.

# Department of Computer Science Undergraduate Events

## More details @
https://my.cs.ubc.ca/students/development/events

**Industry Panel: Cultivating a Career**

**In Vancouver's Thriving Tech Sector**

**Wed. Mar 25**

**5:30 pm – 6:30 pm**

**DMP 310**

**CSSS Boat Cruise**

**Sat., Apr 4**

**6 pm – 10:30 pm**

**501 Denman St.**

**Tickets on sale in CSSS Office,**

**ICCS 021**

# CPSC 221 Administrative Notes

- Programming project #1 is marked
  - Feedback on your mark will be emailed to you

- Programming Project #2 is posted
  - Due date: Tue, 07 Apr @ 21.00
  - No Milestones but make sure you start early

- Sample solution for Assignment #2 is posted

- PeerWise #4

# So, Where Were We?

- Benefits of Hashing
  - O(1) Access

- Good Hash functions
  - is easy (fast) to compute
  - distributes data evenly
  - Uses the whole hash table

- The Pigeonhole Principle
  - Examples of how to do proofs with it
  - How it is related to hashing

- The Birthday Paradox

# Collision Resolution

- What do we do when two keys hash to the same entry?
    - chaining: put little dictionaries in each entry

        *shove extra pigeons in one hole!*

    - open addressing: pick a next entry to try

# Hashing with Chaining

- Put a little dictionary at each entry
  - choose type as appropriate
  - common case is unordered move-to-front linked list (chain)

- Properties

  - $\lambda$ can be greater than 1
  - performance degrades with length of chains

$h(a) = h(d)$
$h(e) = h(b)$
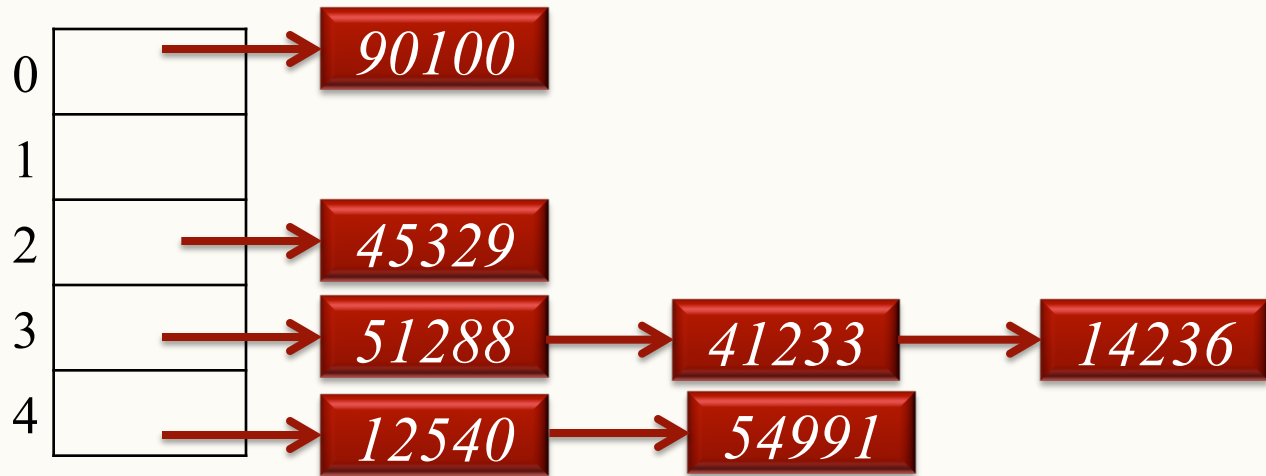


$$load\ factor\ \lambda = \frac{\#\ of\ entries\ in\ table}{tableSize}$$

# In-class exercise

Example:  Suppose $h(x) = \lfloor x/10 \rfloor \bmod 5$

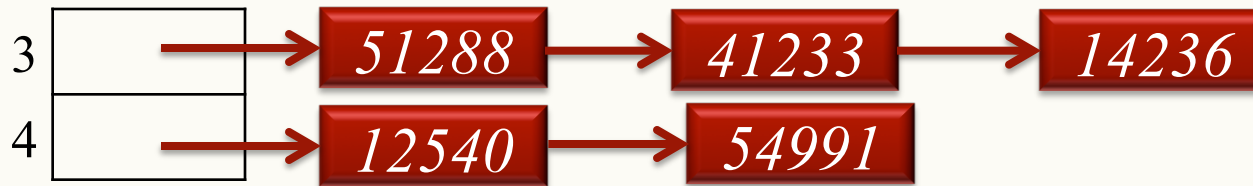Hash: 12540, 51288, 90100, 41233, 54991, 45329, 14236



Example: find node with key 14236

# Deleting when using chaining

Example: Suppose $h(x) = \lfloor x/10 \rfloor \bmod 5$

Hash: 12540, 51288, 41233, 54991, 14236

| 3 | → | 51288 | → | 41233 | → | 14236 |
| 4 | → | 12540 | → | 54991 | | |

- Delete 41233

- Remove 41233 from the linked list

| 3 | → | 51288 | → | 14236 |
| 4 | → | 12540 | → | 54991 |

# Chaining Code

```
Dictionary & findBucket(const Key & k) {
  return table[hash(k)%table.size];
}
```

```
void insert(const Key & k, const Value & v){
  findBucket(k).insert(k,v);
}
```

```
void delete(const Key & k){
  findBucket(k).delete(k);
}
```

```
Value & find(const Key & k){
  return findBucket(k).find(k);
}
```

# Load Factor in Chaining

$$\text{load factor } \lambda = \frac{\text{\# of entries in table}}{\text{tableSize}}$$

- Search cost

  – unsuccessful search:

  - On average $\lambda$

  – successful search:

  - On average $\sim\lambda/2 +1$ (what the book says)
  - More precisely $1 + \frac{n-1}{2m} = 1 + \frac{\lambda}{2} - \frac{\lambda}{2n}$

    *(n-1)/m are in this slot*

- Desired load factor:

  - between 1/2 and 1.

# Pros and cons of chaining

**Advantages of Chaining:**

- The size $s$ of the hash table can be smaller than the number of items $n$ hashed. Why is this often a good thing?
    - Fewer blank/wasted cells (especially in the case where the number of cells greatly exceeds the number of keys).
    - Collision handling can be O(1).
    - Can accommodate overflows

**Disadvantages of Chaining:**
- Search time can become O($n$) due to long chains.

# Open Addressing

What if we only allow one Key at each entry?

$h(a) = h(d)$
$h(e) = h(b)$

- two objects that hash to the same spot can't both go there

- first one there gets the spot

- next one must *go in another spot*

- Properties

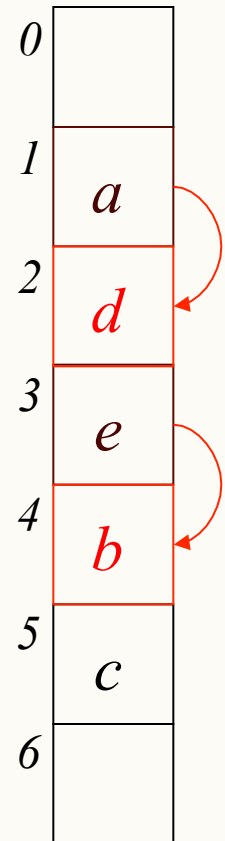$$load\ factor\ \lambda = \frac{\#\ of\ entries\ in\ table}{tableSize}$$

- $\lambda \leq 1$

- performance degrades with difficulty of finding right spot

| 0 | |
| 1 | |
| 2 | *a* |
| | *d* |
| 3 | *e* |
| 4 | *b* |
| 5 | *c* |
| 6 | |

# Probing

- Probing how to:
  - First probe - given a key k, hash to h(k)
  - Second probe - if h(k) is occupied, try h(k) + f(1)
  - Third probe - if h(k) + f(1) is occupied, try h(k) + f(2)
  - And so forth

- Probing properties

  - the $i^{th}$ probe is to (h(k) + f(i)) mod size      where f(0) = 0
  - if i reaches size, the insert has failed
  - depending on f(), the insert may fail sooner
  - long sequences of probes are costly!

# Linear Probing, f(i) = i

- Probe sequence is
  - h(k) mod size
  - h(k) + 1 mod size
  - h(k) + 2 mod size

  - …

- findEntry using linear probing:

```
bool findEntry(const Key & k, Entry *& entry) {
  int probePoint = hash1(k);
  do {
    entry = &table[probePoint];
    probePoint = (probePoint + 1) % size;
  } while (!entry->isEmpty() && entry->key != k);
  return !entry->isEmpty();
}
```

# In-class exercise

- Using the hash function $h(x) = x \% 7$ insert the following values using **linear probing**: *76, 93, 40, 47, 10, 55*

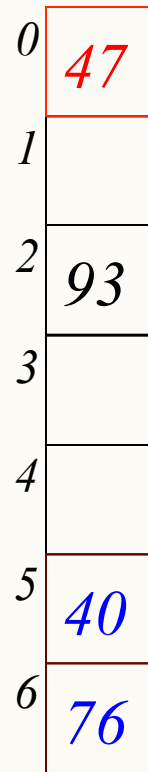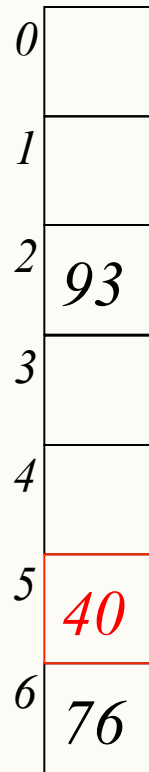| *insert(76)*<br>*76%7 = 6* | *insert(93)*<br>*93%7 = 2* | *insert(40)*<br>*40%7 = 5* | *insert(47)*<br>*47%7 = 5* | *insert(10)*<br>*10%7 = 3* | *insert(55)*<br>*55%7 = 6* |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 *47* | 0 *47* | 0 *47* |
| 1 | 1 | 1 | 1 | 1 | 1 *55* |
| 2 | 2 *93* | 2 *93* | 2 *93* | 2 *93* | 2 *93* |
| 3 | 3 | 3 | 3 | 3 *10* | 3 *10* |
| 4 | 4 | 4 | 4 | 4 | 4 |
| 5 | 5 | 5 *40* | 5 *40* | 5 *40* | 5 *40* |
| 6 *76* | 6 *76* | 6 *76* | 6 *76* | 6 *76* | 6 *76* |

*probes:*  *1*     *1*     *1*     *3*     *1*     *3*

# Load Factor in Linear Probing

$$load\ factor\ \lambda = \frac{\#\ of\ entries\ in\ table}{tableSize}$$

- For *any* $\lambda < 1$, linear probing will find an empty slot
- Search cost (for large table sizes)
  - successful search: $\frac{1}{2}\left(1 + \frac{1}{(1-\lambda)}\right)$

| | λ=0.25 | λ=0.5 | λ=0.75 | λ=0.9 |
|---|---|---|---|---|
| Avg # slots searched | 1.17 | 1.5 | 2.5 | 5.5 |

  - unsuccessful search: $\frac{1}{2}\left(1 + \frac{1}{(1-\lambda)^2}\right)$
    - How performance degrades as $\lambda$ gets bigger

| | λ=0.25 | λ=0.5 | λ=0.75 | λ=0.9 |
|---|---|---|---|---|
| Avg # slots searched | 1.4 | 2.5 | 8.5 | 50.5 |

# Load Factor in Linear Probing

$$load\ factor\ \lambda = \frac{\#\ of\ entries\ in\ table}{tableSize}$$

- For *any* $\lambda < 1$, linear probing will find an empty slot
- Search cost (for large table sizes)
  - successful search: $\dfrac{1}{2}\left(1 + \dfrac{1}{(1-\lambda)}\right)$
  - unsuccessful search: $\dfrac{1}{2}\left(1 + \dfrac{1}{(1-\lambda)^2}\right)$

  *Values hashed close to each other probe the same slots.*

- Linear probing suffers from ***primary clustering***
- Performance quickly degrades for $\lambda > 1/2$

# Quadratic Probing, f(i) = i²

- Probe sequence is
  - h(k) mod size

  - (h(k) + 1) mod size

  - (h(k) + 4) mod size

  - (h(k) + 9) mod size

  - …

- findEntry using quadratic probing:

```cpp
bool findEntry(const Key & k, Entry *& entry) {
  int probePoint = hash1(k), i = 0;
  do {
    entry = &table[(probePoint + i*i) % size];
    i++;
  } while (!entry->isEmpty() && entry->key != key);
  return !entry->isEmpty();
}
```

# Quadratic Probing, (more efficient code)

- Probe sequence is
  - h(k) mod size

  - (h(k) + 1) mod size

  - (h(k) + 4) mod size

  - (h(k) + 9) mod size

- findEntry using quadratic probing:

```
bool findEntry(const Key & k, Entry *& entry) {
  int probePoint = hash1(k), i = 0;
  do {
    entry = &table[probePoint];
    i++;
    probePoint = (probePoint + 2*i – 1) % size;
  } while (!entry–>isEmpty() && entry–>key != key);
  return !entry–>isEmpty();
}
```

# Quadratic Probing Example ☺

- Using the hash function $h(x) = x \% 7$ insert the following values using **quadratic probing**: *76, 40, 48, 5, 55*

| *insert(76)* $76\%7 = 6$ | *insert(40)* $40\%7 = 5$ | *insert(48)* $48\%7 = 6$ | *insert(5)* $5\%7 = 5$ | *insert(55)* $55\%7 = 6$ |
|---|---|---|---|---|
| 0 | 0 | 0 **48** | 0 48 | 0 *48* |
| 1 | 1 | 1 | 1 | 1 |
| 2 | 2 | 2 | 2 **5** | 2 5 |
| 3 | 3 | 3 | 3 | 3 **55** |
| 4 | 4 | 4 | 4 | 4 |
| 5 | 5 **40** | 5 40 | 5 *40* | 5 40 |
| 6 **76** | 6 76 | 6 *76* | 6 *76* | 6 *76* |

| probes: | 1 | 1 | 2 | 3 | 3 |

# Quadratic Probing Example ☹

- Using the hash function $h(x) = x \% 7$ insert the following values using quadratic probing: *76, 93, 40, 35, 47*

*insert(76)*
76%7 = 6

*insert(93)*
93%7 = 2

*insert(40)*
40%7 = 5

*insert(35)*
35%7 = 0

*insert(47)*
47%7 = 5

| | insert(76) | insert(93) | insert(40) | insert(35) | insert(47) |
|---|---|---|---|---|---|
| 0 | | | | 35 | 35 |
| 1 | | | | | |
| 2 | | 93 | 93 | 93 | 93 |
| 3 | | | | | |
| 4 | | | | | |
| 5 | | | 40 | 40 | 40 |
| 6 | 76 | 76 | 76 | 76 | 76 |

*probes:*    1        1        1        1        ∞

# Quadratic Probing Succeeds (for λ ≤ ½)

- Claim: If size is prime, the first size/2 probes are distinct

- Proof (By contradiction) Suppose for some `0 ≤ i< j ≤ size/2`

  `(h(x) + i²) mod size= (h(x) + j²) mod size`

  ` i² mod size = j² mod size`

  `(i² - j²) mod size= 0`

  `[(i + j)(i - j)] mod size= 0`


  but how can `i + j = 0` or `i + j = size` when

  `i ≠ j` and `i,j ≤ size/2`?

  same for `i - j mod size = 0`


- Result: If size is prime and λ ≤ ½, then quadratic probing will find an empty slot in size/2 probes or fewer

# Load Factor in Quadratic Probing

- For *any* λ ≤ ½, quadratic probing will find an empty slot; for greater λ, quadratic probing *may* find a slot

- Quadratic probing does not suffer from primary clustering

- Quadratic probing *does* suffer from **secondary clustering**

  – How could we possibly solve this?

*Values hashed to the SAME index probe the same slots.*

# Double Hashing, $f(i) = i \cdot hash2(x)$

- Probe sequence is
  - $h_1(k)$ mod size
  - $(h_1(k) + 1 \cdot h_2(x))$ mod size
  - $(h_1(k) + 2 \cdot h_2(x))$ mod size
  - …

```
bool findEntry(const Key & k, Entry *& entry) {
  int probePoint = hash1(k), hashIncr = hash2(k);
  do {
    entry = &table[probePoint];
    probePoint = (probePoint + hashIncr) % size;
  } while (!entry->isEmpty() && entry->key != k);
  return !entry->isEmpty();
}
```

# A Good Double Hash Function…

- is quick to evaluate.

- differs from the original hash function.

- never evaluates to 0 (mod size).


- One good choice is to choose a prime R < size
  - $hash_2(x) = R - (x \bmod R)$

# Double Hashing Example

- Using the hash functions $h_1(x) = x \% 7$ and $h_2(x) = 5 - (x \% 5)$ insert the following values using **double hashing** *76, 93, 40, 47, 10, 55*

| *insert(76)* | *insert(93)* | *insert(40)* | *insert(47)* | *insert(10)* | *insert(55)* |
|---|---|---|---|---|---|
| *76%7 = 6* | *93%7 = 2* | *40%7 = 5* | *47%7 = 5* | *10%7 = 3* | *55%7 = 6* |
| | | | *5 - (47%5) = 3* | | *5 - (55%5) = 5* |

|   | insert(76) | insert(93) | insert(40) | insert(47) | insert(10) | insert(55) |
|---|---|---|---|---|---|---|
| 0 | | | | | | |
| 1 | | | | 47 | 47 | 47 |
| 2 | | 93 | 93 | 93 | 93 | 93 |
| 3 | | | | | 10 | 10 |
| 4 | | | | | | 55 |
| 5 | | | 40 | 40 | 40 | 40 |
| 6 | 76 | 76 | 76 | 76 | 76 | 76 |

*probes:*   *1*      *1*      *1*      *2*      *1*      *2*

# Clicker question

The primary hash function is:  $h_1(k) = (2k + 5) \bmod 11$.
The secondary hash function is:  $h_2(k) = 7 - (k \bmod 7)$

Hash these keys, in this order: *12, 44, 13, 88, 23, 94, 11*.
Which cell in the array does key 11 hash to?

**A.0**

**B.2**

**C.3**

**D.4**

**E.10**

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

# Clicker question (answer)

$h_1(k) = (2k + 5) \bmod 11$.  $h_2(k) = 7 - (k \bmod 7)$

*12, 44, 13, 88, 23, 94, 11*. Which cell in the array does key 11 hash to?

*h(12) = (2(12) + 5 ) % 11 = 7*

*h(44) = (2(44) + 5 ) % 11 = 5*

**A.0**

*h(13) = (2(13) + 5 ) % 11 = 9*

**B.2**

**C.3**   *h(88) = (2(88) + 5 ) % 11 = 5 +7 – 88%7 = 8*

**D.4**   *h(23) = (2(23) + 5 ) % 11 = 7 + 7 – 23%7 = 12*

**E.10**  *h(94) = (2(94) + 5 ) % 11 = 6*

*h(11) = (2(11) + 5 ) % 11 = 5 +2( 7 – 11%7) = 11*

| | |
|---|---|
| 0 | *11* |
| 1 | *23* |
| 2 | |
| 3 | |
| 4 | |
| 5 | *44* |
| 6 | *94* |
| 7 | *12* |
| 8 | *88* |
| 9 | *13* |
| 10 | |

# Load Factor in Double Hashing

- For *any* $\lambda < 1$, double hashing will find an empty slot (given appropriate table size and hash$_2$)

- Search cost appears to approach optimal (random hash):

  – successful search: $\dfrac{1}{\lambda} \ln \dfrac{1}{1-\lambda}$

|  | $\lambda$=0.25 | $\lambda$=0.5 | $\lambda$=0.75 | $\lambda$=0.9 |
|---|---|---|---|---|
| Avg # slots searched | 1.5 | 1.4 | 1.8 | 2.6 |

  – unsuccessful search: $\dfrac{1}{1-\lambda}$

|  | $\lambda$=0.25 | $\lambda$=0.5 | $\lambda$=0.75 | $\lambda$=0.9 |
|---|---|---|---|---|
| Avg # slots searched | 1.3 | 2 | 4 | 10 |

- No primary clustering and no secondary clustering
- One extra hash calculation

# Deleting when using probing

**Example:**

- Suppose locations [97] to [101] are occupied in our hash table:

| [0] | [1] | … | [97] | [98] | [99] | [100] | [101] | [102] | … |
|-----|-----|---|------|------|------|-------|-------|-------|---|
|     |     |   | dog  | cat  | goat | owl   | deer  |       |   |

- Suppose a new key hashes to [97]. Assuming a linear collision resolution policy, the key goes to 102.

- Later, suppose we delete the key that was hashed to [98].

- Add a tombstone (i.e., flag, marker) for 2 reasons:
    1. If searching, keep going when you hit a tombstone.

    2. If inserting, stop and add the item here.

This means that table entries can be occupied, deleted, or free

# The Squished Pigeon Principle

- An insert using open addressing *cannot* work with a load factor of 1 or more.

- An insert using open addressing with quadratic probing may not work with a load factor of ½ or more.

- Whether you use chaining or open addressing, large load factors lead to poor performance!

- How can we relieve the pressure on the pigeons?

*Hint: think resizable arrays!*

# Rehashing

- When the load factor gets "too large" (over a constant threshold on $\lambda$), rehash all the elements into a new, larger table:
    - takes $O(n)$, but amortized $O(1)$ as long as we (just about) double table size on the resize
    - spreads keys back out, may drastically improve performance
    - gives us a chance to retune parameterized hash functions
    - avoids failure for open addressing techniques
    - allows arbitrarily large tables starting from a small table
    - clears out lazily deleted items

# Practice: Open Addressing
## (Try linear, quadratic, %7/(1-%5) double hashing.)

Insert 2
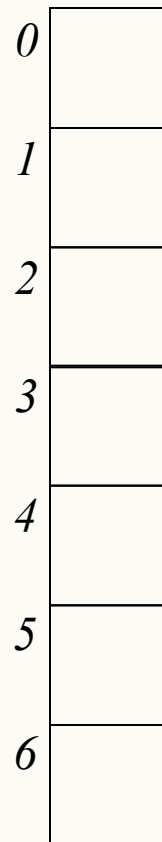
Insert 5

Insert 4

Insert 10

Insert 73

Find 10

Insert 14

Resize/Rehash

Insert -1

Insert 3

| | |
|---|---|
| *0* | |
| *1* | |
| *2* | |
| *3* | |
| *4* | |
| *5* | |
| *6* | |

# Application: De-Duplication

- Given a "stream" of objects
  - Linear scan through a huge file
  - Objects arriving in real time

- Goal: Remove duplicates (keep track of unique objects)
  - Report unique visitors to a web site
  - Avoid duplicates in search results

- Solution: When new object x arrives, look up $h(x)$ and if not found insert.

# Application: The 2-Sum Problem

- Given: Unsorted array of integers $A$, and a target sum $t$
- Goal: Determine whether or not there are two numbers $x$ and $y$ in $A$ such that $x+y=t$

- Naïve solution: *$O(n^2)$* exhaustive search
- Better solution:
  - Sort $A$ *$O(n \lg n)$*
  - For each $x$ in A look for $t-x$ *$O(n \lg n)$*
- Amazing solution:
  - Insert elements of $A$ into hash table $H$ *$O(n)$*
  - For each $x$ in $A$, lookup $t-x$ in $H$ *$O(n)$*

# CPSC Administrative Notes

- Lab 10 Parallelism
  - Starting tomorrow Mar 26 – Apr 2
  - Marking Apr 7 – Apr 10 (Also doing Concept Inventory)
  - Your work on this concept inventory helps us improve our courses! In particular, we use the inventory to track how our approach to teaching foundations of computing concepts affects students' learning over time.

- PeerWise Call #5 due Apr 2 (5pm)
  - The deadline for contributing to your "Answer Score" and "Reputation score" is Monday April 20.

# CPSC Administrative Notes

- Programming project #2 due
  - Apr Tue, 07 Apr @ 21.00
  - How to work on programming projects in pairs?

- We're making the following change to the old scheme, which provides a great opportunity to improve your final grade! The main purpose of this change is to give you more incentive to study for the final.
  - Midterm Exam                                    10%
  - Midterm or Final Exam (your best)      10%

# The Pigeonhole Principle (Full Glory)

- Let X and Y be finite sets with $|X| = n$, $|Y| = m$, and $k = \lceil n/m \rceil$.

If $f : X \rightarrow Y$, then $\exists\ k$ values $x_1, x_2, \ldots, x_k \in X$ such that $f(x_1) = f(x_2) = \ldots f(x_k)$.

Informally: If $n$ pigeons fly into $m$ holes, at least 1 hole contains at least $k = \lceil n/m \rceil$ pigeons.

Proof: Assume there's no such hole. Then, there are at most $(\lceil n/m \rceil - 1)*m$ pigeons in all the holes, which is fewer than $(n/m + 1 - 1)*m = n/m*m = n$, but that is a contradiction. QED

# Pathological Data Sets

- For good hash performance, we need a good hash function

  – Spreads data evenly across buckets

- Ideal: Use super-clever hash function guaranteed to spread every data set out evenly

- Problem: Such a hash function does not exist

  – For every hash function, there is a pathological data set

# Pathological Data Sets

- Reason
  - Fix a hash function $h$
  - Let $U$ be the potential number of keys
  - Let $n$ be the table size

- There exists an array cell $i$, such that at least $U/n$ elements hash to $i$ under $h$

- If data set drawn only from these elements, then everything collides.

- This data set could be quite large since $U >> n$

# Overview of Universal Hashing

- For every deterministic hash function, there is a pathological data set.
  - Solution: Do not commit to a specific hash function

- Use randomization
  - Design a family $H$ of hash functions, such that for every data set $S,$ most functions $h \in H$ spread $S$ out "pretty evenly"

# Learning Goals revisited

After this unit, you should be able to:

- Define various forms of the pigeonhole principle; recognize and solve the specific types of counting and hashing problems to which they apply.

- Provide examples of the types of problems that can benefit from a hash data structure.

- Compare and contrast open addressing and chaining.

- Evaluate collision resolution policies.

- Describe the conditions under which hashing can degenerate from O(1) expected complexity to O(n).

- Identify the types of search problems that do not benefit from hashing (e.g. range searching) and explain why.

- Manipulate data in hash structures both irrespective of implementation and also within a given implementation.