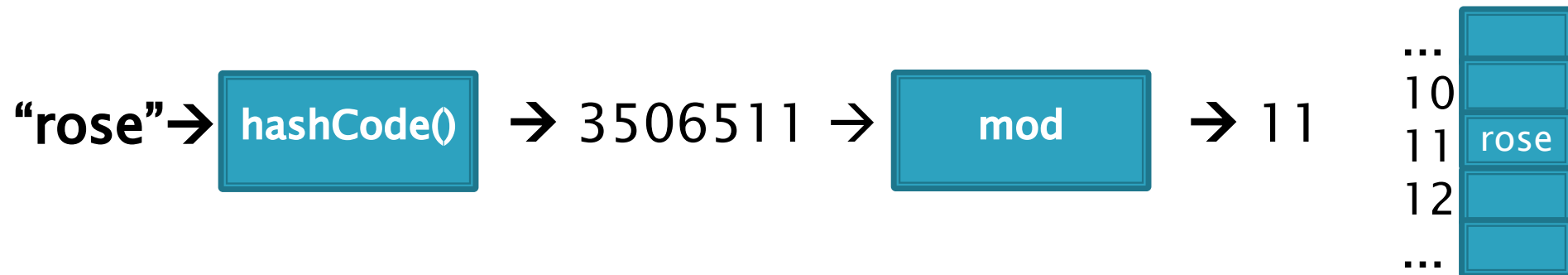


CSSE 230

Hash table basics

How can hash tables perform both **contains()** in $O(1)$ time and **add()** in amortized $O(1)$ time, given enough space?



Midterm feedback

Course – Plus

- Programming assignments help understanding ++++++
- Quizzes help focus lectures ++++++
- Lectures are clear +++++
- In-class coding, examples +++++
- Written homework reinforces material +++++
- Good pace, difficulty of homework. Challenging but manageable +++++
- Everything is clear, smooth ++
- Piazza is helpful +

Course – Delta

- ▶ None, everything good so far +++++
- ▶ Clearer directions on written assignments ++
- ▶ More individual programming ++
- ▶ Would like to choose teammates +
- ▶ Post solutions to quizzes +
- ▶ Get rid of quizzes, or collect them +
- ▶ More small coding questions on written assignments +
- ▶ Lecture sometimes too fast +
- ▶ Sometimes lectures are slow/repetitive +
- ▶ Go over programming assignment solutions, how to do it efficiently

Self – Plus

- Study/work hard, do all assignments & take them seriously ++++++
- Start early / aim to finish assignments early ++++++
- On assignments, solve as much as possible on own +++++
- Reflecting on knowledge +++++
- Taking notes in class ++
- Pay attention in class ++
- Study for exams ++
- Thinking & planning abstractly before starting to code +
- Textbook +
- Reviewing past quizzes +
- Get help, ask questions +
- Help from peers +
- Practice exams +

Self – Delta

- ▶ Start earlier / aim to finish early +++++
- ▶ Keep doing what I'm doing ++
- ▶ Study more for exams ++
- ▶ Ask timely questions about written homework ++
- ▶ More coding practice ++
- ▶ Read textbook more ++
- ▶ Reflect on knowledge, supplement ++
- ▶ Study more of written exam stuff ++
- ▶ Ask more questions when I don't understand +
- ▶ Study more in general +
- ▶ Practice exams +
- ▶ Find more time +

Surprise

- More math / theory than expected ++++_
- Lot of work, had to recalibrate how much effort to devote +++
- More programming than expected +
- Takes a lot of time, especially outside of class +
- Group projects are emphasized +
- Didn't know what to expect +
- Learning a lot +
- Very interesting / cool / fun +
- A lot of trees +
- I'm doing better than expected +
- Most of hard math was first couple weeks
- Lots of recursion
- Not as hard as expected (based on reputation)
- The fact that I really like it in spite of the workload
- Not high volume, but high difficulty

Announcements and questions

- ▶ Questions on HW6?
- ▶ Look at HW7

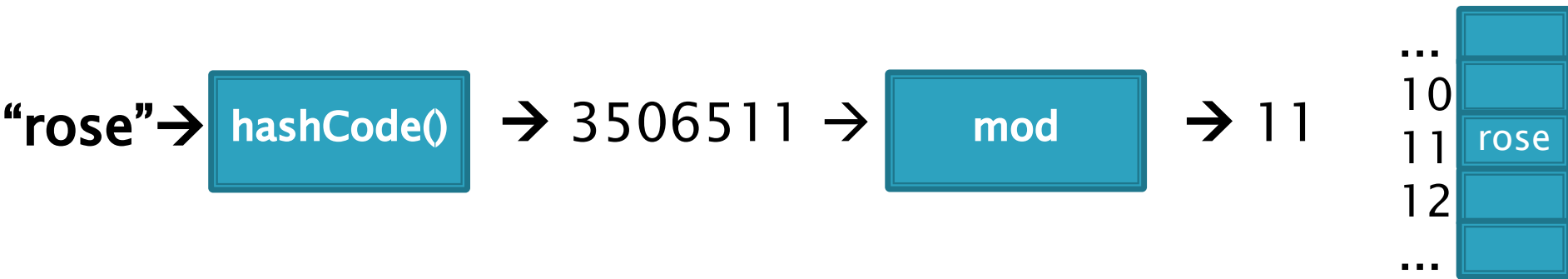
Hashing

Efficiently putting 5 pounds of
data in a 20 pound bag

Reminder: sets hold unique items

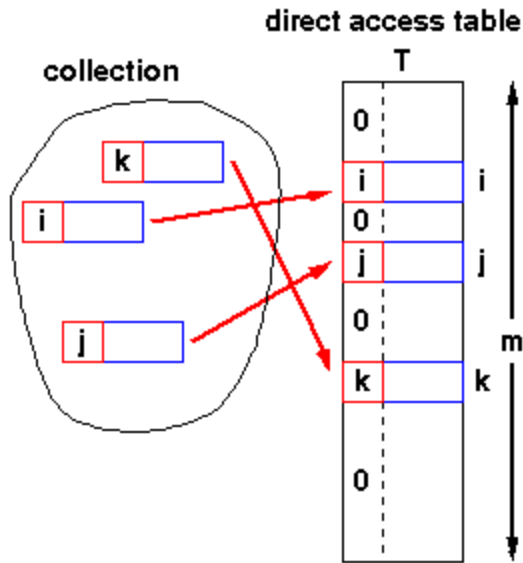
- ▶ **Implementation choices:**
 - **TreeSet** (and TreeMap) uses a balanced tree: $O(\log n)$
 - Uses a red–black tree
 - **HashSet** (and HashMap) uses a hash table: amortized $O(1)$ time
- ▶ Related: maps allow insertion, retrieval, and deletion of items by *key*.
 - Since keys are unique, they form a set.
 - The values just go along for the ride.
 - We'll focus on sets.

Big ideas of hash tables



1. The underlying storage?
Growable array
2. Calculate the index to store an item **from the item itself**. How?
Hashcode. Fast but un-ordered.
3. What if that location is already occupied with another item?
Collision. Two methods to resolve

Direct Address Tables



- ▶ Array of size **m**
- ▶ **n** elements with unique keys
- ▶ If all keys are $\leq m$, then use the key as an array index.
 - Clearly $O(1)$ lookup of keys

▶ Issues?

- Keys must be unique.
- Often the range of potential keys is much larger than the storage we want for an array
 - Example: RHIT student IDs vs. # Rose students

We attempt to create unique keys
by applying a `.hashCode()` function ...



Objects that are `.equals()`
MUST have the same `hashCode` values
A good `hashCode()` also
is **fast** to calculate and
distributes the keys, like:

```
hashCode("rose")= 3506511
```

```
hashCode("hulman")= -1206158341 (can be negative if overflows)
```

```
hashCode("institute") = 36682261
```

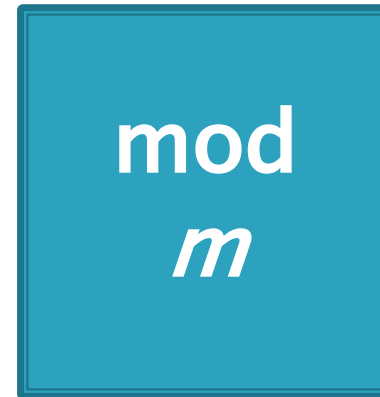
...and then take it mod the table size (m) to get an index into the array.

▶ Example: if $m = 100$:

hashCode("rose") = 3506511

hashCode("hulman") = -1206158341

hashCode("institute") = 36682261



→11

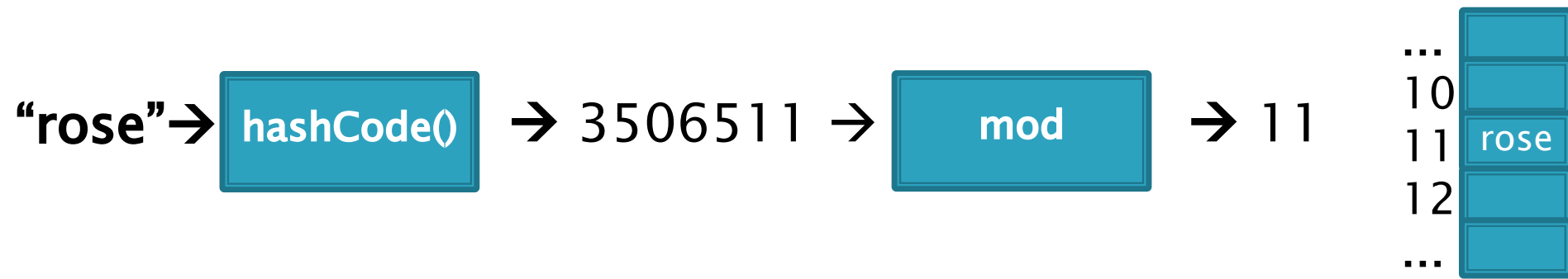
→07*

→61

- * Note: since the hashCode is an integer, it might be negative...
- If it is negative, add `Integer.MAX_VALUE + 1` to make it positive before you mod. (Same as ANDing with `0x7fffffff`, or removing sign bit from two's complement)
 - This mimics what's actually done in practice: when m is a power of 2, say 2^k , we can just truncate, keeping the last k bits (instead of taking mod m). Sign bit is lost.

Index calculated from the object itself, not from a comparison with other objects

- ▶ How Java's `hashCode()` is used:



- Unless this position is already occupied

a “collision”

- ▶ Default if you inherit `Object`'s: memory location (platform-specific, actually)
- ▶ Many JDK classes override `hashCode()`
 - Integer: the value itself
 - Double: XOR first 32 bits with last 32 bits
 - String: we'll see shortly!
 - Date, URL, ...
- ▶ Custom classes should override `hashCode()`
 - Use a combination of **final** fields.
 - If key is based on mutable field, then the hashcode will change and you will lose it!
 - People usually use strings if possible.

A simple hash function for Strings is a function of every character

```
// This could be in the String class
public static int hash(String s) {
    int total = 0;
    for (int i=0; i<s.length(); i++)
        total = total + s.charAt(i);
    return total;
}
```

- ▶ Advantages?
- ▶ Disadvantages?

A better hash function for Strings uses place value

```
// This could be in the String class
public static int hash(String s) {
    int total = 0;
    for (int i=0; i<s.length(); i++)
        total = total*256 + s.charAt(i);
    return total;
}
```

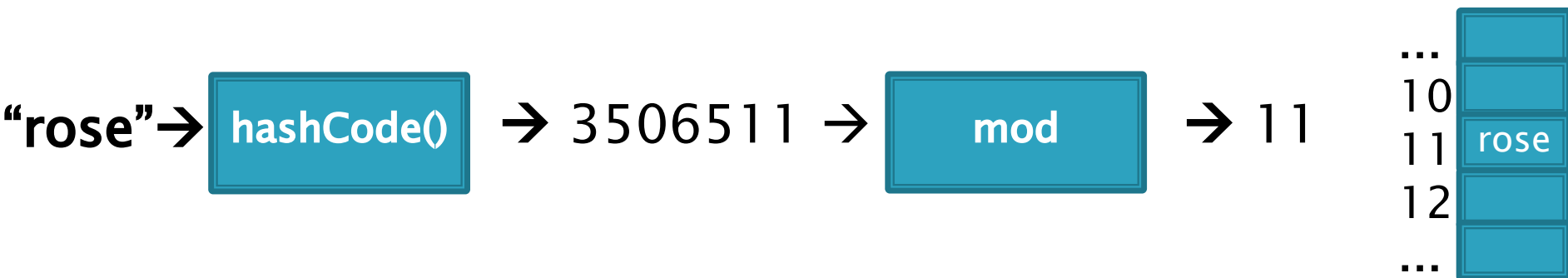
- ▶ Spreads out the values more, and anagrams not an issue.
- ▶ What about overflow during computation?
 - What happens to first characters?

A better hash function for Strings uses place value with a base that's prime

```
// This could be in the String class
public static int hash(String s) {
    int total = 0;
    for (int i=0; i<s.length(); i++)
        total = total*31 + s.charAt(i);
    return total;
}
```

- ▶ Spread out, anagrams OK, overflow OK.
- ▶ This is **String**'s `hashCode()` method.
- ▶ The $(x = 31x + y)$ pattern is a good one to follow.
- ▶ See <https://docs.oracle.com/javase/8/docs/api/java/lang/String.html#hashCode-->

Collisions are inevitable

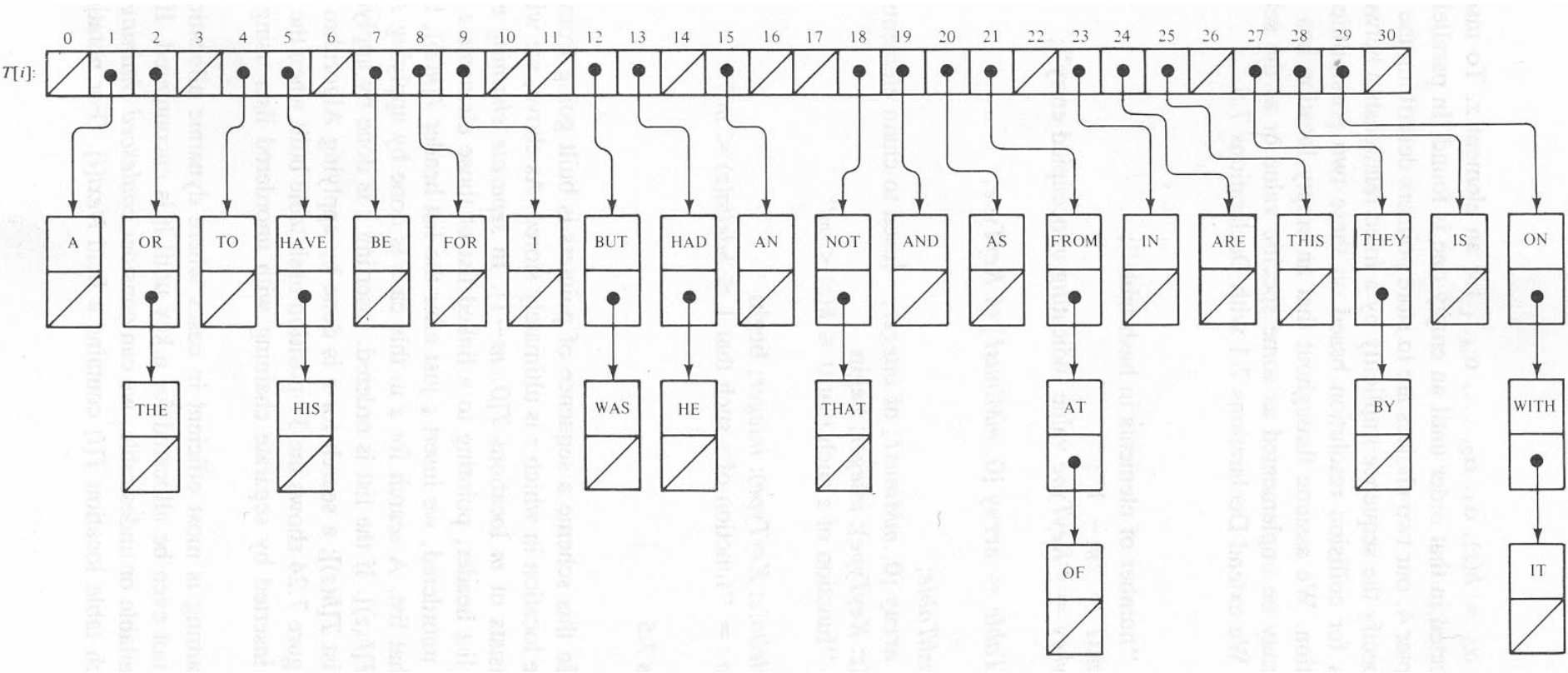


- ▶ A good hashcode distributes keys evenly, but collisions will still happen
- ▶ hashCode() are ints → only ~4 billion unique values.
 - How many 16 character ASCII strings are possible?
- ▶ If n is small, tables should be much smaller
 - mod will cause collisions too!
- ▶ Solutions:
 - Chaining
 - Probing (Linear, Quadratic)

Separate chaining: an array of linked lists

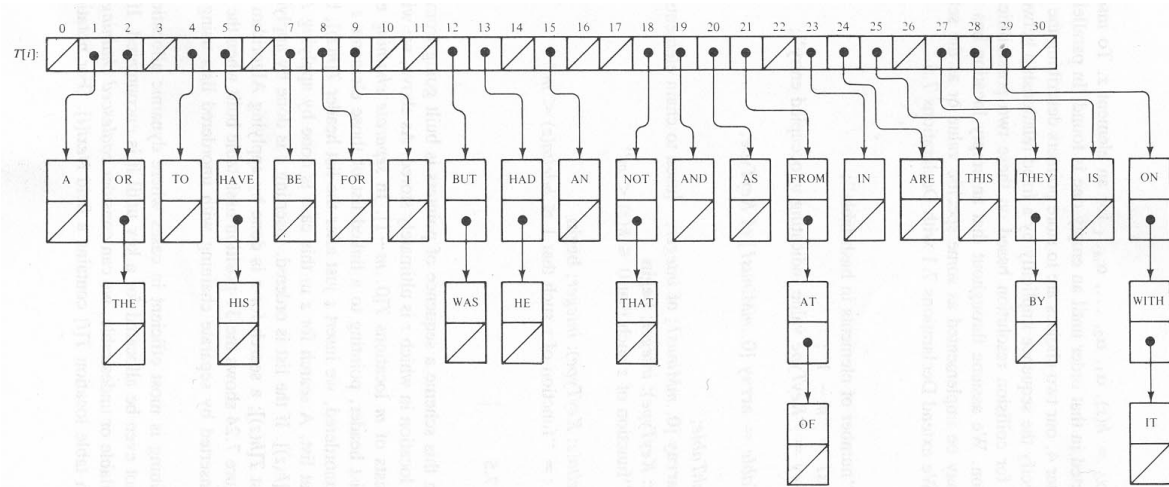
Grow in another direction

Examples: .get("at"), .get("him), (hashcode=18), .add("him"), .delete("with")



Java's **HashMap** uses chaining and a table size that is a power of 2.

Runtime of hashing with chaining depends on the load factor



m array slots,
 n items.

Load factor, $\lambda = n/m$.

Runtime = $O(\lambda)$

Space-time trade-off

1. If m constant, then this is $O(n)$. Why?

2. If keep $m \sim 0.5n$ (by doubling), then this is **amortized $O(1)$** . Why?

Alternative: Store collisions in other array slots.

- ▶ No need to grow in second direction
- ▶ No memory required for pointers
 - Historically, this was important!
 - Still is for some data...
- ▶ Will still need to keep load factor ($\lambda = n/m$) low or else collisions degrade performance
 - We'll grow the array again

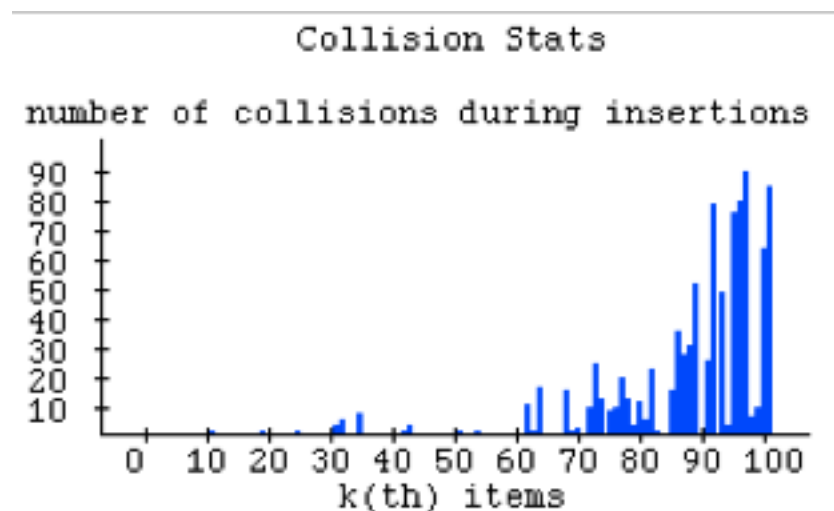
Collision Resolution: Linear Probing

- ▶ Probe H (see if it causes a collision)
- ▶ Collision? Also probe the next available space:
 - Try H , $H+1$, $H+2$, $H+3$, ...
 - Wraparound at the end of the array
- ▶ Example on board: `.add()` and `.get()`

- ▶ Problem: Clustering

- ▶ Animation:
 - http://www.cs.auckland.ac.nz/software/AlgAnim/hash_tables.html
 - Applet deprecated on most browsers.
 - See next slide for a few freeze-frames.

Clustering Example



Linear probing efficiency also depends on load factor, $\lambda = n/m$

- ▶ For probing to work, $0 \leq \lambda \leq 1$.
- ▶ For a given λ , what is the expected number of probes before an empty location is found?

Rough Analysis of Linear Probing

- ▶ Assume all locations are equally likely to be occupied, and equally likely to be the next one we look at.
- ▶ λ is the probability that a given cell is full, $1 - \lambda$ the probability a given cell is empty.
- ▶ What's the expected number?

$$\sum_{p=1}^{\infty} \lambda^{p-1} (1 - \lambda) p = \frac{1}{1 - \lambda}$$

From https://en.wikipedia.org/wiki/List_of_mathematical_series:

$$\sum_{k=1}^n k z^k = z \frac{1 - (n+1)z^n + n z^{n+1}}{(1-z)^2}$$

Better Analysis of Linear Probing

- ▶ **Clustering!**
 - Blocks of occupied cells are formed
 - Any collision in a block makes the block bigger
- ▶ Two sources of collisions:
 - Identical hash values
 - Hash values that hit a cluster
- ▶ Actual average number of probes for large λ :

$$\frac{1}{2} \left(1 + \frac{1}{(1 - \lambda)^2} \right)$$

For a proof, see Knuth, *The Art of Computer Programming, Vol 3: Searching Sorting*, 2nd ed, Addison-Wesley, Reading, MA, 1998.
(1st edition = 1968)

Why consider linear probing?

- ▶ Easy to implement
- ▶ Works well when load factor is low
 - In practice, once $\lambda > 0.5$, we usually **double the size of the array** and rehash
 - This is more efficient than letting the load factor get high
- ▶ Works well with caching

To reduce clustering, probe farther apart

- ▶ **Reminder: Linear probing:**
 - Collision at H ? Try $H, H+1, H+2, H+3, \dots$
- ▶ **New: Quadratic probing:**
 - Collision at H ? Try $H, H+1^2, H+2^2, H+3^2, \dots$
 - Eliminates primary clustering. “Secondary clustering” isn’t as problematic

Quadratic Probing works best with low λ and prime m

- ▶ **Choose a prime number for the array size, m**
- ▶ Then if $\lambda \leq 0.5$:
 - Guaranteed insertion
 - If there is a “hole”, we’ll find it
 - **So no cell is probed twice**
- ▶ **Can show with $m=17$, $H=6$.**

For a proof, see Theorem 20.4:

Suppose the table size is prime, and that we repeat a probe before trying more than half the slots in the table

See that this leads to a contradiction

Quadratic Probing runs quickly if we implement it correctly

Use an algebraic trick to calculate next index

- Difference between successive probes yields:
 - Probe i location, $H_i = (H_{i-1} + 2i - 1) \% M$
- 1. Just use bit shift to multiply i by 2
 - `probeLoc = probeLoc + (i << 1) - 1;`
...faster than multiplication
- 2. Since i is at most $M/2$, can just check:
 - `if (probeLoc >= M)`
 `probeLoc -= M;`
...faster than mod

When growing array, can't double!

- Can use, e.g., `BigInteger.nextProbablePrime()`

Quadratic probing analysis

- ▶ No one has been able to analyze it!
- ▶ Experimental data shows that it works well
 - Provided that the array size is prime, and $\lambda < 0.5$

Summary:

Hash tables are fast for some operations

Structure	insert	Find value	Find max value
Unsorted array			
Sorted array			
Balanced BST			
Hash table			

- ▶ Finish the quiz.
- ▶ Then check your answers with the next slide

Answers:

Structure	insert	Find value	Find max value
Unsorted array	Amortized $\theta(1)$	$\theta(n)$	$\theta(n)$
Sorted array	$\theta(n)$	$\theta(\log n)$	$\theta(1)$
Balanced BST	$\theta(\log n)$	$\theta(\log n)$	$\theta(\log n)$
Hash table	Amortized $\theta(1)$	$\theta(1)$	$\theta(n)$

In practice

- ▶ Constants matter!
- ▶ 727MB data, ~190M elements
 - Many inserts, followed by many finds
 - Microsoft's C++ STL

Structure	build (seconds)	Size (MB)	100k finds (seconds)
Hash map	22	6,150	24
Tree map	114	3,500	127
Sorted array	17	727	25

- ▶ Why?
- ▶ Sorted arrays are nice **if** they don't have to be updated frequently!
- ▶ Trees still nice when interleaved insert/find

Review: discuss with a partner

- ▶ Why use 31 and not 256 as a base in the String hash function?
- ▶ Consider chaining, linear probing, and quadratic probing.
 - What is the purpose of all of these?
 - For which can the load factor go over 1?
 - For which should the table size be prime to avoid probing the same cell twice?
 - For which is the table size a power of 2?
 - For which is clustering a major problem?
 - For which must we grow the array and rehash every element when the load factor is high?

Today's worktime

...is a great time to start HashSet while it's fresh

...is acceptable to use for EditorTrees Milestone 2 group worktime, especially if you have questions for me