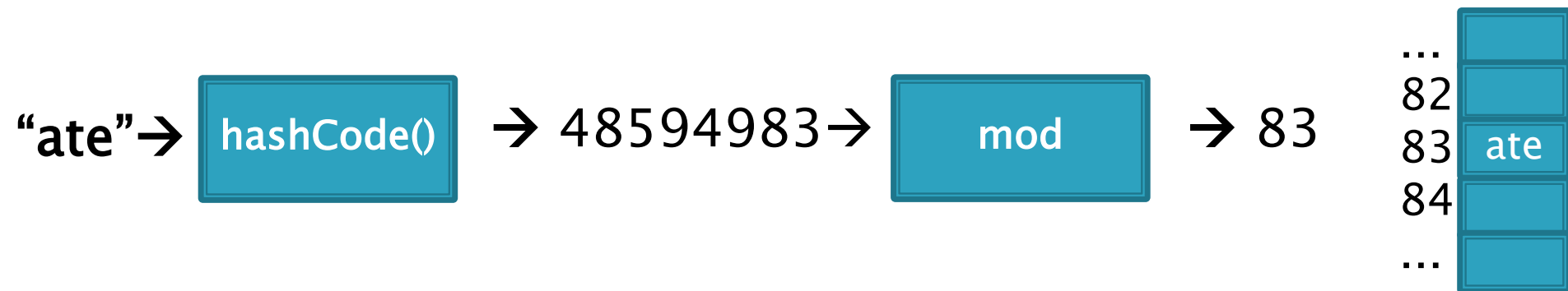


CSSE 230 Day 18

Hash table basics



Reminders / Announcements

- ▶ Reminder from syllabus: EditorTrees worth 10% of term grade
- ▶ See schedule page
 - Exam 2 moved to Friday after break.
- ▶ Short “pop” quiz over AVL rotations now

Exam 2

- ▶ Format same as Exam 1
 - One 8.5x11 sheet of paper (**2-sided**) for written part
 - Same resources as before for programming part
- ▶ Topics: weeks 1–6
 - Reading, programs, in-class, written assignments.
 - Especially
 - Using various data structures (lists, stacks, queues, sets, maps, priority queues)
 - Binary trees, including BST, AVL, and threaded
 - Traversals and iterators, size vs. height, rank
 - Backtracking / Queens problem
 - Hash tables
 - Algorithm analysis in general
- ▶ Through day 19, WA6, and EditorTrees milestone 2

T
F
IDK

Sample exam on Moodle has some good questions (and extras we haven't done, like sorting)
Best practice: written assignments.

Questions

Agenda

- ▶ Hash table basics
- ▶ Collision resolution
- ▶ EditorTrees work time

Hashing

Efficiently putting 5 pounds of
data in a 20 pound bag

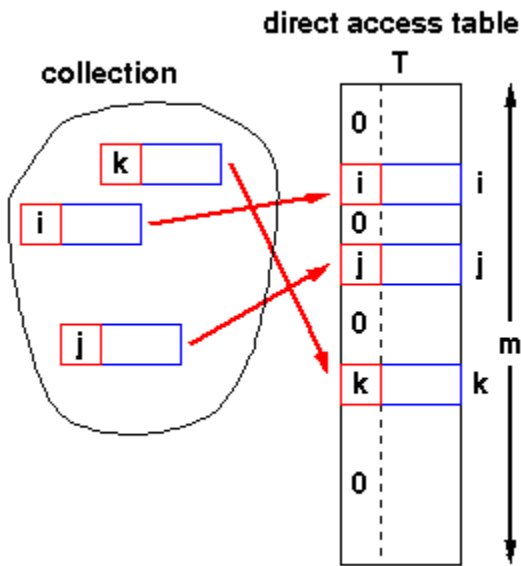
A hash table is a very fast approach to dictionary storage

- ▶ Provides rapid insertion, retrieval, and deletion of items by *key*
- ▶ **HashMap** uses a hash table internally
 - Actual table data is stored in an array
 - **HashSet** uses a **HashMap** internally
- ▶ Insertion and lookup are constant time!
 - With a good “hash function”
 - And large enough storage array



On average

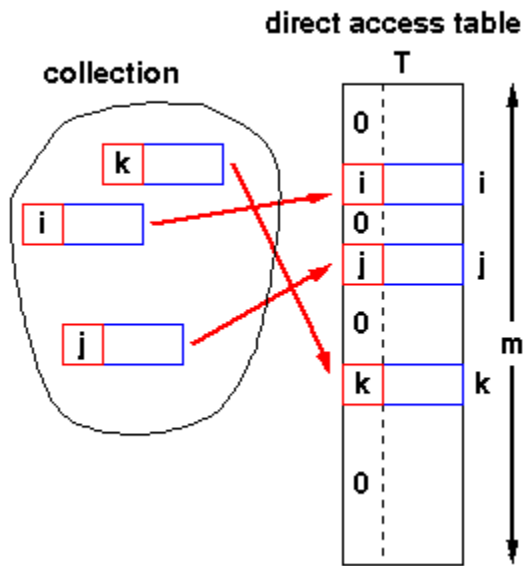
Intro: Direct Address Tables



Contents of this slide are from John Morris, University of Western Australia

- ▶ If we have a collection of n elements whose keys are unique integers in the range $0 \dots m-1$, where $m \geq n$,
- ▶ then we can store the items in a *direct address* table, $T[m]$,
 - where T_i is either empty or contains one of the elements of our collection
 - Searching a direct address table is clearly an $O(1)$ operation:
 - for a key, k , we access T_k ,
 - if it contains an element, return it,
 - if it doesn't, then return a NULL

Intro: Direct Address Tables



- ▶ There are two major constraints:
 1. the keys must be unique
 2. the range of possible keys must be severely bounded

The second constraint is usually impossible to meet

Contents of this slide are from John Morris, University of Western Australia

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



A good hashCode()
distributes the keys, like:

```
hashCode("ate")= 48594983  
hashCode("ape")= 76849201  
hashCode("awe") = 14893202
```

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

- ▶ Example: if $m = 100$:

hashCode("ate") = 48594983
hashCode("ape") = 76849201
hashCode("awe") = 1489036



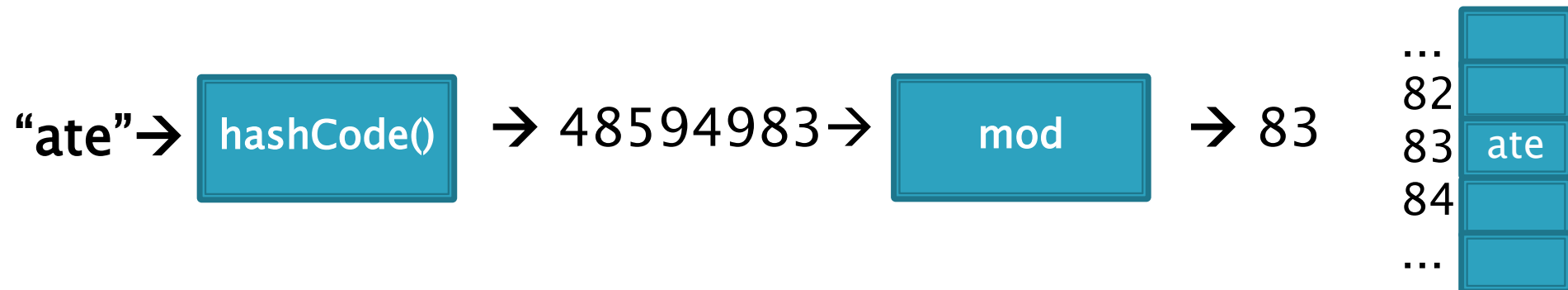
→83

→01

→36

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

- ▶ Every Java object has a **hashCode** method that returns an integer H
 - The hash table uses $H \% m$ as the index into its internal array



- Unless this position is already occupied

a “collision”

Object implements a default *hashCode* method

- ▶ Should we inherit it?
- ▶ JDK classes override the `hashCode()` method
 - Like String
- ▶ If you plan to use instances of your class as keys in a hash table, you probably should too!

hashCode method

- ▶ Should be fast to compute
- ▶ Should distribute keys as evenly as possible
- ▶ These two goals are often contradictory; we need to achieve a balance

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 Math.abs(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 Math.abs(total);
}
```

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

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*23 + s.charAt(i);
    return Math.abs(total);
}
```

- ▶ Spreads out the values more, and anagrams not an issue.
- ▶ We can't entirely avoid collisions. Why?
- ▶ What about overflow during computation?
- ▶ Note: **String** already has a reasonable `hashCode()` method; we don't have to write it ourselves.

Hash Table Caveats

- ▶ Objects that are equal (based on the *equals* method) **MUST** have the same *hashCode* values
- ▶ Different objects should have different *hashCodes* if possible
- ▶ Beware of mutable objects!
- ▶ Hash tables don't maintain sorted order
 - So what's the big-O cost to find min or max element?

Collisions are Inevitable

- ▶ A hash table implementation (like *HashMap*) provides a “collision resolution mechanism”
- ▶ There are a variety of approaches to this
- ▶ Fewer collisions lead to faster performance

Collision Avoidance

▶ Just make hashCode unique?

▶ Impossible!

$|\text{possible key values}| \gg \text{capacity of table}$

- Example: A key may be an array of 16 characters
- How many different values could there be?

▶ So we need to deal with collisions:

- Probing (Linear, Quadratic)
- Chaining

Collision Resolution: Linear Probing

- ▶ Collision? Use the next available space:
 - Try $H+1$, $H+2$, $H+3$, ...
 - Wraparound at the end of the array
- ▶ Problem: Clustering
- ▶ Animation:
 - http://www.cs.auckland.ac.nz/software/AlgAnim/hash_tables.html

```

hash ( 89, 10 ) = 9
hash ( 18, 10 ) = 8
hash ( 49, 10 ) = 9
hash ( 58, 10 ) = 8
hash ( 9, 10 ) = 9

```

After insert 89 After insert 18 After insert 49 After insert 58 After insert 9

0			49	49	49
1				58	58
2					9
3					
4					
5					
6					
7					
8		18	18	18	18
9	89	89	89	89	89

Figure 20.4
 Linear probing hash
 table after each
 insertion

Linear Probing Efficiency

- ▶ Depends on **Load Factor**, λ :
 - Ratio of the number of items stored to table size
 - $0 \leq \lambda \leq 1$.
- ▶ For a given λ , what is the expected number of probes before an empty location is found?

Rough Analysis of Linear Probing

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

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

Better Analysis of Linear Probing

- ▶ “Equally likely” probability is not realistic
- ▶ **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.

Why consider linear probing?

- ▶ Easy to implement
- ▶ Simple code has fast run time per probe
- ▶ Works well when load is low
 - It could be more efficient just to rehash using a bigger table once it starts to fill.
 - And in practice, once $\lambda > 0.5$, we usually **double the size of the array** and rehash

Quadratic Probing

- ▶ Linear probing:
 - Collision at H ? Try $H, H+1, H+2, H+3, \dots$
- ▶ 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 Tricks (1 / 2)

- ▶ **Choose a prime number p for the array size**
- ▶ Then if $\lambda \leq 0.5$:
 - Guaranteed insertion
 - If there is a “hole”, we’ll find it
 - No cell is probed twice
- ▶ See proof of Theorem 20.4:
 - Suppose that we repeat a probe before trying more than half the slots in the table
 - See that this leads to a contradiction
 - Contradicts fact that the table size is prime

Quadratic Probing Tricks (2/2)

- ▶ **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

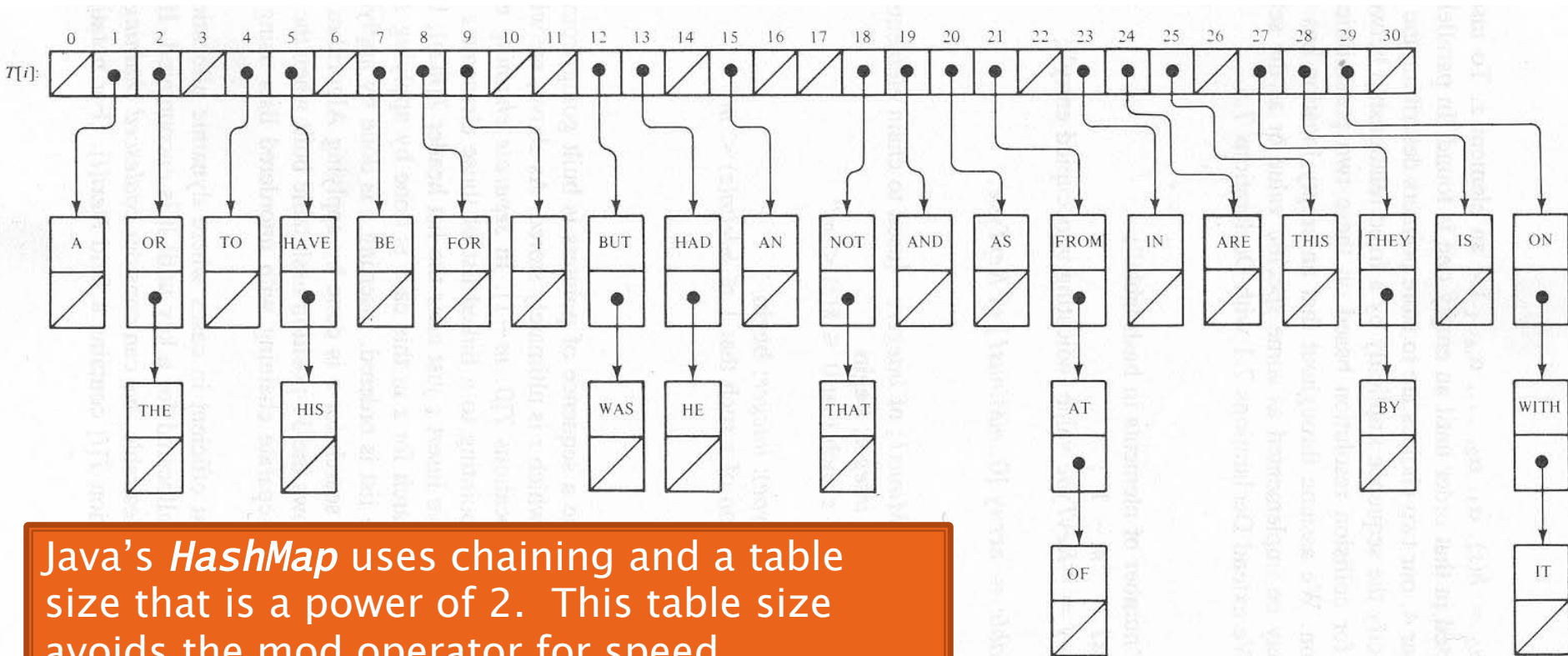
Quadratic probing analysis

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

Another Approach: Separate Chaining

- ▶ Use an array of **linked lists**
- ▶ How would that help resolve collisions?

Hashing with Chaining



Java's *HashMap* uses chaining and a table size that is a power of 2. This table size avoids the mod operator for speed. But since it is susceptible to bad hashes, it always *rehashes* your hash code.

Editor Trees

Immersion in tree
manipulation