

# **CSE 143**

# **Lecture 25**

Hashing

read 11.2

slides created by Marty Stepp

<http://www.cs.washington.edu/143/>

# Recall: ADTs (11.1)

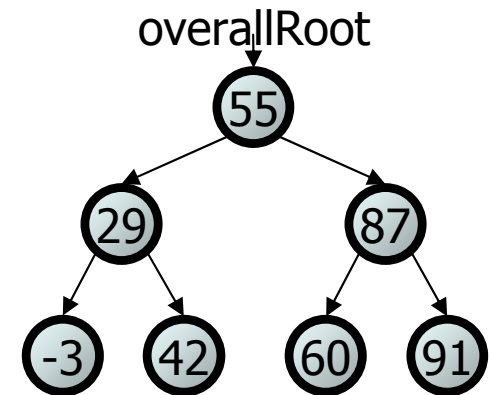
- **abstract data type (ADT):** A specification of a collection of data and the operations that can be performed on it.
  - Describes *what* a collection does, not *how* it does it.
- Java's collection framework describes ADTs with interfaces:
  - Collection, Deque, List, Map, Queue, Set, SortedMap
- An ADT can be implemented in multiple ways by classes:
  - ArrayList and LinkedList            implement List
  - HashSet and TreeSet                implement Set
  - LinkedList , ArrayDeque, etc.    implement Queue

# SearchTree as a set

- We implemented a class `SearchTree` to store a BST of `ints`:
- Our BST is essentially a set of integers.

Operations we support:

- `add`
- `contains`
- `remove`
- ...



- But there are other ways to implement a set...

# How to implement a set?

- Elements of a `TreeSet` (`IntTree`) are in BST sorted order.
  - We need this in order to add or search in  $O(\log N)$  time.
- But it doesn't really matter what order the elements appear in a set, so long as they can be added and searched quickly.
- Consider the task of storing a set in an array.
  - What would make a good ordering for the elements?

|       |   |    |    |    |   |   |   |   |   |   |
|-------|---|----|----|----|---|---|---|---|---|---|
| index | 0 | 1  | 2  | 3  | 4 | 5 | 6 | 7 | 8 | 9 |
| value | 7 | 11 | 24 | 49 | 0 | 0 | 0 | 0 | 0 | 0 |

|       |   |    |   |   |    |   |   |   |   |    |
|-------|---|----|---|---|----|---|---|---|---|----|
| index | 0 | 1  | 2 | 3 | 4  | 5 | 6 | 7 | 8 | 9  |
| value | 0 | 11 | 0 | 0 | 24 | 0 | 0 | 7 | 0 | 49 |

# Hashing

- **hash**: To map a value to an integer index.
  - **hash table**: An array that stores elements via hashing.
- **hash function**: An algorithm that maps values to indexes.
  - one possible hash function for integers:

$$\text{HF}(I) \rightarrow I \% \text{length}$$

```
set.add(11);      // 11 % 10 == 1
set.add(49);      // 49 % 10 == 9
set.add(24);      // 24 % 10 == 4
set.add(7);       // 7 % 10 == 7
```

|       |   |    |   |   |    |   |   |   |   |    |
|-------|---|----|---|---|----|---|---|---|---|----|
| index | 0 | 1  | 2 | 3 | 4  | 5 | 6 | 7 | 8 | 9  |
| value | 0 | 11 | 0 | 0 | 24 | 0 | 0 | 7 | 0 | 49 |

# Efficiency of hashing

```
public static int HF(int i) {           // hash function
    return Math.abs(i) % elementData.length;
}
```

- Add: simply set `elementData[HF(i)] = i;`
- Search: check if `elementData[HF(i)] == i`
- Remove: set `elementData[HF(i)] = 0;`
  
- What is the runtime of `add`, `contains`, and `remove`?
  - **O(1)!** OMGWTFBBQFAST
  
- Are there any problems with this approach?

# Collisions

- **collision:** When a hash function maps two or more elements to the same index.

```
set.add(11);  
set.add(49);  
set.add(24);  
set.add(7);  
set.add(54); // collides with 24!
```

- **collision resolution:** An algorithm for fixing collisions.

|       |   |    |   |   |    |   |   |   |   |    |
|-------|---|----|---|---|----|---|---|---|---|----|
| index | 0 | 1  | 2 | 3 | 4  | 5 | 6 | 7 | 8 | 9  |
| value | 0 | 11 | 0 | 0 | 54 | 0 | 0 | 7 | 0 | 49 |

# Probing

- **probing**: Resolving a collision by moving to another index.
  - **linear probing**: Moves to the next index.

```
set.add(11);  
set.add(49);  
set.add(24);  
set.add(7);  
set.add(54); // collides with 24
```

|       |   |    |   |   |    |           |   |   |   |    |
|-------|---|----|---|---|----|-----------|---|---|---|----|
| index | 0 | 1  | 2 | 3 | 4  | 5         | 6 | 7 | 8 | 9  |
| value | 0 | 11 | 0 | 0 | 24 | <b>54</b> | 0 | 7 | 0 | 49 |

- Is this a good approach?



# Clustering

- **clustering**: Clumps of elements at neighboring indexes.
  - slows down the hash table lookup; you must loop through them.

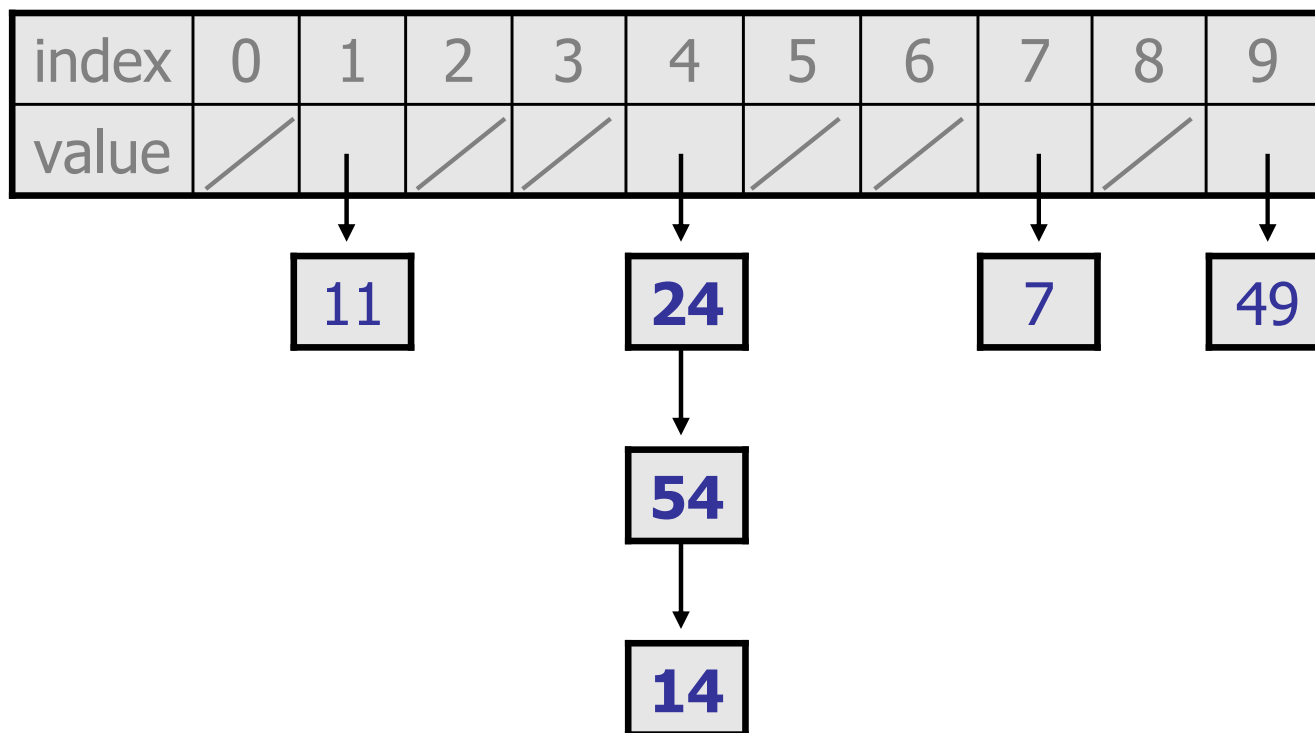
```
set.add(11);  
set.add(49);  
set.add(24);  
set.add(7);  
set.add(54); // collides with 24  
set.add(14); // collides with 24, then 54  
set.add(86); // collides with 14, then 7
```

|       |   |    |   |   |           |           |           |          |           |    |
|-------|---|----|---|---|-----------|-----------|-----------|----------|-----------|----|
| index | 0 | 1  | 2 | 3 | 4         | 5         | 6         | 7        | 8         | 9  |
| value | 0 | 11 | 0 | 0 | <b>24</b> | <b>54</b> | <b>14</b> | <b>7</b> | <b>86</b> | 49 |

- Now a lookup for 94 must look at 7 out of 10 total indexes.

# Chaining

- **chaining:** Resolving collisions by storing a list at each index.
  - add/search/remove must traverse lists, but the lists are short
  - impossible to "run out" of indexes, unlike with probing



# Hash set code

```
import java.util.*;    // for List, LinkedList
// All methods assume value != null; does not rehash
public class HashIntSet {
    private static final int CAPACITY = 137;
    private List<Integer>[] elements;

    // constructs new empty set
    public HashSet() {
        elements = (List<Integer>[]) (new List[CAPACITY]);
    }

    // adds the given value to this hash set
    public void add(int value) {
        int index = HF(value);
        if (elements[index] == null) {
            elements[index] = new LinkedList<Integer>();
        }
        elements[index].add(value);
    }

    // hashing function to convert objects to indexes
    private int HF(int value) {
        return Math.abs(value) % elements.length;
    }
    ...
}
```

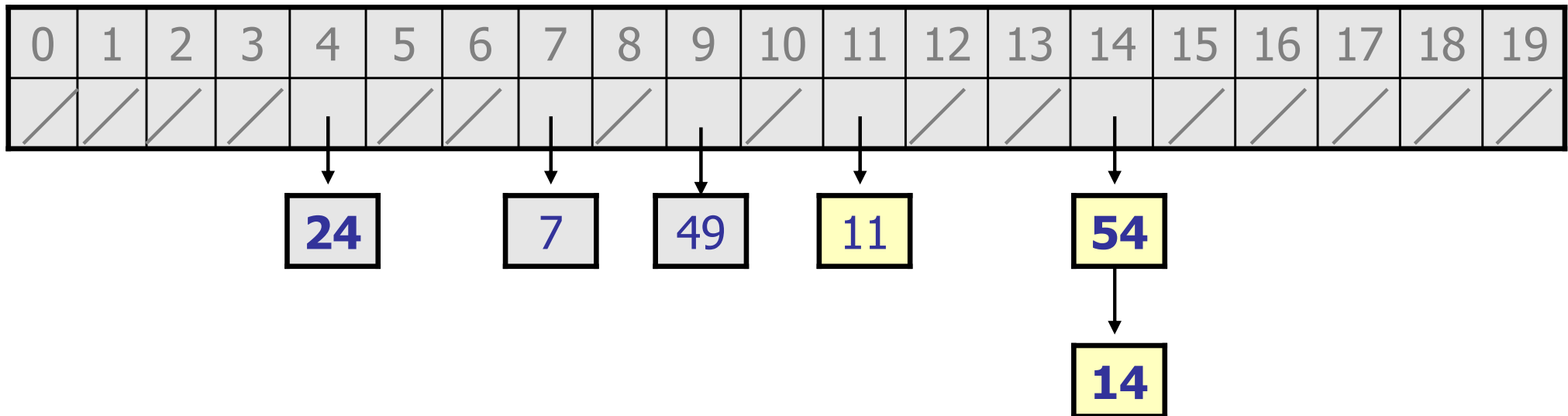
# Final hash set code 2

```
...
// Returns true if this set contains the given value.
public boolean contains(int value) {
    int index = HF(value);
    return elements[index] != null &&
        elements[index].contains(value);
}

// Removes the given value from the set, if it exists.
public void remove(int value) {
    int index = HF(value);
    if (elements[index] != null) {
        elements[index].remove(value);
    }
}
}
```

# Rehashing

- **rehash:** Growing to a larger array when the table is too full.
  - Cannot simply copy the old array to a new one. (Why not?)
- **load factor:** ratio of (*# of elements*) / (*hash table length*)
  - many collections rehash when load factor  $\cong .75$
  - can use big prime numbers as hash table sizes to reduce collisions



# Hashing objects

- It is easy to hash an integer  $I$  (use index  $I \% length$  ).
  - How can we hash other types of values (such as objects)?
- All Java objects contain the following method:

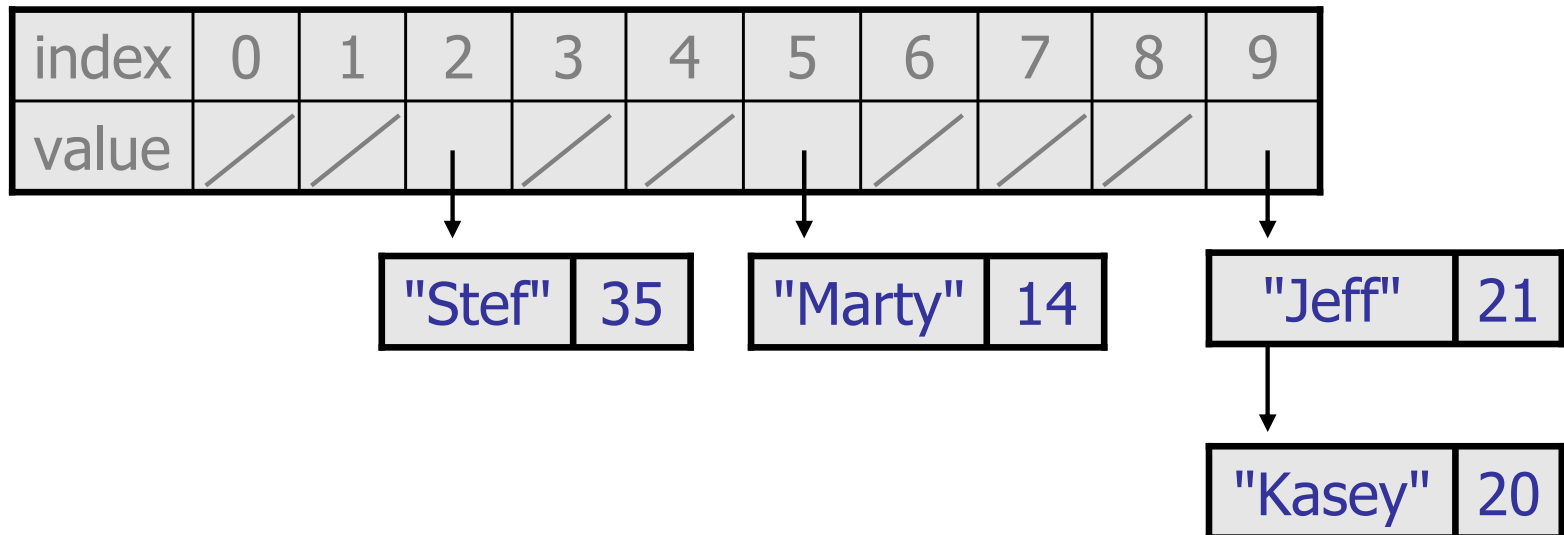
```
public int hashCode()
```

Returns an integer hash code for this object.
  - We can call `hashCode` on any object to find its preferred index.
- How is `hashCode` implemented?
  - Depends on the type of object and its state.
    - Example: a `String`'s `hashCode` adds the ASCII values of its letters.
  - You can write your own `hashCode` methods in classes you write.

# Implementing hash maps

- A hash map is just a set where the lists store key/value pairs:

```
//      key      value
map.put ("Marty", 14);
map.put ("Jeff", 21);
map.put ("Kasey", 20);
map.put ("Stef", 35);
```



- Instead of a `List<Integer>`, write an inner `Entry` node class with `key` and `value` fields; the map stores a `List<Entry>`