# HASHING

## Objectives

- To know what hashing is for (§48.3).

- To obtain the hash code for an object and design the hash function to map a key to an index (§48.4).

- To handle collisions using open addressing (§48.5).

- To know the differences among linear probing, quadratic probing, and double hashing (§48.5).

- To handle collisions using separate chaining (§48.6).

- To understand the load factor and the need for rehashing (§48.7).

- To implement `MyHashMap` using hashing (§48.8).

## 48.1 Introduction

The preceding chapters introduced search trees. An element can be found in $O(\log n)$ time in a well-balanced search tree. Is there a more efficient way to search for an element in a container? This chapter introduces a technique called *hashing*. You can use hashing to implement a map or a set to search, insert, and delete an element in $O(1)$ time.

## 48.2 Map

map
key
value

Recall that a *map* is a data structure that stores entries. Each entry contains two parts: *key* and *value*. The key is also called a *search key*, which is used to search for the corresponding value. For example, a dictionary can be stored in a map, where the words are the keys and the definitions of the words are the values.

dictionary
hash table
associative array

> 📀 **Note**
>
> A map is also called a *dictionary*, a *hash table*, or an associative array.

The Java collections framework defines the `java.util.Map` interface for modeling maps. Three concrete implementations are `java.util.HashMap`, `java.util.LinkedHashMap`, and `java.util.TreeMap`. `java.util.HashMap` is implemented using hashing, `java.util.LinkedHashMap` using `LinkedList`, and `java.util.TreeMap` using red-black trees. You will learn the concept of hashing and use it to implement a map in this chapter. In the chapter exercises, you will implement `LinkedHashMap` and `TreeMap`.

## 48.3 What is Hashing?

hash table
hash function
hashing

If you know the index of an element in the array, you can retrieve the element using the index in $O(1)$ time. So, can we store the values in an array and use the key as the index to find the value? The answer is yes—if you can map a key to an index. The array that stores the values is called a *hash table*. The function that maps a key to an index in the hash table is called a *hash function*. As shown in Figure 48.1, a *hash function* obtains an index from a key and uses the index to retrieve the value for the key. *Hashing* is a technique that retrieves the value using the index obtained from the key without performing a search.
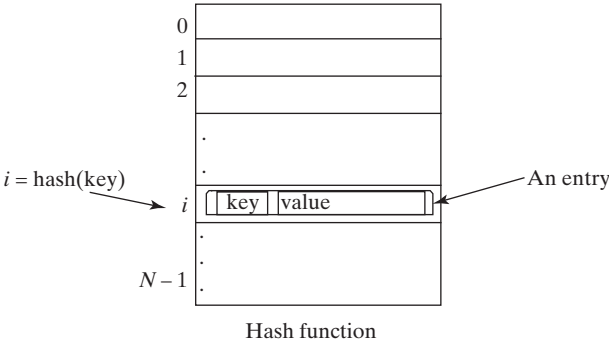


**FIGURE 48.1** A hash function maps a key to an index in the hash table.

perfect hash function

How do you design a hash function that produces an index from a key? Ideally, we would like to design a function that maps each search key to a different index in the hash table. Such a function is called a *perfect hash function*. However, it is difficult to find a perfect hash

function. When two or more keys are mapped to the same hash value, we say that a *collision* has occurred. We will discuss how to deal with collisions later. Although there are ways to deal with them, it is better to avoid collisions in the first place. So, you should design a fast and easy-to-compute hash function that minimizes collisions.

collision

## 48.4 Hash Functions and Hash Codes

A typical hash function first converts a search key to an integer value called a *hash code*, then compresses the hash code into an index to the hash table.

hash code

Java's root class **Object** has the **hashCode** method that returns an integer hash code. By default, the method returns the memory address for the object. The general contract for the **hashCode** is as follows:

- You should override the **hashCode** method whenever the **equals** method is overridden to ensure that two equal objects return the same hash code.

- During the execution of a program, invoking the **hashCode** method multiple times returns the same integer, provided that the object's data are not changed.

- Two unequal objects may have the same hash code, but you should implement the **hashCode** method to avoid too many such cases.

### 48.4.1   Hash Codes for Primitive Types

For a search key of the type **byte**, **short**, **int**, and **char**, simply cast it to **int**. So, two different search keys of any one of these types will have different hash codes.

**byte**, **short**, **int**, **char**

For a search key of the type **float**, use **Float.floatToIntBits(key)** as the hash code. Note that **floatToIntBits(float f)** returns an **int** value whose bit representation is the same as the bit representation for the floating number **f**. So, two different search keys of the **float** type will have different hash codes.

**float**

For a search key of the type **long**, simply casting it to **int** would not be a good choice, because all keys that differ in only the first 32 bits will have the hash code. To take the first 32 bits into consideration, divide the 64 bits into two halves and perform the exclusive-or operation to combine the two halves. This process is called *folding*. So, the hashing code is

**long**

folding

```
int hashCode = (int)(key ^ (key >> 32));
```

Note that `>>` is the right-shift operator that shifts the bits 32 position to the right. For example, `1010110 >> 2` yields `0010101`. The `^` is the bitwise exclusive-or operator. It operates on two corresponding bits of the binary operands. For example, `1010110 ^ 0110111` yields `1100001`.

For a search key of the type **double**, first convert it to a **long** value using **doubleToLongBits**, then perform a folding as follows:

**double**

folding

```
long bits = Double.doubleToLongBits(key);
int hashCode = (int)(bits ^ (bits >> 32));
```

### 48.4.2   Hash Codes for Strings

Search keys are often strings. So, it is important to design a good hash function for strings. An intuitive approach is to sum the Unicode of all characters as the hash code for the string. This approach may work if two search keys in an application don't contain same letters. But it will produce a lot of collisions if the search keys contain the same letters such as **tod** and **dot**.

A better approach is to generate a hash code that takes the position of characters into consideration. Specifically, let the hash code be

$$s_0 * b^{(N-1)} + s_1 * b^{(N-2)} + \cdots + s_{N-1}$$

polynomial hash code

where $s_i$ is **s.charAt(i)**. This expression is a polynomial for some positive $b$. So, this is called a *polynomial hash code*. By Horner's rule, it can be evaluated efficiently as follows:

$$(\ldots((s_0*b + s_1)b + s_2)b + \cdots + s_{N-2})b + s_{N-1}$$

This computation can cause an overflow for long strings. Arithmetic overflow is ignored in Java. You should choose an appropriate value $b$ to minimize collision. Experiments show that the good choices for $b$ are 31, 33, 37, 39, and 41. In the **String** class, the **hashCode** is over-ridden using the polynomial hash code with $b$ being **31**.

### 48.4.3 Compressing Hash Codes

The hash code for a key can be a large integer that is out of the range for the hash-table index. You need to scale it down to fit in the range of the index. Assume the index for a hash table is between **0** and **N-1**. The most common way to scale an integer to between **0** and **N-1** is to use

```
h(hashCode) = hashCode % N
```

To ensure that the indices are spread evenly, choose **N** to be a prime number greater than **2**.

Ideally you should choose a prime number for **N**. However, it is time consuming to find a large prime number. In the Java API implementation for **java.util.HashMap**, **N** is conveniently set to a value of power **2**. To ensure the hashing is evenly distributed, a supplemental hash function is also used along with the primary hash function. The supplemental function is defined as follows:

```
private static int supplementalHash(int h) {
  h ^= (h >>> 20) ^ (h >>> 12);
  return h ^ (h >>> 7) ^ (h >>> 4);
}
```

**^** and **>>>** are bitwise exclusive-or and right-shift operations. See Supplement Part III.D, "Bitwise Operations," on the Companion Website.

The primary hash function is defined as follows:

```
h(hashCode) = supplementalHash(hashCode) % N
```

Note that the function can also be written as

```
h(hashCode) = supplementalHash(hashCode) & (N - 1)
```

since **N** is a power of **2**.

## 48.5 Handling Collisions Using Open Addressing

open addressing
separate chaining
linear probing
quadratic probing
double hashing

A collision occurs when two keys are mapped to the same index in a hash table. Generally, there are two ways for handling collisions: *open addressing* and *separate chaining*.

Open addressing is to find an open location in the hash table in the event of collision. Open addressing has several variations: *linear probing*, *quadratic probing*, and *double hashing*.

### 48.5.1 Linear Probing

add entry

When a collision occurs during the insertion of an entry to a hash table, linear probing finds the next available location sequentially. For example, if a collision occurs at **hashTable[k % N]**, check whether **hashTable[(k+1) % N]** is available. If not, check **hashTable[(k+2) % N]** and so on, until an available cell is found, as shown in Figure 48.2.

New element with
key 26 to be inserted

0 | key: 44
1 |
2 |
3 |
4 | key: 4
5 | key: 16
6 | key: 28
7 |
8 |
9 |
10 | key: 21

For simplicity, only the keys are
shown and the values are not
shown. Here N is 11 and
index = key % N.

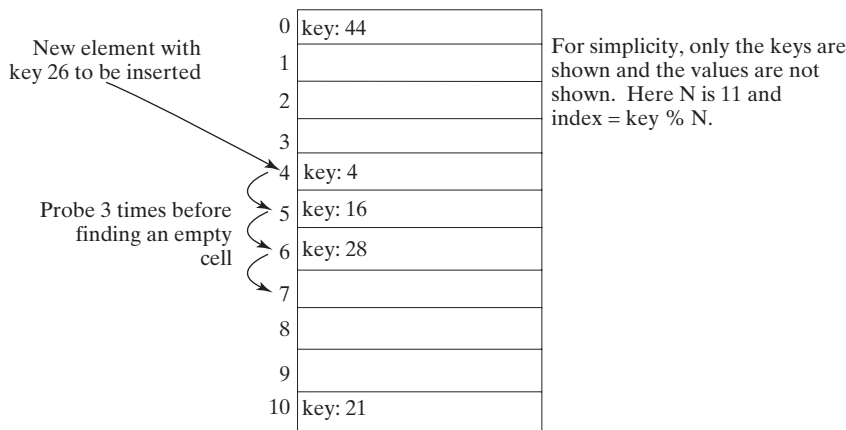Probe 3 times before
finding an empty
cell

**FIGURE 48.2**    Linear probe finds the next available location sequentially.

**Note**

When probing reaches the end of the table, it goes back to the beginning of the table. Thus, the
hash table is treated as if it were circular.

circular hash table

To search for an entry in the hash table, obtain the index, say **k**, from the hash function for the
key. Check whether `hashTable[k % n]` contains the entry. If not, check whether
`hashTable[(k+1) % n]` contains the entry, and so on, until it is found, or an empty cell is
reached.

search entry

To remove an entry from the hash table, search the entry that matches the key. If entry is
found, place a special marker to denote that the entry is available. Each cell in the hash table
has three possible states: occupied, available, or empty. Note that an empty cell is also avail-
able for insertion.

remove entry

Linear probing tends to cause groups of consecutive cells in the hash table to be occupied.
Each group is called a *cluster*. Each cluster is actually a probe sequence that you must search
when retrieving, adding, or removing an entry. As clusters grow in size, they may merge into
even larger clusters, further slowing down the search time. This is a big disadvantage of linear
probing.

cluster

**Pedagogical Note**

Follow the link www.cs.armstrong.edu/liang/animation/HashingLinearProbingAnimation.html to see how to
hashing with linear probing works, as shown in Figure 48.3.

linear probing animation

## 48.5.2   Quadratic Probing

Quadratic probing can avoid the clustering problem in linear probing. Linear probing looks at
the consecutive cells beginning at index **k**. Quadratic probing, on the other hand, looks at the
cells at indices $(k + j^2)$ `% n`, for j ≥ 0, i.e., **k**, **(k + 1) % n**, **(k + 4) % n**, **(k + 9) % n**, ..., 
and so on, as shown in Figure 48.4.

Quadratic probing works in the same way as linear probing except for the change of search
sequence. Quadratic probing avoids the clustering problem in linear probing, but it has its
own clustering problem, called *secondary clustering*; i.e., the entries that collide with an
occupied entry use the same probe sequence.

secondary clustering

Linear probing guarantees that an available cell can be found for insertion as long as the
table is not full. However, there is no such guarantee for quadratic probing.
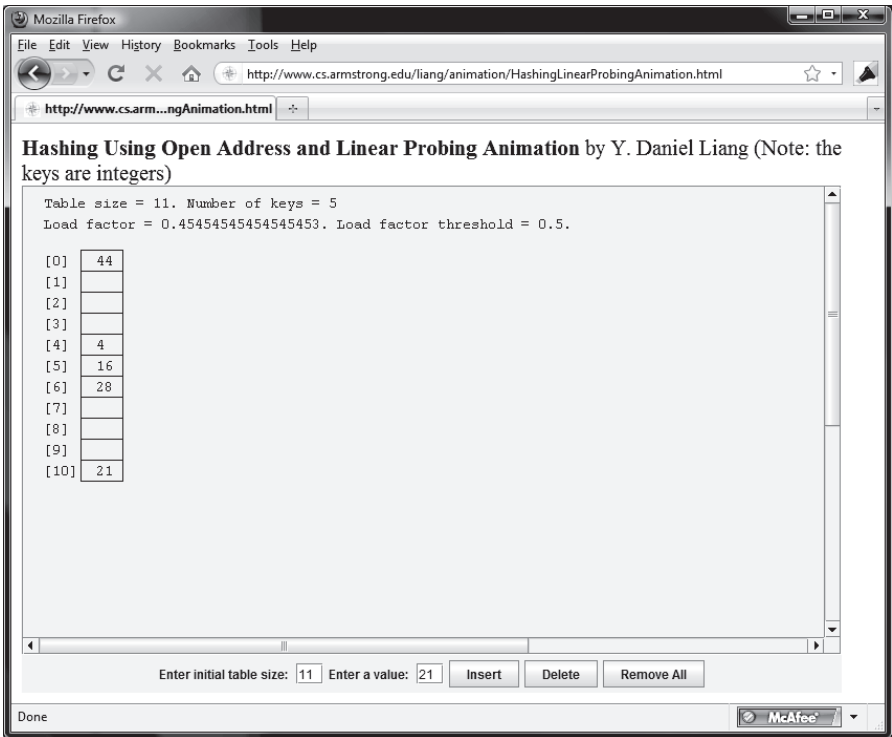
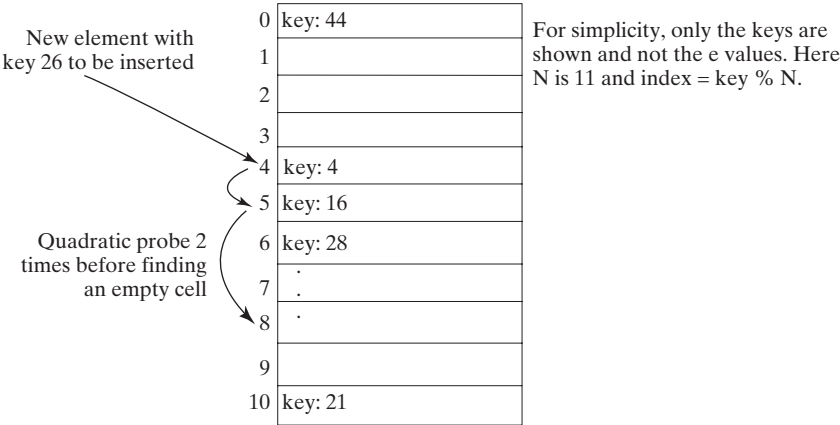**FIGURE 48.3**   The animation tool shows how linear probing works.



**FIGURE 48.4**   Quadratic probe increases the next index in the sequence by $j^2$ for $j = 1, 2, 3, \ldots$.

### 48.5.3   Double Hashing

double hashing

Another open addressing scheme that avoids the clustering problem is known as *double hashing*. Starting from the initial index **k**, both linear probing and quadratic probing add an increment to **k** to define a search sequence. The increment is **1** for linear probing and $j^2$ for quadratic probing. These increments are independent of the keys. Double hashing uses a secondary hash function on the keys to determine the increments to avoid the clustering problem.

For example, let the primary hash function **h** and secondary hash function **h'** on a hash table of size **11** be defined as follows:

```
h(k) = k % 11;
h'(k) = 7 - k % 7;
```

For a search key of **12**, we have

```
h(12) = 12 % 11 = 1;
h'(k) = 7 - 12 % 7 = 2;
```

The probe sequence for key **12** starts at index **1** with an increment **2**, as shown in Figure 48.5.



**FIGURE 48.5**   The secondary hash function in a double hashing determines the increment of the next index in the probe sequence.

The indices of the probe sequence are as follows: 1, 3, 5, 7, 9, 0, 2, 4, 6, 8, 10. This sequence reaches the entire table. You should design your functions to produce the probe sequence that reaches the entire table. Note that the second function should never have a zero value, since zero is not an increment.



**FIGURE 48.6**   Separate chaining chains the entries with the same hash index in a bucket.

## 48.6 Handling Collisions Using Separate Chaining

The preceding section introduced handling collisions using open addressing. The open addressing scheme finds a new location when a collision occurs. This section introduces handling col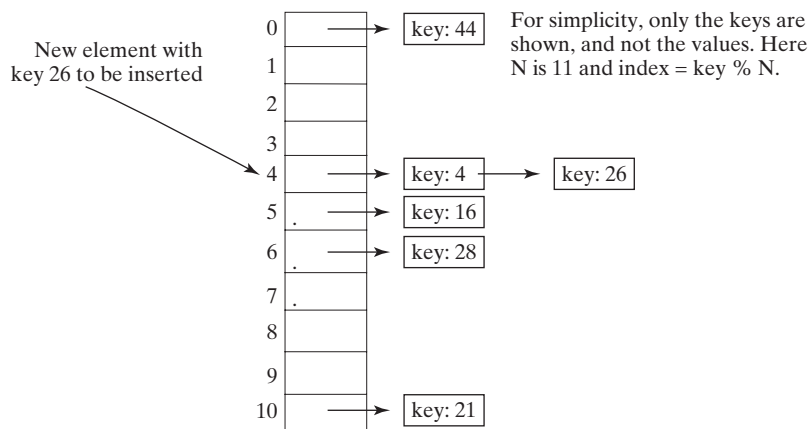lisions using separate chaining. The separate chaining scheme places all entries with the same hash index into the same location, rather than finding new locations. Each location in the separate chaining scheme is called a *bucket*. A bucket is a container that holds multiple entries.

bucket

implementing bucket

You may implement a bucket using an array, **ArrayList**, or **LinkedList**. We will use **LinkedList** for demonstration. You can view each cell in the hash table as the reference to the head of a linked list, and elements in the linked list are chained starting from the head, as shown in Figure 48.6.

## 48.7 Load Factor and Rehashing

load factor

*Load factor* $\lambda$ measures how full the hash table is. It is the ratio of the size of the map to the size of the hash table, i.e., $\lambda = \dfrac{n}{N}$, where $n$ denotes the number of elements and $N$ the number of locations in the hash table.

Note that $\lambda$ is zero if the map is empty. For the open addressing scheme, $\lambda$ is between **0** and **1**; $\lambda$ is **1** if the hash table is full. For the separate chaining scheme, $\lambda$ can be any value. As $\lambda$ increases, the probability of collision increases. Studies show that you should maintain the load factor under **0.5** for the open addressing scheme and under **0.9** for the separate chaining scheme.

threshold

Keeping the load factor under a certain threshold is important for the performance of hashing. In the implementation of **java.util.HashMap** class in the Java API, the threshold **0.75** is used. Whenever the load factor exceeds the threshold, you need to increase the hash-table
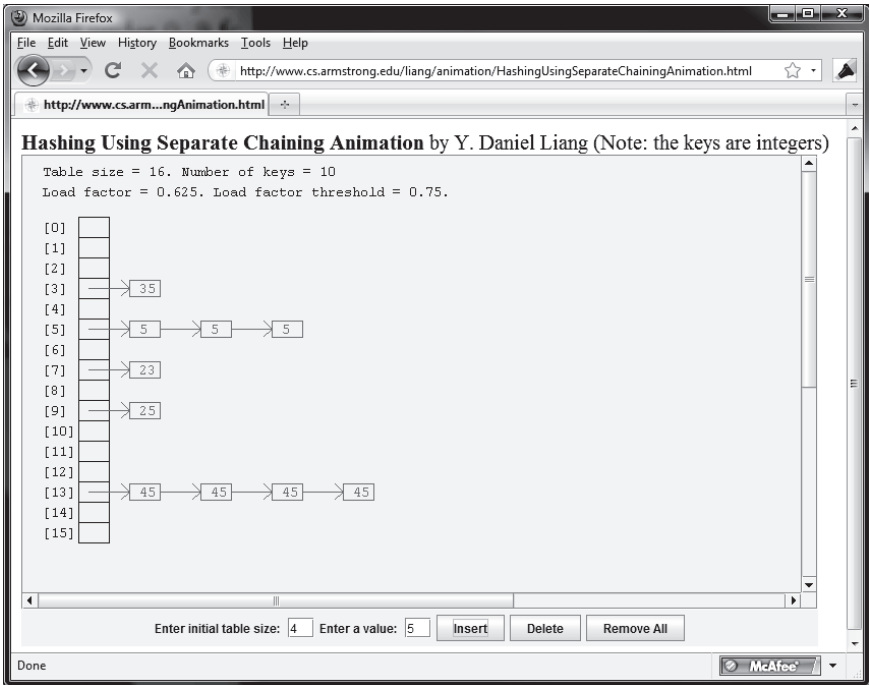


**FIGURE 48.7**   The animation tool shows how separate chaining works.

size and *rehash* all the entries in the map to the new hash table. Notice that you need to change the hash functions, since the hash-table size has been changed. To reduce the likelihood of rehashing, since it is costly, you should at least double the hash-table size. Even with periodic rehashing, hashing is an efficient implementation for map.

*rehash*

**Pedagogical Note**

Follow the link www.cs.armstrong.edu/liang/animation/HashingUsingSeparateChainingAnimation.html to see how to hashing with linear probing works, as shown in Figure 48.7.

separate chaining animation

## 48.8 Implementing a Map Using Hashing

Now you know the concept of hashing. You know how to design a good hash function to map a key to an index in a hash table, how to measure performance using the load factor, and how to increase the table size and rehash to maintain the performance. This section demonstrates how to implement a map using separate chaining.

We design our custom **Map** interface to mirror **java.util.Map** with some minor variations. In the **java.util.Map** interface, the keys are distinct. However, a map may allow duplicate keys. Our map interface allows duplicate keys. We name the interface **MyMap** and a concrete class **MyHashMap**, as shown in Figure 48.8.

duplicate keys



| «interface» *MyMap<K, V>* | |
|---|---|
| +clear(): void | Removes all entries from this map. |
| +containsKey(key: K): boolean | Returns true if this map contains an entry for the specified key. |
| +containsValue(value: V): boolean | Returns true if this map maps one or more keys to the specified value. |
| +entrySet(): Set<Entry<K, V>> | Returns a set consisting of the entries in this map. |
| +get(key: K): V | Returns a value for the specified key in this map. |
| +getAll(key: K): Set<V> | Returns all values for the specified key in this map. |
| +isEmpty(): boolean | Returns true if this map contains no mappings. |
| +keySet(): Set<K> | Returns a set consisting of the keys in this map. |
| +put(key: K, value: V): V | Puts a mapping in this map. |
| +remove(key: K): void | Removes the entries for the specified key. |
| +size(): int | Returns the number of mappings in this map. |
| +values(): Set<V> | Returns a set consisting of the values in this map. |

| **MyHashMap<K, V>** | Concrete class that implements MyMap |
|---|---|

| **MyMap.Entry<K, V>** | |
|---|---|
| -key: K -value: V | |
| +Entry(key: K, value: V) | Constructs an entry with the specified key and value. |
| +getkey(): Key | Returns the key in the entry. |
| +getValue(): Value | Returns the value in the entry. |

**FIGURE 48.8**   **MyHashMap** implements the **MyMap** interface.

[We output the page content]

The **get(key)** method gets one of the values that match the key. The **getAll(key)** method retrieves all values that match the key.

How do you implement **MyHashMap**? If you use an **ArrayList** and store a new entry at the end of the list, the search time will be $O(n)$. If you implement **MyHashMap** using an AVL tree, the search time will be $O(\log n)$. Nevertheless, you can implement **MyHashMap** using hashing to obtain an $O(1)$ time search algorithm. Listing 48.1 shows the **MyMap** interface and Listing 48.2 implements **MyHashMap** using separate chaining.

### LISTING 48.1 MyMap.java

| | |
|---|---|
| interface **MyMap** | |
| | ```
1 public interface MyMap<K, V> {
2   /** Remove all of the entries from this map */
``` |
| **clear** | ```
3   public void clear();
4
``` |
| | ```
5   /** Return true if the specified key is in the map */
``` |
| **containsKey** | ```
6   public boolean containsKey(K key);
7
``` |
| | ```
8   /** Return true if this map contains the specified value */
``` |
| **containsValue** | ```
9   public boolean containsValue(V value);
10
``` |
| | ```
11   /** Return a set of entries in the map */
``` |
| **entrySet** | ```
12   public java.util.Set<Entry<K, V>> entrySet();
13
``` |
| | ```
14   /** Return the first value that matches the specified key */
``` |
| **get** | ```
15   public V get(K key);
16
``` |
| | ```
17   /** Return all values for the specified key in this map */
``` |
| **getAll** | ```
18   public java.util.Set<V> getAll(K key);
19
``` |
| | ```
20   /** Return true if this map contains no entries */
``` |
| **isEmpty** | ```
21   public boolean isEmpty();
22
``` |
| | ```
23   /** Return a set consisting of the keys in this map */
``` |
| **keySet** | ```
24   public java.util.Set<K> keySet();
25
``` |
| | ```
26   /** Add an entry (key, value) into the map */
``` |
| **put** | ```
27   public V put(K key, V value);
28
``` |
| | ```
29   /** Remove an entry for the specified key */
``` |
| **remove** | ```
30   public void remove(K key);
31
``` |
| | ```
32   /** Return the number of mappings in this map */
``` |
| **size** | ```
33   public int size();
34
``` |
| | ```
35   /** Return a set consisting of the values in this map */
``` |
| **values** | ```
36   public java.util.Set<V> values();
37
``` |
| | ```
38   /** Define inner class for Entry */
``` |
| **Entry** inner class | ```
39   public static class Entry<K, V> {
40     K key;
41     V value;
42
43     public Entry(K key, V value) {
44       this.key = key;
45       this.value = value;
46     }
47
48     public K getKey() {
49       return key;
``` |

```
50      }
51
52      public V getValue() {
53        return value;
54      }
55
56      public String toString() {
57        return "[" + key + ", " + value + "]";
58      }
59    }
60 }
```

## LISTING 48.2  MyHashMap.java

```
 1 import java.util.LinkedList;
 2
 3 public class MyHashMap<K, V> implements MyMap<K, V> {          class MyHashMap
 4   // Define the default hash-table size. Must be a power of 2
 5   private static int DEFAULT_INITIAL_CAPACITY = 4;             default initial capacity
 6
 7   // Define the maximum hash-table size. 1 << 30 is same as 2^30
 8   private static int MAXIMUM_CAPACITY = 1 << 30;               maximum capacity
 9
10   // Current hash-table capacity. Capacity is a power of 2
11   private int capacity;                                       current capacity
12
13   // Define default load factor
14   private static float DEFAULT_MAX_LOAD_FACTOR = 0.75f;        default load factor
15
16   // Specify a load factor used in the hash table
17   private float loadFactorThreshold;                          load-factor threshold
18
19   // The number of entries in the map
20   private int size = 0;                                       size
21
22   // Hash table is an array with each cell being a linked list
23   LinkedList<MyMap.Entry<K,V>>[] table;                       hash table
24
25   /** Construct a map with the default capacity and load factor */
26   public MyHashMap() {                                        no-arg constructor
27     this(DEFAULT_INITIAL_CAPACITY, DEFAULT_MAX_LOAD_FACTOR);
28   }
29
30   /** Construct a map with the specified initial capacity and
31    * default load factor */
32   public MyHashMap(int initialCapacity) {                     constructor
33     this(initialCapacity, DEFAULT_MAX_LOAD_FACTOR);
34   }
35
36   /** Construct a map with the specified initial capacity
37    * and load factor */
38   public MyHashMap(int initialCapacity, float loadFactorThreshold) {   constructor
39     if (initialCapacity > MAXIMUM_CAPACITY)
40       this.capacity = MAXIMUM_CAPACITY;
41     else
42       this.capacity = trimToPowerOf2(initialCapacity);
43
44     this.loadFactorThreshold = loadFactorThreshold;
45     table = new LinkedList[capacity];
46   }
```

```
47
48     /** Remove all of the entries from this map */
49     public void clear() {
50       size = 0;
51       removeEntries();
52     }
53
54     /** Return true if the specified key is in the map */
55     public boolean containsKey(K key) {
56       if (get(key) != null)
57         return true;
58        else
59         return false;
60     }
61
62     /** Return true if this map contains the specified value */
63     public boolean containsValue(V value) {
64       for (int i = 0; i < capacity; i++) {
65         if (table[i] != null) {
66           LinkedList<Entry<K, V>> bucket = table[i];
67           for (Entry<K, V> entry: bucket)
68             if (entry.getValue().equals(value))
69               return true;
70         }
71       }
72
73       return false;
74     }
75
76     /** Return a set of entries in the map */
77     public java.util.Set<MyMap.Entry<K,V>> entrySet() {
78       java.util.Set<MyMap.Entry<K, V>> set =
79         new java.util.HashSet<MyMap.Entry<K, V>>();
80
81       for (int i = 0; i < capacity; i++) {
82         if (table[i] != null) {
83           LinkedList<Entry<K, V>> bucket = table[i];
84           for (Entry<K, V> entry: bucket)
85             set.add(entry);
86         }
87       }
88
89       return set;
90     }
91
92     /** Return the first value that matches the specified key */
93     public V get(K key) {
94       int bucketIndex = hash(key.hashCode());
95       if (table[bucketIndex] != null) {
96         LinkedList<Entry<K, V>> bucket = table[bucketIndex];
97         for (Entry<K, V> entry: bucket)
98           if (entry.getKey().equals(key))
99             return entry.getValue();
100      }
101
102      return null;
103    }
104
105    /** Return all values for the specified key in this map */
```

clear

containsKey

containsValue

entrySet

get

```
106    public java.util.Set<V> getAll(K key) {                              getAll
107      java.util.Set<V> set = new java.util.HashSet<V>();
108      int bucketIndex = hash(key.hashCode());
109      if (table[bucketIndex] != null) {
110        LinkedList<Entry<K, V>> bucket = table[bucketIndex];
111        for (Entry<K, V> entry: bucket)
112          if (entry.getKey().equals(key))
113            set.add(entry.getValue());
114      }
115
116      return set;
117    }
118
119    /** Return true if this map contains no entries */
120    public boolean isEmpty() {                                          isEmpty
121      return size == 0;
122    }
123
124    /** Return a set consisting of the keys in this map */
125    public java.util.Set<K> keySet() {                                  keySet
126      java.util.Set<K> set = new java.util.HashSet<K>();
127
128      for (int i = 0; i < capacity; i++) {
129        if (table[i] != null) {
130          LinkedList<Entry<K, V>> bucket = table[i];
131          for (Entry<K, V> entry: bucket)
132            set.add(entry.getKey());
133        }
134      }
135
136      return set;
137    }
138
139    /** Add an entry (key, value) into the map */
140    public V put(K key, V value) {                                      put
141      if (size >= capacity * loadFactorThreshold) {
142        if (capacity == MAXIMUM_CAPACITY)
143          throw new RuntimeException("Exceeding maximum capacity");
144
145        rehash();
146      }
147
148      int bucketIndex = hash(key.hashCode());
149
150      // Create a linked list for the bucket if it is not created
151      if (table[bucketIndex] == null) {
152        table[bucketIndex] = new LinkedList<Entry<K, V>>();
153      }
154
155      // Add an entry (key, value) to hashTable[index]
156      table[bucketIndex].add(new MyMap.Entry<K, V>(key, value));
157
158      size++; // Increase size
159
160      return value;
161    }
162
163    /** Remove the entries for the specified key */
164    public void remove(K key) {                                         remove
165      int bucketIndex = hash(key.hashCode());
```

```
166
167     // Remove the first entry that matches the key from a bucket
168     if (table[bucketIndex] != null) {
169       LinkedList<Entry<K, V>> bucket = table[bucketIndex];
170       for (Entry<K, V> entry: bucket)
171         if (entry.getKey().equals(key)) {
172           bucket.remove(entry);
173           size--; // Decrease size
174           break; // Remove just one entry that matches the key
175         }
176     }
177   }
178
179   /** Return the number of mappings in this map */
180   public int size() {
181     return size;
182   }
183
184   /** Return a set consisting of the values in this map */
185   public java.util.Set<V> values() {
186     java.util.Set<V> set = new java.util.HashSet<V>();
187
188     for (int i = 0; i < capacity; i++) {
189       if (table[i] != null) {
190         LinkedList<Entry<K, V>> bucket = table[i];
191         for (Entry<K, V> entry: bucket)
192           set.add(entry.getValue());
193       }
194     }
195
196     return set;
197   }
198
199   /** Hash function */
200   private int hash(int hashCode) {
201     return supplementalHash(hashCode) & (capacity - 1);
202   }
203
204   /** Ensure the hashing is evenly distributed */
205   private static int supplementalHash(int h) {
206     h ^= (h >>> 20) ^ (h >>> 12);
207     return h ^ (h >>> 7) ^ (h >>> 4);
208   }
209
210   /** Return a power of 2 for initialCapacity */
211   private int trimToPowerOf2(int initialCapacity) {
212     int capacity = 1;
213     while (capacity < initialCapacity) {
214       capacity <<= 1;
215     }
216
217     return capacity;
218   }
219
220   /** Remove all entries from each bucket */
221   private void removeEntries() {
222     for (int i = 0; i < capacity; i++) {
223       if (table[i] != null) {
224         table[i].clear();
225       }
```

Margin labels (left column): size, values, hash, supplementalHash, trimToPowerOf2, removeEntries

```
226       }
227     }
228
229     /** Rehash the map */
230     private void rehash() {                                              rehash
231       java.util.Set<Entry<K, V>> set = entrySet(); // Get entries
232       capacity <<= 1; // Double capacity
233       table = new LinkedList[capacity]; // Create a new hash table
234       size = 0; // Clear size
235
236       for (Entry<K, V> entry: set) {
237         put(entry.getKey(), entry.getValue()); // Store to new table
238       }
229     }
240
241     /** Return a string representation for this map */
242     public String toString() {                                            toString
243       StringBuilder builder = new StringBuilder("[");
244
245       for (int i = 0; i < capacity; i++) {
246         if (table[i] != null && table[i].size() > 0)
247           for (Entry<K, V> entry: table[i])
248             builder.append(entry);
249       }
250
251       builder.append("]");
252       return builder.toString();
253   }
254 }
```

The **MyHashMap** class implements the **MyMap** interface using separate chaining. The parameters that determine the hash-table size and load factors are defined in the class. The default initial capacity is **4** (line 5) and the maximum capacity is $2^{30}$ (line 8). The current hash-table capacity is designed as a power of **2** (line 11). The default load factor threshold is **0.75f** (line 14). You can specify a custom load-factor threshold when constructing a map. The custom load-factor threshold is stored in **loadFactorThreshold** (line 17). The data field **size** denotes the number of entries in the map (line 20). The hash table is an array. Each cell in the array is a linked list (line 23).

hash-table parameters

Three constructors are provided to construct a map. You can construct a default map with the default capacity and load-factor threshold using the no-arg constructor (lines 26–28). You can construct a map with the specified capacity and a default load-factor threshold (lines 32–34). You can construct a map with the specified capacity and load-factor threshold (lines 38–46).

three constructors

The **clear** method removes all entries from the map (lines 49–52). It invokes **removeEntries()** that deletes all entries in the buckets (lines 221–227). This method takes $O(capacity)$ time.

clear

The **get(key)** method returns the value of the first entry with the specified key (lines 93–103). This method takes $O(1)$ time.

get

The **containsKey(key)** method checks whether the specified key is in the map by by invoking the **get** method (lines 55–60). Since **get** method takes $O(1)$ time, the **containsKey(key)** method takes $O(1)$ time.

containsKey

The **containsValue(value)** method checks whether the value is in the map (lines 63–74). This method takes $O(capacity + size)$ time. It is actually $O(capacity)$, since $capacity > size$.

containsValue

The **entrySet()** method returns a set that contains all entries in the map (lines 77–90). This method takes $O(capacity)$ time.

entrySet

**getAll**    The **getAll(key)** method returns the value of all entries with the specified key (lines 106–117). This method takes $O(1)$ time.

**isEmpty**    The **isEmpty()** method simply returns true if the map is empty (lines 120–122). This method takes $O(1)$ time.

**keySet**    The **keySet()** method returns all keys in the map as a set. The method finds the keys from each bucket and add them to a set (lines 125–137). This method takes $O(capacity)$ time.

**put**    The **put(key, value)** method adds a new entry into the map. The method first checks whether the size exceeds the load-factor threshold (line 141). If so, invoke **rehash()** (line 145) to increase the capacity and store entries into the new hash table.

**rehash**    The **rehash()** method first copies all entries in a set (line 231), doubles the capacity (line 232), creates a new hash table (line 233), and clears the size (line 234). The method then copies the entries into the new hash table (lines 236–238). The **rehash** method takes $O(capacity)$ time. If no rehash is performed, the **put** method takes $O(1)$ time to add a new entry.

**remove**    The **remove(key)** method removes all entries with the specified key in the map (lines 164–177). This method takes $O(1)$ time.

**size**    The **size()** method simply returns the size of the map (lines 180–182). This method takes $O(1)$ time.

**values**    The **values()** method returns all values in the map. The method examines each entry from all buckets and add it to a set (lines 185–197). This method takes $O(capacity)$ time.

**hash**    The **hash()** method invokes the **supplementalHash** to ensure that the hashing is evenly distributed to produce an index for the hash table (lines 200–208). This method takes $O(1)$ time.

Table 48.1 summarizes the time complexities of the methods in **MyHashMap**.

Since rehashing does not happen very often, the time complexity for the **put** method is $O(1)$. Note that the complexities of the **clear**, **entrySet**, **keySet**, **values**, and **rehash**

**TABLE 48.1**    Time Complexities for Methods in **MyHashMap**

| Methods | Time |
| --- | --- |
| clear() | $O(capacity)$ |
| containsKey(key: Key) | $O(1)$ |
| containsValue(value: V) | $O(capacity)$ |
| entrySet() | $O(capacity)$ |
| get(key: K) | $O(1)$ |
| getAll(key: K) | $O(1)$ |
| isEmpty() | $O(1)$ |
| keySet() | $O(capacity)$ |
| put(key: K, value: V) | $O(1)$ |
| remove(key: K) | $O(1)$ |
| size() | $O(1)$ |
| values() | $O(capacity)$ |
| rehash() | $O(capacity)$ |

methods depend on `capacity`, so to avoid poor performance for these methods you should choose an initial capacity carefully.

Listing 48.3 gives a test program that uses **MyHashMap**.

**LISTING 48.3** `TestMyHashMap.java`

```
 1 public class TestMyHashMap {
 2   public static void main(String[] args) {
 3     // Create a map
 4     MyMap<String, Integer> map = new MyHashMap<String, Integer>();
 5     map.put("Smith", 30);
 6     map.put("Anderson", 31);
 7     map.put("Lewis", 29);
 8     map.put("Cook", 29);
 9
10     System.out.println("Entries in map: " + map);
11
12     System.out.println("The age for " + "Lewis is " +
13       map.get("Lewis").intValue());
14
15     System.out.println("Is Smith in the map? " +
16       map.containsKey("Smith"));
17     System.out.println("Is age 33 in the map? " +
18         map.containsValue(33));
19
20     map.remove("Smith");
21     System.out.println("Entries in map: " + map);
22
23     map.clear();
24     System.out.println("Entries in map: " + map);
25   }
26 }
```

create a map
put entries

display entries

get value

is key in map?

is value in map?

remove entry

```
Entries in map: [[Anderson, 31][Smith, 30][Lewis, 29][Cook, 29]]
The age for Lewis is 29
Is Smith in the map? true
Is age 33 in the map? false
Entries in map: [[Anderson, 31][Lewis, 29][Cook, 29]]
Entries in map: []
```

The program creates a map using **MyHashMap** (line 4), adds entries to the map (lines 5–8), displays the entries (line 10), gets a value for a key (line 13), checks whether the map contains the key (line 16) and a value (line 18), removes an entry with the key "Smith" (line 20), and redisplays the entries in the map (line 22).

## 48.9 Set

A *set* is a data structure that stores distinct values. The Java collections framework defines the `java.util.Set` interface for modeling sets. Three concrete implementations are `java.util.HashSet`, `java.util.LinkedHashSet`, and `java.util.TreeSet`. `java.util.HashSet` is implemented using hashing, `java.util.LinkedHashSet` using `LinkedList`, and `java.util.TreeSet` using red-black trees.

set

You can implement **MyHashSet** using the same approach for implementing **MyHashMap**. The only difference is that key/value pairs are stored in the map, while elements are stored in the set.

We design our custom **Set** interface to mirror **java.util.Set** with some minor varia-
tions. The **java.util.Set** interface extends **java.util.Collection**. Our set interface is
the root interface. We name the interface **MySet** and a concrete class **MyHashSet**, as shown
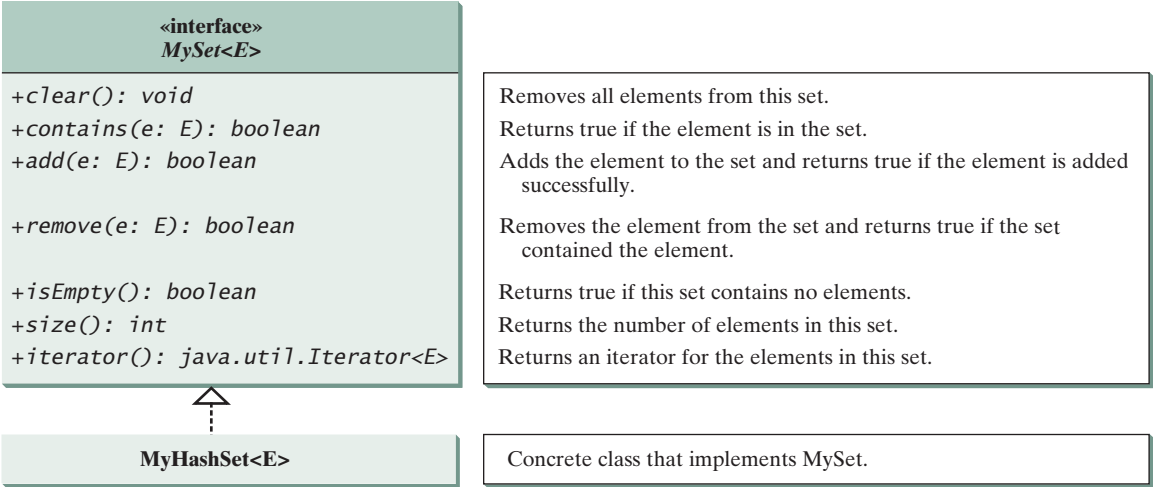in Figure 48.9.

**MySet**
**MyHashSet**

| «interface»<br>*MySet<E>* | |
|---|---|
| +clear(): void | Removes all elements from this set. |
| +contains(e: E): boolean | Returns true if the element is in the set. |
| +add(e: E): boolean | Adds the element to the set and returns true if the element is added successfully. |
| +remove(e: E): boolean | Removes the element from the set and returns true if the set contained the element. |
| +isEmpty(): boolean | Returns true if this set contains no elements. |
| +size(): int | Returns the number of elements in this set. |
| +iterator(): java.util.Iterator<E> | Returns an iterator for the elements in this set. |

| MyHashSet<E> | Concrete class that implements MySet. |
|---|---|

**FIGURE 48.9** **MyHashSet** implements the **MySet** interface.

Listing 48.4 shows the **MySet** interface and Listing 48.5 implements **MyHashSet** using
separate chaining.

## LISTING 48.4 MySet.java

```
1 public interface MySet<E> {
2    /** Remove all elements from this set */
3    public void clear();
4
5    /** Return true if the element is in the set */
6    public boolean contains(E e);
7
8    /** Add an element to the set */
9    public boolean add(E e);
10
11   /** Remove the element from the set */
12   public boolean remove(E e);
13
14   /** Return true if the set contains no elements */
15   public boolean isEmpty();
16
17   /** Return the number of elements in the set */
18   public int size();
19
20   /** Return an iterator for the elements in this set */
21   public java.util.Iterator iterator();
22 }
```

clear

contains

add

remove

isEmpty

size

iterator

## LISTING 48.5 MyHashSet.java

```
1 import java.util.LinkedList;
2
3 public class MyHashSet<E> implements MySet<E> {
```

class **MyHashSet**

```
 4    // Define the default hash-table size. Must be a power of 2
 5    private static int DEFAULT_INITIAL_CAPACITY = 16;                    default initial capacity
 6
 7    // Define the maximum hash-table size. 1 << 30 is same as 2^30
 8    private static int MAXIMUM_CAPACITY = 1 << 30;                       maximum capacity
 9
10    // Current hash-table capacity. Capacity is a power of 2
11    private int capacity;                                               current capacity
12
13    // Define default load factor
14    private static float DEFAULT_MAX_LOAD_FACTOR = 0.75f;               default max load factor
15
16    // Specify a load-factor threshold used in the hash table
17    private float loadFactorThreshold;                                  load-factor threshold
18
19    // The number of entries in the set
20    private int size = 0;                                               size
21
22    // Hash table is an array with each cell that is a linked list
23    private LinkedList<E>[] table;                                      hash table
24
25    /** Construct a set with the default capacity and load factor */
26    public MyHashSet() {                                                no-arg constructor
27      this(DEFAULT_INITIAL_CAPACITY, DEFAULT_MAX_LOAD_FACTOR);
28    }
29
30    /** Construct a set with the specified initial capacity and
31     * default load factor */
32    public MyHashSet(int initialCapacity) {                             constructor
33      this(initialCapacity, DEFAULT_MAX_LOAD_FACTOR);
34    }
35
36    /** Construct a set with the specified initial capacity
37     * and load factor */
38    public MyHashSet(int initialCapacity, float loadFactorThreshold) {  constructor
39      if (initialCapacity > MAXIMUM_CAPACITY)
40        this.capacity = MAXIMUM_CAPACITY;
41      else
42        this.capacity = trimToPowerOf2(initialCapacity);
43
44      this.loadFactorThreshold = loadFactorThreshold;
45      table = new LinkedList[capacity];
46    }
47
48    /** Remove all elements from this set */
49    public void clear() {                                               clear
50      size = 0;
51      removeElements();
52    }
53
54    /** Return true if the element is in the set */
55    public boolean contains(E e) {                                      contains
56      int bucketIndex = hash(e.hashCode());
57      if (table[bucketIndex] != null) {
58        LinkedList<E> bucket = table[bucketIndex];
59        for (E element: bucket)
60          if (element.equals(e))
61            return true;
62      }
63
```

```
64        return false;
65      }
66
67      /** Add an element to the set */
68      public boolean add(E e) {
69        if (contains(e))
70          return false;
71
72        if (size > capacity * loadFactorThreshold) {
73          if (capacity == MAXIMUM_CAPACITY)
74            throw new RuntimeException("Exceeding maximum capacity");
75
76          rehash();
77        }
78
79        int bucketIndex = hash(e.hashCode());
80
81        // Create a linked list for the bucket if it is not created
82        if (table[bucketIndex] == null) {
83          table[bucketIndex] = new LinkedList<E>();
84        }
85
86        // Add e to hashTable[index]
87        table[bucketIndex].add(e);
88
89        size++; // Increase size
90
91        return true;
92      }
93
94      /** Remove the element from the set */
95      public boolean remove(E e) {
96        if (!contains(e))
97          return false;
98
99        int bucketIndex = hash(e.hashCode());
100
101        // Create a linked list for the bucket if it is not created
102        if (table[bucketIndex] != null) {
103          LinkedList<E> bucket = table[bucketIndex];
104          for (E element: bucket)
105            if (e.equals(element)) {
106              bucket.remove(element);
107              break;
108            }
109        }
110
111        size--; // Decrease size
112
113        return true;
114      }
115
116      /** Return true if the set contains no elements */
117      public boolean isEmpty() {
118        return size == 0;
119      }
120
121      /** Return the number of elements in the set */
122      public int size() {
123        return size;
124      }
```

add

remove

isEmpty

size

```
125
126    /** Return an iterator for the elements in this set */
127    public java.util.Iterator<E> iterator() {                          iterator
128      return new MyHashSetIterator(this);
129    }
130
131    /** Inner class for iterator */
132    private class MyHashSetIterator implements java.util.Iterator<E> {   inner class
133      // Store the elements in a list
134      private java.util.ArrayList<E> list;
135      private int current = 0; // Point to the current element in list
136      MyHashSet<E> set;
137
138      /** Create a list from the set */
139      public MyHashSetIterator(MyHashSet<E> set) {
140        this.set = set;
141        list = setToList();
142      }
143
144      /** Next element for traversing? */
145      public boolean hasNext() {
146        if (current < list.size())
147          return true;
148
149        return false;
150      }
151
152      /** Get the current element and move cursor to the next */
153      public E next() {
154        return list.get(current++);
155      }
156
157     /** Remove the current element and refresh the list */
158      public void remove() {
159        // Delete the current element from the hash set
160        set.remove(list.get(current));
161       list.remove(current); // Remove the current element from the list
162      }
163    }
164
165    /** Hash function */
166    private int hash(int hashCode) {                                    hash
167      return supplementalHash(hashCode) & (capacity - 1);
168    }
169
170    /** Ensure the hashing is evenly distributed */
171    private static int supplementalHash(int h) {                        supplementalHash
172      h ^= (h >>> 20) ^ (h >>> 12);
173      return h ^ (h >>> 7) ^ (h >>> 4);
174    }
175
176    /** Return a power of 2 for initialCapacity */
177    private int trimToPowerOf2(int initialCapacity) {                   trimToPowerOf2
178      int capacity = 1;
179      while (capacity < initialCapacity) {
180        capacity <<= 1;
181      }
182
183      return capacity;
184    }
```

```
185
186    /** Remove all e from each bucket */
187    private void removeElements() {
188      for (int i = 0; i < capacity; i++) {
189        if (table[i] != null) {
190          table[i].clear();
191        }
192      }
193    }
194
195    /** Rehash the set */
196    private void rehash() {
197      java.util.ArrayList<E> list = setToList(); // Copy to a list
198      capacity <<= 1; // Double capacity
199      table = new LinkedList[capacity]; // Create a new hash table
200      size = 0;
201
202      for (E element: list) {
203        add(element); // Add from the old table to the new table
204      }
205    }
206
207    /** Copy elements in the hash set to an array list */
208    private java.util.ArrayList<E> setToList() {
209      java.util.ArrayList<E> list = new java.util.ArrayList<E>();
210
211      for (int i = 0; i < capacity; i++) {
212        if (table[i] != null) {
213          for (E e: table[i]) {
214            list.add(e);
215          }
216        }
217      }
218
219      return list;
220    }
221
222    /** Return a string representation for this set */
223    public String toString() {
224      java.util.ArrayList<E> list = setToList();
225      StringBuilder builder = new StringBuilder("[");
226
227      // Add the elements except the last one to the string builder
228      for (int i = 0; i < list.size() - 1; i++) {
229        builder.append(list.get(i) + ", ");
230      }
231
232      // Add the last element in the list to the string builder
233      if (list.size() == 0)
234        builder.append("]");
235      else
236        builder.append(list.get(list.size() - 1) + "]");
237
238      return builder.toString();
239    }
240 }
```

rehash *(margin label at line 196)*

setToList *(margin label at line 208)*

toString *(margin label at line 223)*

The **MyHashSet** class implements the **MySet** interface using separate chaining. Implementing **MyHashSet** is very similar to implementing **MyHashMap** except for the following differences:

1. The elements are stored in the hash table for **MyHashSet**, but the entries (key/value pairs) are stored in the hash table for **MyHashMap**.

2. The elements are all distinct in **MyHashSet**, but two entries may have the same keys in **MyHashMap**.

Three constructors are provided to construct a set. You can construct a default set with the default capacity and load factor using the no-arg constructor (lines 26–28). You can construct a set with the specified capacity and a default load factor (lines 32–34). You can construct a set with the specified capacity and load factor (lines 38–46).

The **clear** method removes all entries from the map (lines 49–52). It invokes **removeElements()** that deletes all elements in the buckets (lines 187–193). This method takes $O(capacity)$ time.

The **contains(element)** method checks whether the specified element is in the set by examining whether the designated bucket contains the element (lines 55–65). This method takes $O(1)$ time.

The **add(element)** method adds a new element into the set. The method first checks whether the size exceeds the load-factor threshold (line 72). If so, invoke **rehash()** (line 76) to increase the capacity and store entries into the new hash table.

The **rehash()** method first copies all elements in a list (line 197), doubles the capacity (line 198), obtains a new threshold (line 198), and creates a new hash table (line 199), and clears the size (line 200). The method then copies the entries into the new hash table (lines 202–203). The **rehash** method takes $O(capacity)$ time. If no rehash is performed, the **add** method takes $O(1)$ time to add a new element.

The **remove(element)** method removes the specified element in the set (lines 95–114). This method takes $O(1)$ time.

The **size()** method simply returns the size of the set (lines 122–124). This method takes $O(1)$ time.

The **iterator()** method returns an instance of **java.util.Iterator**. The **MyHashSetIterator** class implements **java.util.Iterator** to create a forward iterator. When a **MyHashSetIterator** is constructed, it copies all the elements in the set to a list (line 141). The variable **current** points to the element in the list. Initially, **current** is **0** (line 135), which points to the first element in the list. **MyHashSetIterator** implements the methods **hasNext()**, **next()**, and **remove()** in **java.util.Iterator**. Invoking **hasNext()** returns true if **current < list.size()**. Invoking **next()** returns the current element and moves **current** to point to the next element (line 153). Invoking **remove()** removes the current element in the iterator from the set.

The **hash()** method invokes the **supplementalHash** to ensure that the hashing is evenly distributed to produce an index for the hash table (lines 166–174). This method takes $O(1)$ time.

Table 48.2 summarizes the time complexity of the methods in **MyHashSet**.

Listing 48.6 gives a test program that uses **MyHashSet**.

## LISTING 48.6  TestMyHashSet.java

```
1 public class TestMyHashSet {
2   public static void main(String[] args) {
3     // Create a MyHashSet
4     MySet<String> set = new MyHashSet<String>();
5     set.add("Smith");
6     set.add("Anderson");
```

```
display elements
set size

remove element
display elements

clear set
```

```
 7       set.add("Lewis");
 8       set.add("Cook");
 9
10       System.out.println("Elements in set: " + set);
11       System.out.println("Number of elements in set: " + set.size());
12       System.out.println("Is Smith in set? " + set.contains("Smith"));
13
14       set.remove("Smith");
15       System.out.println("Elements in set: " + set);
16
17       set.clear();
18       System.out.println("Elements in set: " + set);
19   }
20 }
```

```
Elements in set: [Smith, Lewis, Anderson, Cook]
Number of elements in set: 4
Is Smith in set? true
Elements in set: [Lewis, Anderson, Cook]
Elements in set: []
```

**TABLE 48.2** Time Complexities for Methods in **MyHashMap**

| Methods | Time |
|---|---|
| clear() | $O(capacity)$ |
| contains(e: E) | $O(1)$ |
| add(e: E) | $O(1)$ |
| remove(e: E) | $O(1)$ |
| isEmpty() | $O(1)$ |
| size() | $O(1)$ |
| iterator() | $O(capacity)$ |
| rehash() | $O(capacity)$ |

The program creates a set using **MyHashSet** (line 4), adds elements to the set (lines 5–8), displays the elements (line 10), gets the size (line 11), checks whether the set contains the element (line 12), removes an element (line 14), and clears the set (line 17).

## KEY TERMS

associative array    48–2
clustering    48–5
dictionary    48–2
double hashing    48–4
hash code    48–3
hash function    48–2
hash map    48–2
hash set    48–18
hash table    48–2

linear probing    48–4
load factor    48–8
open addressing    48–4
perfect hash function    48–2
polynomial hash code    48–4
rehashing    48–8
secondary clustering    48–5
separate chaining    48–4

## CHAPTER SUMMARY

**1.** A *map* is a data structure that stores entries. Each entry contains two parts: *key* and *value*. The key is also called a *search key*, which is used to search for the corresponding value. You can implement a map to obtain $O(1)$ time complexity on search, retrieval, insertion, and deletion, using the hashing technique.

**2.** A set is a data structure that stores elements. You can use the hashing technique to implement a set to achieve $O(1)$ time complexity on search, insertion, and deletion for a set.

**3.** *Hashing* is a technique that retrieves the value using the index obtained from key without performing a search. A typical hash function first converts a search key to an integer value called a *hash code*, then compresses the hash code into an index to the hash table.

**4.** A collision occurs when two keys are mapped to the same index in a hash table. Generally, there are two ways for handling collisions: *open addressing* and *separate chaining*.

**5.** Open addressing is finding an open location in the hash table in the event of collision. Open addressing has several variations: *linear probing*, *quadratic probing*, and *double hashing*.

**6.** The separate chaining scheme places all entries with the same hash index into the same location, rather than finding new locations. Each location in the separate chaining scheme is called a *bucket*. A bucket is a container that holds multiple entries.

## REVIEW QUESTIONS

**Sections 48.1–48.5**

**48.1** What is a hash function? What is a perfect hash function? What is a collision?

**48.2** What is a hash code? What is the hash code for **Byte**, **Short**, **Integer**, and **Character**?

**48.3** How is the hash code for a **Float** object computed?

**48.4** How is the hash code for a **Long** object computed?

**48.5** How is the hash code for a **Double** object computed?

**48.6** How is the hash code for a **String** object computed?

**48.7** How is a hash code compressed to an integer representing the index in a hash table?

**48.8** What is open addressing? What is linear probing? What is quadratic probing? What is double hashing?

**48.9** Describe the clustering problem for linear probing.

**48.10** What is the secondary clustering?

**48.11** Show the hash table of size 11 after inserting entries with keys 34, 29, 53, 44, 120, 39, 45, and 40, using linear probing.

**48.12** Show the hash table of size 11 after inserting entries with keys 34, 29, 53, 44, 120, 39, 45, and 40, using quadratic probing.

**48.13** Show the hash table of size 11 after inserting entries with keys 34, 29, 53, 44, 120, 39, 45, and 40, using double hashing with the following functions:

```
h(k) = k % 11;
h'(k) = 7 - k % 7;
```

**48.14** Suppose the size of the table is 10. What is the probe sequence for a key 12 using the following double hashing functions?

```
h(k) = k % 10;
h'(k) = 7 - k % 7;
```

### Sections 48.6–48.8

**48.15** Show the hash table of size 11 after inserting entries with keys 34, 29, 53, 44, 120, 39, 45, and 40, using separate chaining.

**48.16** In Listing 48.5, the **remove** method in the iterator removes the current element from the set. It also removes the current element from the internal list (line 165):

```
// Remove the current element from the list
list.remove(current);
```

Why is it necessary?

## PROGRAMMING EXERCISES

**48.1\*\*** (*Implementing MyMap using open addressing with linear probing*) Create a new concrete class that implements **MyMap** using open addressing with linear probing. For simplicity, use **f(key) = key % size** as the hash function, where size is the hash-table size. Initially, the hash-table size is **4**. The table size is doubled whenever the load factor exceeds the threshold (**0.5**).

**48.2\*\*** (*Implementing MyMap using open addressing with quadratic probing*) Create a new concrete class that implements **MyMap** using open addressing with quadratic probing. For simplicity, use **f(key) = key % size** as the hash function, where **size** is the hash-table size. Initially, the hash-table size is **4**. The table size is doubled whenever the load factor exceeds the threshold (**0.5**).

**48.3\*\*** (*Implementing MyMap using open addressing with double hashing*) Create a new concrete class that implements **MyMap** using open addressing with double probing. For simplicity, use **f(key) = key % size** as the hash function, where **size** is the hash-table size. Initially, the hash-table size is **4**. The table size is doubled whenever the load factor exceeds the threshold (**0.5**).

**48.4\*\*** (*Modifying MyHashMap with distinct keys*) Modify **MyHashMap** so that all entries in it have different keys.

**48.5\*\*** (*Implementing MyHashSet using MyHashMap*) Implement **MyHashSet** using **MyHashMap**. Note that you can create entries with (key, key), rather than (key, value).

**48.6\*\*** (*Animating linear probing*) Write a Java applet that animates linear probing as shown in Figure 48.3. You can change the initial size of the hash-table in the applet. Assume the load-factor threshold is **0.75**.

**48.7\*\*** (*Animating separate chaining*) Write a Java applet that animates **MyHashSet** as shown in Figure 48.7. You can change the initial size of the hash table. Assume the load-factor threshold is **0.75**.