

SETS AND MAPS

Chapter 7

Introduction

2

- Java Collection Framework (`ArrayList` and `LinkedList`)
- The classes that implement the `List` interface are all *indexed* collections
 - ▣ An index or subscript is associated with each element
 - ▣ The element's index often reflects the relative order of its insertion into the list
 - ▣ Searching for a particular value in a list is generally $O(n)$
 - ▣ An exception is a binary search of a sorted object, which is $O(\log n)$

Introduction (cont.)

3

- Next, we consider another part of the `Collection` hierarchy: the `Set` interface and the classes that implement it
- `Set` objects
 - ▣ are not indexed
 - ▣ do not reveal the order of insertion of items
 - ▣ enable efficient search and retrieval of information
 - ▣ allow removal of elements without moving other elements around

Introduction (cont.)

4

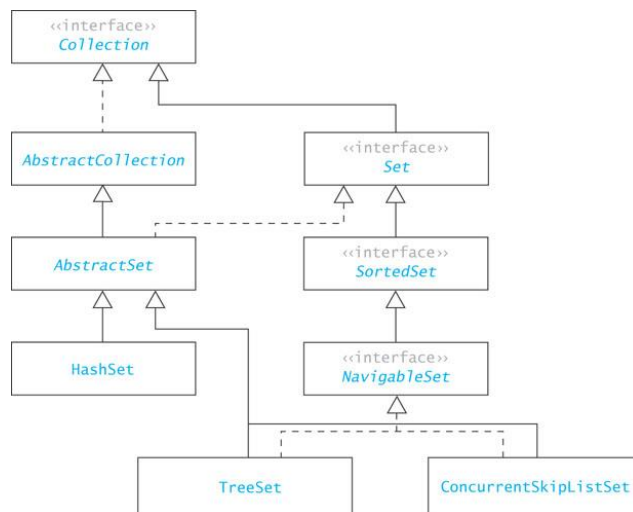
- Relative to a `Set`, `Map` objects provide efficient search and retrieval of entries that contain pairs of objects (a unique key and the information)
- Hash tables (implemented by a `Map` or `Set`) store objects at arbitrary locations and offer an average constant time for insertion, removal, and searching

5 Sets and the Set Interface

Section 7.1

Sets and the Set Interface

6



The Set Abstraction

7

- A set is a collection that contains no duplicate elements and at most one null element
 - adding "apples" to the set {"apples", "oranges", "pineapples"} results in the same set (no change)
- Operations on sets include:
 - testing for membership
 - adding elements
 - removing elements
 - union $A \cup B$
 - intersection $A \cap B$
 - difference $A - B$
 - subset $A \subset B$

The Set Abstraction(cont.)

8

- The union of two sets A, B is a set whose elements belong either to A or B or to both A and B.
Example: $\{1, 3, 5, 7\} \cup \{2, 3, 4, 5\}$ is $\{1, 2, 3, 4, 5, 7\}$
- The intersection of sets A, B is the set whose elements belong to both A and B.
Example: $\{1, 3, 5, 7\} \cap \{2, 3, 4, 5\}$ is $\{3, 5\}$
- The difference of sets A, B is the set whose elements belong to A but not to B.
Examples: $\{1, 3, 5, 7\} - \{2, 3, 4, 5\}$ is $\{1, 7\}$; $\{2, 3, 4, 5\} - \{1, 3, 5, 7\}$ is $\{2, 4\}$
- Set A is a subset of set B if every element of set A is also an element of set B.
Example: $\{1, 3, 5, 7\} \subset \{1, 2, 3, 4, 5, 7\}$ is true

The Set Interface and Methods(cont.)

9

Method	Behavior
<code>boolean add(E obj)</code>	Adds item <code>obj</code> to this set if it is not already present (optional operation) and returns true . Returns false if <code>obj</code> is already in the set.
<code>boolean addAll(Collection<E> coll)</code>	Adds all of the elements in collection <code>coll</code> to this set if they're not already present (optional operation). Returns true if the set is changed. Implements <i>set union</i> if <code>coll</code> is a <code>Set</code> .
<code>boolean contains(Object obj)</code>	Returns true if this set contains an element that is equal to <code>obj</code> . Implements a test for <i>set membership</i> .
<code>boolean containsAll(Collection<E> coll)</code>	Returns true if this set contains all of the elements of collection <code>coll</code> . If <code>coll</code> is a set, returns true if this set is a subset of <code>coll</code> .
<code>boolean isEmpty()</code>	Returns true if this set contains no elements.
<code>Iterator<E> iterator()</code>	Returns an iterator over the elements in this set.
<code>boolean remove(Object obj)</code>	Removes the set element equal to <code>obj</code> if it is present (optional operation). Returns true if the object was removed.
<code>boolean removeAll(Collection<E> coll)</code>	Removes from this set all of its elements that are contained in collection <code>coll</code> (optional operation). Returns true if this set is changed. If <code>coll</code> is a set, performs the <i>set difference</i> operation.
<code>boolean retainAll(Collection<E> coll)</code>	Retains only the elements in this set that are contained in collection <code>coll</code> (optional operation). Returns true if this set is changed. If <code>coll</code> is a set, performs the <i>set intersection</i> operation.
<code>int size()</code>	Returns the number of elements in this set (its cardinality).

Using sets in Java

10

```
import java.util.Set;

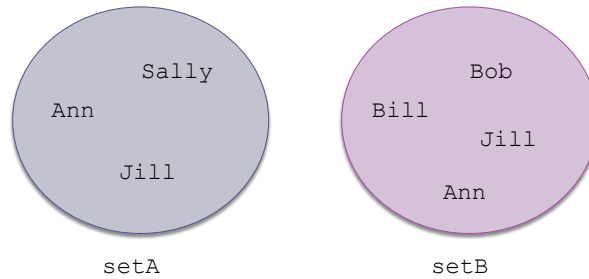
Set<String> setA = new HashSet<String>();
Set<String> setB = new TreeSet<String>();
```

HashSet is implemented using Hash Table (coming next)

TreeSet is implemented using a special kind of Binary Search Tree – Red-Black Trees.

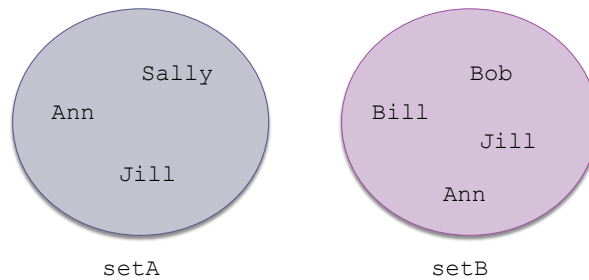
The Set Interface and Methods(cont.)

11



The Set Interface and Methods(cont.)

12

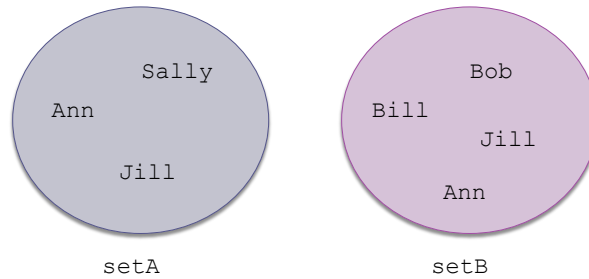


```
setA.addAll(setB); // Set Union  
System.out.println(setA);
```

Outputs:
[Bill, Jill, Ann, Sally, Bob]

The Set Interface and Methods(cont.)

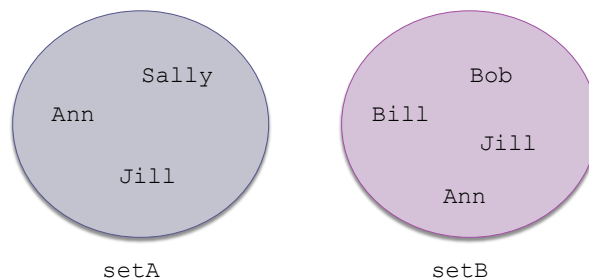
13



If a copy of original setA is in setACopy, then ...

The Set Interface and Methods(cont.)

14



```
setACopy.retainAll(setB); // Set Intersection
```

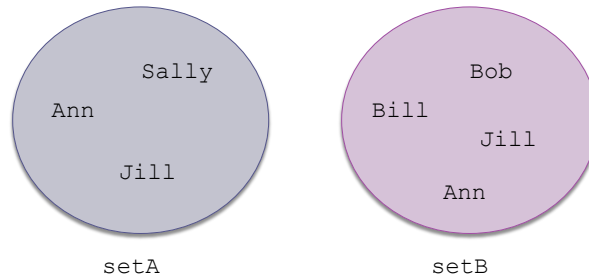
```
System.out.println(setACopy);
```

Outputs:

```
[Jill, Ann]
```

The Set Interface and Methods(cont.)

15



```
setACopy.removeAll(setB); // Set Difference
```

```
System.out.println(setACopy);
```

Outputs:
[Sally]

Comparison of Lists and Sets

16

- Collections implementing the `Set` interface may contain only unique elements
- Unlike the `List.add` method, the `Set.add` method returns `false` if you attempt to insert a duplicate item
- Unlike a `List`, a `Set` does not have a `get` method—elements cannot be accessed by index

Comparison of Lists and Sets (cont.)

17

- You can iterate through all elements in a `Set` using an `Iterator` object, but the elements will be accessed in arbitrary order

```
for (String nextItem : setA) {  
    //Do something with nextItem  
    ...  
}
```

18

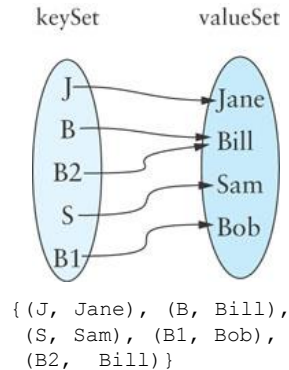
Maps and the `Map` Interface

Section 7.2

Maps and the Map Interface

19

- The Map is related to the Set
- Mathematically, a Map is a set of ordered pairs whose elements are known as the **key** and the **value**
- Keys must be unique, but values need not be unique
- You can think of each key as a “mapping” to a particular value
- A map provides efficient storage and retrieval of information in a table
- A map can have *many-to-one* mapping: (B, Bill), (B2, Bill)



Maps and the Map Interface(cont.)

20

- In an *onto* mapping, all the elements of valueSet have a corresponding member in keySet
- The Map interface should have methods of the form
 - V.get (Object key)
 - V.put (K key, V value)

Maps and the Map Interface(cont.)

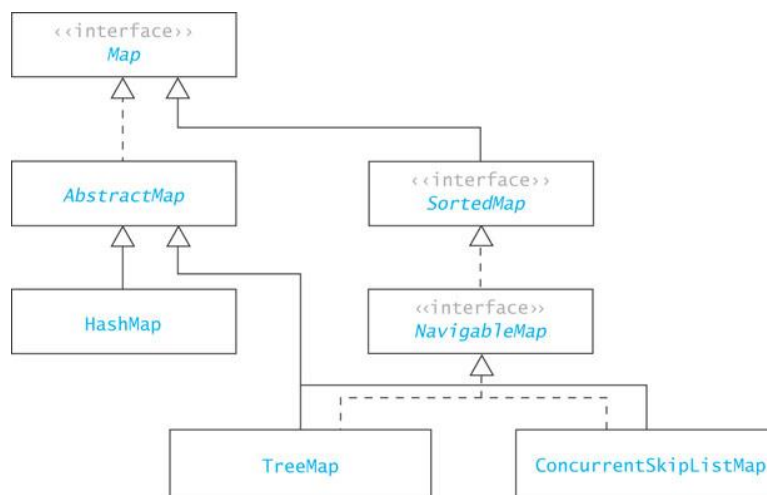
21

- When information about an item is stored in a table, the information should have a unique ID
- A unique ID may or may not be a number
- This unique ID is equivalent to a key

Type of item	Key	Value
University student	Student ID number	Student name, address, major, grade point average
Online store customer	E-mail address	Customer name, address, credit card information, shopping cart
Inventory item	Part ID	Description, quantity, manufacturer, cost, price

Map Hierarchy

22



Map Interface

23

Method	Behavior
V get(Object key)	Returns the value associated with the specified key. Returns null if the key is not present.
boolean isEmpty()	Returns true if this map contains no key-value mappings.
V put(K key, V value)	Associates the specified value with the specified key in this map (optional operation). Returns the previous value associated with the specified key, or null if there was no mapping for the key.
V remove(Object key)	Removes the mapping for this key from this map if it is present (optional operation). Returns the previous value associated with the specified key, or null if there was no mapping for the key.
int size()	Returns the number of key-value mappings in this map.

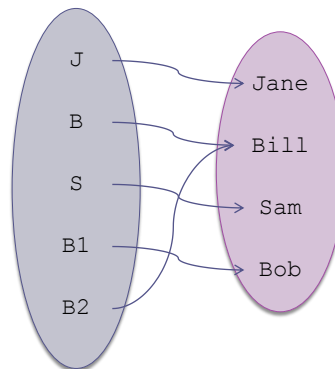
Map Interface (cont.)

24

- The following statements build a Map object:

```
Map<String, String> aMap =
    new HashMap<String,
        String>();
```

```
aMap.put("J", "Jane");
aMap.put("B", "Bill");
aMap.put("S", "Sam");
aMap.put("B1", "Bob");
aMap.put("B2", "Bill");
```



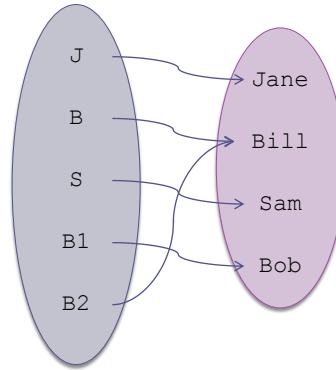
Map Interface (cont.)

25

```
aMap.get("B1")
```

returns:

"Bob"



Map Interface (cont.)

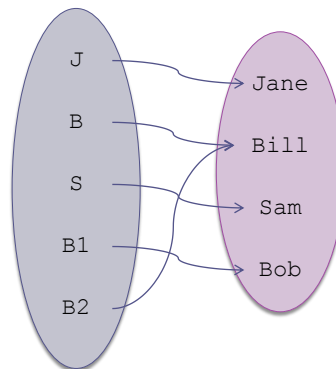
26

```
aMap.get("Bill")
```

returns:

null

("Bill" is a value, not a key)



Map **Interface** (cont.)

27

```
Map<String, String> places = new HashMap<String, Place>();

places.put("Bryn MawrPA", new Place("Bryn Mawr", "PA",
"19010"));

Places.get("Bryn MawrPA");

returns

<Bryn Mawr, PA, 19010>
```

28

Hash Tables

Section 7.3

Hash Tables

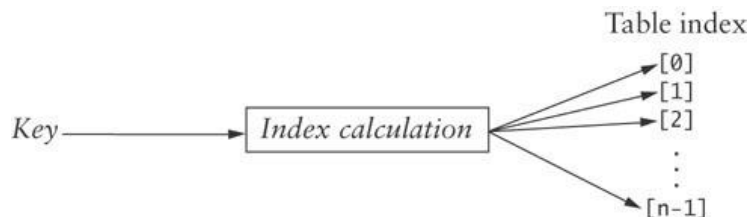
29

- The goal of hash table is to be able to access an entry based on its key value, not its location
- We want to be able to access an entry directly through its key value, rather than by having to determine its location first by searching for the key value in an array
- Using a hash table enables us to retrieve an entry in constant time (on average, $O(1)$)

Hash Codes and Index Calculation

30

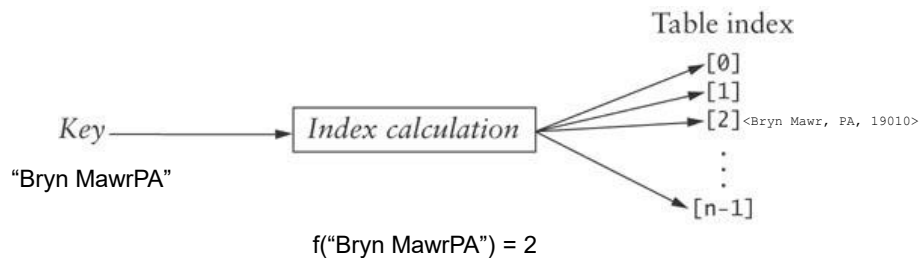
- The basis of hashing is to transform the item's key value into an integer value (its *hash code*) which is then transformed into a table index



Hash Codes and Index Calculation

31

- The basis of hashing is to transform the item's key value into an integer value (its *hash code*) which is then transformed into a table index



Hash Codes and Index Calculation

32

- $F(\langle \text{key} \rangle)$ returns an index in range $[0..n-1]$
- Goal: Similar keys map to different locations in an array.

$f(\text{"Richards"}) = 53$

$f(\text{"Richardson"}) = 417$

- When two or more keys map to the same location, it is called a **collision**.

E.g. $f(\text{"Richards"}) = 53$

$f(\text{"Deepak"}) = 53$

Methods for Generating Hash Codes

33

- In most applications, a key will consist of strings of letters or digits (such as a social security number, an email address, or a partial ID) rather than a single character
- The number of possible key values is much larger than the table size. E.g. 10-letter strings have 26^{10} keys!
- Generating good hash codes typically is an experimental process
- The goal is a *random distribution of values*
- Simple algorithms sometimes generate lots of **collisions**

Java hashCode Method

34

- Example hash function: sum the `int` values of all

$$f1 = (\text{code}(S_0) + \text{code}(S_1) + \dots + \text{code}(S_{N-1})) \% \text{tableSize}$$
 But, this returns the same hash code for "sign" and "sing"!
 - Example hash function: use position + code

$$f2 = [1 * \text{code}(S_0) + 2 * \text{code}(S_1) + \dots + N * \text{code}(S_{N-1})] \% \text{tableSize}$$
 - Java uses the following:

$$f3 = S_0 \times 31^{(n-1)} + S_1 \times 31^{(n-2)} + \dots + S_{N-1}$$
- <string>.hashCode() method in Java...

Java hashCode Method

35

- The Java API algorithm accounts for position of the characters as well
- `<string>.hashCode()` returns the integer calculated by the formula:

$$s_0 \times 31^{(n-1)} + s_1 \times 31^{(n-2)} + \dots + s_{n-1}$$

where s_i is the i th character of the string, and n is the length of the string

- “Cat” has a hash code of:

$$'C' \times 31^2 + 'a' \times 31 + 't' = 67,510$$

- 31 is a prime number, and prime numbers generate relatively few collisions

Java hashCode Method

36

String	hashCode()
“Tom”	84274
“Dick”	2129869
“Harry”	69496448
“Sam”	82879
“Pete”	2484038

Java hashCode Method (cont.)

37

- Because there are too many possible strings, the integer value returned by `hashCode()` can't be unique
- However, because the `hashCode()` method distributes the hash code values fairly evenly throughout the range, the probability of two strings having the same hash code is low
- The probability of a collision with

```
s.hashCode() % table.length
```

 is proportional to how full the table is

Methods for Generating Hash Codes (cont.)

38

- A good hash function should be relatively simple and efficient to compute
- It doesn't make sense to use an $O(n)$ hash function to avoid doing an $O(n)$ search

Open Addressing

39

- We now consider two ways to organize hash tables:
 - ▣ open addressing
 - ▣ chaining
- In open addressing, *linear probing* can be used to access an item in a hash table
 - ▣ If the index calculated for an item's key is occupied by an item with that key, we have found the item
 - ▣ If that element contains an item with a different key, increment the index by one
 - ▣ Keep incrementing until you find the key or a `null` entry (assuming the table is not full)

Open Addressing (cont.)

40

Algorithm for Accessing an Item in a Hash Table

1. Compute the index by taking the item's `hashCode()` % `table.length`.
2. **if** `table[index]` is `null`
3. The item is not in the table.
4. **else if** `table[index]` is equal to the item
5. The item is in the table.
6. **else**
6. Continue to search the table by incrementing the index until either the item is found or a `null` entry is found.

Table Wraparound and Search Termination

41

- As you increment the table index, your table should wrap around as in a circular array
- This enables you to search the part of the table before the hash code value in addition to the part of the table after the hash code value
- But it could lead to an infinite loop
- How do you know when to stop searching if the table is full and you have not found the correct value?
 - ▣ Stop when the index value for the next probe is the same as the hash code value for the object
 - ▣ Ensure that the table is never full by increasing its size after an insertion when its load factor exceeds a specified threshold

Hash Code Insertion Example (cont.)

42

	[0]	Dick
Pete	[1]	Sam
	[2]	Pete
	[3]	Harry
	[4]	Tom

Name	hashCode()	hashCode()%5
"Tom"	84274	4
"Dick"	2129869	4
"Harry"	69496448	3
"Sam"	82879	4
"Pete"	2484038	3

Retrieval of "Tom" or "Harry" takes one step,
 $O(1)$

Because of collisions, retrieval of the others
requires a linear search

Hash Code Insertion Example (cont.)

43

Name	hashCode()	hashCode()%11
"Tom"	84274	3
"Dick"	2129869	5
"Harry"	69496448	10
"Sam"	82879	5
"Pete"	2484038	7

[0]	
[1]	
[2]	
[3]	
[4]	
[5]	
[6]	
[7]	
[8]	
[9]	
[10]	

Hash Code Insertion Example (cont.)

44

Name	hashCode()	hashCode()%11
"Tom"	84274	3
"Dick"	2129869	5
"Harry"	69496448	10
"Sam"	82879	5
"Pete"	2484038	7

[0]	
[1]	
[2]	
[3]	Tom
[4]	
[5]	Dick
[6]	Sam
[7]	Pete
[8]	
[9]	
[10]	Harry

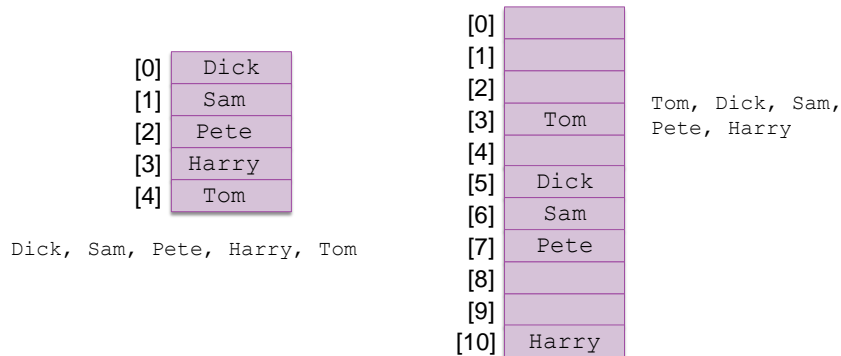
Only one
collision
occurred

The best way to reduce the possibility of collision (and reduce linear search retrieval time because of collisions) is to increase the table size

Traversing a Hash Table

45

- You cannot traverse a hash table in a meaningful way since the sequence of stored values is arbitrary



Deleting an Item Using Open Addressing

46

- When an item is deleted, you cannot simply set its table entry to null
- If we search for an item that may have collided with the deleted item, we may conclude incorrectly that it is not in the table.
- Instead, store a dummy value or mark the location as available, but previously occupied
- Deleted items reduce search efficiency which is partially mitigated if they are marked as available
- You cannot simply replace a deleted item with a new item until you verify that the new item is not in the table

Reducing Collisions by Expanding the Table Size

47

- Use a prime number for the size of the table to reduce collisions
- A fuller table results in more collisions, so, when a hash table becomes sufficiently full, a larger table should be allocated and the entries reinserted
- You must reinsert (*rehash*) values into the new table; do not copy values as some search chains which were wrapped may break
- Deleted items are not reinserted, which saves space and reduces the length of some search chains

Reducing Collisions Using Quadratic Probing

48

- Linear probing tends to form clusters of keys in the hash table, causing longer search chains
- *Quadratic probing* can reduce the effect of clustering
 - ▣ Increments form a quadratic series ($1 + 2^2 + 3^2 + \dots$)

```
probeNum++;
index = (startIndex + probeNum * probeNum) % table.length
```

- If an item has a hash code of 5, successive values of index will be 6 ($5+1$), 9 ($5+4$), 14 ($5+9$), ...

[0]		[0]	
[1]		[1]	
[2]		[2]	
[3]		[3]	
[4]		[4]	
[5]	1 st item with hash code 5	[5]	1 st item with hash code 5
[6]	1 st item with hash code 6	[6]	1 st item with hash code 6
[7]	2 nd item with hash code 5	[7]	2 nd item with hash code 6
[8]	2 nd item with hash code 6	[8]	1 st item with hash code 7
[9]	1 st item with hash code 7	[9]	2 nd item with hash code 5
[10]		[10]	

Problems with Quadratic Probing

49

- The disadvantage of quadratic probing is that the next index calculation is time-consuming, involving multiplication, addition, and modulo division
- A more efficient way to calculate the next index is:

```
k += 2;  
index = (index + k) % table.length;
```

Problems with Quadratic Probing

(cont.)

50

- Examples:
 - If the initial value of k is -1 , successive values of k will be $1, 3, 5, \dots$
 - If the initial value of index is 5 , successive value of index will be $6 (= 5 + 1), 9 (= 5 + 1 + 3), 14 (= 5 + 1 + 3 + 5), \dots$

Problems with Quadratic Probing (cont.)

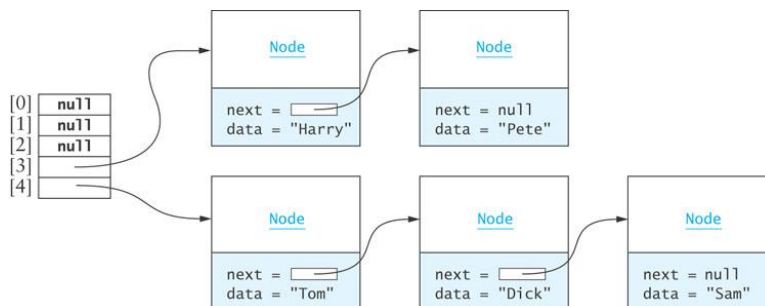
51

- A more serious problem is that not all table elements are examined when looking for an insertion index; this may mean that
 - ▣ an item can't be inserted even when the table is not full
 - ▣ the program will get stuck in an infinite loop searching for an empty slot
- If the table size is a prime number and it is never more than half full, this won't happen
- However, requiring a half empty table wastes a lot of memory

Chaining

52

- *Chaining* is an alternative to open addressing
- Each table element references a linked list that contains all of the items that hash to the same table index
 - ▣ The linked list often is called a *bucket*
 - ▣ The approach sometimes is called *bucket hashing*



Chaining (cont.)

53

- Advantages relative to open addressing:
 - ▣ Only items that have the same value for their hash codes are examined when looking for an object
 - ▣ You can store more elements in the table than the number of table slots (indices)
 - ▣ Once you determine an item is not present, you can insert it at the beginning or end of the list
 - ▣ To remove an item, you simply delete it; you do not need to replace it with a dummy item or mark it as deleted

Performance of Hash Tables

54

- *Load factor*

$$\text{load factor} = \frac{\text{\#filled cells}}{\text{table size}}$$

- *Lower load factor -> better performance*
- *Higher load factor -> worse performance*
- *If there are no collisions, performance for search and retrieval is $O(1)$ regardless of table size*

Performance of Open Addressing versus Chaining

55

- Donald E. Knuth derived the following formula for the expected number of comparisons, c , required for finding an item that is in a hash table using open addressing with linear probing and a load factor L

$$c = \frac{1}{2} \left(1 + \frac{1}{1 - L} \right)$$

Performance of Open Addressing versus Chaining (cont.)

56

- Using chaining, if an item is in the table, on average we must examine the table element corresponding to the item's hash code and then half of the items in each list

$$c = 1 + \frac{L}{2}$$

where L is the average number of items in a list (the number of items divided by the table size)

Performance of Hash Tables versus Sorted Array and Binary Search Tree

57

- The number of comparisons required for a binary search of a sorted array is $O(\log n)$
 - ▣ A sorted array of size 128 requires up to 7 probes (2^7 is 128) which is more than for a hash table of any size that is 90% full
 - ▣ A binary search tree performs similarly
- Insertion or removal

hash table	$O(1)$ expected; worst case $O(n)$
unsorted array	$O(n)$
binary search tree	$O(\log n)$; worst case $O(n)$