

### ★ 11.3.3 Universal hashing

If a malicious adversary chooses the keys to be hashed by some fixed hash function, then the adversary can choose  $n$  keys that all hash to the same slot, yielding an average retrieval time of  $\Theta(n)$ . Any fixed hash function is vulnerable to such terrible worst-case behavior; the only effective way to improve the situation is to choose the hash function *randomly* in a way that is *independent* of the keys that are actually going to be stored. This approach, called **universal hashing**, can yield provably good performance on average, no matter which keys the adversary chooses.

In universal hashing, at the beginning of execution we select the hash function at random from a carefully designed class of functions. As in the case of quick-sort, randomization guarantees that no single input will always evoke worst-case behavior. Because we randomly select the hash function, the algorithm can behave differently on each execution, even for the same input, guaranteeing good average-case performance for any input. Returning to the example of a compiler's symbol table, we find that the programmer's choice of identifiers cannot now cause consistently poor hashing performance. Poor performance occurs only when the compiler chooses a random hash function that causes the set of identifiers to hash poorly, but the probability of this situation occurring is small and is the same for any set of identifiers of the same size.

Let  $\mathcal{H}$  be a finite collection of hash functions that map a given universe  $U$  of keys into the range  $\{0, 1, \dots, m-1\}$ . Such a collection is said to be **universal** if for each pair of distinct keys  $k, l \in U$ , the number of hash functions  $h \in \mathcal{H}$  for which  $h(k) = h(l)$  is at most  $|\mathcal{H}|/m$ . In other words, with a hash function randomly chosen from  $\mathcal{H}$ , the chance of a collision between distinct keys  $k$  and  $l$  is no more than the chance  $1/m$  of a collision if  $h(k)$  and  $h(l)$  were randomly and independently chosen from the set  $\{0, 1, \dots, m-1\}$ .

The following theorem shows that a universal class of hash functions gives good average-case behavior. Recall that  $n_i$  denotes the length of list  $T[i]$ .

#### **Theorem 11.3**

Suppose that a hash function  $h$  is chosen randomly from a universal collection of hash functions and has been used to hash  $n$  keys into a table  $T$  of size  $m$ , using chaining to resolve collisions. If key  $k$  is not in the table, then the expected length  $E[n_{h(k)}]$  of the list that key  $k$  hashes to is at most the load factor  $\alpha = n/m$ . If key  $k$  is in the table, then the expected length  $E[n_{h(k)}]$  of the list containing key  $k$  is at most  $1 + \alpha$ .

**Proof** We note that the expectations here are over the choice of the hash function and do not depend on any assumptions about the distribution of the keys. For each pair  $k$  and  $l$  of distinct keys, define the indicator random variable

$X_{kl} = \mathbb{I}\{h(k) = h(l)\}$ . Since by the definition of a universal collection of hash functions, a single pair of keys collides with probability at most  $1/m$ , we have  $\Pr\{h(k) = h(l)\} \leq 1/m$ . By Lemma 5.1, therefore, we have  $\mathbb{E}[X_{kl}] \leq 1/m$ .

Next we define, for each key  $k$ , the random variable  $Y_k$  that equals the number of keys other than  $k$  that hash to the same slot as  $k$ , so that

$$Y_k = \sum_{\substack{l \in T \\ l \neq k}} X_{kl}.$$

Thus we have

$$\begin{aligned} \mathbb{E}[Y_k] &= \mathbb{E}\left[\sum_{\substack{l \in T \\ l \neq k}} X_{kl}\right] \\ &= \sum_{\substack{l \in T \\ l \neq k}} \mathbb{E}[X_{kl}] \quad (\text{by linearity of expectation}) \\ &\leq \sum_{\substack{l \in T \\ l \neq k}} \frac{1}{m}. \end{aligned}$$

The remainder of the proof depends on whether key  $k$  is in table  $T$ .

- If  $k \notin T$ , then  $n_{h(k)} = Y_k$  and  $|\{l : l \in T \text{ and } l \neq k\}| = n$ . Thus  $\mathbb{E}[n_{h(k)}] = \mathbb{E}[Y_k] \leq n/m = \alpha$ .
- If  $k \in T$ , then because key  $k$  appears in list  $T[h(k)]$  and the count  $Y_k$  does not include key  $k$ , we have  $n_{h(k)} = Y_k + 1$  and  $|\{l : l \in T \text{ and } l \neq k\}| = n - 1$ . Thus  $\mathbb{E}[n_{h(k)}] = \mathbb{E}[Y_k] + 1 \leq (n - 1)/m + 1 = 1 + \alpha - 1/m < 1 + \alpha$ . ■

The following corollary says universal hashing provides the desired payoff: it has now become impossible for an adversary to pick a sequence of operations that forces the worst-case running time. By cleverly randomizing the choice of hash function at run time, we guarantee that we can process every sequence of operations with a good average-case running time.

#### **Corollary 11.4**

Using universal hashing and collision resolution by chaining in an initially empty table with  $m$  slots, it takes expected time  $\Theta(n)$  to handle any sequence of  $n$  INSERT, SEARCH, and DELETE operations containing  $O(m)$  INSERT operations.

**Proof** Since the number of insertions is  $O(m)$ , we have  $n = O(m)$  and so  $\alpha = O(1)$ . The INSERT and DELETE operations take constant time and, by Theorem 11.3, the expected time for each SEARCH operation is  $O(1)$ . By linearity of

expectation, therefore, the expected time for the entire sequence of  $n$  operations is  $O(n)$ . Since each operation takes  $\Omega(1)$  time, the  $\Theta(n)$  bound follows. ■

### Designing a universal class of hash functions

It is quite easy to design a universal class of hash functions, as a little number theory will help us prove. You may wish to consult Chapter 31 first if you are unfamiliar with number theory.

We begin by choosing a prime number  $p$  large enough so that every possible key  $k$  is in the range 0 to  $p - 1$ , inclusive. Let  $\mathbb{Z}_p$  denote the set  $\{0, 1, \dots, p - 1\}$ , and let  $\mathbb{Z}_p^*$  denote the set  $\{1, 2, \dots, p - 1\}$ . Since  $p$  is prime, we can solve equations modulo  $p$  with the methods given in Chapter 31. Because we assume that the size of the universe of keys is greater than the number of slots in the hash table, we have  $p > m$ .

We now define the hash function  $h_{ab}$  for any  $a \in \mathbb{Z}_p^*$  and any  $b \in \mathbb{Z}_p$  using a linear transformation followed by reductions modulo  $p$  and then modulo  $m$ :

$$h_{ab}(k) = ((ak + b) \bmod p) \bmod m. \quad (11.3)$$

For example, with  $p = 17$  and  $m = 6$ , we have  $h_{3,4}(8) = 5$ . The family of all such hash functions is

$$\mathcal{H}_{pm} = \{h_{ab} : a \in \mathbb{Z}_p^* \text{ and } b \in \mathbb{Z}_p\}. \quad (11.4)$$

Each hash function  $h_{ab}$  maps  $\mathbb{Z}_p$  to  $\mathbb{Z}_m$ . This class of hash functions has the nice property that the size  $m$  of the output range is arbitrary—not necessarily prime—a feature which we shall use in Section 11.5. Since we have  $p - 1$  choices for  $a$  and  $p$  choices for  $b$ , the collection  $\mathcal{H}_{pm}$  contains  $p(p - 1)$  hash functions.

#### Theorem 11.5

The class  $\mathcal{H}_{pm}$  of hash functions defined by equations (11.3) and (11.4) is universal.

**Proof** Consider two distinct keys  $k$  and  $l$  from  $\mathbb{Z}_p$ , so that  $k \neq l$ . For a given hash function  $h_{ab}$  we let

$$\begin{aligned} r &= (ak + b) \bmod p, \\ s &= (al + b) \bmod p. \end{aligned}$$

We first note that  $r \neq s$ . Why? Observe that

$$r - s \equiv a(k - l) \pmod{p}.$$

It follows that  $r \neq s$  because  $p$  is prime and both  $a$  and  $(k - l)$  are nonzero modulo  $p$ , and so their product must also be nonzero modulo  $p$  by Theorem 31.6. Therefore, when computing any  $h_{ab} \in \mathcal{H}_{pm}$ , distinct inputs  $k$  and  $l$  map to distinct

values  $r$  and  $s$  modulo  $p$ ; there are no collisions yet at the “mod  $p$  level.” Moreover, each of the possible  $p(p-1)$  choices for the pair  $(a, b)$  with  $a \neq 0$  yields a *different* resulting pair  $(r, s)$  with  $r \neq s$ , since we can solve for  $a$  and  $b$  given  $r$  and  $s$ :

$$\begin{aligned} a &= ((r - s)((k - l)^{-1} \bmod p)) \bmod p, \\ b &= (r - ak) \bmod p, \end{aligned}$$

where  $((k - l)^{-1} \bmod p)$  denotes the unique multiplicative inverse, modulo  $p$ , of  $k - l$ . Since there are only  $p(p - 1)$  possible pairs  $(r, s)$  with  $r \neq s$ , there is a one-to-one correspondence between pairs  $(a, b)$  with  $a \neq 0$  and pairs  $(r, s)$  with  $r \neq s$ . Thus, for any given pair of inputs  $k$  and  $l$ , if we pick  $(a, b)$  uniformly at random from  $\mathbb{Z}_p^* \times \mathbb{Z}_p$ , the resulting pair  $(r, s)$  is equally likely to be any pair of distinct values modulo  $p$ .

Therefore, the probability that distinct keys  $k$  and  $l$  collide is equal to the probability that  $r \equiv s \pmod{m}$  when  $r$  and  $s$  are randomly chosen as distinct values modulo  $p$ . For a given value of  $r$ , of the  $p - 1$  possible remaining values for  $s$ , the number of values  $s$  such that  $s \neq r$  and  $s \equiv r \pmod{m}$  is at most

$$\begin{aligned} \lceil p/m \rceil - 1 &\leq ((p + m - 1)/m) - 1 \quad (\text{by inequality (3.6)}) \\ &= (p - 1)/m. \end{aligned}$$

The probability that  $s$  collides with  $r$  when reduced modulo  $m$  is at most  $((p - 1)/m)/(p - 1) = 1/m$ .

Therefore, for any pair of distinct values  $k, l \in \mathbb{Z}_p$ ,

$$\Pr \{h_{ab}(k) = h_{ab}(l)\} \leq 1/m,$$

so that  $\mathcal{H}_{pm}$  is indeed universal. ■

## Exercises

### 11.3-1

Suppose we wish to search a linked list of length  $n$ , where each element contains a key  $k$  along with a hash value  $h(k)$ . Each key is a long character string. How might we take advantage of the hash values when searching the list for an element with a given key?

### 11.3-2

Suppose that we hash a string of  $r$  characters into  $m$  slots by treating it as a radix-128 number and then using the division method. We can easily represent the number  $m$  as a 32-bit computer word, but the string of  $r$  characters, treated as a radix-128 number, takes many words. How can we apply the division method to compute the hash value of the character string without using more than a constant number of words of storage outside the string itself?

**11.3-3**

Consider a version of the division method in which  $h(k) = k \bmod m$ , where  $m = 2^p - 1$  and  $k$  is a character string interpreted in radix  $2^p$ . Show that if we can derive string  $x$  from string  $y$  by permuting its characters, then  $x$  and  $y$  hash to the same value. Give an example of an application in which this property would be undesirable in a hash function.

**11.3-4**

Consider a hash table of size  $m = 1000$  and a corresponding hash function  $h(k) = \lfloor m(kA \bmod 1) \rfloor$  for  $A = (\sqrt{5} - 1)/2$ . Compute the locations to which the keys 61, 62, 63, 64, and 65 are mapped.

**11.3-5 ★**

Define a family  $\mathcal{H}$  of hash functions from a finite set  $U$  to a finite set  $B$  to be  **$\epsilon$ -universal** if for all pairs of distinct elements  $k$  and  $l$  in  $U$ ,

$$\Pr\{h(k) = h(l)\} \leq \epsilon,$$

where the probability is over the choice of the hash function  $h$  drawn at random from the family  $\mathcal{H}$ . Show that an  $\epsilon$ -universal family of hash functions must have

$$\epsilon \geq \frac{1}{|B|} - \frac{1}{|U|}.$$

**11.3-6 ★**

Let  $U$  be the set of  $n$ -tuples of values drawn from  $\mathbb{Z}_p$ , and let  $B = \mathbb{Z}_p$ , where  $p$  is prime. Define the hash function  $h_b : U \rightarrow B$  for  $b \in \mathbb{Z}_p$  on an input  $n$ -tuple  $\langle a_0, a_1, \dots, a_{n-1} \rangle$  from  $U$  as

$$h_b(\langle a_0, a_1, \dots, a_{n-1} \rangle) = \left( \sum_{j=0}^{n-1} a_j b^j \right) \bmod p,$$

and let  $\mathcal{H} = \{h_b : b \in \mathbb{Z}_p\}$ . Argue that  $\mathcal{H}$  is  $((n-1)/p)$ -universal according to the definition of  $\epsilon$ -universal in Exercise 11.3-5. (*Hint*: See Exercise 31.4-4.)

## 11.4 Open addressing

In **open addressing**, all elements occupy the hash table itself. That is, each table entry contains either an element of the dynamic set or NIL. When searching for an element, we systematically examine table slots until either we find the desired element or we have ascertained that the element is not in the table. No lists and

no elements are stored outside the table, unlike in chaining. Thus, in open addressing, the hash table can “fill up” so that no further insertions can be made; one consequence is that the load factor  $\alpha$  can never exceed 1.

Of course, we could store the linked lists for chaining inside the hash table, in the otherwise unused hash-table slots (see Exercise 11.2-4), but the advantage of open addressing is that it avoids pointers altogether. Instead of following pointers, we *compute* the sequence of slots to be examined. The extra memory freed by not storing pointers provides the hash table with a larger number of slots for the same amount of memory, potentially yielding fewer collisions and faster retrieval.

To perform insertion using open addressing, we successively examine, or *probe*, the hash table until we find an empty slot in which to put the key. Instead of being fixed in the order  $0, 1, \dots, m-1$  (which requires  $\Theta(n)$  search time), the sequence of positions probed *depends upon the key being inserted*. To determine which slots to probe, we extend the hash function to include the probe number (starting from 0) as a second input. Thus, the hash function becomes

$$h : U \times \{0, 1, \dots, m-1\} \rightarrow \{0, 1, \dots, m-1\} .$$

With open addressing, we require that for every key  $k$ , the *probe sequence*

$$\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$$

be a permutation of  $\{0, 1, \dots, m-1\}$ , so that every hash-table position is eventually considered as a slot for a new key as the table fills up. In the following pseudocode, we assume that the elements in the hash table  $T$  are keys with no satellite information; the key  $k$  is identical to the element containing key  $k$ . Each slot contains either a key or NIL (if the slot is empty). The HASH-INSERT procedure takes as input a hash table  $T$  and a key  $k$ . It either returns the slot number where it stores key  $k$  or flags an error because the hash table is already full.

HASH-INSERT( $T, k$ )

```

1   $i = 0$ 
2  repeat
3       $j = h(k, i)$ 
4      if  $T[j] == \text{NIL}$ 
5           $T[j] = k$ 
6          return  $j$ 
7      else  $i = i + 1$ 
8  until  $i == m$ 
9  error “hash table overflow”
```

The algorithm for searching for key  $k$  probes the same sequence of slots that the insertion algorithm examined when key  $k$  was inserted. Therefore, the search can

terminate (unsuccessfully) when it finds an empty slot, since  $k$  would have been inserted there and not later in its probe sequence. (This argument assumes that keys are not deleted from the hash table.) The procedure **HASH-SEARCH** takes as input a hash table  $T$  and a key  $k$ , returning  $j$  if it finds that slot  $j$  contains key  $k$ , or **NIL** if key  $k$  is not present in table  $T$ .

**HASH-SEARCH**( $T, k$ )

```

1   $i = 0$ 
2  repeat
3       $j = h(k, i)$ 
4      if  $T[j] == k$ 
5          return  $j$ 
6       $i = i + 1$ 
7  until  $T[j] == \text{NIL}$  or  $i == m$ 
8  return NIL
```

Deletion from an open-address hash table is difficult. When we delete a key from slot  $i$ , we cannot simply mark that slot as empty by storing **NIL** in it. If we did, we might be unable to retrieve any key  $k$  during whose insertion we had probed slot  $i$  and found it occupied. We can solve this problem by marking the slot, storing in it the special value **DELETED** instead of **NIL**. We would then modify the procedure **HASH-INSERT** to treat such a slot as if it were empty so that we can insert a new key there. We do not need to modify **HASH-SEARCH**, since it will pass over **DELETED** values while searching. When we use the special value **DELETED**, however, search times no longer depend on the load factor  $\alpha$ , and for this reason chaining is more commonly selected as a collision resolution technique when keys must be deleted.

In our analysis, we assume *uniform hashing*: the probe sequence of each key is equally likely to be any of the  $m!$  permutations of  $\{0, 1, \dots, m - 1\}$ . Uniform hashing generalizes the notion of simple uniform hashing defined earlier to a hash function that produces not just a single number, but a whole probe sequence. True uniform hashing is difficult to implement, however, and in practice suitable approximations (such as double hashing, defined below) are used.

We will examine three commonly used techniques to compute the probe sequences required for open addressing: linear probing, quadratic probing, and double hashing. These techniques all guarantee that  $\langle h(k, 0), h(k, 1), \dots, h(k, m - 1) \rangle$  is a permutation of  $\{0, 1, \dots, m - 1\}$  for each key  $k$ . None of these techniques fulfills the assumption of uniform hashing, however, since none of them is capable of generating more than  $m^2$  different probe sequences (instead of the  $m!$  that uniform hashing requires). Double hashing has the greatest number of probe sequences and, as one might expect, seems to give the best results.

### Linear probing

Given an ordinary hash function  $h' : U \rightarrow \{0, 1, \dots, m-1\}$ , which we refer to as an **auxiliary hash function**, the method of **linear probing** uses the hash function

$$h(k, i) = (h'(k) + i) \bmod m$$

for  $i = 0, 1, \dots, m-1$ . Given key  $k$ , we first probe  $T[h'(k)]$ , i.e., the slot given by the auxiliary hash function. We next probe slot  $T[h'(k) + 1]$ , and so on up to slot  $T[m-1]$ . Then we wrap around to slots  $T[0], T[1], \dots$  until we finally probe slot  $T[h'(k) - 1]$ . Because the initial probe determines the entire probe sequence, there are only  $m$  distinct probe sequences.

Linear probing is easy to implement, but it suffers from a problem known as **primary clustering**. Long runs of occupied slots build up, increasing the average search time. Clusters arise because an empty slot preceded by  $i$  full slots gets filled next with probability  $(i+1)/m$ . Long runs of occupied slots tend to get longer, and the average search time increases.

### Quadratic probing

**Quadratic probing** uses a hash function of the form

$$h(k, i) = (h'(k) + c_1i + c_2i^2) \bmod m, \quad (11.5)$$

where  $h'$  is an auxiliary hash function,  $c_1$  and  $c_2$  are positive auxiliary constants, and  $i = 0, 1, \dots, m-1$ . The initial position probed is  $T[h'(k)]$ ; later positions probed are offset by amounts that depend in a quadratic manner on the probe number  $i$ . This method works much better than linear probing, but to make full use of the hash table, the values of  $c_1$ ,  $c_2$ , and  $m$  are constrained. Problem 11-3 shows one way to select these parameters. Also, if two keys have the same initial probe position, then their probe sequences are the same, since  $h(k_1, 0) = h(k_2, 0)$  implies  $h(k_1, i) = h(k_2, i)$ . This property leads to a milder form of clustering, called **secondary clustering**. As in linear probing, the initial probe determines the entire sequence, and so only  $m$  distinct probe sequences are used.

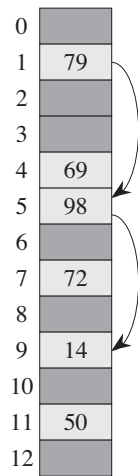
### Double hashing

Double hashing offers one of the best methods available for open addressing because the permutations produced have many of the characteristics of randomly chosen permutations. **Double hashing** uses a hash function of the form

$$h(k, i) = (h_1(k) + ih_2(k)) \bmod m,$$

where both  $h_1$  and  $h_2$  are auxiliary hash functions. The initial probe goes to position  $T[h_1(k)]$ ; successive probe positions are offset from previous positions by the





**Figure 11.5** Insertion by double hashing. Here we have a hash table of size 13 with  $h_1(k) = k \bmod 13$  and  $h_2(k) = 1 + (k \bmod 11)$ . Since  $14 \equiv 1 \pmod{13}$  and  $14 \equiv 3 \pmod{11}$ , we insert the key 14 into empty slot 9, after examining slots 1 and 5 and finding them to be occupied.

amount  $h_2(k)$ , modulo  $m$ . Thus, unlike the case of linear or quadratic probing, the probe sequence here depends in two ways upon the key  $k$ , since the initial probe position, the offset, or both, may vary. Figure 11.5 gives an example of insertion by double hashing.

The value  $h_2(k)$  must be relatively prime to the hash-table size  $m$  for the entire hash table to be searched. (See Exercise 11.4-4.) A convenient way to ensure this condition is to let  $m$  be a power of 2 and to design  $h_2$  so that it always produces an odd number. Another way is to let  $m$  be prime and to design  $h_2$  so that it always returns a positive integer less than  $m$ . For example, we could choose  $m$  prime and let

$$\begin{aligned} h_1(k) &= k \bmod m, \\ h_2(k) &= 1 + (k \bmod m'), \end{aligned}$$

where  $m'$  is chosen to be slightly less than  $m$  (say,  $m - 1$ ). For example, if  $k = 123456$ ,  $m = 701$ , and  $m' = 700$ , we have  $h_1(k) = 80$  and  $h_2(k) = 257$ , so that we first probe position 80, and then we examine every 257th slot (modulo  $m$ ) until we find the key or have examined every slot.

When  $m$  is prime or a power of 2, double hashing improves over linear or quadratic probing in that  $\Theta(m^2)$  probe sequences are used, rather than  $\Theta(m)$ , since each possible  $(h_1(k), h_2(k))$  pair yields a distinct probe sequence. As a result, for

such values of  $m$ , the performance of double hashing appears to be very close to the performance of the “ideal” scheme of uniform hashing.

Although values of  $m$  other than primes or powers of 2 could in principle be used with double hashing, in practice it becomes more difficult to efficiently generate  $h_2(k)$  in a way that ensures that it is relatively prime to  $m$ , in part because the relative density  $\phi(m)/m$  of such numbers may be small (see equation (31.24)).

### Analysis of open-address hashing

As in our analysis of chaining, we express our analysis of open addressing in terms of the load factor  $\alpha = n/m$  of the hash table. Of course, with open addressing, at most one element occupies each slot, and thus  $n \leq m$ , which implies  $\alpha \leq 1$ .

We assume that we are using uniform hashing. In this idealized scheme, the probe sequence  $\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$  used to insert or search for each key  $k$  is equally likely to be any permutation of  $\{0, 1, \dots, m-1\}$ . Of course, a given key has a unique fixed probe sequence associated with it; what we mean here is that, considering the probability distribution on the space of keys and the operation of the hash function on the keys, each possible probe sequence is equally likely.

We now analyze the expected number of probes for hashing with open addressing under the assumption of uniform hashing, beginning with an analysis of the number of probes made in an unsuccessful search.

#### **Theorem 11.6**

Given an open-address hash table with load factor  $\alpha = n/m < 1$ , the expected number of probes in an unsuccessful search is at most  $1/(1-\alpha)$ , assuming uniform hashing.

**Proof** In an unsuccessful search, every probe but the last accesses an occupied slot that does not contain the desired key, and the last slot probed is empty. Let us define the random variable  $X$  to be the number of probes made in an unsuccessful search, and let us also define the event  $A_i$ , for  $i = 1, 2, \dots$ , to be the event that an  $i$ th probe occurs and it is to an occupied slot. Then the event  $\{X \geq i\}$  is the intersection of events  $A_1 \cap A_2 \cap \dots \cap A_{i-1}$ . We will bound  $\Pr\{X \geq i\}$  by bounding  $\Pr\{A_1 \cap A_2 \cap \dots \cap A_{i-1}\}$ . By Exercise C.2-5,

$$\Pr\{A_1 \cap A_2 \cap \dots \cap A_{i-1}\} = \Pr\{A_1\} \cdot \Pr\{A_2 \mid A_1\} \cdot \Pr\{A_3 \mid A_1 \cap A_2\} \cdots \Pr\{A_{i-1} \mid A_1 \cap A_2 \cap \dots \cap A_{i-2}\}.$$

Since there are  $n$  elements and  $m$  slots,  $\Pr\{A_1\} = n/m$ . For  $j > 1$ , the probability that there is a  $j$ th probe and it is to an occupied slot, given that the first  $j-1$  probes were to occupied slots, is  $(n-j+1)/(m-j+1)$ . This probability follows

because we would be finding one of the remaining  $(n - (j - 1))$  elements in one of the  $(m - (j - 1))$  unexamined slots, and by the assumption of uniform hashing, the probability is the ratio of these quantities. Observing that  $n < m$  implies that  $(n - j)/(m - j) \leq n/m$  for all  $j$  such that  $0 \leq j < m$ , we have for all  $i$  such that  $1 \leq i \leq m$ ,

$$\begin{aligned} \Pr\{X \geq i\} &= \frac{n}{m} \cdot \frac{n-1}{m-1} \cdot \frac{n-2}{m-2} \cdots \frac{n-i+2}{m-i+2} \\ &\leq \left(\frac{n}{m}\right)^{i-1} \\ &= \alpha^{i-1}. \end{aligned}$$

Now, we use equation (C.25) to bound the expected number of probes:

$$\begin{aligned} E[X] &= \sum_{i=1}^{\infty} \Pr\{X \geq i\} \\ &\leq \sum_{i=1}^{\infty} \alpha^{i-1} \\ &= \sum_{i=0}^{\infty} \alpha^i \\ &= \frac{1}{1-\alpha}. \end{aligned}$$

■

This bound of  $1/(1-\alpha) = 1 + \alpha + \alpha^2 + \alpha^3 + \cdots$  has an intuitive interpretation. We always make the first probe. With probability approximately  $\alpha$ , the first probe finds an occupied slot, so that we need to probe a second time. With probability approximately  $\alpha^2$ , the first two slots are occupied so that we make a third probe, and so on.

If  $\alpha$  is a constant, Theorem 11.6 predicts that an unsuccessful search runs in  $O(1)$  time. For example, if the hash table is half full, the average number of probes in an unsuccessful search is at most  $1/(1 - .5) = 2$ . If it is 90 percent full, the average number of probes is at most  $1/(1 - .9) = 10$ .

Theorem 11.6 gives us the performance of the HASH-INSERT procedure almost immediately.

### **Corollary 11.7**

Inserting an element into an open-address hash table with load factor  $\alpha$  requires at most  $1/(1 - \alpha)$  probes on average, assuming uniform hashing.

**Proof** An element is inserted only if there is room in the table, and thus  $\alpha < 1$ . Inserting a key requires an unsuccessful search followed by placing the key into the first empty slot found. Thus, the expected number of probes is at most  $1/(1-\alpha)$ . ■

We have to do a little more work to compute the expected number of probes for a successful search.

**Theorem 11.8**

Given an open-address hash table with load factor  $\alpha < 1$ , the expected number of probes in a successful search is at most

$$\frac{1}{\alpha} \ln \frac{1}{1-\alpha},$$

assuming uniform hashing and assuming that each key in the table is equally likely to be searched for.

**Proof** A search for a key  $k$  reproduces the same probe sequence as when the element with key  $k$  was inserted. By Corollary 11.7, if  $k$  was the  $(i+1)$ st key inserted into the hash table, the expected number of probes made in a search for  $k$  is at most  $1/(1-i/m) = m/(m-i)$ . Averaging over all  $n$  keys in the hash table gives us the expected number of probes in a successful search:

$$\begin{aligned} \frac{1}{n} \sum_{i=0}^{n-1} \frac{m}{m-i} &= \frac{m}{n} \sum_{i=0}^{n-1} \frac{1}{m-i} \\ &= \frac{1}{\alpha} \sum_{k=m-n+1}^m \frac{1}{k} \\ &\leq \frac{1}{\alpha} \int_{m-n}^m (1/x) dx \quad (\text{by inequality (A.12)}) \\ &= \frac{1}{\alpha} \ln \frac{m}{m-n} \\ &= \frac{1}{\alpha} \ln \frac{1}{1-\alpha}. \quad \blacksquare \end{aligned}$$

If the hash table is half full, the expected number of probes in a successful search is less than 1.387. If the hash table is 90 percent full, the expected number of probes is less than 2.559.

## Exercises

### 11.4-1

Consider inserting the keys 10, 22, 31, 4, 15, 28, 17, 88, 59 into a hash table of length  $m = 11$  using open addressing with the auxiliary hash function  $h'(k) = k$ . Illustrate the result of inserting these keys using linear probing, using quadratic probing with  $c_1 = 1$  and  $c_2 = 3$ , and using double hashing with  $h_1(k) = k$  and  $h_2(k) = 1 + (k \bmod (m - 1))$ .

### 11.4-2

Write pseudocode for HASH-DELETE as outlined in the text, and modify HASH-INSERT to handle the special value DELETED.

### 11.4-3

Consider an open-address hash table with uniform hashing. Give upper bounds on the expected number of probes in an unsuccessful search and on the expected number of probes in a successful search when the load factor is  $3/4$  and when it is  $7/8$ .

### 11.4-4 ★

Suppose that we use double hashing to resolve collisions—that is, we use the hash function  $h(k, i) = (h_1(k) + ih_2(k)) \bmod m$ . Show that if  $m$  and  $h_2(k)$  have greatest common divisor  $d \geq 1$  for some key  $k$ , then an unsuccessful search for key  $k$  examines  $(1/d)$ th of the hash table before returning to slot  $h_1(k)$ . Thus, when  $d = 1$ , so that  $m$  and  $h_2(k)$  are relatively prime, the search may examine the entire hash table. (*Hint:* See Chapter 31.)

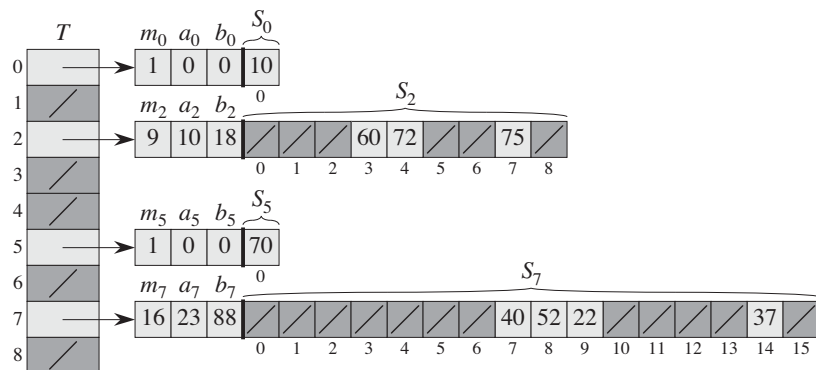
### 11.4-5 ★

Consider an open-address hash table with a load factor  $\alpha$ . Find the nonzero value  $\alpha$  for which the expected number of probes in an unsuccessful search equals twice the expected number of probes in a successful search. Use the upper bounds given by Theorems 11.6 and 11.8 for these expected numbers of probes.

---

## ★ 11.5 Perfect hashing

Although hashing is often a good choice for its excellent average-case performance, hashing can also provide excellent *worst-case* performance when the set of keys is *static*: once the keys are stored in the table, the set of keys never changes. Some applications naturally have static sets of keys: consider the set of reserved words in a programming language, or the set of file names on a CD-ROM. We



**Figure 11.6** Using perfect hashing to store the set  $K = \{10, 22, 37, 40, 52, 60, 70, 72, 75\}$ . The outer hash function is  $h(k) = ((ak + b) \bmod p) \bmod m$ , where  $a = 3$ ,  $b = 42$ ,  $p = 101$ , and  $m = 9$ . For example,  $h(75) = 2$ , and so key 75 hashes to slot 2 of table  $T$ . A secondary hash table  $S_j$  stores all keys hashing to slot  $j$ . The size of hash table  $S_j$  is  $m_j = n_j^2$ , and the associated hash function is  $h_j(k) = ((a_jk + b_j) \bmod p) \bmod m_j$ . Since  $h_2(75) = 7$ , key 75 is stored in slot 7 of secondary hash table  $S_2$ . No collisions occur in any of the secondary hash tables, and so searching takes constant time in the worst case.

call a hashing technique **perfect hashing** if  $O(1)$  memory accesses are required to perform a search in the worst case.

To create a perfect hashing scheme, we use two levels of hashing, with universal hashing at each level. Figure 11.6 illustrates the approach.

The first level is essentially the same as for hashing with chaining: we hash the  $n$  keys into  $m$  slots using a hash function  $h$  carefully selected from a family of universal hash functions.

Instead of making a linked list of the keys hashing to slot  $j$ , however, we use a small **secondary hash table**  $S_j$  with an associated hash function  $h_j$ . By choosing the hash functions  $h_j$  carefully, we can guarantee that there are no collisions at the secondary level.

In order to guarantee that there are no collisions at the secondary level, however, we will need to let the size  $m_j$  of hash table  $S_j$  be the square of the number  $n_j$  of keys hashing to slot  $j$ . Although you might think that the quadratic dependence of  $m_j$  on  $n_j$  may seem likely to cause the overall storage requirement to be excessive, we shall show that by choosing the first-level hash function well, we can limit the expected total amount of space used to  $O(n)$ .

We use hash functions chosen from the universal classes of hash functions of Section 11.3.3. The first-level hash function comes from the class  $\mathcal{H}_{pm}$ , where as in Section 11.3.3,  $p$  is a prime number greater than any key value. Those keys

hashing to slot  $j$  are re-hashed into a secondary hash table  $S_j$  of size  $m_j$  using a hash function  $h_j$  chosen from the class  $\mathcal{H}_{p,m_j}$ .<sup>1</sup>

We shall proceed in two steps. First, we shall determine how to ensure that the secondary tables have no collisions. Second, we shall show that the expected amount of memory used overall—for the primary hash table and all the secondary hash tables—is  $O(n)$ .

**Theorem 11.9**

Suppose that we store  $n$  keys in a hash table of size  $m = n^2$  using a hash function  $h$  randomly chosen from a universal class of hash functions. Then, the probability is less than  $1/2$  that there are any collisions.

**Proof** There are  $\binom{n}{2}$  pairs of keys that may collide; each pair collides with probability  $1/m$  if  $h$  is chosen at random from a universal family  $\mathcal{H}$  of hash functions. Let  $X$  be a random variable that counts the number of collisions. When  $m = n^2$ , the expected number of collisions is

$$\begin{aligned} E[X] &= \binom{n}{2} \cdot \frac{1}{n^2} \\ &= \frac{n^2 - n}{2} \cdot \frac{1}{n^2} \\ &< 1/2. \end{aligned}$$

(This analysis is similar to the analysis of the birthday paradox in Section 5.4.1.) Applying Markov's inequality (C.30),  $\Pr\{X \geq t\} \leq E[X]/t$ , with  $t = 1$ , completes the proof. ■

In the situation described in Theorem 11.9, where  $m = n^2$ , it follows that a hash function  $h$  chosen at random from  $\mathcal{H}$  is more likely than not to have *no* collisions. Given the set  $K$  of  $n$  keys to be hashed (remember that  $K$  is static), it is thus easy to find a collision-free hash function  $h$  with a few random trials.

When  $n$  is large, however, a hash table of size  $m = n^2$  is excessive. Therefore, we adopt the two-level hashing approach, and we use the approach of Theorem 11.9 only to hash the entries within each slot. We use an outer, or first-level, hash function  $h$  to hash the keys into  $m = n$  slots. Then, if  $n_j$  keys hash to slot  $j$ , we use a secondary hash table  $S_j$  of size  $m_j = n_j^2$  to provide collision-free constant-time lookup.

---

<sup>1</sup>When  $n_j = m_j = 1$ , we don't really need a hash function for slot  $j$ ; when we choose a hash function  $h_{ab}(k) = ((ak + b) \bmod p) \bmod m_j$  for such a slot, we just use  $a = b = 0$ .

We now turn to the issue of ensuring that the overall memory used is  $O(n)$ . Since the size  $m_j$  of the  $j$ th secondary hash table grows quadratically with the number  $n_j$  of keys stored, we run the risk that the overall amount of storage could be excessive.

If the first-level table size is  $m = n$ , then the amount of memory used is  $O(n)$  for the primary hash table, for the storage of the sizes  $m_j$  of the secondary hash tables, and for the storage of the parameters  $a_j$  and  $b_j$  defining the secondary hash functions  $h_j$  drawn from the class  $\mathcal{H}_{p,m_j}$  of Section 11.3.3 (except when  $n_j = 1$  and we use  $a = b = 0$ ). The following theorem and a corollary provide a bound on the expected combined sizes of all the secondary hash tables. A second corollary bounds the probability that the combined size of all the secondary hash tables is superlinear (actually, that it equals or exceeds  $4n$ ).

**Theorem 11.10**

Suppose that we store  $n$  keys in a hash table of size  $m = n$  using a hash function  $h$  randomly chosen from a universal class of hash functions. Then, we have

$$\mathbb{E} \left[ \sum_{j=0}^{m-1} n_j^2 \right] < 2n ,$$

where  $n_j$  is the number of keys hashing to slot  $j$ .

**Proof** We start with the following identity, which holds for any nonnegative integer  $a$ :

$$a^2 = a + 2 \binom{a}{2} . \tag{11.6}$$

We have

$$\begin{aligned} \mathbb{E} \left[ \sum_{j=0}^{m-1} n_j^2 \right] &= \mathbb{E} \left[ \sum_{j=0}^{m-1} \left( n_j + 2 \binom{n_j}{2} \right) \right] && \text{(by equation (11.6))} \\ &= \mathbb{E} \left[ \sum_{j=0}^{m-1} n_j \right] + 2 \mathbb{E} \left[ \sum_{j=0}^{m-1} \binom{n_j}{2} \right] && \text{(by linearity of expectation)} \\ &= \mathbb{E} [n] + 2 \mathbb{E} \left[ \sum_{j=0}^{m-1} \binom{n_j}{2} \right] && \text{(by equation (11.1))} \end{aligned}$$



$$= n + 2 \mathbb{E} \left[ \sum_{j=0}^{m-1} \binom{n_j}{2} \right] \quad (\text{since } n \text{ is not a random variable}) .$$

To evaluate the summation  $\sum_{j=0}^{m-1} \binom{n_j}{2}$ , we observe that it is just the total number of pairs of keys in the hash table that collide. By the properties of universal hashing, the expected value of this summation is at most

$$\begin{aligned} \binom{n}{2} \frac{1}{m} &= \frac{n(n-1)}{2m} \\ &= \frac{n-1}{2} , \end{aligned}$$

since  $m = n$ . Thus,

$$\begin{aligned} \mathbb{E} \left[ \sum_{j=0}^{m-1} n_j^2 \right] &\leq n + 2 \frac{n-1}{2} \\ &= 2n - 1 \\ &< 2n . \end{aligned} \quad \blacksquare$$

### Corollary 11.11

Suppose that we store  $n$  keys in a hash table of size  $m = n$  using a hash function  $h$  randomly chosen from a universal class of hash functions, and we set the size of each secondary hash table to  $m_j = n_j^2$  for  $j = 0, 1, \dots, m-1$ . Then, the expected amount of storage required for all secondary hash tables in a perfect hashing scheme is less than  $2n$ .

**Proof** Since  $m_j = n_j^2$  for  $j = 0, 1, \dots, m-1$ , Theorem 11.10 gives

$$\begin{aligned} \mathbb{E} \left[ \sum_{j=0}^{m-1} m_j \right] &= \mathbb{E} \left[ \sum_{j=0}^{m-1} n_j^2 \right] \\ &< 2n , \end{aligned} \quad (11.7)$$

which completes the proof. ■

### Corollary 11.12

Suppose that we store  $n$  keys in a hash table of size  $m = n$  using a hash function  $h$  randomly chosen from a universal class of hash functions, and we set the size of each secondary hash table to  $m_j = n_j^2$  for  $j = 0, 1, \dots, m-1$ . Then, the probability is less than  $1/2$  that the total storage used for secondary hash tables equals or exceeds  $4n$ .

**Proof** Again we apply Markov's inequality (C.30),  $\Pr\{X \geq t\} \leq E[X]/t$ , this time to inequality (11.7), with  $X = \sum_{j=0}^{m-1} m_j$  and  $t = 4n$ :

$$\begin{aligned} \Pr\left\{\sum_{j=0}^{m-1} m_j \geq 4n\right\} &\leq \frac{E\left[\sum_{j=0}^{m-1} m_j\right]}{4n} \\ &< \frac{2n}{4n} \\ &= 1/2. \end{aligned}$$

■

From Corollary 11.12, we see that if we test a few randomly chosen hash functions from the universal family, we will quickly find one that uses a reasonable amount of storage.

## Exercises

### 11.5-1 ★

Suppose that we insert  $n$  keys into a hash table of size  $m$  using open addressing and uniform hashing. Let  $p(n, m)$  be the probability that no collisions occur. Show that  $p(n, m) \leq e^{-n(n-1)/2m}$ . (*Hint*: See equation (3.12).) Argue that when  $n$  exceeds  $\sqrt{m}$ , the probability of avoiding collisions goes rapidly to zero.

---

## Problems

### 11-1 Longest-probe bound for hashing

Suppose that we use an open-addressed hash table of size  $m$  to store  $n \leq m/2$  items.

- Assuming uniform hashing, show that for  $i = 1, 2, \dots, n$ , the probability is at most  $2^{-k}$  that the  $i$ th insertion requires strictly more than  $k$  probes.
- Show that for  $i = 1, 2, \dots, n$ , the probability is  $O(1/n^2)$  that the  $i$ th insertion requires more than  $2 \lg n$  probes.

Let the random variable  $X_i$  denote the number of probes required by the  $i$ th insertion. You have shown in part (b) that  $\Pr\{X_i > 2 \lg n\} = O(1/n^2)$ . Let the random variable  $X = \max_{1 \leq i \leq n} X_i$  denote the maximum number of probes required by any of the  $n$  insertions.

- Show that  $\Pr\{X > 2 \lg n\} = O(1/n)$ .
- Show that the expected length  $E[X]$  of the longest probe sequence is  $O(\lg n)$ .

**11-2 Slot-size bound for chaining**

Suppose that we have a hash table with  $n$  slots, with collisions resolved by chaining, and suppose that  $n$  keys are inserted into the table. Each key is equally likely to be hashed to each slot. Let  $M$  be the maximum number of keys in any slot after all the keys have been inserted. Your mission is to prove an  $O(\lg n / \lg \lg n)$  upper bound on  $E[M]$ , the expected value of  $M$ .

- a. Argue that the probability  $Q_k$  that exactly  $k$  keys hash to a particular slot is given by

$$Q_k = \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k} \binom{n}{k}.$$

- b. Let  $P_k$  be the probability that  $M = k$ , that is, the probability that the slot containing the most keys contains  $k$  keys. Show that  $P_k \leq n Q_k$ .
- c. Use Stirling's approximation, equation (3.18), to show that  $Q_k < e^k / k^k$ .
- d. Show that there exists a constant  $c > 1$  such that  $Q_{k_0} < 1/n^3$  for  $k_0 = c \lg n / \lg \lg n$ . Conclude that  $P_k < 1/n^2$  for  $k \geq k_0 = c \lg n / \lg \lg n$ .

- e. Argue that

$$E[M] \leq \Pr \left\{ M > \frac{c \lg n}{\lg \lg n} \right\} \cdot n + \Pr \left\{ M \leq \frac{c \lg n}{\lg \lg n} \right\} \cdot \frac{c \lg n}{\lg \lg n}.$$

Conclude that  $E[M] = O(\lg n / \lg \lg n)$ .

**11-3 Quadratic probing**

Suppose that we are given a key  $k$  to search for in a hash table with positions  $0, 1, \dots, m-1$ , and suppose that we have a hash function  $h$  mapping the key space into the set  $\{0, 1, \dots, m-1\}$ . The search scheme is as follows:

1. Compute the value  $j = h(k)$ , and set  $i = 0$ .
2. Probe in position  $j$  for the desired key  $k$ . If you find it, or if this position is empty, terminate the search.
3. Set  $i = i + 1$ . If  $i$  now equals  $m$ , the table is full, so terminate the search. Otherwise, set  $j = (i + j) \bmod m$ , and return to step 2.

Assume that  $m$  is a power of 2.

- a. Show that this scheme is an instance of the general “quadratic probing” scheme by exhibiting the appropriate constants  $c_1$  and  $c_2$  for equation (11.5).
- b. Prove that this algorithm examines every table position in the worst case.

### 11-4 Hashing and authentication

Let  $\mathcal{H}$  be a class of hash functions in which each hash function  $h \in \mathcal{H}$  maps the universe  $U$  of keys to  $\{0, 1, \dots, m-1\}$ . We say that  $\mathcal{H}$  is ***k-universal*** if, for every fixed sequence of  $k$  distinct keys  $\langle x^{(1)}, x^{(2)}, \dots, x^{(k)} \rangle$  and for any  $h$  chosen at random from  $\mathcal{H}$ , the sequence  $\langle h(x^{(1)}), h(x^{(2)}), \dots, h(x^{(k)}) \rangle$  is equally likely to be any of the  $m^k$  sequences of length  $k$  with elements drawn from  $\{0, 1, \dots, m-1\}$ .

- a. Show that if the family  $\mathcal{H}$  of hash functions is 2-universal, then it is universal.
- b. Suppose that the universe  $U$  is the set of  $n$ -tuples of values drawn from  $\mathbb{Z}_p = \{0, 1, \dots, p-1\}$ , where  $p$  is prime. Consider an element  $x = \langle x_0, x_1, \dots, x_{n-1} \rangle \in U$ . For any  $n$ -tuple  $a = \langle a_0, a_1, \dots, a_{n-1} \rangle \in U$ , define the hash function  $h_a$  by

$$h_a(x) = \left( \sum_{j=0}^{n-1} a_j x_j \right) \bmod p.$$

Let  $\mathcal{H} = \{h_a\}$ . Show that  $\mathcal{H}$  is universal, but not 2-universal. (*Hint*: Find a key for which all hash functions in  $\mathcal{H}$  produce the same value.)

- c. Suppose that we modify  $\mathcal{H}$  slightly from part (b): for any  $a \in U$  and for any  $b \in \mathbb{Z}_p$ , define

$$h'_{ab}(x) = \left( \sum_{j=0}^{n-1} a_j x_j + b \right) \bmod p$$

and  $\mathcal{H}' = \{h'_{ab}\}$ . Argue that  $\mathcal{H}'$  is 2-universal. (*Hint*: Consider fixed  $n$ -tuples  $x \in U$  and  $y \in U$ , with  $x_i \neq y_i$  for some  $i$ . What happens to  $h'_{ab}(x)$  and  $h'_{ab}(y)$  as  $a_i$  and  $b$  range over  $\mathbb{Z}_p$ ?)

- d. Suppose that Alice and Bob secretly agree on a hash function  $h$  from a 2-universal family  $\mathcal{H}$  of hash functions. Each  $h \in \mathcal{H}$  maps from a universe of keys  $U$  to  $\mathbb{Z}_p$ , where  $p$  is prime. Later, Alice sends a message  $m$  to Bob over the Internet, where  $m \in U$ . She authenticates this message to Bob by also sending an authentication tag  $t = h(m)$ , and Bob checks that the pair  $(m, t)$  he receives indeed satisfies  $t = h(m)$ . Suppose that an adversary intercepts  $(m, t)$  en route and tries to fool Bob by replacing the pair  $(m, t)$  with a different pair  $(m', t')$ . Argue that the probability that the adversary succeeds in fooling Bob into accepting  $(m', t')$  is at most  $1/p$ , no matter how much computing power the adversary has, and even if the adversary knows the family  $\mathcal{H}$  of hash functions used.

---

**Chapter notes**

Knuth [211] and Gonnet [145] are excellent references for the analysis of hashing algorithms. Knuth credits H. P. Luhn (1953) for inventing hash tables, along with the chaining method for resolving collisions. At about the same time, G. M. Amdahl originated the idea of open addressing.

Carter and Wegman introduced the notion of universal classes of hash functions in 1979 [58].

Fredman, Komlós, and Szemerédi [112] developed the perfect hashing scheme for static sets presented in Section 11.5. An extension of their method to dynamic sets, handling insertions and deletions in amortized expected time  $O(1)$ , has been given by Dietzfelbinger et al. [86].

---

## 12 Binary Search Trees

The search tree data structure supports many dynamic-set operations, including SEARCH, MINIMUM, MAXIMUM, PREDECESSOR, SUCCESSOR, INSERT, and DELETE. Thus, we can use a search tree both as a dictionary and as a priority queue.

Basic operations on a binary search tree take time proportional to the height of the tree. For a complete binary tree with  $n$  nodes, such operations run in  $\Theta(\lg n)$  worst-case time. If the tree is a linear chain of  $n$  nodes, however, the same operations take  $\Theta(n)$  worst-case time. We shall see in Section 12.4 that the expected height of a randomly built binary search tree is  $O(\lg n)$ , so that basic dynamic-set operations on such a tree take  $\Theta(\lg n)$  time on average.

In practice, we can't always guarantee that binary search trees are built randomly, but we can design variations of binary search trees with good guaranteed worst-case performance on basic operations. Chapter 13 presents one such variation, red-black trees, which have height  $O(\lg n)$ . Chapter 18 introduces B-trees, which are particularly good for maintaining databases on secondary (disk) storage.

After presenting the basic properties of binary search trees, the following sections show how to walk a binary search tree to print its values in sorted order, how to search for a value in a binary search tree, how to find the minimum or maximum element, how to find the predecessor or successor of an element, and how to insert into or delete from a binary search tree. The basic mathematical properties of trees appear in Appendix B.

---

### 12.1 What is a binary search tree?

A binary search tree is organized, as the name suggests, in a binary tree, as shown in Figure 12.1. We can represent such a tree by a linked data structure in which each node is an object. In addition to a *key* and satellite data, each node contains attributes *left*, *right*, and *p* that point to the nodes corresponding to its left child,