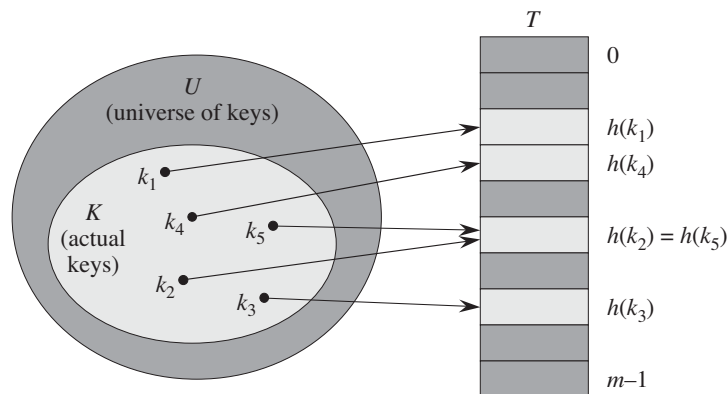## 11.2    Hash tables

The downside of direct addressing is obvious: if the universe $U$ is large, storing a table $T$ of size $|U|$ may be impractical, or even impossible, given the memory available on a typical computer. Furthermore, the set $K$ of keys *actually stored* may be so small relative to $U$ that most of the space allocated for $T$ would be wasted.

When the set $K$ of keys stored in a dictionary is much smaller than the universe $U$ of all possible keys, a hash table requires much less storage than a direct-address table. Specifically, we can reduce the storage requirement to $\Theta(|K|)$ while we maintain the benefit that searching for an element in the hash table still requires only $O(1)$ time. The catch is that this bound is for the *average-case time*, whereas for direct addressing it holds for the *worst-case time*.
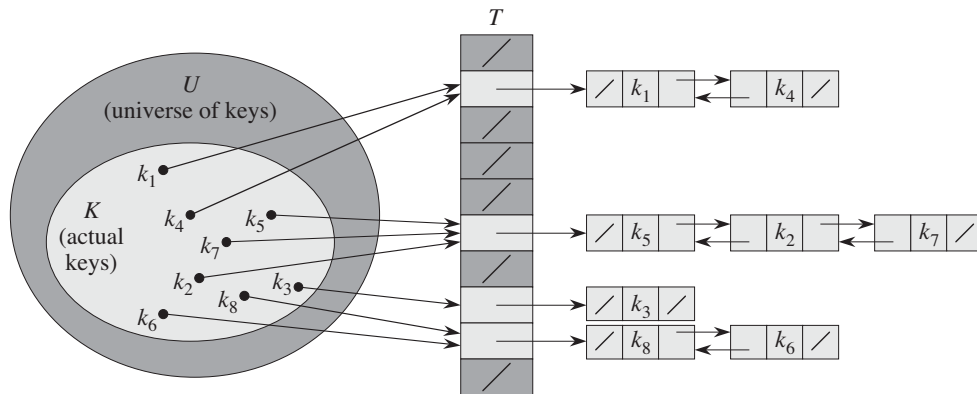
With direct addressing, an element with key $k$ is stored in slot $k$. With hashing, this element is stored in slot $h(k)$; that is, we use a **hash function** $h$ to compute the slot from the key $k$. Here, $h$ maps the universe $U$ of keys into the slots of a **hash table** $T[0 \ldots m-1]$:

$$h : U \rightarrow \{0, 1, \ldots, m-1\} \ ,$$

where the size $m$ of the hash table is typically much less than $|U|$. We say that an element with key $k$ **hashes** to slot $h(k)$; we also say that $h(k)$ is the **hash value** of key $k$. Figure 11.2 illustrates the basic idea. The hash function reduces the range of array indices and hence the size of the array. Instead of a size of $|U|$, the array can have size $m$.



**Figure 11.2**    Using a hash function $h$ to map keys to hash-table slots. Because keys $k_2$ and $k_5$ map to the same slot, they collide.

**Figure 11.3**   Collision resolution by chaining. Each hash-table slot $T[j]$ contains a linked list of all the keys whose hash value is $j$. For example, $h(k_1) = h(k_4)$ and $h(k_5) = h(k_7) = h(k_2)$. The linked list can be either singly or doubly linked; we show it as doubly linked because deletion is faster that way.

There is one hitch: two keys may hash to the same slot. We call this situation a ***collision***. Fortunately, we have effective techniques for resolving the conflict created by collisions.

Of course, the ideal solution would be to avoid collisions altogether. We might try to achieve this goal by choosing a suitable hash function $h$. One idea is to make $h$ appear to be "random," thus avoiding collisions or at least minimizing their number. The very term "to hash," evoking images of random mixing and chopping, captures the spirit of this approach. (Of course, a hash function $h$ must be deterministic in that a given input $k$ should always produce the same output $h(k)$.) Because $|U| > m$, however, there must be at least two keys that have the same hash value; avoiding collisions altogether is therefore impossible. Thus, while a well-designed, "random"-looking hash function can minimize the number of collisions, we still need a method for resolving the collisions that do occur.

The remainder of this section presents the simplest collision resolution technique, called chaining. Section 11.4 introduces an alternative method for resolving collisions, called open addressing.

**Collision resolution by chaining**

In ***chaining***, we place all the elements that hash to the same slot into the same linked list, as Figure 11.3 shows. Slot $j$ contains a pointer to the head of the list of all stored elements that hash to $j$; if there are no such elements, slot $j$ contains NIL.

The dictionary operations on a hash table $T$ are easy to implement when collisions are resolved by chaining:

CHAINED-HASH-INSERT$(T, x)$

1    insert $x$ at the head of list $T[h(x.key)]$

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

1    search for an element with key $k$ in list $T[h(k)]$

CHAINED-HASH-DELETE$(T, x)$

1    delete $x$ from the list $T[h(x.key)]$

The worst-case running time for insertion is $O(1)$. The insertion procedure is fast in part because it assumes that the element $x$ being inserted is not already present in the table; if necessary, we can check this assumption (at additional cost) by searching for an element whose key is $x.key$ before we insert. For searching, the worst-case running time is proportional to the length of the list; we shall analyze this operation more closely below. We can delete an element in $O(1)$ time if the lists are doubly linked, as Figure 11.3 depicts. (Note that CHAINED-HASH-DELETE takes as input an element $x$ and not its key $k$, so that we don't have to search for $x$ first. If the hash table supports deletion, then its linked lists should be doubly linked so that we can delete an item quickly. If the lists were only singly linked, then to delete element $x$, we would first have to find $x$ in the list $T[h(x.key)]$ so that we could update the *next* attribute of $x$'s predecessor. With singly linked lists, both deletion and searching would have the same asymptotic running times.)

### Analysis of hashing with chaining

How well does hashing with chaining perform? In particular, how long does it take to search for an element with a given key?

Given a hash table $T$ with $m$ slots that stores $n$ elements, we define the ***load factor*** $\alpha$ for $T$ as $n/m$, that is, the average number of elements stored in a chain. Our analysis will be in terms of $\alpha$, which can be less than, equal to, or greater than 1.

The worst-case behavior of hashing with chaining is terrible: all $n$ keys hash to the same slot, creating a list of length $n$. The worst-case time for searching is thus $\Theta(n)$ plus the time to compute the hash function—no better than if we used one linked list for all the elements. Clearly, we do not use hash tables for their worst-case performance. (Perfect hashing, described in Section 11.5, does provide good worst-case performance when the set of keys is static, however.)

The average-case performance of hashing depends on how well the hash function $h$ distributes the set of keys to be stored among the $m$ slots, on the average.

Section 11.3 discusses these issues, but for now we shall assume that any given element is equally likely to hash into any of the $m$ slots, independently of where any other element has hashed to. We call this the assumption of *simple uniform hashing*.

For $j = 0, 1, \ldots, m - 1$, let us denote the length of the list $T[j]$ by $n_j$, so that

$$n = n_0 + n_1 + \cdots + n_{m-1} \, , \tag{11.1}$$

and the expected value of $n_j$ is $\mathrm{E}[n_j] = \alpha = n/m$.

We assume that $O(1)$ time suffices to compute the hash value $h(k)$, so that the time required to search for an element with key $k$ depends linearly on the length $n_{h(k)}$ of the list $T[h(k)]$. Setting aside the $O(1)$ time required to compute the hash function and to access slot $h(k)$, let us consider the expected number of elements examined by the search algorithm, that is, the number of elements in the list $T[h(k)]$ that the algorithm checks to see whether any have a key equal to $k$. We shall consider two cases. In the first, the search is unsuccessful: no element in the table has key $k$. In the second, the search successfully finds an element with key $k$.

### Theorem 11.1
In a hash table in which collisions are resolved by chaining, an unsuccessful search takes average-case time $\Theta(1+\alpha)$, under the assumption of simple uniform hashing.

*Proof*   Under the assumption of simple uniform hashing, any key $k$ not already stored in the table is equally likely to hash to any of the $m$ slots. The expected time to search unsuccessfully for a key $k$ is the expected time to search to the end of list $T[h(k)]$, which has expected length $\mathrm{E}[n_{h(k)}] = \alpha$. Thus, the expected number of elements examined in an unsuccessful search is $\alpha$, and the total time required (including the time for computing $h(k)$) is $\Theta(1 + \alpha)$.                                  ∎

The situation for a successful search is slightly different, since each list is not equally likely to be searched. Instead, the probability that a list is searched is proportional to the number of elements it contains. Nonetheless, the expected search time still turns out to be $\Theta(1 + \alpha)$.

### Theorem 11.2
In a hash table in which collisions are resolved by chaining, a successful search takes average-case time $\Theta(1+\alpha)$, under the assumption of simple uniform hashing.

*Proof*   We assume that the element being searched for is equally likely to be any of the $n$ elements stored in the table. The number of elements examined during a successful search for an element $x$ is one more than the number of elements that

appear before $x$ in $x$'s list. Because new elements are placed at the front of the list, elements before $x$ in the list were all inserted after $x$ was inserted. To find the expected number of elements examined, we take the average, over the $n$ elements $x$ in the table, of 1 plus the expected number of elements added to $x$'s list after $x$ was added to the list. Let $x_i$ denote the $i$th element inserted into the table, for $i = 1, 2, \ldots, n$, and let $k_i = x_i.key$. For keys $k_i$ and $k_j$, we define the indicator random variable $X_{ij} = \mathrm{I}\{h(k_i) = h(k_j)\}$. Under the assumption of simple uniform hashing, we have $\Pr\{h(k_i) = h(k_j)\} = 1/m$, and so by Lemma 5.1, $\mathrm{E}[X_{ij}] = 1/m$. Thus, the expected number of elements examined in a successful search is

$$
\begin{aligned}
\mathrm{E}&\left[\frac{1}{n}\sum_{i=1}^{n}\left(1 + \sum_{j=i+1}^{n} X_{ij}\right)\right] \\
&= \frac{1}{n}\sum_{i=1}^{n}\left(1 + \sum_{j=i+1}^{n} \mathrm{E}[X_{ij}]\right) \quad \text{(by linearity of expectation)} \\
&= \frac{1}{n}\sum_{i=1}^{n}\left(1 + \sum_{j=i+1}^{n} \frac{1}{m}\right) \\
&= 1 + \frac{1}{nm}\sum_{i=1}^{n}(n - i) \\
&= 1 + \frac{1}{nm}\left(\sum_{i=1}^{n} n - \sum_{i=1}^{n} i\right) \\
&= 1 + \frac{1}{nm}\left(n^2 - \frac{n(n+1)}{2}\right) \quad \text{(by equation (A.1))} \\
&= 1 + \frac{n-1}{2m} \\
&= 1 + \frac{\alpha}{2} - \frac{\alpha}{2n}\ .
\end{aligned}
$$

Thus, the total time required for a successful search (including the time for computing the hash function) is $\Theta(2 + \alpha/2 - \alpha/2n) = \Theta(1 + \alpha)$. ∎

What does this analysis mean? If the number of hash-table slots is at least proportional to the number of elements in the table, we have $n = O(m)$ and, consequently, $\alpha = n/m = O(m)/m = O(1)$. Thus, searching takes constant time on average. Since insertion takes $O(1)$ worst-case time and deletion takes $O(1)$ worst-case time when the lists are doubly linked, we can support all dictionary operations in $O(1)$ time on average.

## Exercises

### 11.2-1
Suppose we use a hash function $h$ to hash $n$ distinct keys into an array $T$ of length $m$. Assuming simple uniform hashing, what is the expected number of collisions? More precisely, what is the expected cardinality of $\{\{k,l\} : k \neq l$ and $h(k) = h(l)\}$?

### 11.2-2
Demonstrate what happens when we insert the keys $5, 28, 19, 15, 20, 33, 12, 17, 10$ into a hash table with collisions resolved by chaining. Let the table have 9 slots, and let the hash function be $h(k) = k \bmod 9$.

### 11.2-3
Professor Marley hypothesizes that he can obtain substantial performance gains by modifying the chaining scheme to keep each list in sorted order. How does the professor's modification affect the running time for successful searches, unsuccessful searches, insertions, and deletions?

### 11.2-4
Suggest how to allocate and deallocate storage for elements within the hash table itself by linking all unused slots into a free list. Assume that one slot can store a flag and either one element plus a pointer or two pointers. All dictionary and free-list operations should run in $O(1)$ expected time. Does the free list need to be doubly linked, or does a singly linked free list suffice?

### 11.2-5
Suppose that we are storing a set of $n$ keys into a hash table of size $m$. Show that if the keys are drawn from a universe $U$ with $|U| > nm$, then $U$ has a subset of size $n$ consisting of keys that all hash to the same slot, so that the worst-case searching time for hashing with chaining is $\Theta(n)$.

### 11.2-6
Suppose we have stored $n$ keys in a hash table of size $m$, with collisions resolved by chaining, and that we know the length of each chain, including the length $L$ of the longest chain. Describe a procedure that selects a key uniformly at random from among the keys in the hash table and returns it in expected time $O(L \cdot (1 + 1/\alpha))$.