# 11 Hash Tables

Many applications require a dynamic set that supports only the dictionary operations INSERT, SEARCH, and DELETE. For example, a compiler for a computer language maintains a symbol table, in which the keys of elements are arbitrary character strings that correspond to identifiers in the language. A hash table is an effective data structure for implementing dictionaries. Although searching for an element in a hash table can take as long as searching for an element in a linked list—$\Theta(n)$ time in the worst case—in practice, hashing performs extremely well. Under reasonable assumptions, the expected time to search for an element in a hash table is $O(1)$.

A hash table is a generalization of the simpler notion of an ordinary array. Directly addressing into an ordinary array makes effective use of our ability to examine an arbitrary position in an array in $O(1)$ time. Section 11.1 discusses direct addressing in more detail. Direct addressing is applicable when we can afford to allocate an array that has one position for every possible key.

When the number of keys actually stored is small relative to the total number of possible keys, hash tables become an effective alternative to directly addressing an array, since a hash table typically uses an array of size proportional to the number of keys actually stored. Instead of using the key as an array index directly, the array index is *computed* from the key. Section 11.2 presents the main ideas, focusing on "chaining" as a way to handle "collisions" in which more than one key maps to the same array index. Section 11.3 describes how array indices can be computed from keys using hash functions. We present and analyze several variations on the basic theme. Section 11.4 looks at "open addressing," which is another way to deal with collisions. The bottom line is that hashing is an extremely effective and practical technique: the basic dictionary operations require only $O(1)$ time on the average. Section 11.5 explains how "perfect hashing" can support searches in $O(1)$ *worst-case* time, when the set of keys being stored is static (that is, when the set of keys never changes once stored).

## 11.1    Direct-address tables

Direct addressing is a simple technique that works well when the universe $U$ of keys is reasonably small. Suppose that an application needs a dynamic set in which each element has a key drawn from the universe $U = \{0, 1, \ldots, m - 1\}$, where $m$ is not too large. We shall assume that no two elements have the same key.

To represent the dynamic set, we use an array, or *direct-address table*, denoted by $T[0 \ldots m - 1]$, in which each position, or *slot*, corresponds to a key in the universe $U$. Figure 11.1 illustrates the approach; slot $k$ points to an element in the set with key $k$. If the set contains no element with key $k$, then $T[k] = \text{NIL}$.

The dictionary operations are trivial to implement.

DIRECT-ADDRESS-SEARCH$(T, k)$
    **return** $T[k]$

DIRECT-ADDRESS-INSERT$(T, x)$
    $T[key[x]] \leftarrow x$

DIRECT-ADDRESS-DELETE$(T, x)$
    $T[key[x]] \leftarrow \text{NIL}$

Each of these operations is fast: only $O(1)$ time is required.

For some applications, the elements in the dynamic set can be stored in the direct-address table itself. That is, rather than storing an element's key and satellite data in an object external to the direct-address table, with a pointer from a slot in the table to the object, we can store the object in the slot itself, thus saving space. Moreover, it is often unnecessary to store the key field of the object, since if we have the index of an object in the table, we have its key. If keys are not stored, however, we must have some way to tell if the slot is empty.

### Exercises

***11.1-1***
Suppose that a dynamic set $S$ is represented by a direct-address table $T$ of length $m$. Describe a procedure that finds the maximum element of $S$. What is the worst-case performance of your procedure?

***11.1-2***
A *bit vector* is simply an array of bits (0's and 1's). A bit vector of length $m$ takes much less space than an array of $m$ pointers. Describe how to use a bit vector
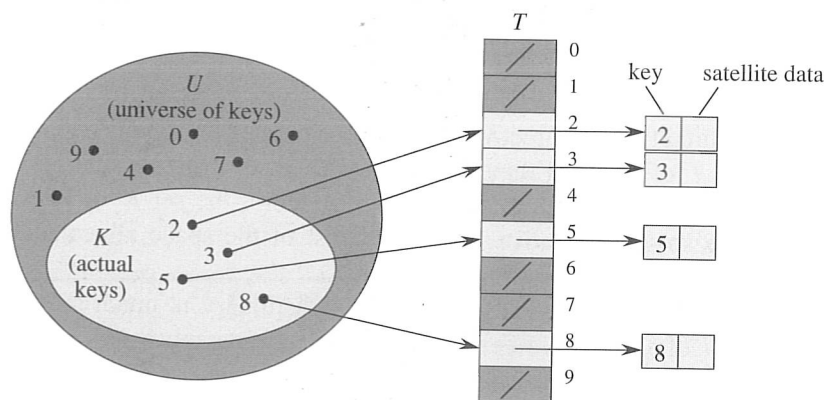
**Figure 11.1**  Implementing a dynamic set by a direct-address table $T$. Each key in the universe $U = \{0, 1, \ldots, 9\}$ corresponds to an index in the table. The set $K = \{2, 3, 5, 8\}$ of actual keys determines the slots in the table that contain pointers to elements. The other slots, heavily shaded, contain NIL.

to represent a dynamic set of distinct elements with no satellite data. Dictionary operations should run in $O(1)$ time.

*11.1-3*

Suggest how to implement a direct-address table in which the keys of stored elements do not need to be distinct and the elements can have satellite data. All three dictionary operations (INSERT, DELETE, and SEARCH) should run in $O(1)$ time. (Don't forget that DELETE takes as an argument a pointer to an object to be deleted, not a key.)

*11.1-4*  ⋆

We wish to implement a dictionary by using direct addressing on a *huge* array. At the start, the array entries may contain garbage, and initializing the entire array is impractical because of its size. Describe a scheme for implementing a direct-address dictionary on a huge array. Each stored object should use $O(1)$ space; the operations SEARCH, INSERT, and DELETE should take $O(1)$ time each; and the initialization of the data structure should take $O(1)$ time. (*Hint:* Use an additional stack, whose size is the number of keys actually stored in the dictionary, to help determine whether a given entry in the huge array is valid or not.)

## 11.2    Hash tables

The difficulty with 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, the storage requirements can be reduced to $\Theta(|K|)$ while we maintain the benefit that searching for an element in the hash table still requires only $O(1)$ time. The only catch is that this bound is for the *average 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** $T[0 \ldots m-1]$:

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

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 point of the hash function is to reduce the range of array indices that need to be handled. Instead of $|U|$ values, we need to handle only $m$ values. Storage requirements are correspondingly reduced.

There is one hitch: two keys may hash to the same slot. We call this situation a **collision**. Fortunately, there are 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 $h(k)$.) Since $|U| > m$, 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.
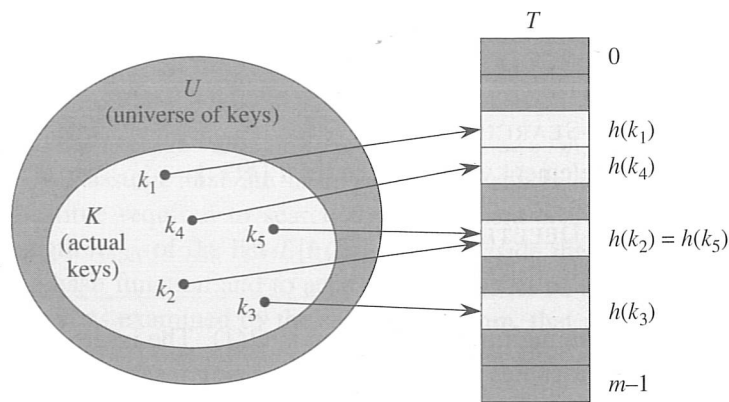
**Figure 11.2**   Using a hash function $h$ to map keys to hash-table slots. Keys $k_2$ and $k_5$ map to the same slot, so they collide.
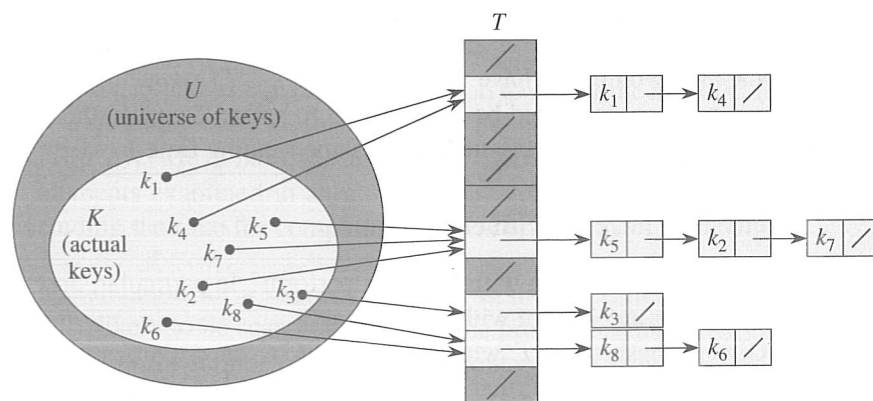


**Figure 11.3**   Collision resolution by chaining. Each hash-table slot $T[j]$ contains a linked list of all the keys whose hash value is $j$. For example, $h(k_1) = h(k_4)$ and $h(k_5) = h(k_2) = h(k_7)$.

## Collision resolution by chaining

In ***chaining***, we put all the elements that hash to the same slot in a linked list, as shown in Figure 11.3. 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($T, x$)
    insert $x$ at the head of list $T[h(key[x])]$

CHAINED-HASH-SEARCH($T, k$)
    search for an element with key $k$ in list $T[h(k)]$

CHAINED-HASH-DELETE($T, x$)
    delete $x$ from the list $T[h(key[x])]$

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; this assumption can be checked if necessary (at additional cost) by performing a search before insertion. For searching, the worst-case running time is proportional to the length of the list; we shall analyze this operation more closely below. Deletion of an element $x$ can be accomplished in $O(1)$ time if the lists are doubly linked. (Note that CHAINED-HASH-DELETE takes as input an element $x$ and not its key $k$, so we don't have to search for $x$ first. If the lists were singly linked, it would not be of great help to take as input the element $x$ rather than the key $k$. We would still have to find $x$ in the list $T[h(key[x])]$, so that the *next* link of $x$'s predecessor could be properly set to splice $x$ out. In this case, deletion and searching would have essentially the same running time.)

### 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, hash tables are not used for their worst-case performance. (Perfect hashing, described in Section 11.5, does however provide good worst-case performance when the set of keys is static.)

The average 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 average value of $n_j$ is $\mathrm{E}[n_j] = \alpha = n/m$.

We assume that the hash value $h(k)$ can be computed in $O(1)$ time, 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 are checked to see if their keys are 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 expected 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 is still $\Theta(1 + \alpha)$.

### Theorem 11.2

In a hash table in which collisions are resolved by chaining, a successful search takes time $\Theta(1 + \alpha)$, on the average, 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 1 more than the number of elements that appear before $x$ in $x$'s list. Elements before $x$ in the list were all inserted after $x$ was inserted, because new elements are placed at the front of the list. 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 ta-

ble, for $i = 1, 2, \ldots, n$, and let $k_i = key[x_i]$. 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, all dictionary operations can be supported 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 the insertion of 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 substantial performance gains can be obtained if we modify the chaining scheme so that each list is kept 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 storage for elements can be allocated and deallocated 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*

Show that if $|U| > nm$, there is a subset of $U$ 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 $\Theta(n)$.

## 11.3    Hash functions

In this section, we discuss some issues regarding the design of good hash functions and then present three schemes for their creation. Two of the schemes, hashing by division and hashing by multiplication, are heuristic in nature, whereas the third scheme, universal hashing, uses randomization to provide provably good performance.

### What makes a good hash function?

A good hash function satisfies (approximately) the assumption of simple uniform hashing: each key is equally likely to hash to any of the $m$ slots, independently of where any other key has hashed to. Unfortunately, it is typically not possible to check this condition, since one rarely knows the probability distribution according to which the keys are drawn, and the keys may not be drawn independently.

Occasionally we do know the distribution. For example, if the keys are known to be random real numbers $k$ independently and uniformly distributed in the range $0 \leq k < 1$, the hash function

$$h(k) = \lfloor km \rfloor$$

satisfies the condition of simple uniform hashing.

In practice, heuristic techniques can often be used to create a hash function that performs well. Qualitative information about distribution of keys may be useful in this design process. For example, consider a compiler's symbol table, in which the keys are character strings representing identifiers in a program. Closely related symbols, such as pt and pts, often occur in the same program. A good hash function would minimize the chance that such variants hash to the same slot.

A good approach is to derive the hash value in a way that is expected to be independent of any patterns that might exist in the data. For example, the "division method" (discussed in Section 11.3.1) computes the hash value as the remainder when the key is divided by a specified prime number. This method frequently gives good results, assuming that the prime number is chosen to be unrelated to any patterns in the distribution of keys.

Finally, we note that some applications of hash functions might require stronger properties than are provided by simple uniform hashing. For example, we might want keys that are "close" in some sense to yield hash values that are far apart. (This property is especially desirable when we are using linear probing, defined in Section 11.4.) Universal hashing, described in Section 11.3.3, often provides the desired properties.

### Interpreting keys as natural numbers

Most hash functions assume that the universe of keys is the set $\mathbf{N} = \{0, 1, 2, \ldots\}$ of natural numbers. Thus, if the keys are not natural numbers, a way is found to interpret them as natural numbers. For example, a character string can be interpreted as an integer expressed in suitable radix notation. Thus, the identifier pt might be interpreted as the pair of decimal integers (112, 116), since $p = 112$ and $t = 116$ in the ASCII character set; then, expressed as a radix-128 integer, pt becomes $(112 \cdot 128) + 116 = 14452$. It is usually straightforward in an application to devise some such method for interpreting each key as a (possibly large) natural number. In what follows, we assume that the keys are natural numbers.

### 11.3.1   The division method

In the *division method* for creating hash functions, we map a key $k$ into one of $m$ slots by taking the remainder of $k$ divided by $m$. That is, the hash function is

$$h(k) = k \bmod m .$$

For example, if the hash table has size $m = 12$ and the key is $k = 100$, then $h(k) = 4$. Since it requires only a single division operation, hashing by division is quite fast.

When using the division method, we usually avoid certain values of $m$. For example, $m$ should not be a power of 2, since if $m = 2^p$, then $h(k)$ is just the $p$ lowest-order bits of $k$. Unless it is known that all low-order $p$-bit patterns are equally likely, it is better to make the hash function depend on all the bits of the key. As Exercise 11.3-3 asks you to show, choosing $m = 2^p - 1$ when $k$ is a character string interpreted in radix $2^p$ may be a poor choice, because permuting the characters of $k$ does not change its hash value.

A prime not too close to an exact power of 2 is often a good choice for $m$. For example, suppose we wish to allocate a hash table, with collisions resolved by chaining, to hold roughly $n = 2000$ character strings, where a character has 8 bits. We don't mind examining an average of 3 elements in an unsuccessful search, so we allocate a hash table of size $m = 701$. The number 701 is chosen because it is a prime near $2000/3$ but not near any power of 2. Treating each key $k$ as an integer, our hash function would be

$$h(k) = k \bmod 701 \ .$$

### 11.3.2   The multiplication method

The ***multiplication method*** for creating hash functions operates in two steps. First, we multiply the key $k$ by a constant $A$ in the range $0 < A < 1$ and extract the fractional part of $kA$. Then, we multiply this value by $m$ and take the floor of the result. In short, the hash function is

$$h(k) = \lfloor m \, (k \, A \bmod 1) \rfloor \ ,$$

where "$k \, A \bmod 1$" means the fractional part of $kA$, that is, $kA - \lfloor kA \rfloor$.

An advantage of the multiplication method is that the value of $m$ is not critical. We typically choose it to be a power of 2 ($m = 2^p$ for some integer $p$) since we can then easily implement the function on most computers as follows. Suppose that the word size of the machine is $w$ bits and that $k$ fits into a single word. We restrict $A$ to be a fraction of the form $s/2^w$, where $s$ is an integer in the range $0 < s < 2^w$. Referring to Figure 11.4, we first multiply $k$ by the $w$-bit integer $s = A \cdot 2^w$. The result is a $2w$-bit value $r_1 2^w + r_0$, where $r_1$ is the high-order word of the product and $r_0$ is the low-order word of the product. The desired $p$-bit hash value consists of the $p$ most significant bits of $r_0$.

Although this method works with any value of the constant $A$, it works better with some values than with others. The optimal choice depends on the characteristics of the data being hashed. Knuth [185] suggests that
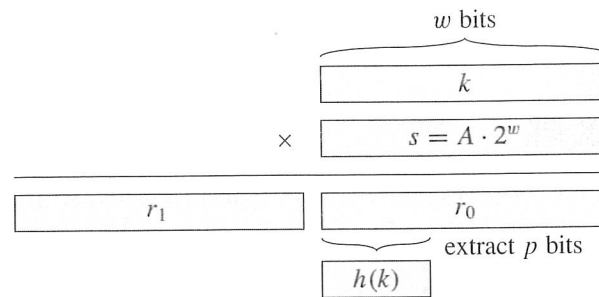
**Figure 11.4**   The multiplication method of hashing. The $w$-bit representation of the key $k$ is multiplied by the $w$-bit value $s = A \cdot 2^w$. The $p$ highest-order bits of the lower $w$-bit half of the product form the desired hash value $h(k)$.

$$A \approx (\sqrt{5} - 1)/2 = 0.6180339887\ldots \tag{11.2}$$

is likely to work reasonably well.

As an example, suppose we have $k = 123456$, $p = 14$, $m = 2^{14} = 16384$, and $w = 32$. Adapting Knuth's suggestion, we choose $A$ to be the fraction of the form $s/2^{32}$ that is closest to $(\sqrt{5} - 1)/2$, so that $A = 2654435769/2^{32}$. Then $k \cdot s = 327706022297664 = (76300 \cdot 2^{32}) + 17612864$, and so $r_1 = 76300$ and $r_0 = 17612864$. The 14 most significant bits of $r_0$ yield the value $h(k) = 67$.

★    ### 11.3.3   Universal hashing

If a malicious adversary chooses the keys to be hashed by some fixed hash function, then he 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 what keys are chosen by the adversary.

The main idea behind universal hashing is to select the hash function at random from a carefully designed class of functions at the beginning of execution. As in the case of quicksort, randomization guarantees that no single input will always evoke worst-case behavior. Because of the randomization, 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