

---

## 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, we typically have no way to check this condition, since we rarely know the probability distribution from which the keys are drawn. Moreover, the keys might not be drawn independently.

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

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

satisfies the condition of simple uniform hashing.

In practice, we can often employ heuristic techniques to create a hash function that performs well. Qualitative information about the 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 derives the hash value in a way that we expect 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 we choose a prime number that is 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  $\mathbb{N} = \{0, 1, 2, \dots\}$  of natural numbers. Thus, if the keys are not natural numbers, we find a way to interpret them as natural numbers. For example, we can interpret a character string as an integer expressed in suitable radix notation. Thus, we might interpret the identifier `pt` 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  $(112 \cdot 128) + 116 = 14452$ . In the context of a given application, we can usually 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.

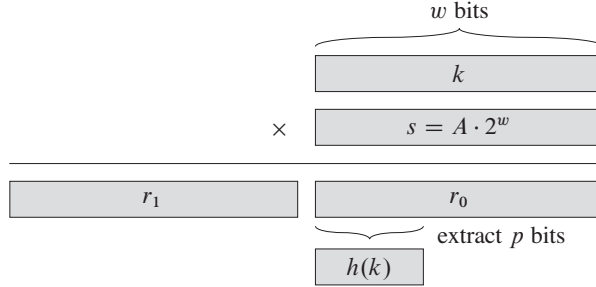
When using the division method, we usually avoid certain values of  $m$ . For example,  $m$  should not be a power of 2, since if  $m = 2^p$ , then  $h(k)$  is just the  $p$  lowest-order bits of  $k$ . Unless we know that all low-order  $p$ -bit patterns are equally likely, we are better off designing the hash function to 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, and so we allocate a hash table of size  $m = 701$ . We could choose  $m = 701$  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



**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  $h(k)$ .

fractional part of  $kA$ . Then, we multiply this value by  $m$  and take the floor of the result. In short, the hash function is

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

where “ $kA \bmod 1$ ” 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 [211] suggests that

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

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 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

$X_{kl} = \mathbf{I}\{h(k) = h(l)\}$ . 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  $E[X_{kl}] \leq 1/m$ .

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

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

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  $E[n_{h(k)}] = E[Y_k] \leq n/m = \alpha$ .
- 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  $E[n_{h(k)}] = E[Y_k] + 1 \leq (n - 1)/m + 1 = 1 + \alpha - 1/m < 1 + \alpha$ . ■

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?