

## 10.4. Skip Lists 437

### 10.4 Skip Lists

An interesting data structure for realizing the sorted map ADT is the **skip list**. In Section 10.3.1, we saw that a sorted array will allow  $O(\log n)$ -time searches via the binary search algorithm. Unfortunately, update operations on a sorted array have  $O(n)$  worst-case running time because of the need to shift elements. In Chapter 7 we demonstrated that linked **lists** support very efficient update operations, as long as the position within the **list** is identified. Unfortunately, we cannot perform fast searches on a standard linked **list**; an efficient means for direct accessing an element of a sequence by index.

**Skip lists** provide a clever compromise to efficiently support search and update operations. A **skip list**  $S$  for a map  $M$  consists of a series of **lists**  $\{S_0, S_1, \dots, S_h\}$ .

Each **list**  $S_i$  stores a subset of the items of  $M$  sorted by increasing keys, plus items with two sentinel keys denoted  $-\infty$  and  $+\infty$ , where  $-\infty$  is smaller than every possible key that can be inserted in  $M$  and  $+\infty$  is larger than every possible key that can be inserted in  $M$ . In addition, the **lists** in  $S$  satisfy the following:

- **List**  $S_0$  contains every item of the map  $M$  (plus sentinels  $-\infty$  and  $+\infty$ ).
- For  $i=1, \dots, h-1$ , **list**  $S_i$  contains (in addition to  $-\infty$  and  $+\infty$ ) a randomly generated subset of the items in **list**  $S_{i-1}$ .
- **List**  $S_h$  contains only  $-\infty$  and  $+\infty$ .

An example of a **skip list** is shown in Figure 10.10. It is customary to visualize a **skip list**  $S$  with **list**  $S_0$  at the bottom and **lists**  $S_1, \dots, S_h$  above it. Also, we refer to  $h$  as the height of **skip list**  $S$ .

Intuitively, the **lists** are set up so that  $S_{i+1}$  contains more or less alternate items of  $S_i$ . As we shall see in the details of the insertion method, the items in  $S_{i+1}$  are chosen at random from the items in  $S_i$  by picking each item from  $S_i$  to also be in  $S_{i+1}$  with probability  $1/2$ . That is, in essence, we flip a coin for each item in  $S_i$ .

3125

25

-----

--

--1717

17

17 12 S5

S4

S3

S2

S1

S0 55

55

55

55 12 17 20 25 31 38 39 44 50 +--+--+--+

44 38 31 25

Figure 10.10: Example of a **skip list** storing 10 items. For simplicity, we show only the items' keys, not their associated values.

## 438 Chapter 10. Maps, Hash Tables, and Skip Lists

and place that item in  $S_{i+1}$  if the coin comes up ·heads·. Thus, we expect  $S_1$  to have about  $n/2$  items,  $S_2$  to have about  $n/4$  items, and, in general,  $S_i$  to have about  $n/2^i$  items. In other words, we expect the height of  $S$  to be about  $\log n$ . The halving of the number of items from one list to the next is not enforced as an explicit property of skip lists, however. Instead, randomization is used.

Functions that generate numbers that can be viewed as random numbers are built into most modern computers, because they are used extensively in computer games, cryptography, and computer simulations. Random number generators, generate random-like numbers, starting with an initial seed. (See discussion of random module in Section 1.11.1.) Other methods use hardware devices to extract ·true· random numbers from nature. In any case, we will assume that our computer has a good random number generator. The main advantage of using randomization in data structure and algorithm design is that the structures and functions that result are usually simple and efficient.

The skip list has the same logarithmic time bounds for searching as is achieved by the binary search algorithm, yet it extends that performance to update methods when inserting or deleting items. skip list, while binary search has a worst-case bound with a sorted table.

A skip list makes random choices in arranging its structure in such a way that search and update times are  $O(\log n)$  on average, where  $n$  is the number of items in the map. Interestingly, the notion of average time complexity used here does not depend on the probability distribution used to help decide where to place the new item. The running time is averaged over all possible outcomes of the random numbers used when inserting entries.

Using the position abstraction used for lists and trees, we view a skip list as a two-dimensional collection of positions arranged horizontally into levels and vertically into towers. Each level is a list  $S$

and each tower contains positions storing

the same item across consecutive lists. The positions in a skip list can be traversed using the following operations:

$\text{next}(p)$  : Return the position following  $p$  on the same level.

$\text{prev}(p)$  : Return the position preceding  $p$  on the same level.

$\text{below}(p)$  : Return the position below  $p$  in the same tower.

$\text{above}(p)$  : Return the position above  $p$  in the same tower.

We conventionally assume that the above operations return `None` if the position

requested does not exist. Without going into the details, we note that we can easily implement a skip list by means of a doubly linked list structure. The structure is essentially a collection of doubly linked lists aligned at towers, which are also doubly linked lists.

## Removal in a Skip List

Like the search and insertion algorithms, the removal algorithm for a skip list is quite simple. In fact, it is even easier than the insertion algorithm. That is, to perform the map operation `del M[k]` we begin by executing method `SkipSearch(k)`. If the position `p` stores an entry with key different from `k`, we raise a `KeyError`.

Otherwise, we remove `p` and all the positions above `p`, which are easily accessed by using above operations to climb up the tower of this entry in `S` starting at position `p`. While removing levels of the tower, we reestablish links between the horizontal neighbors of each removed position. The removal algorithm is illustrated in Figure 10.13 and a detailed de-

$O(\log n)$  expected running time.

Before we give this analysis, however, there are some minor improvements to the skip-list data structure we would like to discuss. First, we do not actually need to store references to values at more efficiently represent a tower as a single object, storing the key-value pair, and maintaining `jprevious` references and `jnext` references if the tower reaches level `S`.

j. Second, for the horizontal axes, it is possible to keep the list singly linked,

storing only the next references. We can perform insertions and removals in strictly a top-down, scan-forward fashion.

Exercise C-10.44. Neither of these optimizations improve the asymptotic performance of skip lists by more than a constant factor, but these improvements can, nevertheless, be meaningful in practice.

search trees, which are discussed in Chapter 11.

31 S5

S4

S3

S2

S1 ---

-- 1212 --

1717 25

25 31

3142

55 5055+

+..+

+

+.....

17

38

38 39 424242

44

445555+

17

17

20 2525

S0

Figure 10.13: Removal of the entry with key 25 from the skip list of Figure 10.12.

The positions visited after the search for the position of `S0` holding the entry are highlighted. The positions removed are drawn with dashed lines.

Algorithm `SkipSearch(k)` :

Input: A search key  $k$

Output: Position  $p$  in the bottom list  $S_0$  with the largest key such that  $\text{key}(p) \leq k$

$p = \text{start}$  {begin at start position }

while  $\text{below}(p) \neq \text{None}$  do

$p = \text{below}(p)$  {drop down }

while  $k < \text{key}(\text{next}(p))$  do

$p = \text{next}(p)$  {scan forward }

return  $p$ .

Code Fragment 10.12: Algorithm to search a skip list for key  $k$ .

As it turns out, the expected running time of algorithm `SkipSearch` on a skip list with  $n$  entries is  $O(\log n)$ . We postpone the justification of this fact, however, until after we discuss the implementation of the update methods for skip lists. Navigation starting at the position identified by `SkipSearch(k)` can be easily used to provide the additional forms of searches in the sorted map ADT (e.g., `find`, `findRange`).

Insertion into a Skip List

The execution of the map operation  $M[k] = v$  begins with a call to `SkipSearch(k)`.

This gives us the position  $p$  of the bottom-level item with the largest key less than or equal to  $k$  (note that  $p$  may hold the special item with key  $-\infty$ ). If  $\text{key}(p) = k$ , the

associated value is overwritten with  $v$ . Otherwise, we need to create a new tower for item  $(k, v)$ . We insert  $(k, v)$  immediately after position  $p$  within  $S_0$ . After inserting

the new item at the bottom level, we use randomization to decide the height of the tower for the new item. We flip a coin, and if the flip comes up tails, then we stop here. Else (the flip comes up heads), we backtrack to the previous (next higher)

level and insert  $(k, v)$  in this level at the appropriate position. We again flip a coin;

if it comes up heads, we go to the next higher level and repeat. Thus, we continue to insert the new item  $(k, v)$  in list

We link together all the references to the new item  $(k, v)$  created in this process to

create its tower. A coin flip can be simulated with Python's built-in pseudo-random number generator from the range 0 or 1, each with probability  $1/2$ .

We give the insertion algorithm for a skip list in Code Fragment 10.13 and

we illustrate it in Figure 10.12. The algorithm uses an `insertAfterAbove(p, q, (k, v))`

method that inserts a position storing the item  $(k, v)$  after position  $p$  (on the same level as  $p$ ) and above position  $q$ , returning the new position  $r$  (and setting internal references so that `next`, `prev`, `above`, and `below` methods will work correctly for  $p$ ,

$q$ , and  $r$ ). The expected running time of the insertion algorithm on a skip list with  $n$  entries is  $O(\log n)$ , which we show in Section 10.4.2.

#### 444 Chapter 10. Maps, Hash Tables, and Skip Lists

##### Bounding the Height of a Skip List

Because the insertion step involves randomization, a more accurate analysis of skip lists involves a bit of probability. At first, this might seem like a major undertaking, for a complete and thorough probabilistic analysis could require deep mathematics (and, indeed, there are several ways to understand the expected asymptotic behavior of skip lists). The informal and intuitive probabilistic analysis we give below uses only basic concepts of probability theory. Let us begin by determining the expected value of the height of a skip list  $S$  with  $n$  entries (assuming that we do not terminate insertions early). The probability that a given entry has a tower of height  $i+1$  is equal to the probability of getting  $i$  consecutive heads when flipping a coin, that is, this probability is  $1/2^i$ .

i. Hence, the

probability that level  $i$  has at least one position is at most

$$P_i \leq n/2^i,$$

for the probability that any one of  $n$  different events occurs is at most the sum of the probabilities that each occurs.

The probability that the height of  $S$  is larger than  $i$  is equal to the probability that level  $i$  has at least one position, that is, it is no more than  $P_i$ . This means that  $h(S)$  is larger than, say,  $3 \log n$  with probability at most

$P_{3 \log n} \leq n/2^{3 \log n} = n/n^3 = 1/n^2$ .

For example, if  $n=1000$ , this probability is a one-in-a-million long shot. More generally, given a constant  $c > 1$ ,  $h(S)$  is larger than  $c \log n$  with probability at most  $1/n^{c-1}$ . That is, the probability that  $h(S)$  is smaller than  $c \log n$  is at least  $1 - 1/n^{c-1}$ .

Thus, with high probability, the height of  $S$  is  $O(\log n)$ .

Analyzing Search Time in a Skip List

Next, consider the running time of a search in skip list  $S$ , and recall that such a search involves two nested while loops. The inner loop performs a scan forward on level  $i$  of  $S$  as long as the next key is no greater than the search key  $k$ , and the outer loop drops down to the next level and repeats the scan forward iteration. Since the height of  $S$  is  $O(\log n)$  with high probability, the search time is  $O(\log n)$  with high probability.

## 10.4. Skip Lists 443

### Maintaining the Topmost Level

A skip list must maintain a reference to the start position (the topmost, left position in  $S$ ) as an instance variable, and must have a policy for any insertion that wishes to continue inserting a new entry past the top level of  $S$ . There are two possible courses of action we can take, both of which have their merits.

One possibility is to restrict the top level,  $h$ , to be kept at some fixed value that is a function of  $n$ , the number of entries currently in the map (from the analysis we will see that  $h = \max\{10, 2 \cdot \log n\}$  is a reasonable choice, and picking  $h = 3 \cdot \log n$  is even safer). Implementing this choice means that we must modify the insertion algorithm to stop inserting a new position once we reach the topmost level (unless  $\log n < \log(n+1)$ , in which case we can now go at least one more level, since the bound on the height is increasing).

The other possibility is to let an insertion continue inserting a new position as long as heads keeps getting returned from the random number generator. This is the approach taken by algorithm in the analysis of skip lists, the probability that an insertion will go to a level that is more than  $O(\log n)$  is very low, so. Either choice will still result in the expected  $O(\log n)$  time to perform search, insertion, and removal, however, which we show in the next section.

### 10.4.2 Probabilistic Analysis of Skip Lists

As we have shown above, skip lists provide a simple implementation of a sorted map. In terms of worst-case performance, the probability of having a fair coin repeatedly come up heads forever is 0. Moreover, we cannot infinitely add positions to a list without eventually running out of memory. In any case, if we terminate positions, the worst-case running time for performing the

getitem

,

setitem

, and

delitem

map operations in a skip list with  $n$  entries and height  $h$  is  $O(n+h)$ . This worst-case performance occurs when the tower of every entry reaches level  $h-1$ , where  $h$  is the height of  $S$ . However, this event has very low probability. Judging from this worst case, we might conclude that the skip-list structure is strictly inferior to the other map implementations. But, as we will see, this is not a fair analysis, for this worst-case behavior is a gross overestimate.

#### 10.4. Skip Lists 439

##### 10.4.1 Search and Update Operations in a Skip List

The skip-list structure affords simple map search and update algorithms. In fact, all of the skip-list search and update algorithms are based on an elegant SkipSearch method that takes a key  $k$  and finds the position  $p$  of the item in list  $S$  that has the largest key less than or equal to  $k$  (which is possibly  $..$ ).

##### Searching in a Skip List

Suppose we are given a search key  $k$ . We begin the SkipSearch method by setting a position variable  $p$  to the topmost, left position in the skip list  $S$ , called the start position of  $S$ . That is, the start position is the position of  $S$  storing the special entry with key  $..$ . We then perform the following steps (see Figure 10.11), where  $\text{key}(p)$  denotes the key of the item at position  $p$ :

1. If  $S.\text{below}(p)$  is  $\text{None}$ , then the search terminates—we are at the bottom and have located the item in  $S$  with the largest key less than or equal to the search key  $k$ . Otherwise, we drop down to the next lower level in the present tower by setting  $p = S.\text{below}(p)$ .

2. Starting at position  $p$ , we move  $p$  forward until it is at the rightmost position on the present level such that  $\text{key}(p) \leq k$ . We call this the scan forward step.

Note that such a position always exists, since each level contains the keys  $+$  and  $..$ . It may be that  $p$  remains where it was when we started the scan forward for this level.

3. Return to step 1.

55S1S2S3S4S5

+..

+..+

+..

+..+

---

-- 1212 --17

17 25

25 20 17 31 38 39--

-- 1717 25

25 31

31 38 44

44 50555555

S0

Figure 10.11: Example of a search in a skip list. The positions examined when searching for key 50 are highlighted.

We give a pseudo-code description of the skip-list search algorithm, SkipSearch, in Code Fragment 10.12. Given this method, the map operation  $M[k]$  is performed by computing  $p = \text{SkipSearch}(k)$  and testing whether or not  $\text{key}(p) = k$ . If these two keys are equal, we return the associated value; otherwise, we raise a `KeyError`.

#### 10.4. Skip Lists 441

Algorithm SkipInsert(k,v) :

Input: Key and value v

Output: Topmost position of the item inserted in the skip list

p = SkipSearch (k)

q = None {q will represent top node in new item's tower }

i = 1

repeat

i = i + 1

if i = h then

h = h + 1 {add a new level to the skip list }

t = next(s)

s = insertAfterAbove (None, s, (∞, None )) {grow leftmost tower }

insertAfterAbove (s, t, (∞, None )) {grow rightmost tower }

while above (p) is None do

p = prev (p) {scan backward }

p = above (p) {jump up to higher level }

q = insertAfterAbove (p, q, (k, v)) {increase height of new item's tower }

until coinFlip () == tails

n = n + 1

return q

Code Fragment 10.13: Insertion in a skip list. Method coinFlip () returns 'heads' or 'tails', each with probability 1 / 2. Instance variables n, h, and s hold the number of entries, the height, and the start node of the skip list.

55 S1S2S3S4S5

+.

+..+

+.

+..+

---

-- 1212 --17

17 25

25 20 17 31--

-- 1717 25

25 31

31 38 44

44424242

5555

55 38 39 42 50 S0

Figure 10.12: Insertion of an entry with key 42 into the skip list of Figure 10.10. We assume that the random 'coin' 'ips' for the new entry came up heads three times in a row, followed by tails. The positions visited are highlighted. The positions inserted to hold the new entry are drawn



718 Chapter 15. Memory Management and B-Trees

C-15.11 Describe an external-memory data structure to implement the queue ADT so that the total number of disk transfers needed to process a sequence of enqueue and dequeue operations is  $O(k/B)$ .

C-15.12 Describe an external-memory version of the Positional List ADT (Section 7.4), with block size  $B$ , such that an iteration of a list of length  $n$  is completed using  $O(n/B)$  transfers in the worst case, and all other methods of the ADT require only  $O(1)$  transfers.

C-15.13 Change the rules that define red-black trees so that each red-black tree  $T$  has a corresponding  $(4,8)$  tree, and vice versa.

C-15.14 Describe a modified version of the B-tree insertion algorithm so that each time we create an overflow block, we split the keys among all of  $w$ 's siblings, so that each sibling holds roughly the same number of keys (possibly cascading the split up to the parent of  $w$ ). What is the minimum fraction of each block that will always be filled using this scheme?

C-15.15 Another possible external-memory map implementation is to use a skip list, but to collect consecutive groups of  $O(B)$  nodes, in individual blocks, on any level in the skip list. In particular, we define an order- $d$  B-skip list to be such a representation of a skip list structure, where each block contains at least  $d/2$  list nodes and at most  $d$  list nodes. Let us also choose  $d$  in this case to be the maximum number of list nodes from a level of a skip list that can fit into one block. Describe how we should modify the skip-list insertion and removal algorithms for a B-skip list so that the expected height of the structure is  $O(\log n / \log B)$ .

C-15.16 Describe how to use a B-tree to implement the partition (union-find) ADT (from Section 14.7.3) so that the union and find operations each use at most  $O(\log n / \log B)$  disk transfers.

C-15.17 Suppose we are given a sequence  $S$  of  $n$  elements with integer keys such that some elements in  $S$  are colored blue and some elements in  $S$  are colored red. In addition, say that a red element pairs with a blue element if they have the same key value. Describe an efficient external-memory algorithm for finding all the red-blue pairs in  $S$ . How many disk transfers does your algorithm perform?

C-15.18 Consider the page caching problem where the memory cache can hold  $m$  pages, and we are given a sequence  $P$  of  $n$  requests taken from a pool of  $m+1$  possible pages. Describe the optimal strategy for the offline algorithm and show that it causes at most  $m + n/m$  page misses in total, starting from an empty cache.

C-15.19 Describe an efficient external-memory algorithm that determines whether an array of  $n$  integers contains a

#### 10.4. Skip Lists 445

So we have yet to bound the number of scan-forward steps we make. Let  $n_i$  be the number of keys examined while scanning forward at level  $i$ . Observe that, after the key at the starting position, each additional key examined in a scan-forward at level  $i$  cannot also belong to level  $i+1$ . If any of these keys were on the previous level, we would have encountered them in the previous scan-forward step. Thus, the probability that any key is co-  
 $n_i$  is  $1/2$ . Therefore, the expected value of  $n_i$  is exactly equal to the expected number of times we must flip a fair coin before it comes up heads. This expected value is 2. Hence, the expected amount of time spent scanning forward at any probability, a search in Stakes expected time  $O(\log n)$ . By a similar analysis, we can show that the expected running time of an insertion or a removal is  $O(\log n)$ .

#### Space Usage in a Skip List

Finally, let us turn to the space requirement of a skip list with  $n$  entries. As we observed above, the expected number of positions at level  $i$  is  $n/2^i$ , which means that the expected total number of positions in  $S$  is

$$\sum_{i=0}^h n/2^i$$

$$\leq \sum_{i=0}^{\infty} n/2^i$$

Using Proposition 3.5 on geometric summations, we have

$$\sum_{i=0}^h n/2^i \leq \sum_{i=0}^{\infty} n/2^i = n \sum_{i=0}^{\infty} 1/2^i = n \cdot 2 = 2n$$

Hence, the expected space requirement of  $S$  is  $O(n)$ .

Table 10.4 summarizes the performance of a sorted map realized by a skip list.

Operation	Running Time
$\text{len}(M)$	$O(1)$
$\text{kin}M$	$O(\log n)$ expected
$M[k] = v$	$O(\log n)$ expected
$\text{del } M[k]$	$O(\log n)$ expected
$M.\text{nd}$	$\min(), M.\text{nd}$
$\text{max}()$	$O(1)$
$M.\text{nd}$	$\min(), M.\text{nd}$

## 10.6. Exercises 453

R-10.12 What is the result of Exercise R-10.9 when collisions are handled by double hashing using the secondary hash function  $h_{\text{prime}}(k) = 7 \cdot (k \bmod 7)$ ?

R-10.13 What is the worst-case time for putting  $n$  entries in an initially empty hash table, with collisions resolved by chaining? What is the best case?

R-10.14 Show the result of rehashing the hash table shown in Figure 10.6 into a table of size 19 using the new hash function  $h(k) = 3k \bmod 17$ .

R-10.15 Our `HashMapBase` class maintains a load factor  $\approx 0.5$ . Reimplement that class to allow the user to specify the maximum load, and adjust the concrete subclasses accordingly.

R-10.16 Give a pseudo-code description of an insertion into a hash table that uses quadratic probing to resolve collisions, assuming we also use the trick of replacing deleted entries with a special marker.

R-10.17 Modify our `ProbeHashMap` to use quadratic probing.

R-10.18 Explain why a hash table is not suited to implement a sorted map.

R-10.19 Describe how a sorted `list` implemented as a doubly linked `list` could be used to implement the sorted map.

R-10.20 What is the worst-case asymptotic running time for performing  $n$  deletions from a `SortedTableMap` instance that initially contains  $2n$  entries?

R-10.21 Consider the following variant of the

```
.nd
index method from Code Fragment 10.8, in the context of the SortedTableMap class:
```

```
def
.nd
index(self, k, low, high):
    if high < low:
        return high + 1
    else:
        mid = (low + high) // 2
        if self.table[mid].key < k:
            return self.index(k, mid + 1, high)
        else:
            return self.index(k, low, mid - 1)
```

Does this always produce the same result as the original version? Justify your answer.

R-10.22 What is the expected running time of the methods for maintaining a max-ima set if we insert  $n$  pairs such that each pair has lower cost and performance than one before it? What is contained in the sorted map at the end of this series of operations? What if each pair had a lower cost and higher performance than the one before it?

R-10.23 Draw an example `skip list` that results from performing the following series of operations on the `skip list` shown in Figure 10.13: `del S[38]`, `S[48] =`

```
x
, S[24] =
y
```

456 Chapter 10. Maps, Hash Tables, and Skip Lists

C-10.39 Although keys in a map are distinct, the binary search algorithm can be applied in a more general setting in which an array stores possibly duplicative elements in nondecreasing order. Consider the goal of identifying the index of the leftmost element with key  $g$ .

Does the

`find`

index method as given in Code Fragment 10.8 guarantee

such a result? Does the

`find`

index method as given in Exercise R-10.21

guarantee such a result? Justify your answers.

C-10.40 Suppose we are given two sorted search tables  $S$  and  $T$ , each with  $n$  entries (with  $S$  and  $T$  being implemented with arrays). Describe an  $O(\log^2 n)$ -time algorithm for finding the  $k$ th smallest key in the union of the keys from  $S$  and  $T$  (assuming no duplicates).

C-10.41 Give an  $O(\log n)$ -time solution for the previous problem.

C-10.42 Suppose that each row of an  $n \times n$  array  $A$  consists of 1-s and 0-s such that, in any row of  $A$ , all the 1-s come before any 0-s in that row. Assuming  $A$  is already in memory, describe a method running in  $O(n \log n)$  time (not  $O(n^2)$  time!) for counting the number of 1-s in  $A$ .

C-10.43 Given a collection  $C$  of cost-performance pairs  $(c, p)$ , describe an algorithm for finding the maxima pairs of  $C$  in  $O(n \log n)$  time.

C-10.44 Show that the methods `above(p)` and `prev(p)` are not actually needed to efficiently implement a map using a skip list. That is, we can implement insertions and deletions in a skip list using only `find` and `insert`. In the insertion algorithm, first repeatedly flip the coin to determine the level where you should start inserting the new entry.

C-10.45 Describe how to modify a skip-list representation so that index-based operations, such as retrieving the item at index  $j$ , can be performed in  $O(\log n)$  expected time.

C-10.46 For sets  $S$  and  $T$ , the syntax  $S \cdot T$  returns a new set that is the symmetric difference, that is, a set of elements that are in precisely one of  $S$  or  $T$ . This syntax is supported by the special

`xor`

method. Provide an

implementation of that method in the context of the `MutableSet` abstract base class, relying only on the `remove` primary abstract methods of that class.

C-10.47 In the context of the `MutableSet` abstract base class, describe a concrete implementation of the

`intersect`

method, which supports the syntax  $S \& T$

for computing the intersection of two existing sets.

C-10.48 An inverted file is a critical data structure for implementing a search engine or the index of a book. Given a document  $D$ , which can be viewed as an unordered, numbered list of words, an inverted file is an ordered list of words,  $L$ , such that, for each word  $w$  in  $L$ , we store the indices of the places in  $D$  where  $w$  appears. Design an efficient algorithm for constructing  $L$  from  $D$ .

## 10.6. Exercises 457

C-10.49 Python's collections module provides an `OrderedDict` class that is unrelated to our sorted map abstraction. An `OrderedDict` is a subclass of the standard hash-based dictclass that retains the expected  $O(1)$  performance for the primary map operations, but that also guarantees that the `iter`

method reports items of the map according to `first-in, first-out` (FIFO) order. That is, the key that has been in the dictionary the longest is reported `first`. (The order is unaffected when the value for an existing key is overwritten.) Describe an algorithmic approach for achieving such performance.

### Projects

P-10.50 Perform a comparative analysis that studies the collision rates for various hash codes for character strings, such as various polynomial hash codes for different values of the parameter  $a$ . Use a hash table to determine collisions, but only count collisions where different strings map to the same hash code (not if they map to the same value). Test these hash codes on text files found on the Internet.

P-10.51 Perform a comparative analysis as in the previous exercise, but for 10-digit telephone numbers instead of character strings.

P-10.52 Implement an `OrderedDict` class, as described in Exercise C-10.49, ensuring that the primary map operations run in  $O(1)$  expected time.

P-10.53 Design a Python class that implements the `skip-list` data structure. Use this class to create a complete implementation of the sorted map ADT.

P-10.54 Extend the previous project by providing a graphical animation of the `skip-list` operations. Visualize how entries move up the `skip list` during insertions and are linked out of the `skip list` during deletions.

P-10.55 Write a spell-checker class that stores a lexicon of words,  $W$ , in a Python set, and implements a method, `check(s)`, which performs a spell check on the string  $s$  with respect to the set of words,  $W$ . If  $s \in W$ , then the call to `check(s)` returns a `list` containing only  $s$ , as it is assumed to be spelled correctly in this case. If  $s$  is not in  $W$ , then the call to `check(s)` returns a `list` of every word in  $W$  that might be a correct spelling of  $s$ . Your program should be able to handle all the common ways that  $s$  might be a misspelling of a word in  $W$ , including swapping adjacent characters in a word, inserting a single character in between two adjacent characters in a word, deleting a single character from a word with another character. For an extra challenge, consider phonetic substitutions as well.

### Chapter Notes

Hashing is a well-studied technique. The reader interested in further study is encouraged to explore the book by Knuth [65], as well as the book by Vitter and Chen [100]. Skip lists were introduced by Pugh [86]. Our analysis of skip lists is a simplification of a presentation given by Motwani and Raghavan [80]. For a more in-depth analysis of skip lists, please see the various research papers on skip lists that have appeared in the data structures literature [59, 81, 84]. Exercise C-10.36 was contributed by James Lee.

## Chapter

### 10 Maps, Hash Tables, and Skip Lists

#### Contents

10.1 Maps and Dictionaries .....	402
10.1.1 The Map ADT .....	403
10.1.2 Application: Counting Word Frequencies .....	405
10.1.3 Python's MutableMapping Abstract Base Class .....	406
10.1.4 Our Map Base Class .....	407
10.1.5 Simple Unsorted Map Implementation .....	408
10.2 Hash Tables .....	410
10.2.1 Hash Functions .....	411
10.2.2 Collision-Handling Schemes .....	412
10.2.3 Load Factors, Rehashing, and Efficiency .....	420
10.2.4 Python Hash Table Implementation .....	422
10.3 Sorted Maps .....	427
10.3.1 Sorted Search Tables .....	428
10.3.2 Two Applications of Sorted Maps .....	434
10.4 Skip Lists .....	437
10.4.1 Search and Update Operations in a Skip List .....	439
10.4.2 Probabilistic Analysis of Skip Lists .....	443
10.5 Sets, Multisets, and Multimaps .....	446
10.5.1 The Set ADT .....	446
10.5.2 Python's MutableSet Abstract Base Class .....	448
10.5.3 Implementing Sets, Multisets, and Multimaps .....	450
10.6 Exercises .....	452

## 408 Chapter 10. Maps, Hash Tables, and Skip Lists

1classMapBase(MutableMapping):

2...Our own abstract base class that includes a nonpublic  
Item class...

3

4#----- nested

Item class -----

5class

Item:

6 ...Lightweight composite to store key-value pairs as map items...

7

slots

=

\_key

,

\_value

89 def

init

(self,k ,v ) :

10 self.

key = k

11 self.

value = v

1213 def

eq

(self,o t h e r ) :

14 return self .

key == other.

key # compare items based on their keys

1516 def

ne

(self,o t h e r ) :

17 return not (self== other) #o p p o s i t e o f

eq

18

19 def

lt

(self,o t h e r ) :

20 return self .

key<other.

key # compare items based on their keys

Code Fragment 10.2: Extending the MutableMapping abstract base class to provide  
a nonpublic

Item class for use in our various map implementations.

### 10.1.5 Simple Unsorted Map Implementation

We demonstrate the use of the MapBase class with a very simple concrete implementation of the map ADT. Code Fragment 10.3 presents an UnsortedTableMap class that relies on storing key-value pairs in arbitrary order within a Python list.

An empty table is initialized as self.

table within the constructor for our map.



#### 414 Chapter 10. Maps, Hash Tables, and Skip Lists

An implementation of a cyclic-shift hash code computation for a character string in Python appears as follows:

```
defhash
code(s):
mask = (1 <<32)-1 # limit to 32-bit integers
h=0
forcharacter ins:
h=( h <<5&m a s k ) |(h>>27) # 5-bit cyclic shift of running sum
h += ord(character) # add in value of next character
return h
```

As with the traditional polynomial hash code, fine-tuning is required when using a cyclic-shift hash code, as we must wisely choose the amount to shift by for each new character. Our choice of a 5 shift amounts (see Table 10.1).

Collisions

Shift

Total

Max

0

234735

623

1

165076

43

2

38471

13

3

7174

5

4

1379

3

5

190

3

6

502

2

7

560

2

8

5546

4

9

393

3

10

5194

### Sets

Although sets and maps have very different public interfaces, they are really quite similar. A set is simply a map in which keys do not have associated values. Any data structure used to implement a set is storing set elements as keys, and using None as an irrelevant value, but such an implementation is unnecessarily wasteful. An efficient set implementation should abandon the

Item composite that we use in our MapBase class and instead store set elements directly in a data structure.

### Multisets

The same element may occur several times in a multiset. All of the data structures we have seen can be reimplemented in which the map key is a (distinct) element of the multiset, and the associated value is a count of the number of occurrences.

Python's standard collections module includes a definition for a class named Counter that is in essence a multiset. Formally, the Counter class is a subclass of dict, with the expectation that values are integers, and with additional functionality like a most

common(n) method that returns a list of the n most common elements.

The standard

iter

reports each element only once (since those are formally the keys of the dictionary). There is another method named elements() that iterates through the multiset with each element being repeated according to its count.

### Multimaps

Although there is no multimap in Python's standard libraries, a common implementation approach is to use a standard dictionary. This uses the standard dict class as the map, and a list of values as a composite value in the dictionary. We have designed the class so that a different map implementation can easily be substituted by overriding the MapType attribute at line 3.

454 Chapter 10. Maps, Hash Tables, and Skip Lists

R-10.24 Give a pseudo-code description of the `delitem`

`map` operation when

using a skip list.

R-10.25 Give a concrete implementation of the `pop` method, in the context of a `MutableSet` abstract base class, that relies only on the core set behaviors described in Section 10.5.2.

R-10.26 Give a concrete implementation of the `isdisjoint` method in the context of the `MutableSet` abstract base class, relying only on the primary abstract methods of that class. Your algorithm should run in  $O(\min(n, m))$  where  $n$  and  $m$  denote the respective cardinalities of the two sets.

R-10.27 What abstraction would you use to manage a database of friends' birthdays in order to support efficient queries such as "find all friends whose birthday is today" and "find the friend who was born on this day"? Creativity

C-10.28 On page 406 of Section 10.1.3, we give an implementation of the `methodsetdefault` as it might appear in `UnsortedTableMap`. While that method accomplishes the goal in a general fashion, its efficiency is less than ideal. In particular, when `getitem`

`getitem`, and then a subsequent insertion via

`setitem`

`setitem`

`setitem`. For a concrete implementation, such as

`theUnsortedTableMap`, this is twice the work because a complete scan

of the table will take place during the failed

`getitem`

`getitem`, and then another

complete scan of the table takes place due to the implementation of

`setitem`

`setitem`. A better solution is for the `UnsortedTableMap` class to override

`setdefault` to provide a direct solution that performs a single search.

Give such an implementation of `UnsortedTableMap.setdefault`.

C-10.29 Repeat Exercise C-10.28 for the `ProbeHashMap` class.

C-10.30 Repeat Exercise C-10.28 for the `ChainHashMap` class.

C-10.31 For an ideal compression function, the capacity of the bucket array for a hash table should be a prime number.

Finding such a prime by using the sieve algorithm. In this algorithm, we

allocate a 2M-cell Boolean array  $A$ , such that cell  $i$  is associated with the

integer  $i$ . We then initialize the array cells to all be `true` and we mark

off all the cells that are multiples of 2, 3, 5, 7, and so on. This process can stop after it reaches a number larger than

$2M$ . (Hint: Consider a

bootstrapping method for finding the primes up to

$2M$ .)

## 15.3. External Searching and B-Trees 711

### 15.3 External Searching and B-Trees

Consider the problem of maintaining a large collection of items that does not fit in main memory, such as a typical database. In this context, we refer to the secondary-memory blocks as disk blocks. The difference between secondary memory and primary memory is a disk transfer. Recalling the great time difference that exists between main memory accesses and disk accesses, the main goal of maintaining such a collection is to minimize the number of disk transfers. The number of disk transfers that count as the I/O complexity of the algorithm involved.

#### Some Inefficient External-Memory Representations

A typical operation we would like to support is the search for a key in a map. If we were to store items unordered, a search for a key within the linked list requires  $n$  transfers in the worst case, since each link hop we perform on the linked list might access a different block of memory.

We can reduce the number of block transfers by using an array-based sequence.

A sequential search of an array can be performed using only  $O(n/B)$  block transfers because of spatial locality of reference, where  $B$  denotes the number of elements that fit into a block. This is because the block transfer when accessing the  $i$ -th element of the array actually retrieves the  $i$ -th block, and so on with each successive block. It is worth noting that the bound of  $O(n/B)$  transfers is only achieved when using a compact array representation (see Section 5.2.2). The standard Python list class is a referential container, and so even though the sequence of references are stored in an array, the actual elements that must be examined during a search are not generally contiguous, resulting in  $O(n)$  transfers in the worst case.

We could alternately store a sequence using a sorted array. In this case, a search performs  $O(\log n)$  transfers.

via binary search, which is a nice improvement. But we do not get significant benefit from block transfers because each query during a binary search is likely in a different block. Block transfers are expensive for a sorted array.

Since these simple implementations are I/O inefficient, we should consider the logarithmic-time internal-memory strategies that use balanced binary trees (for example, AVL trees or red-black trees) or other search structures with logarithmic average-case query and update time. Typically, each node accessed for a query or update in one of these structures will be in a different block. Thus, these methods all require  $O(\log n)$  transfers in the worst case to perform a query or update operation. But we can do better! We can perform map queries and updates using  $O(\log n / \log B)$  transfers.

## xvi Contents

10.2.2 Collision-Handling Schemes . . . . .	417
10.2.3 Load Factors, Rehashing, and Efficiency . . . . .	420
10.2.4 Python Hash Table Implementation . . . . .	422
10.3 Sorted Maps . . . . .	427
10.3.1 Sorted Search Tables . . . . .	428
10.3.2 Two Applications of Sorted Maps . . . . .	434
10.4 Skip Lists . . . . .	437
10.4.1 Search and Update Operations in a Skip List . . . . .	439
10.4.2 Probabilistic Analysis of Skip Lists . . . . .	443
10.5 Sets, Multisets, and Multimaps . . . . .	446
10.5.1 The Set ADT . . . . .	446
10.5.2 Python's Mutable Set Abstract Base Class . . . . .	448
10.5.3 Implementing Sets, Multisets, and Multimaps . . . . .	450
10.6 Exercises . . . . .	452
11 Search Trees . . . . .	459
11.1 Binary Search Trees . . . . .	460
11.1.1 Navigating a Binary Search Tree . . . . .	461
11.1.2 Searches . . . . .	463
11.1.3 Insertions and Deletions . . . . .	465
11.1.4 Python Implementation . . . . .	468
11.1.5 Performance of a Binary Search Tree . . . . .	473
11.2 Balanced Search Trees . . . . .	475
11.2.1 Python Framework for Balancing Search Trees . . . . .	478
11.3 AVL Trees . . . . .	481
11.3.1 Update Operations . . . . .	483
11.3.2 Python Implementation . . . . .	488
11.4 Splay Trees . . . . .	490
11.4.1 Splaying . . . . .	490
11.4.2 When to Splay . . . . .	494
11.4.3 Python Implementation . . . . .	496
11.4.4 Amortized Analysis of Splaying . . . . .	497
11.5 (2,4) Trees . . . . .	502
11.5.1 Multiway Search Trees . . . . .	502
11.5.2 (2,4)-Tree Operations . . . . .	505
11.6 Red-Black Trees . . . . .	512
11.6.1 Red-Black Tree Operations . . . . .	514
11.6.2 Python Implementation . . . . .	525
11.7 Exercises . . . . .	528

## 10.3.1 Sorted Search Tables

Several data structures can efficiently support the sorted map ADT, and we will examine some advanced techniques in Section 10.4 and Chapter 11. In this section, we begin by exploring a simple realization of the sorted map ADT, assuming the keys have a naturally defined order. (See Figure 10.8.) We refer to this implementation of a map as a sorted search table.

9 2 4 5 7 8 12 14 17 19 22 25 27 28 33 5

3 7 0 1 2 3 4 6 7 8 9 1 1 1 1 2 1 3 1 4 1 5

Figure 10.8: Realization of a map by means of a sorted search table. We show only the keys for this map, so as to highlight their ordering.

As was the case with the unsorted table map of Section 10.1.5, the sorted search table has a space requirement that is  $O(n)$ , assuming we grow and shrink the array to keep its size proportional to the number of items in the map. The primary advantage of this representation, and our reason for insisting that it be array-based, is that it allows us to use the binary search algorithm for a variety of efficient operations.

## Binary Search and Inexact Searches

We originally presented the binary search algorithm in Section 4.1.3, as a means for detecting whether a given target is stored within a sorted sequence. In our original presentation (Code Fragment 4.3 on page 156), a binary search function returned

True or False to designate whether the desired target was found. While such an approach could be used to implement the

contains

method of the map ADT, we can adapt the binary search algorithm to provide far more useful information when performing forms of inexact search in support of the sorted map ADT.

The important realization is that while performing a binary search, we can determine the index at or near where a target might be found. During a successful search, the standard implementation determines the precise index at which the target is found. During an unsuccessful search, although the target is not found, the algorithm will effectively determine a pair of indices designating elements of the collection that are just less than and just greater than the target.

As a motivating example, our original simulation from Figure 4.5 on page 156 shows a successful binary search for a target of 22, using the same data we portray in Figure 10.8. Had we instead been searching for a target of 20, we would have determined that the missing target lies in the gap between values 19 and 22 in that example.

## 436 Chapter 10. Maps, Hash Tables, and Skip Lists

### Maintaining a Maxima Set with a Sorted Map

We can store the set of maxima pairs in a sorted map,  $M$ , so that the cost is the key -eld and performance (speed) is the value -eld. We can then implement operations `add(c,p)`, which adds a new cost-performance pair  $(c,p)$ , and `best(c)`, which returns the best pair with cost at most  $c$ , as shown in Code Fragment 10.11.

```
1 class CostPerformanceDatabase:
2     ... Maintain a database of maximal (cost, performance) pairs. ...
3
4     def
5         init
6         (self):
7             ... Create an empty database. ...
8             self.
9             M = SortedTableMap( ) # or a more efficient sorted map
10
11     def best(self, c) :
12         ... Return (cost, performance) pair with largest cost not exceeding c.
13         Return None if there is no such pair.
14
15     def add(self, c, p) :
16         ... Add new entry with cost c and performance p. ...
17         # determine if (c,p) is dominated by an existing pair
18         other = self.
19         M.nd
20         le(c) # other is at least as cheap as c
21         if other is not None and other[1] >= p : # if its performance is as good,
22             return # (c,p) is dominated, so ignore
23         self.
24         M[c] = p # else, add (c,p) to database
25         # and now remove any pairs that are dominated by (c,p)
26         other = self.
27         M.nd
28         gt(c) # other more expensive than c
29         while other is not None and other[1] <= p :
30             del self.
31             M[other[0]]
32             other = self.
33         M.nd
34         gt(c)
```

Code Fragment 10.11: An implementation of a class maintaining a set of maxima cost-performance pairs using a sorted map.

Unfortunately, if we implement `Musing` the `SortedTableMap`, the `add` behavior has  $O(n)$  worst-case running time. If, on the other hand, we implement `Musing` a **skip list**, which we next describe, we can perform `best(c)` queries in  $O(\log n)$  expected time and `add(c,p)` updates in  $O((1+r)\log n)$  expected time, where  $r$  is the number of points removed.

## Bibliography 735

- [63] D. E. Knuth, 'Big omicron and big omega and big theta,' in SIGACT News ,v o l .8 , pp. 18-24, 1976.
- [64] D. E. Knuth, Fundamental Algorithms , vol. 1 of The Art of Computer Programming . Reading, MA: Addison-Wesley, 3rd ed., 1997.
- [65] D. E. Knuth, Sorting and Searching ,v o l .3o f The Art of Computer Programming . Reading, MA: Addison-Wesley, 2nd ed., 1998.
- [66] D. E. Knuth, J. H. Morris, Jr., and V . R. Pratt, 'Fast pattern matching in strings,' SIAM J. Comput. , vol. 6, no. 1, pp. 323-350, 1977.
- [67] J. B. Kruskal, Jr., 'On the shortest spanning subtree of a graph and the traveling salesman problem,' Proc. Amer . Math. Soc. , vol. 7, pp. 48-50, 1956.
- [68] R. Lesuisse, 'Some lessons drawn from the history of the binary search algorithm,' The Computer Journal , vol. 26, pp. 154-163, 1983.
- [69] N. G. Leveson and C. S. Turner, 'An investigation of the Therac-25 accidents,' IEEE Computer , vol. 26, no. 7, pp. 18-41, 1993.
- [70] A. Levitin, 'Do we teach the right algorithm design techniques?,' in 30th ACM SIGCSE Symp. on Computer Science Education , pp. 179-183, 1999.
- [71] B. Liskov and J. Guttag, Abstraction and Speci-cation in Program Development . Cambridge, MA/New York: The MIT Press/McGraw-Hill, 1986.
- [72] M. Lutz, Programming Python . O-Reilly Media, 4th ed., 2011.
- [73] E. M. McCreight, 'A space-economical suf-x tree construction algorithm,' Journal of Algorithms , vol. 23, no. 2, pp. 262-272, 1976.
- [74] C. J. H. McDiarmid and B. A. Reed, 'Building heaps fast,' Journal of Algorithms , vol. 10, no. 3, pp. 352-365, 1989.
- [75] N. Megiddo, 'Linear programming in linear time when the dimension is -xed,' J. ACM , vol. 31, pp. 114-127, 1984.
- [76] K. Mehlhorn, Data Structures and Algorithms 1: Sorting and Searching ,v o l .1 of EATCS Monographs on Theoretical Computer Science . Heidelberg, Germany: Springer-Verlag, 1984.
- [77] K. Mehlhorn, Data Structures and Algorithms 2: Graph Algorithms and NP-Completeness , vol. 2 of EATCS Monographs on Theoretical Computer Science .H e i - delberg, Germany: Springer-Verlag, 1984.
- [78] K. Mehlhorn and A. Tsakalidis, 'Data structures,' in Handbook of Theoretical Computer Science (J. van Leeuwen, ed.), vol. A. Algorithms and Complexity, pp. 301-341, Amsterdam: Elsevier, 1990.
- [79] D. R. Morrison, 'PATRICIA-practical algorithm to retrieve information coded in alphanumeric,' Journal of the ACM , vol. 15, no. 4, pp. 514-534, 1968.
- [80] R. Motwani and P. Raghavan, Randomized Algorithms . New York, NY: Cambridge University Press, 1995.
- [81] T. Papadakis, J. I. Munro, and P. V . Poblete, 'Average search and update costs in skip lists,' BIT, vol. 32, pp. 316-332, 1992.
- [82] L. Perkovic, Introduction to Computing Using Python: An Application Development Fo c u s . Wiley, 2011.
- [83] D. Phillips, Python 3: Object Oriented Programming . Packt Publishing, 2010.
- [84] P. V . Poblete, J. I. Munro, and T. Papadaki s, 'The binomial transform and its application to the analysis of skip lists,' in Proceedings of the European Symposium on Algorithms (ESA) , pp. 554-569, 1995.



Index 743  
Karger, David, 696  
Karp, Richard, 361  
KeyboardInterrupt , 33, 83, 303  
KeyError , 33, 34, 83, 303, 403, 404, 422, 460  
keyword parameter, 27  
Klein, Philip, 696  
Kleinberg, Jon, 580  
Knuth, Donald, 147, 227, 298, 361, 400, 458, 535, 580, 618, 696, 719  
Knuth-Morris-Pratt algorithm, 590-593  
Kosaraju, S. Rao, 696  
Kruskal-s algorithm, 676-684  
Kruskal, Joseph, 696  
L-H -opital-s rule, 731  
Landis, Evgenii, 481, 535  
Langston, Michael, 580  
last-in, -rst-out (LIFO), 229  
lazy evaluation, 39, 80  
LCS, see longest common subsequence  
leaves, 302  
Lecroq, Thierry, 618  
Leiserson, Charles, 535, 696  
len function, 29  
Lesuisse, R., 182  
Letscher, David, 55, 108  
level in a tree, 315  
level numbering, 325, 371  
lexicographic order, 15, 203, 385, 565  
LIFO, 229  
linear exponential, 728  
linear function, 117  
linear probing, 418  
linearity of expectation, 573, 730  
linked **list**, 256-293  
doubly linked, 260, 270-276, 281  
singly linked, 256-260  
linked structure, 317  
LinkedBinaryTree class, 303, 318-324 , 335, 348  
LinkedDeque class, 275-276  
LinkedQueue class, 264-265 , 271, 306, 335  
LinkedStack class, 261-263  
Lins, Rafael, 719  
Liotta, Giuseppe, 361, 696  
Liskov, Barbara, 108, 254, 298 **list**  
of favorites, 286-291  
positional, 277-285  
**list** class, 7, 9, 202-207  
sort method, 23, 569  
**list** comprehension, 43, 207, 209, 221  
literal, 6  
Littman, Michael, 580  
live objects, 700  
load factor, 417, 420-421  
local scope, 23-25, 46, 96

Preface vii

## Contents and Organization

The chapters for this book are organized to provide a pedagogical path that starts with the basics of Python programming and object-oriented design. We then add foundational techniques like algorithm analysis and recursion. In the main portion of the book, we present fundam

1. Python Primer
2. Object-Oriented Programming
3. Algorithm Analysis
4. Recursion
5. Array-Based Sequences
6. Stacks, Queues, and Deques
7. Linked Lists
8. Trees
9. Priority Queues
10. Maps, Hash Tables, and Skip Lists
11. Search Trees
12. Sorting and Selection
13. Text Processing
14. Graph Algorithms
15. Memory Management and B-Trees
- A. Character Strings in Python
- B. Useful Mathematical Facts

A more detailed table of contents follows this preface, beginning on page xi.

## Prerequisites

We assume that the reader is at least vaguely familiar with a high-level programming language, such as C, C++

- Variables and expressions.
- Decision structures (such as if-statements and switch-statements).
- Iteration structures (for loops and while loops).
- Functions (whether stand-alone or object-oriented methods).

For readers who are familiar with these concepts, but not with how they are expressed in Python, we provide a p  
give a comprehensive treatment of Python.

## Contents xv

8.4.1 Preorder and Postorder Traversals of General Trees . . .	328
8.4.2 Breadth-First Tree Traversal .....	330
8.4.3 Inorder Traversal of a Binary Tree .....	331
8.4.4 Implementing Tree Traversals in Python .....	333
8.4.5 Applications of Tree Traversals .....	337
8.4.6 Euler Tours and the Template Method Pattern .....	341
8.5 Case Study: An Expression Tree .....	348
8.6 Exercises .....	352
9 Priority Queues	362
9.1 The Priority Queue Abstract Data Type .....	363
9.1.1 Priorities .....	363
9.1.2 The Priority Queue ADT .....	364
9.2 Implementing a Priority Queue .....	365
9.2.1 The Composition Design Pattern .....	365
9.2.2 Implementation with an Unsorted List .....	366
9.2.3 Implementation with a Sorted List .....	368
9.3 Heaps .....	370
9.3.1 The Heap Data Structure .....	370
9.3.2 Implementing a Priority Queue with a Heap .....	372
9.3.3 Array-Based Representation of a Complete Binary Tree .	376
9.3.4 Python Heap Implementation .....	376
9.3.5 Analysis of a Heap-Based Priority Queue .....	379
9.3.6 Bottom-Up Heap Construction .....	380
9.3.7 Python's <code>heapq</code> Module .....	384
9.4 Sorting with a Priority Queue .....	385
9.4.1 Selection-Sort and Insertion-Sort .....	386
9.4.2 Heap-Sort .....	388
9.5 Adaptable Priority Queues .....	390
9.5.1 Locators .....	390
9.5.2 Implementing an Adaptable Priority Queue .....	391
9.6 Exercises .....	395
10 Maps, Hash Tables, and Skip Lists	401
10.1 Maps and Dictionaries .....	402
10.1.1 The Map ADT .....	403
10.1.2 Application: Counting Word Frequencies .....	405
10.1.3 Python's <code>MutableMapping</code> Abstract Base Class .....	406
10.1.4 Our <code>Map</code> Base Class .....	407
10.1.5 Simple Unsorted Map Implementation .....	408
10.2 Hash Tables .....	410
10.2.1 Hash Functions .....	411

## 10.1 Maps and Dictionaries

Python's dict class is arguably the most significant data structure in the language. It represents an abstraction known as a dictionary in which unique keys are mapped to associated values. Because of the relationship they express between keys and values, dictionaries are commonly known as associative arrays or maps. In this book, we use the term dictionary when specifically discussing Python's dict class, and the term map when discussing the more general notion of the abstract data type. As a simple example, Figure 10.1 illustrates a map from the names of countries to their associated units of currency.

Rupee	Turkey	Spain	China	United States	India	Greece
Lira	Euro	Yuan	Dollar			

Figure 10.1: A map from countries (the keys) to their units of currency (the values).

We note that the keys (the country names) are assumed to be unique, but the values (the currency units) are not necessarily unique. For example, we note that Spain and Greece both use the euro for currency. Maps use an array-like syntax for indexing, such as `currency[Greece]` to access a value associated with a given key.

or `currency[Greece]` to access a value associated with a given key.

to remap it to a new value. Unlike a standard array, indices for a map need not be consecutive nor even numeric. Common applications of maps include the following.

• A university's information system relies on some form of a student ID as a key that is mapped to that student's associated record (such as the student's name, address, and course grades) serving as the value.

• The domain-name system (DNS) maps a host name, such as `www.wiley.com`, to an Internet-Protocol (IP) address, such as `208.215.179.146`.

• A social media site typically relies on a (nonnumeric) username as a key that can be efficiently mapped to a particular user's profile.

• A computer graphics system may map a color name, such as `turquoise`, to the triple of numbers that describes the color's RGB (red-green-blue) representation, such as `(64,224,208)`.

• Python uses a dictionary to represent each namespace, mapping an identifying string, such as `pi`, to an associated object, such as `3.14159`.

In this chapter and the next we demonstrate that a map may be implemented so that a search for a key, and its associated value, can be performed very efficiently, thereby supporting fast lookup.

to the triple of numbers that describes the color's RGB (red-green-blue) representation, such as `(64,224,208)`.

• Python uses a dictionary to represent each namespace, mapping an identifying string, such as `pi`, to an associated object, such as `3.14159`.

In this chapter and the next we demonstrate that a map may be implemented so that a search for a key, and its associated value, can be performed very efficiently, thereby supporting fast lookup.

to the triple of numbers that describes the color's RGB (red-green-blue) representation, such as `(64,224,208)`.

• Python uses a dictionary to represent each namespace, mapping an identifying string, such as `pi`, to an associated object, such as `3.14159`.

#### 404 Chapter 10. Maps, Hash Tables, and Skip Lists

`M.popitem()` :Remove an arbitrary key-value pair from the map, and return a (k,v) tuple representing the removed pair. If map is empty, raise a `KeyError` .

`M.clear()` :Remove all key-value pairs from the map.

`M.keys()` :Return a set-like view of all keys of M.

`M.values()` :Return a set-like view of all values of M.

`M.items()` :Return a set-like view of (k,v) tuples for all entries of M.

`M.update(M2)` :Assign `M[k] = v` for every (k,v) pair in map M2.

`M == M2` :Return True if maps M and M2 have identical key-value associations.

`M != M2` :Return True if maps M and M2 do not have identical key-value associations.

Example 10.1: In the following, we show the effect of a series of operations on an initially empty map storing items with integer keys and single-character values. We use the literal syntax for Python's dictclass to describe the map contents.

Operation

Return Value

Map

`len(M)`

0

`{}`

`M[`

`K`

`]=2`

`.`

`{`

`K`

`:2}`

`M[`

`B`

`]=4`

`.`

`{`

`K`

`:2 ,`

`B`

`:4}`

`M[`

`U`

`]=2`

`.`

`{`

`K`

`:2 ,`

`B`

`:4 ,`

`U`

`:2}`

`M[`

## 10.1.3 Python's MutableMapping Abstract Base Class

Section 2.4.3 provides an introduction to the concept of an abstract base class and the role of such classes in Python's collections module. Methods that are declared to be abstract in such a base class must be implemented by concrete subclasses. However, an abstract base class may provide concrete implementation of other methods that depend upon use of the presumed abstract methods. (This is an example of the template method design pattern.)

The collections module provides two abstract base classes that are relevant to our current discussion: the Mapping and MutableMapping classes. The Mapping class includes all nonmutating methods supported by Python's dict class, while the MutableMapping class extends that to include the mutating methods. What we define as the map ADT in Section 10.1.1 is akin to the MutableMapping abstract base class in Python's collections module.

The significance of these abstract base classes is that they provide a framework to assist in creating a user-defined map class. In particular, the MutableMapping class provides concrete implementations for all behaviors other than the `__getitem__` outlined in Section 10.1.1:

```
__getitem__
,
__setitem__
,
__delitem__
,
__len__
, and
__iter__
```

. As we implement the map abstraction with various data structures, as

long as we provide the `__getitem__` core behaviors, we can inherit all other derived behaviors by simply declaring MutableMapping.

To better understand the MutableMapping class, we provide a few examples of

how concrete behaviors can be derived from the `__getitem__` core abstractions. For example, the `contains`

method, supporting the syntax `key in M`, could be implemented by making a guarded attempt to retrieve `self[k]` to determine if the key exists.

```
def contains( self, k ):
    try:
        self[k] # access via
        __getitem__
        (ignore result)
    return True
    except KeyError:
        return False # attempt failed
```

A similar approach might be used to provide the logic of the `setdefault` method.

```
def setdefault( self, k, d ) :
```

```
    try:
        return self [k] # if
        __getitem__
        succeeds, return value
```

## 10.2 Hash Tables

In this section, we introduce one of the most practical data structures for implementing a map, and the one that is used by Python's own implementation of the dict class. This structure is known as a hash table.

Intuitively, a map  $M$  supports the abstraction of using keys as indices with a syntax such as  $M[k]$ . As a mental warm-up, consider a restricted setting in which a map with  $n$  items uses keys that are known to be integers in a range from 0 to  $N-1$  for some  $N$ . In this case, we can represent the map using a lookup table of length  $N$ , as diagrammed in Figure 10.3.

```
0 1 2 3 4 5 6 7 8 9 10
D Z C Q
```

Figure 10.3: A lookup table with length 11 for a map containing items (1,D), (3,Z), (6,C), and (7,Q).

In this representation, we store the value associated with key  $k$  at index  $k$  of the table (presuming that we have a distinct way to represent an empty slot). Basic map operations of

```
,
getitem
, and
delitem
can be implemented in
O(1) worst-case time.
```

There are two challenges in extending this framework to the more general setting of a map. First, we may not wish to devote an array of length  $N$  if it is the case that  $N \gg n$ . Second, we do not in general require that a map's keys be integers.

The novel concept for a hash table is the use of a hash function to map general keys to corresponding indices in a table. Ideally, keys will be well distributed in the range from 0 to  $N-1$  by a hash function, so that distinct keys that get mapped to the same index. As a result, we will conceptualize our table as a bucket array, as shown in Figure 10.4, in which each bucket may manage a collection of items that are sent to a specific index by the hash function.

(To save space, an empty bucket may be replaced by `None`.)

```
0 1 2 3 4 5 6 7 8 9 10
(1,D) (25,C)
(3,F)
(14,Z) (39,C) (6,A) (7,Q)
```

Figure 10.4: A bucket array of capacity 11 with items (1,D), (25,C), (3,F), (14,Z), (6,A), (39,C), and (7,Q), using a simple hash function.

## Hash Codes

The first action that a hash function performs is to take an arbitrary key  $k$  in our map and compute an integer that is called the hash code for  $k$ ; this integer need not be in the range  $[0, N-1]$ , and may even be negative. We desire that the set of hash codes assigned to our keys should avoid collisions as much as possible. For if the hash codes of our keys cause collisions, then there is no hope for our compression function to avoid them. In this subsection, we begin by discussing the theory of hash codes. Following that, we discuss

## Treating the Bit Representation as an Integer

To begin, we note that, for any data type  $X$  that is represented using at most as many bits as our integer hash codes, we can simply take as a hash code for  $x$  an integer interpretation of its bits. For example, the hash code for key 314 could simply be 314. The hash code for a floating-point number such as 3.14 could be based upon an interpretation of the bits of the floating-point representation as an integer.

For a type whose bit representation is longer than a desired hash code, the above scheme is not immediately applicable. For example, Python relies on 32-bit hash codes. If a floating-point number uses a 64-bit representation, its bits cannot be viewed directly as a hash code. One possibility is to use only the high-order 32 bits (or the low-order 32 bits). This hash code, of course, ignores half of the information present in the original key, and

A better approach is to combine in some way the high-order and low-order portions of a 64-bit key to form a 32-bit hash code, which takes all the original bits into consideration. A simple implementation

binary representation can be viewed as an  $n$ -tuple  $(x_0, x_1, \dots, x_{n-1})$  of 32-bit integers,

for example, by forming a hash code for  $x$  as

$$h(x) = x_0 \oplus x_1 \oplus \dots \oplus x_{n-1},$$

where the  $\oplus$ -symbol represents the bitwise exclusive-or operation (which is `^` in Python).

## Polynomial Hash Codes

The summation and exclusive-or hash codes, described above, are not good choices for character strings or other sequences of characters. For example, consider a 16-bit hash code for a character string  $s$  that sums the Unicode values of the characters in  $s$ . This hash code unfortunately produces lots of unwanted



### Compression Functions

The hash code for a key  $k$  will typically not be suitable for immediate use with a bucket array, because the integer hash code may be negative or may exceed the capacity of the bucket array. Thus, once we have determined an integer hash code for a key object  $k$ , there is still the issue of mapping that integer into the range  $[0, N-1]$ .

This computation, known as a compression function, is the second action performed as part of an overall hash function. A good compression function is one that minimizes the number of collisions for a given set of distinct hash codes.

### The Division Method

A simple compression function is the division method, which maps an integer  $h$  to  $h \bmod N$ ,

where  $N$ , the size of the bucket array, is a fixed positive integer. Additionally, if we take  $N$  to be a prime number, then this compression function helps spread out the distribution of hashed values. Indeed, if  $N$  is not prime, then there is greater risk

that patterns in the distribution of hash codes will be repeated in the distribution of hash values, thereby causing collisions. For example, if we map the hash codes  $\{200, 205, 210, 215, 220, \dots, 600\}$  into a bucket array of size 100, then each hash code will collide with three others. But if we use a bucket array of size 101, then there will be no collisions. If a hash function is chosen well, it should ensure that the probability of two different hash codes colliding is small.

Choosing  $N$  to be a prime number is not always enough, however, for if there is a repeated pattern of hash codes of the form  $pN + q$  for several different  $p$ 's, then there will still be collisions.

### The MAD Method

A more sophisticated compression function, which helps eliminate repeated patterns in a set of integer keys, is the MAD method.

This method maps an integer  $h$  to

$$[(a \cdot h + b) \bmod p] \bmod N,$$

where  $N$  is the size of the bucket array,  $p$  is a prime number larger than  $N$ ,  $a$  and  $b$  are integers chosen at random from the interval  $[0, p-1]$ , with  $a > 0$ . This

compression function is chosen in order to eliminate repeated patterns in the set of hash codes and get us closer to the same as we would have if these keys were thrown into a uniformly at random.

## 418 Chapter 10. Maps, Hash Tables, and Skip Lists

### Open Addressing

The separate chaining rule has many nice properties, such as affording simple implementations of map operations, but it nevertheless has one slight disadvantage:

It requires the use of an auxiliary data structure—a **list**—to hold items with colliding keys. If space is at a premium (as it often is), these structures are employed, but it requires a bit more complexity to deal with collisions. There are several variants of this approach, collectively referred to as open addressing schemes, which we discuss next. Open addressing requires that the load factor is always at most 1 and that items are stored directly in the cells of the bucket array itself.

### Linear Probing and Its Variants

A simple method for collision handling with open addressing is linear probing.

With this approach, if we try to insert an item  $(k, v)$  into a bucket  $A[j]$  that is already occupied, where  $j = h(k)$ , then we next try  $A[(j+1) \bmod N]$ . If  $A[(j+1) \bmod N]$

is also occupied, then we try  $A[(j+2) \bmod N]$ , and so on, until we find an empty

bucket that can accept the new item. Once this bucket is located, we simply in-

sert the item there. Of course, this collision resolution strategy requires that we change the implementation when

getitem

,

setitem

, or

delitem

operations. In particular, to attempt

to locate an item with key equal to  $k$ , we must examine consecutive slots, starting

from  $A[h(k)]$ , until we either find an item with that key or we find an empty bucket.

(See Figure 10.7.) The name “linear probing” comes from the fact that accessing a cell of the bucket array can be

26123456789 1 0 0 New element with

key = 15 to be inserted Must probe 4 times

before finding empty slot

53 7 1 6 2 1 13

Figure 10.7: Insertion into a hash table with integer keys using linear probing. The

hash function is  $h(k) = k \bmod 11$ . Values associated with keys are not shown.

## 10.2.3 Load Factors, Rehashing, and Efficiency

In the hash table schemes described thus far, it is important that the load factor,  $\alpha = n/N$ , be kept below 1. With separate chaining, as  $\alpha$  gets very close to 1, the probability of a collision greatly increases, which adds overhead to our operations, since we must revert to linear-time list-based methods in buckets that have collisions.

Experiments and average-case analyses suggest that we should maintain  $\alpha < 0.9$  for hash tables with separate chaining.

With open addressing, on the other hand, as the load factor  $\alpha$  grows beyond 0.5

and starts approaching 1, clusters of entries in the bucket array start to grow as well. These clusters cause the probability of a collision to increase significantly.

We suggest that we should maintain  $\alpha < 0.5$  for an open addressing scheme with linear

probing, and perhaps only a bit higher for other open addressing schemes (for example, Python's implementation).

If an insertion causes the load factor of a hash table to go above the specified

threshold, then it is common to resize the table (to regain the specified load factor)

and to reinsert all objects into this new table. Although we need not define a new

hash code for each object, we do need to reapply a new compression function that takes into consideration the size of the new bucket array.

When rehashing to a new table, it is a good requirement for the new hash function to

scatter the items throughout the new bucket array. Indeed, if we always double the size of the table with each rehashing operation, then

we can amortize the cost of rehashing all the entries in the table against the time used to insert them in the first place.

## Efficiency of Hash Tables

Although the details of the average-case analysis of hashing are beyond the scope of this book, its probabilistic behavior can be analyzed.

Thus, to store  $n$  entries, the expected number of keys in a bucket would

be  $n/N$ , which is  $O(1)$  if  $n$  is  $O(N)$ .

The costs associated with a periodic rehashing, to resize a table after occasional

insertions or deletions can be accounted for separately, leading to an additional

$O(1)$  amortized cost for

setitem

and

getitem

.

In the worst case, a poor hash function could map every item to the same bucket.

This would result in linear-time performance for the core map operations with separate

chaining, or with any open addressing model in which the secondary sequence of probes depends only on the primary hash value.

Table 10.2.

#### 10.2.4 Python Hash Table Implementation

In this section, we develop two implementations of a hash table, one using separate chaining and the other using open addressing with linear probing. While these approaches to collision resolution are different, we can define a common base class (from Code Fragment 10.2), to define a new HashMapBase class (see Code Fragment 10.4), providing much of the common functionality to our two hash table implementations. The main design of the base class is as follows:

- The bucket array is represented as a Python list, named self.

- The bucket array, with all

- entries initialized to None.

- We maintain an instance variable self.

- n that represents the number of dis-

- tinct items that are currently stored in the hash table.

- If the load factor of the table increases beyond 0.5, we double the size of the table and rehash all items into the new table.

- We define a

- hash

- function utility method that relies on Python's built-in

- hash function to produce hash codes for keys, and a randomized Multiply-

- Add-and-Divide (MAD) formula for the compression function.

What is not implemented in the base class is any notion of how a bucket

should be represented. With separate chaining, each bucket will be an independent structure. With open addressing

In our design, the HashMapBase class presumes the following to be abstract

methods, which must be implemented by each concrete subclass:

- 

- bucket

- getitem(j, k)

This method should search bucket j for an item having key k, returning the associated value, if found, or else raising a KeyError.

- 

- bucket

- setitem(j, k, v)

This method should modify bucket j so that key k becomes associated with value v. If the key already exists, the new value overwrites the existing value.

Otherwise, a new item is inserted and this method is responsible for incrementing self.

- n.

- 

- bucket

- delitem(j, k)

This method should remove the item from bucket j having key k, or raise a KeyError if no such item exists. ( self.

is decremented after this method.)

- 

- iter

This is the standard map method to iterate through all keys of the map. Our base class does not delegate this on

## 424 Chapter 10. Maps, Hash Tables, and Skip Lists

### Separate Chaining

Code Fragment 10.5 provides a concrete implementation of a hash table with separate chaining, in the form of the ChainHashMap class. To represent a single bucket, it relies on an instance of the UnsortedTableMap class from Code Fragment 10.3.

The first three methods in the class use index j to access the potential bucket in the bucket array, and a check for the special case in which that table entry is None.

The only time we need a new bucket structure is when

bucket

setitem is called on

an otherwise empty slot. The remaining functionality relies on map behaviors that are already supported by the inbuilt dict class. A small bit of forethought to determine whether the application of

setitem

on the chain

causes a net increase in the size of the map (that is, whether the given key is new).

1 class ChainHashMap(HashMapBase):

2 """Hash map implemented with separate chaining for collision resolution. ...

34 def

bucket

getitem( self, j, k ) :

5 bucket = self.

table[j]

6 if bucket is None :

7 raise KeyError(

Key Error:

+ repr( k ) ) # no match found

8 return bucket[k] # may raise KeyError

9

10 def

bucket

setitem( self, j, k, v ) :

11 if self.

table[j] is None :

12 self.

table[j] = UnsortedTableMap( ) # bucket is new to the table

13 oldsize = len( self.

table[j])

14 self.

table[j][k] = v

15 if len(self.

table[j]) > oldsize: # key was new to the table

16 self.

n += 1 # increase overall map size

17 18 def

bucket

delitem( self, j, k ) :

19 bucket = self.

table[j]

20 if bucket is None :

21 raise KeyError(

426 Chapter 10. Maps, Hash Tables, and Skip Lists

```
26def
bucket
    getitem( self,j ,k ) :
27 found, s = self.
    .nd
    slot(j, k)
28 if not found:
29 raiseKeyError(
    Key Error:
    +r e p r ( k ) ) # no match found
30 return self .
table[s].
    value
31
32def
bucket
    setitem( self,j ,k ,v ) :
33 found, s = self.
    .nd
    slot(j, k)
34 if not found:
35 self.
    table[s] = self.
    Item(k,v) # insert new item
36 self.
    n+ =1 # size has increased
37 else:
38 self.
    table[s].
        value = v # overwrite existing
3940def
bucket
    delitem( self,j ,k ) :
41 found, s = self.
    .nd
    slot(j, k)
42 if not found:
43 raiseKeyError(
    Key Error:
    +r e p r ( k ) ) # no match found
44 self.
    table[s] = ProbeHashMap.
    AVAIL # mark as vacated
4546def
iter
(self):
47 forjinrange(len( self.
table)): # scan entire table
48 if not self .
```

#### 430 Chapter 10. Maps, Hash Tables, and Skip Lists

```

1class SortedTableMap(MapBase):
2    ...Map implementation using a sorted table...
3
4    #----- nonpublic behaviors -----
5    def
6        .nd
7        index(self, k, low, high):
8            ...Return index of the leftmost item with key greater than or equal to k.
9            Return high + 1 if no such item qualifies.
10           That is, j will be returned such that:
11           all items of slice table[low:j] have key < k
12           all items of slice table[j:high+1] have key >= k
13           ...
14           if high < low:
15               return high + 1 # no element qualifies
16           else:
17               mid = (low + high) // 2
18               if k == self.table[mid].key:
19                   return mid # found exact match
20               elif k < self.table[mid].key:
21                   return self.index(k, low, mid - 1) # Note: may return mid
22           else:
23               return self.index(k, mid + 1, high) # answer is right of mid
24
25    #----- public behaviors -----
26    def
27        .init
28        (self):
29            ...Create an empty map...
30            self.table = []
31
32    def
33        .len
34        (self):
35            ...Return number of items in the map...
36            return len(self.table)
37
38    def
39        .getitem
40        (self, k):
41            ...Return value associated with key k (raise KeyError if not found)...

```

432 Chapter 10. Maps, Hash Tables, and Skip Lists

```
78def nd
ge(self,k) :
79 ...Return (key,value) pair with least key greater than or equal to k...
80 j=self.
.nd
index(k, 0, len( self.
table) - 1) #j
sk e y>=k
81 if j<len(self.
table):
82 return (self.
table[j].
key,self.
table[j].
value)
83 else:
84 return None
85
86def nd
lt(self,k) :
87 ...Return (key,value) pair with greatest key strictly less than k...
88 j=self.
.nd
index(k, 0, len( self.
table) - 1) #j
sk e y>=k
89 if j>0:
90 return (self.
table[j - 1].
key,self.
table[j - 1].
value) # Note use of j-1
91 else:
92 return None
9394def nd
gt(self,k) :
95 ...Return (key,value) pair with least key strictly greater than k...
96 j=self.
.nd
index(k, 0, len( self.
table) - 1) #j
sk e y>=k
97 if j<len(self.
table) and self .
table[j].
key == k:
98 j+=1 # advanced past match
99 if j<len(self.
table):
```



## 10.3.2 Two Applications of Sorted Maps

In this section, we explore applications in which there is particular advantage to using a sorted map rather than a traditional (unsorted) map. To apply a sorted map, keys must come from a domain that is totally ordered. Furthermore, to take advantage of the inexact or range searches afforded by a sorted map, there should be some reason why nearby keys have relevance to a search.

## Flight Databases

There are several Web sites on the Internet that allow users to perform queries on flight databases to find flights between two cities and to buy a ticket. To make a query, a user specifies origin and destination cities, a departure date, and a departure time. To support such queries, we can model the flight database as a map, where keys are tuples of these four parameters. That is, a key is a tuple

$k = (\text{origin}, \text{destination}, \text{date}, \text{time})$ .

Additional information about a flight, such as the flight number, the number of seats still available in first (F) and coach (C) class, can be stored in the value object.

Finding a requested flight is not simply a matter of finding an exact match

for a requested query. Although a user typically wants to exactly match the origin and destination cities, he or she may have flexibility for the departure date, and certainly will have some flexibility for the departure time.

We can handle such a query by ordering our keys lexicographically. Then, an efficient implementation for a sorted map would be a good way to satisfy users' queries. For instance, given a user query  $q$ , we can find the first flight between the desired cities, having a departure date and time matching the desired query or later. Because the keys are ordered lexicographically, we can use `range(k1, k2)` to find all flights within a given range of times. For example, if  $k1 = (\text{ORD}, \text{PVD}, \text{05May}, \text{09:30})$ , and  $k2 = (\text{ORD}, \text{PVD}, \text{05May}, \text{20:00})$ , a respective call to `range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

`range(k1, k2)` might result in the following sequence of key-value pairs:

## 10.5 Sets, Multisets, and Multimaps

We conclude this chapter by examining several additional abstractions that are closely related to the map ADT, and that can be implemented using data structures similar to those for a map.

- A set is an unordered collection of elements, without duplicates, that typically supports efficient membership tests. In essence, elements of a set are like keys of a map, but without any associated values.
- A multiset (also known as a bag) is a set-like container that allows duplicates.
- A multimap is similar to a traditional map, in that it associates values with keys; however, in a multimap the same key can be mapped to multiple values. For example, the index of this book maps a given term to one or more locations at which the term occurs elsewhere in the book.

### 10.5.1 The Set ADT

Python provides support for representing the mathematical notion of a set through the built-in classes `frozenset` and `set`, with `frozenset` being an immutable form. Both of those classes are implemented using hash tables in Python.

Python's `collections` module defines abstract base classes that essentially mirror these built-in classes. Although the choice of names is counterintuitive, the abstract base class `collections.Set` mirrors the behavior of the built-in `set` class. `collections.MutableSet` is akin to the concrete `set` class.

In our own discussion, we equate the "set ADT" with the behavior of the built-in `set` class (and thus, the `collections.MutableSet` base class). We begin by listing what we consider to be the "most fundamental behaviors for a set `S`":

`S.add(e)`: Add element `e` to the set. This has no effect if the set already contains `e`.

`S.discard(e)`: Remove element `e` from the set, if present. This has no effect if the set does not contain `e`.

`e in S`: Return `True` if the set contains element `e`. In Python, this is implemented with the special `__contains__` method.

`len(S)`: Return the number of elements in set `S`. In Python, this is implemented with the special method `__len__`.

`iter(S)`: Generate an iteration of all elements of the set. In Python, this is implemented with the special method `__iter__`.

## 10.5.2 Python's MutableSet Abstract Base Class

To aid in the creation of user-defined set classes, Python's collections module provides a MutableSet abstract base class (just as it provides the MutableMapping abstract base class discussed in Section 10.1.3). The MutableSet base class provides concrete implementations for all methods described in Section 10.5.1, except for the core behaviors (add, discard, contains

```
,
len
, and
iter
) that must
```

be implemented by any concrete subclass. This design is an example of what is known as the template method pattern, as the concrete methods of the MutableSet class rely on the presumed abstract methods that will subsequently be provided by a subclass.

For the purpose of illustration, we examine algorithms for implementing several of the derived methods of the MutableSet base class. For example, to determine if one set is a proper subset of another, we must verify two conditions: a proper subset must have size strictly smaller than that of its superset, and each element of a subset must be contained in the superset. An implementation of the corresponding

method based on this logic is given in Code Fragment 10.14.

```
def
It
(self, other) : # supports syntax S < T
...Return true if this set is a proper subset of other.
...if len(self) >= len(other):
return False # proper subset must have strictly smaller size
for element in self :
    if element not in other:
return False # not a subset since element missing from other
return True # success; all conditions are met
```

Code Fragment 10.14: A possible implementation of the MutableSet.

It  
method, which tests if one set is a proper subset of another.

As another example, we consider the computation of the union of two sets.

The set ADT includes two forms for computing a union. The syntax  $S \mid T$  should produce a new set that has contents equal to the union of existing sets  $S$  and  $T$ . This operation is implemented through the special method

or  
in Python. Another

syntax,  $S \mid= T$  is used to update existing set  $S$  to become the union of itself and set  $T$ . Therefore, all elements of  $T$  that are not already contained in  $S$  should be added to  $S$ . We note that this 'in-place' operation may be implemented more efficiently than if we were to rely on the first form, using the syntax  $S = S \mid T$ , in

which identifier  $S$  is reassigned to a new set instance that represents the union. For

convenience, Python's built-in set class supports named version of these behaviors, with  $S.union(T)$  equivalent to those named versions are not formally provided by the MutableSet abstract base class).

## 452 Chapter 10. Maps, Hash Tables, and Skip Lists

### 10.6 Exercises

For help with exercises, please visit the site, [www.wiley.com/college/goodrich](http://www.wiley.com/college/goodrich).

#### Reinforcement

R-10.1 Give a concrete implementation of the `pop` method in the context of the `MutableMapping` class, relying only on the `__delitem__` primary abstract methods of that class.

R-10.2 Give a concrete implementation of the `items()` method in the context of the `MutableMapping` class, relying only on the `__getitem__` primary abstract methods of that class. What would its running time be if directly applied to the `UnsortedTableMap` subclass?

R-10.3 Give a concrete implementation of the `items()` method directly within the `UnsortedTableMap` class, ensuring that the entire iteration runs in  $O(n)$  time.

R-10.4 What is the worst-case running time for inserting  $n$  key-value pairs into an initially empty map  $M$  that is implemented with the `UnsortedTableMap` class?

R-10.5 Reimplement the `UnsortedTableMap` class from Section 10.1.5, using the `PositionalList` class from Section 7.4 rather than a Python `list`.

R-10.6 Which of the hash table collision-handling schemes could tolerate a load factor above 1 and which could not?

R-10.7 Our `Position` classes for `lists` and trees support the `eq` method so that

two distinct position instances are considered equivalent if they refer to the same underlying node in a structure. Implement a `hash`

method that is consistent with this notion of equivalence. Provide such a `hash` method.

R-10.8 What would be a good hash code for a vehicle identification number that is a string of numbers and letters?

R-10.9 Draw the 11-entry hash table that results from using the hash function,  $h(i) = (3i + 5) \bmod 11$ , to hash the keys 16, and 5, assuming collisions are handled by chaining.

R-10.10 What is the result of the previous exercise, assuming collisions are handled by linear probing?

R-10.11 Show the result of Exercise R-10.9, assuming collisions are handled by quadratic probing, up to the point where the method fails.

## 11.7. Exercises 533

C-11.56 The standard splaying step requires two passes, one downward pass to find the node  $x$  to splay, followed by an upward pass to splay the node

$x$ . Describe a method for splaying and searching for  $x$  in one downward pass. Each substep now requires that you consider the next two nodes in the path down to  $x$ , with a possible zig substep performed at the end.

Describe how to perform the zig-zig, zig-zag, and zig steps.

C-11.57 Consider a variation of splay trees, called half-splay trees, where splaying a node at depth  $d$  stops as soon as the node reaches depth  $\lceil d/2 \rceil$ . Perform an amortized analysis of half-splay trees.

C-11.58 Describe a sequence of accesses to an  $n$ -node splay tree  $T$ , where  $n$  is odd, that results in  $T$  consisting of a single chain of nodes such that the path down  $T$  alternates between left children and right children.

C-11.59 As a positional structure, our TreeMap implementation has a subtle flaw.

A position instance associated with an key-value pair  $(k, v)$  should remain valid as long as that item remains in the map. In particular, that position should be unaffected by calls to insert. However, our implementation may fail to provide such a guarantee, in particular because of our rule for using the inorder predecessor of a key as a replacement when deleting a key that is located in a node with two children.

C-11.60 How might the TreeMap implementation be changed to avoid the flaw described in the previous problem?

### Projects

P-11.61 Perform an experimental study to compare the speed of our AVL tree, splay tree, and red-black tree implementations for various sequences of operations.

P-11.62 Redo the previous exercise, including an implementation of skip lists. (See Exercise P-10.53.)

P-11.63 Implement the Map ADT using a  $(2,4)$  tree. (See Section 10.1.1.)

P-11.64 Redo the previous exercise, including all methods of the Sorted Map ADT. (See Section 10.3.)

P-11.65 Redo Exercise P-11.63 providing positional support, as we did for binary search trees (Section 11.1.1), so as to include methods `first()`, `last()`, `before(p)`, `after(p)`, and

`position(k)`. Each item should have a dis-

tinct position in this abstraction, even though several items may be stored at a single node of a tree.

## 734 Bibliography

- [41] E. Gamma, R. Helm, R. Johnson, and J. Vlissides, *Design Patterns: Elements of Reusable Object-Oriented Software* . Reading, MA: Addison-Wesley, 1995.
- [42] A. Goldberg and D. Robson, *Smalltalk-80: The Language* . Reading, MA: Addison-Wesley, 1989.
- [43] M. H. Goldwasser and D. Letscher, *Object-Oriented Programming in Python* . Upper Saddle River, NJ: Prentice Hall, 2008.
- [44] G. H. Gonnet and R. Baeza-Yates, *Handbook of Algorithms and Data Structures in Pascal and C* . Reading, MA: Addison-Wesley, 1991.
- [45] G. H. Gonnet and J. I. Munro, 'Heaps on heaps,' *SIAM J. Comput.* , vol. 15, no. 4, pp. 964-971, 1986.
- [46] M. T. Goodrich, J.-J. Tsay, D. E. Vengroff, and J. S. Vitter, 'External-memory computational geometry,' in *Proc. 34th Annu. IEEE Sympos. Found. Comput. Sci.* , pp. 714-723, 1993.
- [47] R. L. Graham and P. Hell, 'On the history of the minimum spanning tree problem,' *Annals of the History of Computing* , vol. 7, no. 1, pp. 43-57, 1985.
- [48] L. J. Guibas and R. Sedgwick, 'A dichromatic framework for balanced trees,' in *Proc. 19th Annu. IEEE Sympos. Found. Comput. Sci.* , *Lecture Notes Comput. Sci.*, pp. 8-21, Springer-Verlag, 1978.
- [49] Y. Gurevich, 'What does  $O(n)$  mean?,' *SIGACT News* , vol. 17, no. 4, pp. 61-63, 1986.
- [50] J. Hennessy and D. Patterson, *Computer Architecture: A Quantitative Approach* . San Francisco: Morgan Kaufmann, 2nd ed., 1996.
- [51] C. A. R. Hoare, 'Quicksort,' *The Computer Journal* , vol. 5, pp. 10-15, 1962.
- [52] J. E. Hopcroft and R. E. Tarjan, 'Efficient algorithms for graph manipulation,' *Communications of the ACM* , vol. 16, no. 6, pp. 372-378, 1973.
- [53] B.-C. Huang and M. Langston, 'Practical in-place merging,' *Communications of the ACM* , vol. 31, no. 3, pp. 348-352, 1988.
- [54] J. Jája, *An Introduction to Parallel Algorithms* . Reading, MA: Addison-Wesley, 1992.
- [55] V. Jarník, 'O jistém problému minimálního,' *Praca Moravske Prirodovedecké Společnosti* , vol. 6, pp. 57-63, 1930. (in Czech).
- [56] R. Jones and R. Lins, *Garbage Collection: Algorithms for Automatic Dynamic Memory Management* . John Wiley and Sons, 1996.
- [57] D. R. Karger, P. Klein, and R. E. Tarjan, 'A randomized linear-time algorithm to find minimum spanning trees,' *Journal of the ACM* , vol. 42, pp. 321-328, 1995.
- [58] R. M. Karp and V. Ramachandran, 'Parallel algorithms for shared memory machines,' in *Handbook of Theoretical Computer Science* (J. van Leeuwen, ed.), pp. 869-941, Amsterdam: Elsevier/The MIT Press, 1990.
- [59] P. Kirschenhofer and H. Prodinger, 'The path length of random skip lists,' *Acta Informatica* , vol. 31, pp. 775-792, 1994.
- [60] J. Kleinberg and É. Tardos, *Algorithm Design* . Reading, MA: Addison-Wesley, 2006.
- [61] A. Klink and J. Wälde, 'Efficient denial of service attacks on web application platforms,' 2011.
- [62] D. E. Knuth, *Sorting and Searching* , vol. 3 of *The Art of Computer Programming* . Reading, MA: Addison-Wesley, 1973.