

# CSCI 230 — Data Structures and Algorithms

## Project 3 — Jog Phone Company

March 17, 1999

### **Submission Deadline: Wednesday, April 7 at 11:59:59 pm**

Many applications require fast storage and retrieval of data of some type `T` based on *keys* of some other type `Key`. If the set of possible keys is a reasonably small range of integers then the table could simply be an array indexed by the keys. Often though the number of possible keys is much too large for an array, as for example when the keys are people's names or social security numbers. (Since social security numbers are 9 digits long, the array would have to have room for one billion entries.) So instead we must store both keys and their associated data values, in an abstract data type called a `table` or `dictionary`. When we need to retrieve the value associated with a given key  $k$ , we have to search for  $k$  among the stored keys. We want to arrange the table so that such searches can be done quickly, and we assume that it's also required that inserting or deleting (key, data) pairs must be efficient operations also. The first solution we might think of is an array or vector of (key, data) pairs, sorted based on the keys, so that keys can be looked up quickly using a binary search. (Why is a vector to be preferred over a plain array?) In order to keep the table sorted, a new entry must be made at the correct place in the table (we can't just put it at the beginning or end). In an array, such insertions take linear time (time proportional to the size of the table), which can be very slow for a large table. Thus, using an array or vector is satisfactory only if the number of insertions and deletions is very small compared to the number of lookup operations.

The STL library provides *sorted associative containers*, which are a solution to this problem. With these containers, insertion, deletion, and lookup are all fast operations—they work in time proportional to the *logarithm* of the size of the container. These containers are implemented with *balanced binary search trees*, which, as we've seen in Chapter 4 of the textbook, does support these operations with a logarithmic bound on their computing time. In this project, you are to write a couple of short programs that construct tables of (key, data) pairs. One program should use of the STL associative containers, `map`. However, it's possible to do better than logarithmic time for insertion, deletion, and lookup, on average, by using a *hash table*, which is another type of associative container. Hash tables have average case *constant* time bounds on insertion, deletion, and lookup. Their disadvantage is that their worst case behavior can be very bad compared to balanced binary search trees—linear time instead of logarithmic. Achieving the optimal constant time behavior requires careful choice of the hash function and other parameters. In this project you will experiment with these factors and make some measurements of the performance of your programs using timing functions that are provided in another of the standard C/C++ libraries, `time.h`.

### **An application of associative containers**

The application to be programmed and timed is as follows. The Jog Phone Company keeps track of phone connections in a table. The key of an entry in the table is the calling phone number and the associated data is the called number and the time the call started.

Insertions and deletions in this table are made in response to *call events*. A call event is either a connection or a disconnection.

The program operates in either a verbose mode, in which it outputs a commentary on the events and its actions, or a silent mode, in which it only prints overall statistics at the end of its operation.

In verbose mode, the program's response to a connection event is either "*n2* is busy" or "At time *t* connected call from *n1* to *n2*", where *n1* is the calling number, *n2* is the called number and *t* is the time at which the event was processed. The "busy" response occurs if there is already a connection involving the called number. Otherwise, in either mode, an entry is made in the table with *n1* as the key and (*n2*, *t*) as the associated data, and another entry with *n2* as the key and (*n1*, *t*) as the associated data.

The response to a disconnection event with number *n* is "At time *t* disconnected call between *n1* and *n2*, duration: *d*" where *n1* or *n2* is *n*, *t* is the time at which the event was processed, and *d* is the difference between *t* and the start time of the connection. There is no action taken if there is no call in progress with *n1* or *n2* equal to *n* (i.e., the event is ignored). But if there is such a call in progress, it is removed from the table by deleting both of the entries involving phone number *n*.

## Also a discrete event simulation

The program you are to write (each version) performs a simulation of the phone company's operation in handling a specified number of call events and entering them or deleting them from the table. The type of simulation it should do is called a *discrete event simulation*. (See also Weiss, Section 6.4.3.) Such simulations are an important application of priority queues. A discrete event simulation consists of processing of events, such as the connection and disconnection events of the phone company. Events, including the times at which they occur, might be generated randomly. These times do not have to be exact time-of-day times; instead we can just use multiples of some quantum unit, called a tick. Similarly, times between events can also be generated randomly. A simple form of simulation would be to start the clock at zero ticks, then advance it one tick at a time, checking to see if there is an event. This is a *discrete time-driven simulation*. The problem with a discrete time-driven simulation is that it can waste a lot of computing time just advancing the clock tick by tick to the next event. The problem is worse when the events occur infrequently compared to the granularity of the clock. A better way: after completing each event, advance the clock instantly to the time for the next event. This is called a discrete event-driven simulation. Since we need to find the nearest event in the future among those waiting to be processed, a priority queue is the natural data structure to use to hold the events.

## Program parameters and output

The program should take four command line parameters. The first parameter *N* indicates the number of events to generate, the second parameter *p* is a probability that determines the rate at which connection events occur, the third parameter is a seed for the random number generator, and the fourth is a flag that controls whether the program operates in verbose or silent mode. The value `log` means verbose mode, while `no log` means silent mode.

A random connection event is generated by using a random number generator to generate two 7-digit numbers and a start time for the call. This information is packaged in the event and it is placed in the event queue, which is set up as a priority queue in which earliest (smallest) times have the highest priority. When the connection event is later removed from the queue, a disconnection event is constructed for it with its time set to the current time plus a randomly generated duration, and this event is entered in the queue.

After the prescribed number of events has been processed, the program should display (on `cout`) a line that tells the number of calls completed and the average duration of the

completed calls and the average number of active connections (calls that are still in progress are ignored in this summary). Finally, it should display the amount of time taken to process all of the events. This (real, not simulated) time can be computed using the function `clock` from `time.h`:

```
clock_t start, duration;
...
start = clock();
... code to be timed
duration = (double)(clock() - start)/CLOCKS_PER_SEC;
cout << "Time used: " << duration << " seconds\n";
```

The type `clock_t` and the factor `CLOCKS_PER_SEC` are also defined in `time.h`.

## Two versions required: maps versus hash tables

Write two versions of the program, one that keeps the connected calls in an STL `map`, and the other that keeps them in a hash table that you implement. Represent the phone numbers as strings (C++ strings, using the `<string>` header), and the times or time differences as floats. Thus the `map` should be declared as

```
map<string, pair<string, float> >
```

That is, it associates a pair consisting of a `string` and a `float` with each key of type `string`, and the ordering of the keys is computed using the usual `<` operator on strings. In the `map`, each entry is of the form `pair(n1, pair(n2, t))`.

For the hash table version, you must also implement a hash table class. In

</dept/cs/cs230/public/projects/project3><sup>1</sup>

there is already a hash table class, implemented using separate chaining. (It is similar to the hash table class in the Chapter 5 online notes, but the interface is extended and modified some to make it more similar to the the STL `map` interface.) You must implement your hash table with *exactly the same public interface*, but using *open addressing with quadratic probing* instead of separate chaining. (See Section 5.4.2 in the textbook, where some of the implementation is given.)

As mentioned, the public interface of this hash table class is designed to be very similar to that of the STL `map` class—more precisely, it is very similar to the partial specialization

```
template <class ElementType> map<string, ElementType>
```

i.e., the keys must be strings. This similarity means that very few changes to your phone company simulation program should be required to change it from the `map` version to the hash table version. For simplicity, the given hash table class omits some features of the `map` class; for example, it defines `iterator` and `const_iterator` types but doesn't provide them with all of the usual iterator operations. This isn't a problem for the phone company simulation, because it doesn't need full-blown iterators. Note that GNU C++ does provide a set of hashed associative containers (`hash_set`, `hash_multiset`, `hash_map`, `hash_multimap`) that have full-blown iterators and are as fully compatible with the sorted associated containers as possible. (These containers were actually developed by the generic programming project headed by Alexander Stepanov at Silicon Graphics; they are not part of the C++ standard library but are part of SGI's version of STL. For purposes of this project you are not permitted to use any of SGI's hash tables. In fact, they don't meet the project specification, because they are implemented with separate chaining.)

---

<sup>1</sup></dept/cs/cs230/public/projects/project3>

## Performance comparisons

Since insertions and deletions in the connections table are a large part of the computation in the the phone company simulation, it is a good benchmark for comparing the performance of hashed associative containers versus sorted associative containers (hash tables versus balanced binary search trees). The hash table version of your program should be somewhat faster than the map version. Is it? How much faster (or slower)? Include in the README file you submit with your project a brief report on the performance results you obtained. You should report the times for the two versions on several (silent) runs with at least the values 1000, 10000, 20000, 40000, 80000 for the number of connections.

## Examples of operation

To get a better understanding of what the Job Phone Company simulation is supposed to do, look at some sample runs.<sup>2</sup>

You can also find two compiled versions of the program (both using a map) in

`/dept/cs/cs230/public/projects/project33`

called `jog` (SUN version) or `jog.exe` (Windows PC version). Try them out with various command-line parameters to give you still more information about what is expected.

## Additional notes

Some additional notes and advice are available.<sup>4</sup>

---

<sup>2</sup><http://www.rpi.edu/dept/cs/cs230/www/projdir/project3/examples.html>

<sup>3</sup>[/dept/cs/cs230/public/projects/project3](http://dept/cs/cs230/public/projects/project3)

<sup>4</sup><http://www.rpi.edu/dept/cs/cs230/www/projdir/project3/notes.html>