

CSCI 6620 Spring 2016
Program 4: A Radix Sort Using Queues

1 Goals.

- To implement and use a composite data structure: an array of queues.
- To build your own queue class, not use the C++ STL queue.
- To master an excellent sorting algorithm that is at least half a century old. Wikipedia does not even give its history.
- To gain some skill with manipulating linked lists.

2 Overview.

A Radix sort¹ is the fastest sort for a vary large number of data items². It is an $O(n)$ sort that works by using groups of bits from the sort keys rather than by comparing key values to each other. It is the fastest way to sort a large number of items. However, it *does* involve considerable setup time and data structures, and works only with fixed-length sort keys. A data item could be any type of object (or pointer to object) that has one data member that acts as a sort key.

The Radix sort described here is applicable to any large data set in which the key is an unsigned number, or can be treated like an unsigned number. However, the presentation given here is specific to the data set you will be using.

2.1 The Data to Sort

- We will be sorting objects with two data fields: (1) the time (type `time_t`), a 4-byte unsigned integer. (2) the temperature, a 4-byte float.
- The data file will consist of a very large number of time/temperature measurements, recorded in January at an apartment in Tuckahoe, New York, arranged by time. Here is the first line from one of the files: 2016-01-01 00:00:04.039251 1451624404 01948 4.9

As you can see, the fields are (1) date, (2) time, (3) date and time again, as a Unix `time_T` value, (4) device number, and (5) temperature in Celsius. So there are five fields on each line, and you only need two of those fields: (3) the `time_t` value and (5) the temperature. When you read the file, you only need to save those two fields in your data structure.

Further, to sort the temperatures, the Celsius temperature should be converted to Kelvin by adding 273.15 to the value in the file. Store the Kelvin temperature in your data objects. This eliminates negative numbers. You will need to subtract the constant before outputting the answers.

- There will be three files, one from each of three sensors on different sides of the building. You only need to process one of the files, but different people should choose different files. Handle file opening and EOF properly.

2.2 Overview of the Sort

The algorithm divides the sort key into k fields, based on its binary representation and makes k passes through the entire data list. On pass 0, the least significant (rightmost) field of the key is isolated and used to assign the data item to a bucket-queue. On pass p , the p th field is used. At the end of each

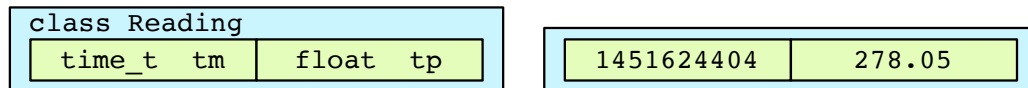
¹This technique has been in use since the days of punched cards and card-handling machines.

²Google's map-reduce is based on a Radix sort.

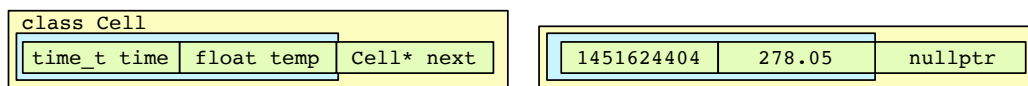
pass, all the data is collected again into a single queue, ready for re-distribution. After k passes, the data is sorted.

This program will make 4 passes, each time using 256 buckets.

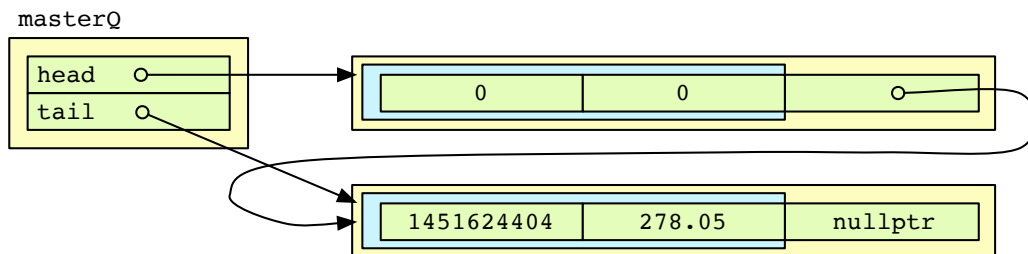
- First, all of the data must be read from a file and stored in a Reading object. The diagram shows the structure of the class and an actual reading object after initialization.



Then each reading must be installed in a new cell. Again, there are two diagrams. The first is the Cell structure, the second of an initialized Cell object.



Then each cell must each be linked onto the end of a queue called the master queue. The diagram shows the master queue after inserting the first data cell. (Remember – the queue has a dummy header cell).



- Finally, data cells will be distributed into buckets. For this purpose, an array of 256 buckets must be created. Each bucket is a linked queue, initially empty.
- When the data and buckets are ready, sorting begins. The process will work in four passes (numbered 0, 1, 2, and 3), where each pass consists of a distribute operation (next slide) and a collect operation.
- After distributing all the data into the buckets. the queues from the buckets must be collected and joined back into a single master queue.
- At the end of four passes, the data will again be in the master queue, and will be sorted. Output the sorted data to a file.
- The number of fields in the key determines the number of passes.
- The number of bits in each field determines the number of buckets needed: for k bits, we need 2^k buckets.
- The algorithm is simpler to implement if each field is a half-byte, a byte, or two bytes, resulting in 16 or 256 or 35536 queues.

Illustration of a radix sort process. In these diagrams, we show a bucket sort of 14 short integers, written here in hex. We make 4 passes, sorting on one nybble each time. The original data list is shown on the left:

| Data: | BUCKETS for distribution pass 0: (rightmost hex digit) | | | | | | | | | | | | | | | | Collect 0 | Collect 1 | Collect 2 | Collect3 |
|-------|--|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|-----------|-----------|-----------|----------|
| 9123 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | Distr.1 | Distr.2 | Distr.3 | Sorted |
| 438B | FA10 | 4341 | | 9123 | | 84C5 | | C437 | 63A8 | | | 438B | 973C | A18D | 245E | | FA10 | F00D | F00D | 1743 |
| 1743 | | | | 1743 | | | | | | | | | | F00D | | | 4341 | FA10 | 9123 | 245E |
| C437 | | | | | | | | | | | | | | BEAD | | | 9123 | 9123 | A18D | 4341 |
| A18D | | | | | | | | | | | | | | DEAD | | | 1743 | C437 | 4341 | 438B |
| F00D | BUCKETS for distribution pass 1: (thrid hex digit) | | | | | | | | | | | | | | | | 84C5 | 4341 | 438B | 63A8 |
| BEAD | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | C437 | 1743 | 63A8 | 84C5 |
| FA10 | F00D | FA10 | 9123 | C437 | 4341 | 245E | | | 438B | | 63A8 | A18D | 84C5 | | | | 63A8 | 245E | C437 | 9123 |
| 245E | | | | | 1743 | | | | | | BEAD | | 973C | | | | 438B | 438B | 245E | 973C |
| 63A8 | | | | | | | | | | | DEAD | | | | | | 973C | 63A8 | 84C5 | A18D |
| DEAD | BUCKETS for distribution pass 2: (second hex digit) | | | | | | | | | | | | | | | | A18D | BEAD | 1743 | BEAD |
| 84C5 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | F00D | DEAD | 973C | C437 |
| 973C | F00D | 9123 | | 4341 | C437 | | | 1743 | | | FA10 | | | | BEAD | | BEAD | A18D | FA10 | DEAD |
| 4341 | | A18D | | 438B | 245E | | | 973C | | | | | | | DEAD | | DEAD | 84C5 | BEAD | F00D |
| | | | | 63A8 | 84C5 | | | | | | | | | | | | 245E | 973C | DEAD | FA10 |
| | BUCKETS for distribution pass 3: (first hex digit) | | | | | | | | | | | | | | | | | | | |
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | | | | |
| | | 1743 | 245E | | 4341 | | 63A8 | | 84C5 | 9123 | A18D | BEAD | C437 | DEAD | | F00D | | | | |
| | | | | 438B | | | | | | 973C | | | | | | FA10 | | | | |

3 The Radix Sort Algorithm.

3.1 One distribution pass.

Distribute each cell from the master list into a bucket as follows:

- Remove the cell from the front of the master queue (call it currCell).
- Get the temperature out of the Reading that is in the Cell, call it currTemp. Since the temperature is private in the Reading, and the data field of the Cell is private in the Cell, you will need to use get functions.
- CurrTemp is a float, but you need to do integer operations on it. To do this, you must use a *reinterpret cast*. This is a kind of cast that works on pointers, not values. It happens at compile time, not runtime, and simply changes the type of the object in the compiler's internal table. It allows you to copy an object into a variable of the wrong type. It does not change any bits in the object itself. This function does the conversion:

```
unsigned long getTempAsUnsigned( float* pCurrTemp ){
    return *(unsigned long*) pCurrTemp;
}
```

A proper call on it is:

```
unsigned long currTempAsUns = getTempAsUnsigned( &currTemp);.
```

- Access one byte of the temperature and use the result as a subscript to select a bucket. Isolate that byte using shifting and masking. Define `mask = 0xff;` to isolate one byte (8 bits).
 - Compute `subscript = mask & (currTempAsUns >> n)` where `n` is $8 * k$, and `k` is the pass number.
 - The computed subscript is the number of the bucket to use for this data item.
 - Push currCell onto the selected bucket queue.

- In pass k , use the k th byte from the right.
- Attach the cell to the end of the queue in that bucket.
- In this process, the actual sort key is reinterpreted, copied, shifted, and masked, but the changed value must not ever be stored back into the Reading object.

3.2 Collecting the data into the master queue.

At the end of each distribution pass, your master queue is empty. Now collect all the cells, as follows:

- Find the first non-empty bucket, (k), in the bucket array.
- Start at bucket k and process buckets, in order, through number 255.
- If the bucket is not empty,
 - Attach the last cell of the master queue to the first non-dummy cell in the bucket.
 - Set the back pointer of the master list equal to the back pointer of the bucket.
 - All the cells have now been removed from the bucket, so reset it to empty.

3.3 Finishing the Sort.

At the end of the fourth pass, the data is sorted. Now you need to write out the sorted file. Write out the data in the Reading, but remember that you need to subtract 273.15 from every temperature (to convert back from Kelvin to Celsius) before writing it out.

Extra credit, 1 point. Write out the data in the same form as when you read it in: with the date and time in a readable format, then the time_t value, then the device ID number, then the temperature. You will need the `ctime()` function from the C time library.

4 Turn in...

Your source code with part of your output. It is often difficult or impossible to email very large files. For that reason, please make a special file for submitting your output. Use a .txt file. It should start with the first 50 lines of output from the sorted list and end with the last 50 lines of the sorted list. Put a dashed divider line between the two parts of the output.