

## P2: uMessage

**Checkpoint 1:** Due Tue, Oct 20th

**Checkpoint 2:** Due Mon, Nov 2nd

**P2 Due Date:** Due Tue, Nov 10th

The purpose of this project is to implement various data structures and algorithms described in class. You will also implement the back-end for a chat application called “uMessage”.

## Overview

One of the most important ADTs is the Dictionary and one of the most studied problems is sorting. In this assignment, you will write multiple implementations (AVLTree, HashTable, etc.) of Dictionary and multiple sorting algorithms.

All of these implementations will be used to drive *word suggestion*, *spelling correction*, and *autocompletion* in a chat application called uMessage. These algorithms are very similar to the ones smartphones use for these problems, and you will see that they do relatively well with a small effort. Since uMessage has many components and is difficult to test, we will ask you to test your code by writing another client for WordSuggestor.

We have provided the boring pieces of these programs (e.g., GUIs, printing code, etc.), but you will write the data structures that back all of the code we’ve written.

## Project Restrictions

- You *must* work in a group of two unless you successfully petition to work by yourself.
- You may not use **any** of the built-in Java data structures. One of the main learning outcomes is to write everything yourself.
- You may use the `math` package.
- You may not edit any file in the `cse332.*` packages.
- The *design and architecture* of your code are a *substantial* part of your grade.
- The Write-Up is a *substantial* part of your grade; do **not** leave it to the last minute.
- Make sure to not duplicate fields that are in super-classes (e.g., `size`). This will lead to unexpected behavior and failures of tests.

## P1 and Beyond

This project actually extends on p1 a lot! You will need to port over (i.e., put them in the same packages) the following:

- `datastructures.worklists`: All your simple WorkLists: `ArrayStack`, `ListFIFOQueue`, `CircularArrayFIFOQueue`
- `datastructures.dictionaries`: Your `HashTrieSet` and your `HashTrieMap`

Be sure you do NOT place these in `cse332.datastructures.worklists`. After you port these files over, `CircularArrayFIFOQueue` won’t compile. It defines a type parameter `E` in `CircularArrayFIFOQueue<E>` at the top of the class, but you should replace this `E` with “`E extends Comparable<E>`”.

## Provided Code

Several of the interfaces and implementations from p1 also appear in p2. We will only describe the *new* classes in an attempt to be less overwhelming.

- `cse332.interfaces.misc`
  - `DeletelessDictionary.java`: Like a dictionary, but the `delete` method is unsupported.
  - `ComparableDictionary.java`: A `DeletelessDictionary` that requires comparable keys.
  - `SimpleIterator.java`: A simplification of Java's `Iterator` that has no `remove` method.
- `cse332.datastructures.*`
  - `Item.java`: A simple container for a key and a value. This is intended to be used as the object stored in your dictionaries.
  - `BinarySearchTree.java`: An implementation of `Dictionary` using a binary search tree. It is provided as an example of how to use function objects and iterators. The iterators you write will not be as difficult.
- `cse332.*`
  - `WordReader.java`: Standardizes inputs into lower case without punctuation.
  - `LargeValueFirstItemComparator.java`: A comparator that considers larger values as "smaller", and breaks ties by considering the keys.
  - `InsertionSort.java`: A provided implementation of `InsertionSort`.
  - `AlphabeticString.java`: This type is a `BString` that is just a wrapper for a standard `String`.
  - `NGram.java`: This type is a `BString` that represents an  $n$ -gram.
- `p2.wordcorrector`
  - `AutocompleteTrie.java`: This is the trie used by `uMessage`; it is backed by `HashTrieMap`.
  - `SpellingCorrector.java`: This is the spelling corrector used by `uMessage`.
- `p2.wordsuggestor`
  - `NGramToNextChoicesMap.java`: Client data structure that will be used to drive `WordSuggestor`.
  - `ParseFBMessages.java`: This program downloads your facebook messages. It is intended to be used as a way of generating a personal corpus for the `WordSuggestor`. There are more instructions for using this in the writeup spec
  - `WordSuggestor.java`: This is the word suggestor used by `uMessage`.
- `chat`
  - `uMessage.java`: This is the main driver program for `uMessage`.

You will implement data structures `MinFourHeap` (in `datastructures.dictionaries`), `MoveToFrontList`, `AVLTree`, and `ChainingHashTable` (in `datastructures.dictionaries`) and sorting algorithms `HeapSort`, `QuickSort`, and `TopKSort` (in `p2.sorts`).

## `uMessage`

After you have finished all the implementations, you will be ready to try out `uMessage`. We expect you to actually play with the application, and the Write-Up will ask you to do several things with it. Importantly, there are configuration settings ( $n$  and the corpus) at the top of `uMessage.java` which you will want to edit.

## Project Checkpoints

This project will have **two** checkpoints (and a final due date). A *checkpoint* is a check-in on a certain date to make sure you are making reasonable progress on the project. For each checkpoint, you (and your partner) will turn in a Gradescope survey individually.

*As long as you turn in the checkpoint survey, the checkpoint will not affect your grade in any way.*

**Checkpoint 1: (1), (2)**

**Tue, Oct 20th**

**Checkpoint 2: (3), (4), (5), (6)**

**Mon, Nov 2nd**

**P2 Due Date: (7), (8)**

**Tue, Nov 10th**

## Part 0: Importing .jar files

P2 uses a number of .jar files which you will need to import. Make sure to follow the steps below for all of the files in the directory that end in .jar as well as the javafx directory.

To import the jar file or directory

- (1) In the main menu, go to **File | Project Structure**
- (2) Go to **Project Settings | Libraries**
- (3) Click on the + button
- (4) Locate the file or folder in the main directory
- (5) Apply settings and click **Ok**

## Part 1: Another WorkList

First, implement one more WorkList:

### (1) MinFourHeapComparable

Your MinFourHeapComparable should be an implementation of the heap data structure we've discussed in class. It should be an array-based implementation which starts at index 0. Unlike the implementation discussed in lecture, it should be a *four*-heap (not a two-heap). In other words, each node should have *four* children, not two. All operations should have the efficiencies we've discussed in class.

## Part 2: Implementing Dictionary Classes and Sorts

### (2) MoveToFrontList: Another Dictionary

In this part, you will implement MoveToFrontList, a new type of Dictionary.

For the remainder of the Dictionary classes you will implement, we will not ask you to write delete—it is possible (and you can do it for fun), but it's not educational enough to be part of the actual project. As a result, your Dictionary classes will inherit from DeletelessDictionary which is the same as Dictionary except it does not require that you implement a delete method.

MoveToFrontList is a type of linked list where new items are inserted at the front of the list, and an existing item gets moved to the front whenever it is referenced. Although it has  $\mathcal{O}(n)$  worst-case time operations, it has a very good amortized analysis.

You will also be implementing an Iterator for this dictionary. The runtime for all Iterator operations should run in  $\mathcal{O}(1)$ . We will not be discussing Iterators in class so if you need, you can reference the Iterator for BinarySearchTree.

`MoveToFrontList` relies on *equality* testing of elements. In Java, we deal with this by defining an `equals` method. If you look in `BString`, it relies on `CircularArrayFIFOQueue` having a reasonable definition of equality. Before `MoveToFrontList` will work, you will need to define the `equals` method for `CircularArrayFIFOQueue`. You may not use `toString` to implement `equals`; we expect you to build it from scratch. You might be wondering how to figure out the type of the parameter for `equals`; in Java, the `equals` method takes an `Object`. You will want to do research on the Java instance of operator, as it will be a part of your solution.

In addition to equality testing, we also need to be able to *compare* two `Objects`. To do this, you should complete the `compareTo` method in `CircularArrayFIFOQueue`. You may not use `toString` to implement `compareTo`; we expect you to build it from scratch.

The reason we implement this is that our *tree* dictionaries in the next part will need to be able to do comparisons instead of equality testing.

Remember, in any `Dictionary` implementation, you may use any of your `WorkList` implementations.

### (3) AVLTree: Another Another Dictionary

In this part, you will implement `AVLTree`. We recommend waiting to do this until we have discussed it in lecture. Just like before, you do not have to implement `delete`. Your `AVLTree` should be a sub-class of `BinarySearchTree` which we have written for you. Be careful to not duplicate code in rotation. You should use an array implementation of left and right children as in `BinarySearchTree`. Your `insert(K key, V value)` should run in  $\mathcal{O}(\log(n))$ . If your rotation code is repetitive or does not run in  $\mathcal{O}(\log(n))$ , you will lose a substantial amount of points.

A note on `AVLTree` Inheritance. `AVLTree` extends `BinarySearchTree`, and `BinarySearchTree` has a couple methods we might think could be useful: `find(K key, V value)` and `find(K key)`. Some of you may be trying to use the former (`find(K key, V value)`) to access the appropriate spot in your tree without duplicating code, but there's actually an issue with this: `find(K key, V value)` puts `BSTNodes` in your `AVLTree` and returns them to you. These nodes can't be cast to `AVLNode` (because they were initialized as `BSTNodes`), and, since they are `BSTNodes`, they don't have that all-important height field, so you can't use them.

In other words, you should not call the `find(key, value)` method (with a non-null second argument) in `BinarySearchTree` as part of your `insert` method. It's okay if you end up duplicating some of the `find(K key, V value)` logic in your `insert()` method.

You will not need to write a separate `find(key)` method, though, since the behavior of that method will be the same for both tree types, meaning that the inherited method already behaves correctly.

Recall that all BSTs rely on a reasonable definition of comparison. Our BST and your `AVLTree` will both rely on the `compareTo` that you wrote in the previous part.

A note on debugging. You can "fail fast" by adding your `verify-avl` code as a private helper method, checking validity after every modification to the tree, and throwing an exception if the check fails. This will help you identify which sections of the code are breaking the tree. These checks will be expensive and should be disabled in the final version, but can be helpful when debugging.

### (4) ChainingHashTable: Another Another Another Dictionary

In this part, you will implement `ChainingHashTable`. We recommend waiting to do this until we have discussed it in lecture. Just like before, you do not have to implement `delete`. Your hash table must use separate chaining—not probing. Furthermore, you must make the type of chain generic. In particular, you should be able to use *any* dictionary implementation as the type inside the buckets. Your `HashTable` should rehash as appropriate (use an appropriate load factor as discussed in the class), and its capacity should always be a prime number. Your `HashTable` should be able to work with `uMessage` which means

there shouldn't be a hard cap on how much it can grow; though, it doesn't have to use primes past 200,000.

Pick a reasonable starting size for your `HashTable`. You should use a hardcoded list of primes to resize up to 200,000. Do not hard code every prime up to 200,000 - just pick a reasonable range of prime numbers. After this point, you should continue to resize your table using some other mechanism. Note that you **MUST GROW** the table past 200,000. It's ok if you just double the size of your table when you re-hash!

Recall that all Hash Tables rely on a reasonable definition of hash code. Just like you needed to define `equals` and `compareTo` for various other data structures, you will need to define `hashCode` in `CircularArrayFIFOQueue` for `ChainingHashTable`. You may not use `toString` to implement `hashCode`; we expect you to build it from scratch.

At some point, you will want to test various types of chains in your `ChainingHashTable`. It is confusing to do this initially; so, we have provided some examples in the `NGramTester` class.

### (5) `HashTrieMap`: Full Circle!

Now that you have written your own hash map, replace the dependency on Java's `HashMap` with your `ChainingHashTable`! This is not only okay, it's a great example of unexpected refactoring. Refactoring will usually set off a chain reaction where you also have to edit other code.

You will want to look at [the `SimpleEntry` javadoc](#). Remember that you may edit any class that is not in a `cse332.*` package.

Here is a general guide on what to change:

- You will need to fix `AutocompleteTrie.java` as part of your refactor
- Some methods might become impossible to implement with `ChainingHashTable`. In this case, it is okay to throw a `UnsupportedOperationException`.
- You will notice a mismatch between the type of iterator returned from `ChainingHashTable` and the one that you need in `HashTrieNode`. This is an example of a common issue you run into while refactoring code.
  - (a) You'll need to add a (small) bit of code to `HashTrieNode/HashTrieMap` to work around this type mismatch. You can do it using what you've already learned about Iterators.
  - (b) Note that you shouldn't modify the `ChainingHashTable.iterator()` return type, because then it wouldn't match the `Dictionary` interface, and you also shouldn't add superfluous iterator methods to `ChainingHashTable` to solve a problem in `HashTrieMap`.

You have now written pretty much all of the data structures that you've used from Java's library! You now understand all the magic under the hood! Take a minute to bask in the glory that is data structures nirvana.

### (6) `MinFourHeap` and The Sorts

The `MinFourHeap` you wrote before was only able to compare elements in a single way (based on the `compareTo`). There is a more general idea called a `Comparator` which allows the user to specify a comparison *function*. The first thing you should do in this part is implement `MinFourHeap` to use a `Comparator`. You can copy the logic from `MinFourHeapComparable` and modify the logic to use a `Comparator` instead of `compareTo`. Make sure to modify the array type such that it no longer extend `Comparable`. This is necessary to make the sorts (below) work.

After you've edited `MinFourHeap`, you will be ready to write the following sorting algorithms:

- `HeapSort`: Consists of two steps:
  - (1) Insert each element to be sorted into a heap (`MinFourHeap`)
  - (2) Remove each element from the heap, storing them in order in the original array.
- `QuickSort`: Implement quicksort. As with the other sorts, your code should be generic. Your sorting algorithm should meet its expected runtime bound.
- `TopKSort`: An easy way to implement this would be to sort the input as usual and then just print  $k$  largest of them. This approach finds the  $k$  largest items in time  $\mathcal{O}(n \lg n)$ . However, your implementation should have  $\mathcal{O}(n \lg k)$  runtime, assuming  $k$  is less than or equal to  $n$ . Efficiently tracking the  $k$  largest will require a different comparator than what you used in `HeapSort`. `TopKSort` should put the top  $k$  elements in the first  $k$  spots in the array, and **all the other indices should be null**. In other words, if  $A = \text{quicksort}(B)$  for some array  $B$ , then:  $\text{topKSort}(k, A) = [A[n - k], A[n - (k - 1)], \dots, A[n - 1], \text{null}, \text{null}, \dots, \text{null}]$ .  
**(Hint:** Use a heap, but never put more than  $k$  elements into it. Think about why this gives  $\mathcal{O}(n \lg k)$  runtime bound).

## Part 3: The Write-Up

### (7) Write-Up

Approximately half of your grade will be based on your write-up. The analysis part of this project is incredibly important, and we expect you to spend an entire week's worth of work on it. You will find the write-up questions here in the [P2 Write-up Template](#). Remember to follow the instruction on the first and second page to receive full credit!

Some of the write-up questions will ask you to design and run some experiments to determine which implementations are faster for various inputs. Answering these questions will require writing additional code to run the experiments, collecting operation counting or timing information and producing result tables and graphs, together with relatively long answers. Do not wait until the last minute! We will post more information about the difference between operation counting or timing.

Insert tables and graphs into your write-up as appropriate, and be sure to give each one a title and label the axes for the graphs. **IMPORTANT: Place all your operation counting or timing code into the `package experiment`.** Be careful not to leave any write-up related code in the normal files. To prevent losing points due to the modifications made for the write-up experiments, you should copy all files that need to be modified for the experiments into the `package experiment`, and start working from there. Files in different packages can have the same name, but when editing be sure to check if you are using the correct file! If your code does not compile because you did not follow these instructions, you will receive a 0 for all automated tests.

You will need to write a second hashing function. To exaggerate the difference between the two hash functions, you will want to compare a very simple hash function with a decent one (the one used in Part 2). For all experimental results, we would like to see a detailed interpretation, especially when the results do not match your expectations.

### (8) `uMessage` - Do not wait until the last minute for this!

Now that you are done with all of the coding (and most of the write-up) for the project, you are ready to attempt to run `uMessage`. As many folks saw when they ran `zip` on P1, this may expose problems with code you wrote earlier. Do not wait until the last minute for this step!

Note: When using `uMessage`, our course policy requires that you use your CSE or UWNnetID as your username. Using the system with a pseudonym or anything other than your netID is grounds for failing the assignment. This is a fun program to play around with (please do!) but anyone found using the system to annoy or harass others will be referred to the appropriate university authorities.

Before you run `uMessage`, you will want to do the following:

- Increase the allowed heap size in IntelliJ. In particular, `uMessage` runs significantly more smoothly if you give it **6GBs** of memory.
  - (1) On the **Help** menu, click **Edit Custom VM Options**.
  - (2) Set the `-Xmx` option to 3G without the quotation marks, so the final line should be `-Xmx3G`
  - (3) Restart IntelliJ
- Make sure your computer is plugged in. (Yes, this will make a difference.)
- Finish the `getWordsAfter` method in `NGramToNextChoicesMap`. You should replace `InsertionSort` with a faster, standard sort, and if  $k \geq 0$ , you should run `TopKSort`. You might have to do something more than *just* run `TopKSort` to get the most frequent words out. Figuring out exactly what to do here is part of the challenge.

There are several variables at the top of `uMessage` which you will have to edit: the corpus, the “n”, the “inner dictionary” and the “outer dictionary”. If you leave the corpus as `eggs.txt`, the suggestions will be garbage. If you leave the inner and outer dictionaries as tries, `uMessage` will probably be too slow. The point of `uMessage` is that it is a cool application that uses all of the code you wrote. Just like `Zip` was a good stress test for P1, `uMessage` is a good stress test for P2.

Once you start working on `uMessage`, if you’ve implemented `getWordsAfter` correctly, the word suggestions you get in `uMessage` should be sorted by frequency (conditioned on the previous words), with highest frequency on the left and lowest frequency on the right. Note that inputs with apostrophes may not work, (ej. `can't`, `wouldn't`), and throw an `SSLPeerUnverifiedException`. This is a bug in `uMessage` and you can just ignore it :)

As a simple example, with `irc.corpus`, the words suggested as first words on a newly-opened chat with nothing typed should be [`"i"`, `"and"`, `"yeah"`, `"well"`], in that order, since those are the four words with the highest frequency at the start of a line, with `"i"` being the most frequent of the four.

Trying to debug issues with your ordering code on `irc.corpus` will take a long time (since this corpus takes a while to load), so it might be a good idea to make a simple test corpus with only a few sentences where you can work out what the suggested words should be, and using that to quickly figure out `getWordsAfter`.