# Spell Checker Implementation Using Custom HashMap and HashSet

Hasan Can İstekli

210104004058

CSE222 - Data Structures and Algorithms

HW6

May 13, 2025

# Contents

# Abstract

This project builds a high-performance spell checker using custom `GTUHashMap` and `GTUHashSet`. It detects misspellings and suggests corrections via edit distance, avoiding Java's built-in collections and following object-oriented principles. The report covers methodology, results, challenges, and solutions.

# 1.   Introduction

## 1.1.   Literature Review

Spell checking is a key text processing task used in word processors and search engines. Traditional Java implementations use built-in data structures like HashMap and HashSet, which hide underlying mechanisms and hinder optimization. This project builds custom data structures to offer similar functionality with clearer design and performance insights.

## 1.2.   Objectives

The primary objectives of this project are:

- Implement a custom `GTUHashMap` and `GTUHashSet` using open addressing with linear probing.

- Develop a spell checker that suggests corrections based on edit distance.

- Ensure the implementation is efficient and meets performance requirements (response time < 100ms).

- Avoid the use of Java's built-in collection libraries.

## 1.3.   Scope of the Project

This project focuses on implementing a spell checker that can handle a dictionary of over 120,000 words. The spell checker must efficiently identify misspelled words and suggest corrections within an edit distance of 2. The implementation must be robust, handling edge cases such as empty inputs, special characters, and large dictionaries.

# 2. Methodology

## 2.1. Overview

The project is divided into three main components:

1. Implementation of `GTUHashMap` using open addressing.

2. Implementation of `GTUHashSet` built on top of `GTUHashMap`.

3. Development of a spell checker using `GTUHashSet`.

## 2.2. GTUHashMap Implementation

### 2.2.1. Open Addressing and Probing

The `GTUHashMap` uses open addressing with linear probing to resolve collisions. Each entry in the hash table is represented by an `Entry<K, V>` object, which includes a tombstone flag to handle deletions.

### 2.2.2. Rehashing

To maintain performance, the hash table is rehashed when the load factor exceeds a threshold. The new table size is the next prime number greater than twice the current capacity.

### 2.2.3. Methods Implemented

The following methods were implemented:

- `put(K key, V value)`: Inserts a key-value pair into the hash table.

- `get(K key)`: Retrieves the value associated with a key.

- `remove(K key)`: Marks an entry as deleted using a tombstone.

- `containsKey(K key)`: Checks if a key exists in the hash table.

- `size()`: Returns the number of active entries in the hash table.

## 2.3.  GTUHashSet Implementation

The `GTUHashSet` is built on top of `GTUHashMap`, reusing its efficient hashing mechanism. The `add`, `remove`, `contains`, and `size` methods delegate their operations to the underlying `GTUHashMap`.

## 2.4.  Spell Checker Implementation

### 2.4.1.  Dictionary Loading

The dictionary is loaded from a text file containing over 50,000 words. Each word is added to a `GTUHashSet` for efficient lookup.

### 2.4.2.  Edit Distance Calculation

The spell checker generates all possible variants of a word within an edit distance of 2 using character-level operations (insertion, deletion, substitution). Each variant is checked against the dictionary using `GTUHashSet.contains()`.

### 2.4.3.  Performance Optimization

To meet the performance requirement, the implementation avoids scanning the entire dictionary. Instead, it generates and tests only relevant variants of the input word.

# 3. Results and Discussion

## 3.1. Correctness

The spell checker was tested with various inputs, and the results were as expected. For example:

- Input: `speling`

- Suggestions: `spelling, spewing, spieling, spending, opening ...`

- Input: `wrld`

- Suggestions: `weld, wold, wild, world, welt, well ...`

## 3.2. Performance

The average response time for a query was measured at 85ms, meeting the performance requirement. The rehashing mechanism ensured that the hash table maintained efficient lookup times even with a large dictionary.

## 3.3. Challenges and Solutions

### 3.3.1. Handling Collisions

Collisions were handled using linear probing. However, clustering was observed in some cases. This was mitigated by rehashing to a larger capacity.

### 3.3.2. Edge Cases

Edge cases such as empty inputs and special characters were handled by preprocessing the input before generating variants.

### 3.3.3. High Initial Operation Time

During the implementation and testing of the spell checker, it was observed that the initial operations, such as loading the dictionary and performing the first word lookup, took significantly more time compared to subsequent operations. This behavior was analyzed, and the following reasons were identified:

- **Disk I/O Overhead:** The dictionary file (`dictionary.txt`) contains 84,099 words, and reading this large file from disk introduces a significant delay during the first load. However, subsequent loads benefit from the operating system's disk caching mechanism, which reduces the time required for file access.

- **JVM Warm-Up and JIT Compilation:** The Java Virtual Machine (JVM) requires a warm-up period during the first execution. During this time, the Just-In-Time (JIT) compiler optimizes frequently executed code paths, which results in improved performance for subsequent operations.

### 3.3.4. Solutions Implemented

To address the high time for initial operations, the following solutions were implemented:

1. **Increased Buffer Size for File Reading:** The buffer size for the `BufferedReader` was increased from the default size (8 KB) to 64 KB. This allows the program to read larger chunks of data at once, reducing the number of I/O operations and improving file reading performance.

```
reader = new BufferedReader(new FileReader("dictionary.txt"), 65536);
```

2. **Optimized Initial Capacity of `GTUHashSet`:** The initial capacity of the `GTUHashSet` was increased to 120,000 to accommodate the 84,099 words in the dictionary without triggering rehashing. This reduces the overhead of resizing the hash table during the dictionary loading process.

```
GTUHashSet<String> dictionary = new GTUHashSet<>(120000);
```

### 3.3.5. Results After Optimization

After applying these optimizations, the dictionary loading time was reduced significantly. The following table shows the performance improvement:

| Operation | Before Optimization (ms) | After Optimization (ms) |
|---|---|---|
| Dictionary Loading | 50 ms | 30 ms |
| First Word Lookup | 150 ms | 100 ms |

Table 3.1: Performance improvement after optimizations.

### 3.3.6. Are Loading and First Lookup Independent?

The dictionary loading process and the first word lookup are **not completely independent**. The first lookup may still be affected by:

- JVM warm-up and JIT optimizations, which occur during the initial execution of the program.

- Memory allocation and initialization of the `GTUHashSet` during the dictionary loading phase.

### 3.3.7. Further Improvements

To further reduce the time for the first operation, the following approaches could be considered:

- **Preloading the Dictionary into Memory:** Serialize the dictionary into a binary format and load it directly into memory during program startup. This avoids the overhead of parsing the text file.

- **Warm-Up Phase:** Perform a dummy lookup operation after loading the dictionary to allow the JVM to optimize the code before the first real lookup.

# 4.  Conclusions

## 4.1.  Summary

The project successfully implemented a high-performance spell checker using custom data structures. The objectives were met, and the implementation is both efficient and robust.

## 4.2.  Future Work

Future improvements could include:

- Implementing quadratic probing for collision resolution.

- Adding support for multi-language dictionaries.

- Tracking and reporting metrics such as collision counts and memory usage.

# References

- Assignment Description Document

- Data Structures and Algorithms Textbook

- Online Resources on Open Addressing and Edit Distance

# A. Appendices

## A.1. Code Snippets

### A.1.1. GTUHashMap.java

Below is a screenshot of the GTUHashMap class, which implements a custom HashMap using open addressing and linear probing. The put() function handles key-value insertion, while resolving collisions efficiently.

```java
public void put(K key, V value) {
    if (key == null) {
        // Decide on null key handling: throw exception or handle specially
        // For simplicity in this homework, let's assume keys are not null.
        // If allowed, a specific strategy for null key hashing/storage would be needed.
        System.err.println("Null keys are not supported in this GTUHashMap implementation.");
        return;
    }

    if ((double) (size + tombstones) / capacity >= LOAD_FACTOR_THRESHOLD) {
        rehash();
    }

    int index = findIndex(key, forInsert:true);

    if (table[index] == null || table[index] == Entry.TOMBSTONE) {
        if (table[index] == Entry.TOMBSTONE) {
            tombstones--;
        }
        table[index] = new Entry<>(key, value);
        size++;
    } else { // Key already exists, update value
        table[index].value = value;
    }
}
```

Figure A.1: Screenshot of the put() function in the GTUHashMap class.

The get() function retrieves the value associated with a given key. Below is a screenshot of its implementation.

```
/**
 * Retrieves the value associated with the given key.
 *
 * @param key The key to search for.
 * @return The value associated with the key, or null if not found.
 */
public V get(K key) {
    if (key == null)
        return null; // Consistent with not supporting null keys in put
    int index = findIndex(key, forInsert:false);
    if (index != -1 && table[index] != null && table[index] != Entry.TOMBSTONE && table[index].key.equals(key)) {
        return table[index].value;
    }
    return null;
}
```

Figure A.2: Screenshot of the `get()` function in the `GTUHashMap` class.

---

### A.1.2. `GTUHashSet.java`

The `GTUHashSet` class is built on top of `GTUHashMap` and provides set functionality. Below is a screenshot of the `add()` function, which adds an element to the set.

```
/**
 * Adds an element to the set.
 *
 * @param element The element to be added.
 * @return true if the element was added, false if it already exists.
 */
public boolean add(E element) {
    if (map.containsKey(element)) {
        return false; // Element already present
    }
    map.put(element, PRESENT);
    return true; // Element added
}
```

Figure A.3: Screenshot of the `add()` function in the `GTUHashSet` class.

The `contains()` function checks whether an element exists in the set. Below is a screenshot of its implementation.

Figure A.4: Screenshot of the `contains()` function in the `GTUHashSet` class.

—

### A.1.3. `SpellChecker.java`

The `SpellChecker` class is responsible for loading the dictionary, checking the correctness of input words, and generating suggestions based on edit distance. Below is a screenshot of the main function.

```java
public class SpellChecker {
    public static void main(String[] args) throws IOException {
        GTUHashSet<String> dictionary = new GTUHashSet<>(initialCapacity:70000); // Estimate dictionary size for initial capacity

        // Load dictionary from file
        System.out.println("Loading dictionary...");
        long loadStartTime = System.nanoTime();
        BufferedReader reader = new BufferedReader(new FileReader("dictionary.txt"));
        String word;
        int count = 0;
        while ((word = reader.readLine()) != null) {
            dictionary.add(word.trim().toLowerCase()); // Store words in lowercase
            count++;
        }
        reader.close();
        long loadEndTime = System.nanoTime();
        System.out.printf("Dictionary loaded with %d words in %.2f ms.\n", count, (loadEndTime - loadStartTime) / 1e6);

        Scanner scanner = new Scanner(System.in); // Scanner is allowed

        // Main loop for user input
        while (true) {
            System.out.print("\nEnter a word (or type 'exit' to quit): ");
            String input = scanner.nextLine().trim().toLowerCase();

            if (input.equalsIgnoreCase("exit")) {
                break;
            }
            if (input.isEmpty()) {
                continue;
            }

            long startTime = System.nanoTime();

            // Check if the word is in the dictionary
            if (dictionary.contains(input)) {
                System.out.println("'" + input + "' is spelled correctly.");
            } else {
                System.out.println("'" + input + "' is misspelled.");
                System.out.print("Suggestions: ");

                GTUArrayList<String> suggestions = new GTUArrayList<>();
                GTUHashSet<String> uniqueSuggestions = new GTUHashSet<>(); // To ensure suggestions are unique

                // Generate and check edit distance 1
                GTUArrayList<String> edits1 = EditDistanceHelper.getEdits1List(input);
                for (int i = 0; i < edits1.size(); ++i) {
                    String variant1 = edits1.get(i);
                    if (dictionary.contains(variant1)) {
                        if (uniqueSuggestions.add(variant1)) { // add returns true if new
                            suggestions.add(variant1);
                        }
                    }
                }

                // Generate and check edit distance 2
                // For each variant from edits1, generate its edits1 list (these are edits2 from
                // original)
                for (int i = 0; i < edits1.size(); ++i) {
                    String variant1 = edits1.get(i);
                    // If variant1 itself is a valid word, we already added it.
                    // Now, generate edits from variant1 to find edit distance 2 words from original
                    // input.
                    GTUArrayList<String> edits2_from_variant1 = EditDistanceHelper.getEdits1List(variant1);
                    for (int j = 0; j < edits2_from_variant1.size(); ++j) {
                        String variant2 = edits2_from_variant1.get(j);
                        if (dictionary.contains(variant2)) {
                            if (uniqueSuggestions.add(variant2)) {
                                suggestions.add(variant2);
                            }
                        }
                    }
                }

                if (suggestions.isEmpty()) {
                    System.out.println("No suggestions found.");
                } else {
                    System.out.println(suggestions.toString());
                }
            }

            long endTime = System.nanoTime();
            System.out.printf("Lookup and suggestion generation took %.2f ms.\n", (endTime - startTime) / 1e6);
        }

        // Close resources
        scanner.close();
        System.out.println("Spell checker terminated.");
```

Figure A.5: Screenshot of the main function in the `SpellChecker` class.

Below is an example output of the `SpellChecker` for the input word `"speling"`.

Figure A.6: Example output of the `SpellChecker` for the input `"speling"`.

Another example output for the input word `"wrld"` is shown below.



Figure A.7: Example output of the `SpellChecker` for the input `"wrld"`.

—

### A.1.4.  Test Results

The following test cases were executed using the `TestRunner` class to verify the correctness and functionality of the `GTUHashMap` and `GTUHashSet` implementations.

### A.1.5.  Test Case 1: `GTUHashMap` Operations

This test verifies the basic operations of the `GTUHashMap` class, including:

- Adding key-value pairs using the `put()` method.

- Retrieving values using the `get()` method.

- Removing a key-value pair using the `remove()` method.

- Checking the existence of a key using the `containsKey()` method.

16

Figure A.8: Test Case 1: `GTUHashMap` operations (`put()`, `get()`, `remove()`, and `containsKey()`).

### A.1.6. Test Case 2: `GTUHashSet` Operations

This test verifies the basic operations of the `GTUHashSet` class, including:

- Adding elements using the `add()` method.

- Checking the existence of an element using the `contains()` method.

- Removing an element using the `remove()` method.



Figure A.9: Test Case 2: `GTUHashSet` operations (`add()`, `contains()`, and `remove()`).

### A.1.7. Summary of Test Results

The following table summarizes the results of the test cases executed using the `TestRunner` class:

| Test Case | Description | Result |
|---|---|---|
| Test Case 1 | `GTUHashMap` operations | Passed |
| Test Case 2 | `GTUHashSet` operations | Passed |

Table A.1: Summary of test results for `GTUHashMap` and `GTUHashSet`.