

# High-Performance Spell Checker with Custom HashMap and Set

**Course:** CSE 222 - Data Structures and Algorithms

**Assignment:** Homework 6

**Student Name:** Emirhan Çalışkan

**Date:** May 2025

# Abstract

This project presents a high-performance spell checker developed in Java, utilizing fully custom implementations of HashMap and HashSet without relying on the Java Collections Framework. The spell checker reads from a dictionary of over 50,000 words and provides fast suggestions for misspelled words using edit-distance logic, specifically limited to an edit distance of two. The focus of the project is both algorithmic and structural performance, avoiding full-dictionary scans and responding within strict time constraints (under 100ms). The design, testing, and evaluation phases prove that the custom data structures can achieve reliable and scalable performance in real-world tasks.

---

## 1.0 INTRODUCTION

Spell checkers are integral components of modern software systems, whether embedded in word processors, search engines, mobile applications, or developer tools. They help identify and suggest corrections for misspelled words, improving user experience, language correctness, and information retrieval quality.

The challenge in building efficient spell checkers lies not only in their linguistic accuracy but also in the performance and responsiveness they offer to users. Most production-grade spell checkers use advanced data structures like tries, Levenshtein automata, or hash maps built into modern programming languages. However, in this project, we set a unique constraint: **no Java Collections Framework usage.**

Instead, we explore how **custom object-oriented hash structures** can be implemented and optimized to support a high-performance spell checker — a design choice that grants us low-level control over collision handling, memory, and load factors.

## 1.1 Literature Review

The academic and industrial literature on spell checkers has covered many algorithmic models — from probabilistic models (e.g., noisy channel models) to statistical learning, to classical edit-distance-based approaches.

Hashing and open addressing in hash tables are well-studied topics in algorithm design. Knuth, Cormen et al., and others have documented their efficiency in average-case constant-time lookup scenarios. Their simplicity and speed make them an ideal choice when exact matching is the primary operation.

In many competitive programming and constrained environments (e.g., embedded systems or security-restricted software), using low-level or custom data structures becomes necessary. Our approach was inspired by such environments.

## 1.2 Objectives

The project was guided by the following objectives:

- Develop a working spell checker from scratch using object-oriented principles.
- Implement `GTUHashMap<K, V>` with open addressing and linear probing.
- Build `GTUHashSet<E>` on top of the map to simulate standard set behavior.
- Use the custom set to build an efficient suggestion mechanism based on edit distance  $\leq 2$ .
- Avoid scanning the entire dictionary for suggestions — instead generate possible variants of the input.
- Ensure system responds to queries within 100 milliseconds.
- Write a structured project report following academic guidelines.

---

## 2.0 METHODOLOGY

This section explains the technical foundation and design decisions that governed our implementation.

## 2.1 Analytical Work

We first analyzed the behavior of open-addressed hash tables. In such systems, instead of chaining collided keys into linked lists, all entries reside in a single array. When a collision occurs (i.e., two keys hash to the same slot), linear probing searches for the next available empty slot. This method offers cache efficiency but requires careful handling of deletions.

In our GTUHashMap, we designed an `Entry<K, V>` class with an `isDeleted` flag, also known as a **tombstone**. This allows the probing chain to remain intact after deletions.

The rehashing mechanism triggers when the load factor exceeds 0.6. The new capacity is chosen as the next prime number greater than twice the current capacity to reduce clustering.

## 2.2 Computer Simulation

We wrote a full console-based application in Java, structured as follows:

- GTUHashMap.java:
  - Maintains an internal array of Entry objects.
  - Implements `put()`, `get()`, `delete()`, `containsKey()`, and `resize()` for dynamic capacity increase.
  - Hash function: Uses `key.hashCode()` modulo array size.

- Load factor threshold set to 0.75.
- GTUHashSet.java:
  - Uses GTUHashMap to store unique elements.
  - Overrides basic Set operations such as add, remove, contains.
- EditDistanceHelper.java:
  - Contains a static method generateEditDistance1(String word).
  - Returns a MyList<String> containing all possible words one edit away.
- SpellChecker.java:
  - Loads dictionary at startup.
  - Accepts a word from user.
  - Checks if the word exists in the dictionary set.
  - If not, generates candidates and displays valid suggestions.
- MyList.java
  - Custom dynamic array for storing unknown-sized results.
  - Supports add, get, and size operations.
- MyTestFramework.java
  - Performs basic testing for GTUHashMap and GTUHashSet functionalities.

- Ensures all core operations behave as expected.

To simulate realistic usage, we used the dictionary.txt provided and tested inputs of varying lengths and correctness.

## **2.3 Experimental Work**

We ran manual test cases by entering common misspellings ("helo", "aple", "watter") and checked whether our system detected them as incorrect and returned appropriate suggestions.

Performance was measured using `System.nanoTime()`, and logging was added to monitor memory usage and probe collisions (optional bonus feature).

Additionally, we stressed-tested the spell checker with longer non-existent inputs to test worst-case probing scenarios.

---

## 3.0 RESULTS AND DISCUSSION

The spell checker performed accurately in both detection and suggestion modes. For every misspelled word, suggestions were generated by producing all edit-distance-1 variants and checking if they exist in the dictionary using our GTUHashSet.

### **Performance:**

The system consistently responded in under 100ms (average: 45ms), even with 50,000+ dictionary entries. Rehashing occurred smoothly at high load factors.

### **Accuracy:**

We verified that:

- Valid words were accepted.
- Invalid words triggered suggestions.
- Suggestions included real dictionary words within one insertion, deletion, or substitution operation.



## MemoryUsage:

As entries grow, memory expands via rehashing. The use of tombstones slightly increases table size, but ensures correctness.

### Example Results:

```
Memory used while loading dictionary: 10.67 MB
Enter a word: helo
Incorrect. Suggestions:
[rely, hells, silo, lo, polo, veto, bell, hold, hole, kilo, belt, holm, helms, self, hugo, sell, hmsa, heck, reno, kelp, homo, ilo, tel, cell, gel, herb, herd, celt, geo, here, te
le, elf, hero, tell, elk, helot, elm, hers, helots, elt, peso, heed, ely, repo, deli, keno, dell, hess, helps, cello, milo, dels, leo, bolo, meld, hobo, heft, demo, zero, melt, he
els, solo, redo, hem, hen, hula, her, eels, hew, hex, hey, vela, hulk, hull, memo, malo, heals, del, helen, nell, hecto, hallo, el, eo, fell, felo, felt, bello, mel, weld, hewn, h
eir, nemo, hems, well, heron, welt, helga, geld, deco, hypo, gels, gelt, hill, hilt, reo, eel, kelso, pelf, he, felon, halos, ho, sego, held, pelt, hell, helm, help, below, neo, n
olo, yell, palo, yelp, hale, oslo, half, wheelk, hall, hemp, whelp, halo, hems, ego, shelf, halt, shell, hemi, melon, hullo, hens, nero, gebo, head, heal, heap, hear, hello, heat]
Total collisions: 185
Completed in 59.61 ms
```

---

## 4.0 CONCLUSIONS / RECOMMENDED FUTURE WORK

We conclude that it is entirely feasible to implement efficient and scalable spell-checking systems using fully custom object-oriented data structures in Java. Our approach avoids standard libraries, handles collisions and deletions efficiently, and adheres to strict performance goals.

### Future Work Recommendations:

- Implement **quadratic probing** or **double hashing** to reduce primary clustering.
  - Develop a **graphical user interface (GUI)** for broader usability.
  - Extend the dictionary to **multilingual support**.
  - Use **frequency-based ranking** for suggestions.
  - Parallelize variant generation using **multi-threading**.
- 

## REFERENCES

1. Oracle Java Documentation. *Java SE API Specification*.  
<https://docs.oracle.com/javase/8/docs/api/>
2. Cormen, T.H., Leiserson, C.E., Rivest, R.L., Stein, C. *Introduction to Algorithms*. MIT Press, 3rd Edition.
3. Knuth, D.E. *The Art of Computer Programming, Volume 3: Sorting and Searching*. Addison-Wesley.
4. Norvig, P. (2007). *How to Write a Spelling Corrector*. <https://norvig.com/spell-correct.html>

---

# APPENDICES

## Appendix A – Sample Input and Output

```
Enter a word: aple
Incorrect. Suggestions:
[aptly, le, alec, pole, ague, acte, ave, mule, apple, ales, hole, ogle, alex, apnea, sample, table, awe, apply, ail, isle, awl, awls, maple, bale, tile, eagle, anne, marple, axe,
allen, sale, pre, ache, alley, habile, kale, adler, tele, caplet, apia, ae, ampler, able, aye, pee, al, aide, apis, apses, ably, axle, mile, ala, alb, ale, pe, all, alp, alt, bole,
rape, sable, agile, dale, vile, jape, sole, pule, nile, arse, epee, file, male, cole, gable, haply, staple, amble, yule, wile, cable, pie, cape, acme, vale, ally, abe, tape, ape,
dole, auld, acne, apt, axles, ladle, rule, pile, ace, wale, idle, mole, basle, ample, kyle, fable, gale, are, amply, plc, angle, aloe, gayle, vole, maples, ailed, ply, ase, pale,
plea, pled, naples, addle, hale, nape, ails, ate, duple, dapple, lapse, rile, acre, aisle, yale, arles, aprs, ipse, age, aped, apse, apace, bile, poe, abler, apes, apex, apples,
gape]
Total collisions: 246
```

```
Enter a word: sceres
Incorrect. Suggestions:
[cere, secures, scents, sere, scire, scales, sces, screws, scores, scorer, stereo, stereos, scites, scree, scarps, serbs, scored, sees, scopes, scenes, secedes, screw, schemes, sc
ene, seers, scorers, sabres, scones, sires, spares, sceptres, shares, score, steles, acres, sperms, cores, snores, sores, sieves, scarves, serves, cures, swerves, scarfs, scare, s
cars, shores, serfs, scares, sterns, series, scorn, cedes, cares, scared, sieges, spires, shires, sexes, spheres, heres]
Total collisions: 48
Completed in 12.97 ms
```

Enter a word: blasetae

Incorrect. Suggestions:

No suggestions found.

Completed in 54.28 ms

Enter a word: detaw

Incorrect. Suggestions:

[seta, decal, petal, deter, decay, degas, deltas, setae, betas, fetal, delta, dental, dotal, demar, data, eta, debar, dew, taw, draw, beta, metal, dead, delay, deaf, detach, deal, dean, dear, detain, detail]

Total collisions: 13

Completed in 7.18 ms

Enter a word: sceres

Incorrect. Suggestions:

[cere, secures, scents, sere, scire, scales, sces, screws, scores, scorer, stereo, stereos, scites, scree, scarps, serbs, scored, sees, scopes, scenes, secedes, screw, schemes, scene, seers, scorers, sabres, scones, sires, spares, sceptres, shares, score, steles, acres, sperms, cores, snores, sores, sieves, scarves, serves, cures, swerves, scarfs, scare, scars, shores, serfs, scares, sterns, series, scorns, cedes, cares, scared, sieges, spires, shires, sexes, spheres, heres]

Total collisions: 48

Completed in 12.97 ms

## Appendix B – Manual Test Plan

Test Case	Input	Expected Output	Result
Valid word	"java"	Correct	Pass
Common typo	"teh"	Suggestions including "the"	Pass
Nonsense input	"asdff"	No suggestions	Pass