

---

# Duplicate Code Detection System

---

**Arun Cheriakara**

Department of Computer Science  
University of Waterloo  
Waterloo, ON N2L 3G1, Canada  
acheriak@uwaterloo.ca

## 1 Introduction

Duplicate code, also known as code duplication or cloning, is a critical problem in software development. It occurs when developers replicate existing code instead of reusing it through modular functions or classes. While it may seem harmless at first, duplicated code can lead to increased maintenance efforts, bugs, and technical debt as projects evolve.

The following paper is a proposal for a Duplicate Code Detection System which is supposed to help in identifying and grouping code clones in a particular directory. The system classifies duplicates into three categories namely; exact duplicates, near duplicates and structural duplicates. Employing sophisticated approaches to string similarity, Abstract Syntax Tree (AST), and hashing, the system tries to enhance the quality of the code as well as the productivity of the developer

## 2 Motivation

The rise of modern software engineering practices has highlighted the need for automated tools that detect and eliminate duplicate code. Duplicate code not only inflates the size of the codebase but also makes debugging and refactoring more error-prone.

The primary motivation behind this project was to create a lightweight, customizable tool that developers can run locally to detect duplicates in a variety of formats, including .js, .ts, .jsx, and .tsx. Additionally, the tool should allow users to categorize duplicates and save results for further analysis.

## 3 Project Details

### 3.1 Objective

The main objective of the project was to develop a command-line tool that:

1. Scans a specified directory for JavaScript and TypeScript files.
2. Identifies three types of code duplicates: exact, near, and structural.
3. Provides a user-friendly interface to choose the type of detection and optionally save the results to text files.

## 25 3.2 Implementation

### 26 3.2.1 File Scanning and Validation

27 The tool begins by prompting the user to input a directory path. Using file scanning utilities, it  
28 recursively searches for all files with supported extensions (**.js**, **.ts**, **.jsx**, **.tsx**). The **validateDirectory**  
29 function ensures the directory exists and is accessible, reducing runtime errors.

### 30 3.2.2 Duplicate Detection

31 The tool categorizes duplicates into the following types:

- 32 • **Exact Duplicates:** Identifies identical blocks of code using hashing algorithms (e.g., MD5).  
33 This is computationally efficient and provides an immediate match for identical code  
34 segments.
- 35 • **Near Duplicates:** Compares code similarity using the string-similarity library. A normaliza-  
36 tion step removes whitespace and comments before comparison to ensure accuracy. The  
37 similarity threshold (default: 0.8) can be adjusted to control sensitivity
- 38 • **Structural Duplicates:** Parses the code into an Abstract Syntax Tree (AST) and analyzes  
39 the structural logic. This enables detection of logically similar but syntactically different  
40 code.

### 41 3.2.3 Exact Duplicates

42 Exact duplicates are identical code blocks replicated across files or functions. These are often the  
43 easiest to detect and correct, as no analysis of syntax or logic is required.

#### 44 Detection Method

- 45 • A hashing algorithm, such as MD5, is used to compute a hash for each function. Identical  
46 code blocks will generate the same hash value.
- 47 • The system compares hashes instead of the actual code, making the process highly efficient.

48 The hash-based detection ensures that exact duplicates are identified in constant time per function,  
49 making it highly scalable.

#### 50 Example for exact duplicate

```
#file1.js and file5.js

function calculateRectangleArea(length, width) {
  return width * length; // Reordered but logically identical
}

The exact duplicate will print out:

{
  "function": "function calculateRectangleArea(length, width) {\n  \n    return width * length;\n}",
  "files": [
    "test\\TestingFilesEndingWithJS\\file2.js",
    "test\\TestingFilesEndingWithJS\\file5.js"
  ]
}
```

Figure 1: Exact Duplicate Example

### 51 3.2.4 Near Duplicates

52 Near duplicates are code segments that are logically identical or highly similar, but with slight  
53 differences in syntax or structure. These differences might include:

- 54 • Variable renaming.
- 55 • Minor formatting changes.
- 56 • Rearrangement of statements without altering functionality.

#### 57 Detection Method

- 58 1. **Normalization:** The tool preprocesses the code by:
  - 59 • Removing comments.
  - 60 • Normalizing whitespace.
  - 61 • Standardizing variable names where possible.
- 62 2. **String Similarity:** The normalized code is compared using a similarity algorithm (e.g.,  
63 string-similarity). A similarity score is computed, and matches exceeding a predefined  
64 threshold (default: 0.8) are considered near duplicates.

```
#File 6
function calculateRectangleArea(length, width) {
  if (length <= 0 || width <= 0) {
    console.error("Invalid dimensions!");
    return 0;
  }
  return length * width;
}

#File 7
function calculateRectangleArea(length, width) {
  if (length > 0 && width > 0) {
    return length * width;
  }
  throw new Error("Invalid dimensions!");
}
```

The near duplicate will print out with a similarity of 0.8203125:

```
{
  "function1": "function calculateRectangleArea(length, width) {if(
    length<=0||width<=0){console.error(\"Invalid dimensions
    !\");return 0;}return length*width;}",
  "function2": "function calculateRectangleArea(length, width) {if(
    length>0&&width>0){return length*width;}throw new
    Error(\"Invalid dimensions!\");}",
  "similarity": 0.8203125,
  "files": [
    "test\\TestingFilesEndingWithJS\\file6.js",
    "test\\TestingFilesEndingWithJS\\file7.js"
  ]
}
```

Figure 2: Near Duplicate Example

### 65 3.2.5 Structural Duplicates

66 Structural duplicates represent logically similar code blocks that may have completely different  
67 syntax. This is the most complex type of duplicate to detect, as it requires an understanding of the  
68 code's structure and semantics.

## 69 Detection Method

- 70 1. **Abstract Syntax Tree (AST) Parsing:** The code is converted into an AST, which represents  
71 the structure of the code
- 72 2. **Structural Comparison:** The tool compares AST nodes to identify similarity in logic. Even  
73 if the code is written differently, identical AST patterns indicate structural duplication.

```
The structural duplicate will normalize the function and print the output
{
  "normalizedFunction": "function_placeholder(placeholder,placeholder)_
    {\n_return_placeholder_placeholder;\n}",
  "originalFunctions": [
    {
      "file": "test\\TestingFilesEndingWithJS\\file4.js",
      "function": "function_subtract(a,b){\n_return_a-b;\n}"
    },
    {
      "file": "test\\TestingFilesEndingWithJS\\file4.js",
      "function": "function_add(a,b){\n_return_a+b;\n}"
    }
  ]
}
```

Figure 3: Structural Duplicate Example

### 74 3.2.6 User Interaction

75 The user interacts with the tool via prompts powered by the **inquirer** library:

- 76 1. Input the directory path to scan.
- 77 2. Choose the type of duplicate detection.
- 78 3. Decide whether to save the results to text files (**exact.txt**, **near.txt**, **structural.txt**).

### 79 3.2.7 Saving Results

80 The results of each detection type are saved in a structured JSON format in a text file. This enables  
81 easy integration with other tools for further analysis or visualization.

## 82 4 Results and Analysis

83 The tool was tested on a set of sample JavaScript and TypeScript files containing varying levels of  
84 duplication. The following observations were made:

- 85 1. **Exact Duplicates:** Accurately detected identical functions across files.
- 86 2. **Near Duplicates:** Effectively identified minor variations in code structure (e.g., variable  
87 renaming or formatting changes)
- 88 3. **Structural Duplicates:** Detected logical similarities, such as identical loops written with  
89 different syntax.

### 90 4.1 Performance

- 91 • For smaller datasets (<100 functions), results were instantaneous.
- 92 • For larger datasets, optimizations like hashing and normalization significantly reduced  
93 runtime, but further improvements are needed for structural comparisons on larger codebases.

## 94 5 Future Improvements

- 95 1. **Parallel Processing:** Utilize worker threads to handle pairwise comparisons for larger  
96 datasets.
- 97 2. **Enhanced Structural Analysis:** Implement more advanced AST analysis techniques to  
98 improve precision.
- 99 3. **Customization:** Provide APIs or CLI options to integrate the tool with CI/CD pipelines for  
100 continuous monitoring of code duplication.
- 101 4. **Extend for other code bases:** Support applications like python, C and so on.
- 102 5. **Extension:** Add a vs code so users can download and test it on their files.

## 103 6 Conclusion

104 This project shows how a tool can be developed to have a high degree of flexibility and efficiency as a  
105 code duplicate detection tool focused on JavaScript and TypeScript projects. With the integration of a  
106 good CLI interface, the application of detection algorithms can be made easy while at the same time  
107 being very powerful. There are however some problems that affect the scalability and the level of  
108 accuracy that the system can achieve, however the groundwork provided in this paper is a solid base  
109 for development. Detection of code duplication at the initial stage is very beneficial for the quality of  
110 software, its further support, and the efficiency of the developers.

## 111 References

- 112 <https://github.com/platisd/duplicate-code-detection-tool>  
113 <https://www.npmjs.com/package/inquirer>  
114 <https://www.npmjs.com/search?q=string20similarity>  
115 <https://ts-ast-viewer.com/>