

Anti-plagiarism Project

Eric Scott Freeman

Abstract—TODO

Index Terms—anti-plagiarism, plagiarism, fingerprinting, stylometry, Moss, JPlag, dupl, Autograder.

I. INTRODUCTION

THE University of Stavanger uses an application called Autograder, written by Heine Furubotten, to automatically grade students' programming assignments. Autograder works with GitHub. Whenever a student pushes their code to GitHub, Autograder will pull a copy of the code and run a series of tests on it. Autograder does not currently have support for plagiarism detection. The goal of this project is to integrate a few anti-plagiarism tools into Autograder, thereby helping the professors save time.

II. RELATED WORK

Unfortunately software plagiarism is a problem both in the classroom and in the workplace. A number of applications have been created to help detect this problem. While these tools can detect similarities in programs, the flagged files must still be manually examined to determine whether or not code was plagiarized.

There are several different general techniques that are used to look for plagiarism. The tools analyzed in this project used either fingerprinting or stylometry.

A. Fingerprinting

Several tools use a technique called fingerprinting to detect plagiarism. In fingerprinting algorithms, hashes of n -grams, substrings that are n characters, are saved and compared to help find plagiarism. Not all hashes are stored due to the large number that would be produced.

1) *Winnowing*: Moss uses a technique called winnowing to select which hashes to save [1]. In the winnowing algorithm, a window, selection of contiguous hashes, is used to help select which hashes to save. The smallest hash from a window is saved, and then the window moves one hash over. The smallest hash from the next window is often the smallest hash from the previous window. If so it is not saved again. Figure 1 shows an example of how winnowing works. The orange box represents the shifting window. The green box shows whenever a new hash is saved.

2) *Running-Karp-Rabin Greedy-String-Tiling*: JPlag uses Running-Karp-Rabin Greedy-String-Tiling (RKS-GST) to compare hashes of code in plagiarism detection [2]. RKS-GST was originally used in YAP3, another plagiarism detection tool. In RKS-GST, the Greedy String half of the algorithm forms pairs of substrings, each from a different string. Then the Karp-Rabin half of the algorithm hashes each substring in the pair [3]. This is done to help detect code reordering.

```
b4 12 03 56 8a 47 43 90 cb 34 5f a1
b4 12 03 56 8a 47 43 90 cb 34 5f a1
b4 12 03 56 8a 47 43 90 cb 34 5f a1
b4 12 03 56 8a 47 43 90 cb 34 5f a1
b4 12 03 56 8a 47 43 90 cb 34 5f a1
b4 12 03 56 8a 47 43 90 cb 34 5f a1
b4 12 03 56 8a 47 43 90 cb 34 5f a1
b4 12 03 56 8a 47 43 90 cb 34 5f a1
```

Fig. 1. Winnowing.

B. Stylometry

Another approach is to use code stylometry, which analyzes the style of writing or coding.

1) *Abstract Syntax Trees*: Caliskan-Islam, et al. use abstract syntax trees (ASTs) to compare the styles of authors [4]. Things that are easily changed in code, such as variable names, become leaves in the AST, while the structure of the tree is harder to change [4]. Figure 2 shows an abstract syntax tree of the code in Figure 3. Note how the leaves, or circular nodes, in Figure 2 are variable names, constants, and a function name.

Michal Bohuslávček's dupl application uses ASTs to find similarities in code [5]. It looks for any copies of code, not just plagiarism. So if a piece of code is duplicated even in the same file, it will test positive.

C. Supported Languages

Fingerprinting and string comparison techniques can be used to analyze source code written in languages other than their officially supported languages. Moss can analyze Go code, even though it is not technically supported. Since ASTs need to parse the code, applications which use them are stricter on which languages they can analyze. For example dupl only supports code written in Go. Table I shows the languages officially supported by several anti-plagiarism tools.

Tool	Java	Go	C	C++	C#	Python	others
Moss	✓		✓	✓	✓	✓	✓
JPlag	✓		✓	✓	✓		✓
Plaggie	✓						
SIM	✓		✓				✓
dupl		✓					

TABLE I
OFFICIALLY SUPPORTED LANGUAGES BY VARIOUS TOOLS

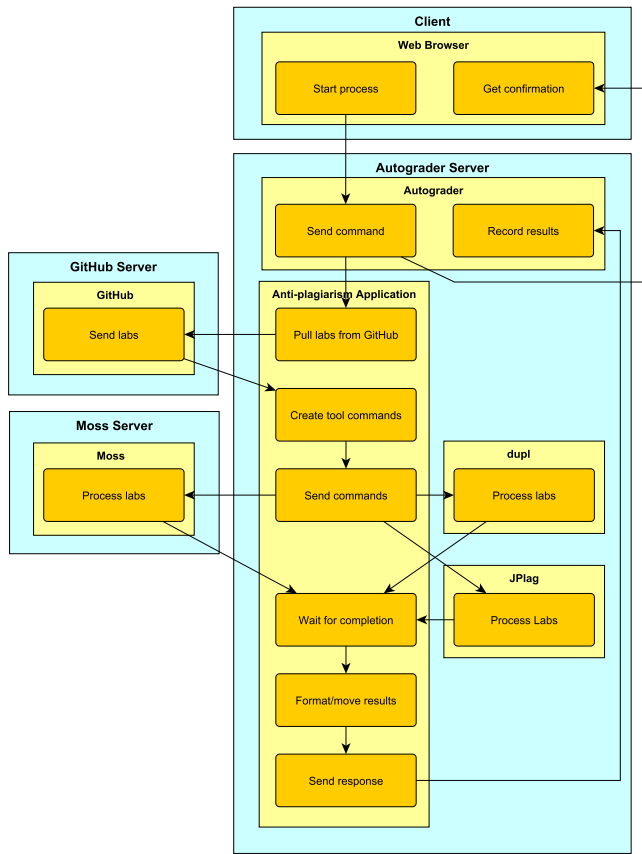


Fig. 5. Process

a gRPC response back to the calling application indicating if it was successful or not.

C. Integrating with Autograder

A few new fields needed to be added to the Autograder's Bolt database. Bolt stores entries as a key/value pair. Each lab assignment entry in each class needed to store the language the assignment is written in. Also each lab assignment entry for each student/group needed to store the results from each of the anti-plagiarism tools.

The user interface was also updated to allow teachers to assign a programming language to each lab (figure 6), to allow teachers to initiate the anti-plagiarism process (figure 7), and to show the results of the process. When the Autograder web service receives a request to start the anti-plagiarism process, it creates the gRPC command. If it was successful it sends the request and informs the user that the request was sent. Since it can take at least several minutes to process the data, it is best to inform the user immediately that things might take a while. Otherwise they could think that it is not working. Then Autograder starts a Go routine (separate thread) to call the anti-plagiarism software and process the results.

The results are stored in the database and also displayed to the teachers through the teacher panel in Autograder. The teacher panel shows a table of all the students in the class and their assignments. If any of the anti-plagiarism tools find evidence of plagiarism, the corresponding cell in the table will

Folder name lab 1

Deadline lab 1

Primary language lab 1

Fig. 6. Selecting language.

Individual Results

Build and test results from the stuc

Test Plagiarism

Fig. 7. Start anti-plagiarism test.

be colored a shade of red. The more tools finding plagiarism, the deeper the color. From the teacher panel, teachers go to a results page to look at more details about each lab. This page shows the results of each anti-plagiarism tool and buttons to show the plagiarized code. See figure 8.

Anti-Plagiarism Results

Moss results:	98%	Show Moss Details
JPlag results:	0%	
Dupl results:	True	Show dupl Details

Fig. 8. Individual results.

While adding the new functionality to Autograder, care was taken not to break any existing functionality. Unit tests were rerun to check that nothing was broken.

IV. RESULTS

TODO

V. ANALYSIS

TODO

Getting three separate tools to behave similarly using a common interface was a bit difficult. Moss was the original anti-plagiarism tool included in the project before it was decided to add other tools. Therefore the results from Moss should be considered more reliable. Moss's threshold says to ignore code that is similar in a certain number of files. So if the threshold is ten and a piece of code appears at least ten times, it is not considered plagiarism. This is very helpful when an instructor gives students a partially completed file or a framework to follow. The other tools produce false positives.

JPlag can produce false positives in two ways in this project. First, JPlag remembers the results from previous executions. So if a student uses a piece of code in one lab and then reuses the same piece of code in another lab, JPlag can detect this as plagiarism. The second way is JPlag is not ignoring code provided by the instructor. While it is possible to tell JPlag to ignore certain files, this would break the common interface shared among the tools. Also it would require the instructors to

format the assignments in such a ways that all of the instructor provided code is separate from all of the student code.

dupl was not designed to detect plagiarism, only duplicate pieces of code. So it will detect all the instructor provided code as plagiarism. It can even detect two pieces of code in the same file as plagiarism. For example the two functions in figure 9 are from the same file and are shown to be duplicates. Increasing the threshold can help lessen this problem.

#1 found 2 clones

example.go:3

```
func sample1(int x) int {
    return 2 * x
}
```

example.go:7

```
func sample2(int x) int {
    return 3 * x
}
```

Fig. 9. Sample dupl results from same file.

Since dupl is only looking for duplicate code and not plagiarism, it does not provide a number indicating the percentage of code it thinks is plagiarized. It also places all the results in a single HTML file rather than splitting them up into separate files.

The results from Moss should be considered first, and the results from the other two tools should be considered as supplemental.

VI. FUTURE WORK

Currently the anti-plagiarism application is only intended to be called locally, so there is no encryption enabled in gRPC. If it is to be called remotely, then transport security needs to be added in the anti-plagiarism application and in Autograder. If the application was called remotely, the result files would need to be sent to the other machine. Additional anti-plagiarism tools can be added later.

VII. CONCLUSION

TODO

REFERENCES

- [1] S. Schleimer, D. S. Wilkerson, and A. Aiken, "Winnowing: Local algorithms for document fingerprinting," in *Proceedings of the 2003 ACM SIGMOD international conference on Management of data*, pp. 76–85, ACM, 2003.
- [2] L. Prechelt, G. Malpohl, and M. Philippsen, "Jplag: Finding plagiarisms among a set of programs," *Journal of Universal Computer Science*, vol. 8, 2002.
- [3] M. Wise, "Yap3: Improved detection of similarities in computer program and other texts," in *SIGCSE '96 Proceedings of the twenty-seventh SIGCSE technical symposium on Computer science education*, pp. 130–134, 1996.
- [4] A. Caliskan-Islam, R. Harang, A. Liu, A. Narayanan, C. Voss, F. Yamaguchi, and R. Greenstadt, "De-anonymizing programmers via code stylometry," in *24th USENIX Security Symposium (USENIX Security 15)*, (Washington, D.C.), pp. 255–270, USENIX Association, 2015.

- [5] M. Bohuslávsek, "dupl." <https://github.com/mibk/dupl>. Accessed: 2015-09-15.

Eric Scott Freeman is currently a master's student at the University of Stavanger. He received his Bachelor of Science degree in Computer Science from Midwestern State University in 2004. He has also worked as a software developer for several companies including RadioShack, Cisco Systems, and TelStrat.