

Cross-Language Source Code Plagiarism Detection using Explicit Semantic Analysis and Scored Greedy String Tiling

Tomáš Foltýnek¹, Richard Všíanský¹, Norman Meuschke^{2,3}, Dita Dlabolová¹, Bela Gipp^{2,3}

¹Mendel University in Brno, Czech Republic, {tomas.foltynek | xvsiansk | dita.dlabolova}@mendelu.cz

²University of Wuppertal, Germany, {meuschke | gipp}@uni-wuppertal.de

³University of Konstanz, Germany

ABSTRACT

We present a method for source code plagiarism detection that is independent of the programming language. Our method EsaGst combines Explicit Semantic Analysis and Greedy String Tiling. Using 25 cases of source code plagiarism in C++, Java, JavaScript, PHP, and Python, we show that EsaGst outperforms a baseline method in identifying plagiarism across programming languages.

CCS CONCEPTS

• Information systems~Information retrieval~Specialized information retrieval

KEYWORDS

Source Code Plagiarism Detection, Explicit Semantic Analysis, Greedy String Tiling

ACM Reference format:

Tomáš Foltýnek, Richard Všíanský, Norman Meuschke, Dita Dlabolová and Bela Gipp. 2020. Cross-Language Source Code Plagiarism Detection via Explicit Semantic Analysis and Scored Greedy String Tiling. In *Proceedings of the ACM/IEEE Joint Conference on Digital Libraries (JCDL'20), August 1 - 5, Virtual event, China*. ACM, New York, NY, USA, 2 pages. <https://doi.org/10.1145/3383583.3398594>

1 Introduction & Related Work

Source code plagiarism detection (SCPD) is an effective deterrent to undue reuse of code in programming assignments, which are common in computer science and related study programs.

Many SCPD methods focus on specific programming languages by employing approximate string matching to identify similar programs [1]. Other methods additionally analyze the structure or semantics of source code [2]. Some methods addressed the cross-language SCPD task using Latent Semantic Analysis [3].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from Permissions@acm.org.
JCDL'20, August 1–5, 2020, Virtual event, China.

© 2020 Copyright is held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-7585-6/20/06...\$15.00.

DOI: <https://doi.org/10.1145/3383583.3398594>

We presume that plagiarists trying to obfuscate reused code preserve the semantics of the identifiers, comments, and other tokens. Thus, we see a semantic analysis as promising for devising a language-independent SCPD method. Therefore, we adapt Explicit Semantic Analysis (ESA) [4], a well-established semantic analysis method, and Greedy String Tiling to the SCPD use case. ESA models text as concept vectors. The concepts are the topics in a knowledge base, which is typically Wikipedia or another encyclopedia. The vector components reflect the relevance of the modeled text for each of the concepts. Greedy String Tiling (GST) is an algorithm with near-linear complexity to find all individually longest substring matches in two strings [5].

2 Method

To perform ESA, we used the EsaPlag system [6] and thirty thousand articles from the categories “Computer programming” and “Fields of mathematics” in the English Wikipedia. Using the title of articles as concepts, we represented each document, i.e., a computer program, by deriving a concept vector for each term in the document. To maintain all semantic information of documents, we only removed line breaks before forming the vectors.

To compute the similarity of documents, we devised the Scored GST algorithm that determines the longest sequence of semantically similar terms. Other than GST, Scored GST matches not only identical elements but all elements whose similarity is above a threshold. Here, the concept vectors for document terms are the elements, whose similarity we computed via the cosine measure. We set the cosine similarity above which we consider concept vectors a match to 50% and the final score above which we report results to 5% as this value maximized the F1 score.

3 Experiments

To evaluate our SCPD method, we created a dataset of simulated source code plagiarism. We implemented a basic programming assignment – a calculator supporting basic arithmetic operations – in the five most common languages on GitHub, i.e., C++, Java, JavaScript, PHP, and Python. For each language-specific implementation, we created four plagiarized versions using the following obfuscation methods: (1) renaming identifiers, (2) renaming identifiers by converting camelCase to snake_case, (3) restructuring the code, and (4) reusing half of the code. To test for false positives, we used unrelated code with no semantic matches.

As a baseline, we used the text matching system Anton [5].



Citation for this Paper

Foltynek, T. & Vsiansky, R. & Meuschke, N. & Dlabolova, D. & Gipp, B., “Cross-Language Source Code Plagiarism Detection using Explicit Semantic Analysis and Scored Greedy String Tilling,” in Proceedings of the ACM/IEEE Joint Conference on Digital Libraries (JCDL), Aug. 2020, DOI: [10.1145/3383583.3398594](https://doi.org/10.1145/3383583.3398594).

BibTeX:

```
@inproceedings{FoltynekVMD20,  
title = {Cross-{Language} {Source} {Code} {Plagiarism} {Detection} using  
{Explicit} {Semantic} {Analysis} and {Scored} {Greedy} {String} {Tilling}},  
doi = {10.1145/3383583.3398594},  
booktitle = {Proceedings of the {ACM}/{IEEE} {Joint} {Conference} on {Digital}  
{Libraries} ({JCDL})},  
author = {Foltynek, Tomas and Vsiansky, Richard and Meuschke, Norman and  
Dlabolova, Dita and Gipp, Bela},  
month = aug,  
year = {2020}  
}
```

RIS:

```
TY - CONF  
TI - Cross-Language Source Code Plagiarism Detection using Explicit Semantic  
Analysis and Scored Greedy String Tilling  
AU - Foltynek, Tomas  
AU - Vsiansky, Richard  
AU - Meuschke, Norman  
AU - Dlabolova, Dita  
AU - Gipp, Bela  
C3 - Proceedings of the ACM/IEEE Joint Conference on Digital Libraries (JCDL)  
DA - 2020/08//  
PY - 2020  
DO - 10.1145/3383583.3398594  
ER -
```

Related Publications: www.gipp.com/pub