

InvAASTCluster: On Applying Invariant-Based Program Clustering to Introductory Programming Assignments

Pedro Orvalho^{a,*}, Mikoláš Janota^b and Vasco Manquinho^a

^aINESC-ID, Instituto Superior Técnico, Universidade de Lisboa, Rua Alves Redol 9, Lisboa, 1000-029, Portugal

^bCIIRC, Czech Technical University in Prague, Dejvice, Prague, 16000, Czechia

ARTICLE INFO

Keywords:

Program Clustering
Program Invariants
Program Equivalence
Program Repair
Programming Education
MOOCs

ABSTRACT

Due to the vast number of students enrolled in programming courses, there has been an increasing number of automated program repair techniques focused on introductory programming assignments (IPAS). Typically, such techniques use program clustering to take advantage of previous correct student implementations to repair a new incorrect submission. These repair techniques use clustering methods since analyzing all available correct submissions to repair a program is not feasible. However, conventional clustering methods rely on program representations based on features such as abstract syntax trees (ASTs), syntax, control flow, and data flow.

This paper proposes INVAASTCLUSTER, a novel approach for program clustering that uses dynamically generated program invariants to cluster semantically equivalent IPAS. INVAASTCLUSTER's program representation uses a combination of the program's semantics, through its invariants, and its structure through its anonymized abstract syntax tree (AASTs). Invariants denote conditions that must remain true during program execution, while AASTs are ASTs devoid of variable and function names, retaining only their types. Our experiments show that the proposed program representation outperforms syntax-based representations when clustering a set of correct IPAS. Furthermore, we integrate INVAASTCLUSTER into a state-of-the-art clustering-based program repair tool. Our results show that INVAASTCLUSTER advances the current state-of-the-art when used by clustering-based repair tools by repairing around 13% more students' programs, in a shorter amount of time.


1. Introduction


Nowadays, thousands of students enroll every year in programming-oriented Massive Open Online Courses (MOOCs) [19]. On top of that, due to the recent pandemic situation, even small-sized programming courses are being taught online. Providing feedback to novice students in introductory programming assignments (IPAS) in these courses requires substantial effort and time by the faculty. Hence, there is an increasing need for automated semantic program repair frameworks [19, 41, 38, 27, 39, 67, 18, 68] capable of providing automated, comprehensive, and personalized feedback for students' incorrect solutions to programming assignments.

Over the last few years, several program repair tools [19, 64, 49, 23] have appeared that use a large number of diverse correct implementations submitted for each IPA by previously enrolled students. Given an incorrect student submission, these frameworks use clustering methods to find the most similar correct submission from previous years to provide a minimal set of repairs to the student. Using the same reference implementation to fix all incorrect programs can potentially generate a large set of repairs. On the other hand, having a similar correct implementation allows computing a smaller set of repairs. However, comparing on the fly all previous correct student submissions against a given incorrect submission is not feasible.

To tackle this problem, different program clustering approaches have been used in program repair tools, which enable focusing only on the representatives of each cluster. CLARA [19] clusters the correct programs based on their dynamic equivalence [40] and control flow, i.e., the order in which program statements, instructions, and function calls are executed. SARFGEN [64] computes program representations based on each program's abstract syntax tree.

*Corresponding author

 pmorvalho@tecnico.ulisboa.pt (P. Orvalho); mikolas.janota@cvut.cz (M. Janota);
vasco.manquinho@tecnico.ulisboa.pt (V. Manquinho)

 <https://arsr.inesc-id.pt/~pmorvalho> (P. Orvalho); <http://people.ciirc.cvut.cz/~janotmik> (M. Janota);
<https://arsr.inesc-id.pt/~vmm> (V. Manquinho)

ORCID(s): 0000-0002-7407-5967 (P. Orvalho); 0000-0003-3487-784X (M. Janota); 0000-0002-4205-2189 (V. Manquinho)

SEMCLUSTER [49] uses each program’s control and data flow. A program’s data flow tracks the number of occurrences of consecutive values a variable takes during its lifetime.

The problem of program equivalence, i.e., deciding if two programs are equivalent, is undecidable [53, 8, 58]. On that account, finding an adequate representation for programs that performs well on program clustering is a challenging problem. The previously-mentioned program representations used in the field of program repair may be brittle. For instance, consider two programs that calculate the sum of natural numbers from 1 to a given number, one using a while-loop and the other a for-loop. Despite producing the same result, their syntactic and structural differences pose challenges for conventional program representations to recognize their semantic equivalence, as we will show in Section 2. To address this problem, we propose to use dynamically-generated program invariants to cluster semantically equivalent programs, overcoming some of the identified weaknesses. A program invariant is a condition that must always be true at a given step of the program during its execution (see Section 3.2). Program invariants are usually used to assert assurances throughout a program (assertions).

This paper proposes to leverage the information of a program’s structure using its *abstract syntax tree* (AST) together with semantic information provided by its invariants. Previous research has been conducted regarding using invariants to promote patch diversity (i.e., diversity in the set of possible repairs to a given incorrect program) on search-based program repair [6, 12, 65]. These works use DAIKON [13] to generate invariant sets for each possible patch. DAIKON is a system that infers likely dynamically generated invariants observed over several program executions. Therefore, these invariants depend on the program executions. Nevertheless, previous work [6] showed promising results in using invariants to semantically cluster patches to provide the user with a semantic reason for a set of similar patches.

This paper presents a novel approach for clustering introductory programming assignments (IPAS) leveraging their sets of invariants. Our approach for clustering IPAS also takes into account each program’s code and *anonymized abstract syntax tree* (AAST). AASTs are essentially ASTs stripped of variable and function identifiers, preserving only their respective types (see Section 4.2). The main contribution of this work is a vector representation of programs based on their invariants and AASTs, bringing together their semantic and syntactic features. The proposed clustering technique has been implemented in a framework INVAASTCLUSTER. This tool has been designed as an independent clustering tool. Therefore, it can be used to help evaluate students’ submissions for IPAS by clustering semantically equivalent solutions for programming exercises. However, INVAASTCLUSTER can also be easily integrated into any clustering-based program repair tool for IPAS. Furthermore, INVAASTCLUSTER can even be used in a plagiarism detection tool, like MOSS [56].

Figure 1 shows the generic architecture of clustering-based program repair frameworks [19, 49, 64]. These frameworks receive an incorrect student submission, a test suite, and a collection of N correct student submissions for the same IPA. For scalability concerns, these frameworks eliminate, through clustering techniques, semantically equivalent solutions, i.e., dynamically equivalent correct programs, given the provided input-output test suite. Those clustering approaches try aggregating the set of N correct solutions into K semantically different clusters ($N \gg K$). Finally, the repair tool uses these K clusters’ representatives to repair the provided incorrect student submission. As Figure 1 shows, INVAASTCLUSTER can be used as the clustering technique of those clustering-based program repair tools. However, some program repair tools [2, 33] use a single reference implementation provided by the lecturer to repair a student’s program. Typically, these tools can only use one correct implementation to repair each program. Therefore, INVAASTCLUSTER was designed to be also capable of finding on a set of correct student submissions which submission is the closest correct solution to the incorrect program. Thus, INVAASTCLUSTER can suggest a specific reference implementation for each incorrect submission that may require fewer changes to fix the program.

We evaluate INVAASTCLUSTER on C-PACK-IPAS [44], a real-world student programs developed during a university introductory programming course. Experimental results show that the proposed invariant-based representation improves upon syntax-based representations when performing program clustering. Additionally, we integrate INVAASTCLUSTER into CLARA, a clustering-based program repair tool, in order to compare our clustering technique against CLARA’s clustering method, which is the current publicly available state-of-the-art method for clustering IPAS.

To summarize, this paper makes the following contributions:

- We propose a novel and efficient approach for clustering submissions for introductory programming assignments (IPAS) based on the submissions’ sets of invariants and AASTs representations.
- We present a study showing the results of using our program clustering tool, INVAASTCLUSTER, on a set of 1620 real-world IPAS correct submissions to show the effectiveness of invariant-based program clustering.

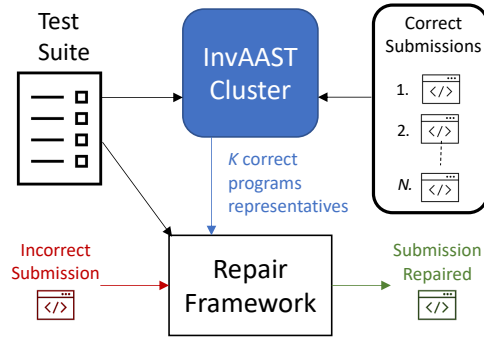


Figure 1: Clustering-based Program Repair.

- We compare INVAASTCLUSTER with the clustering method used by the currently available state-of-the-art program repair tools. Experimental results show that INVAASTCLUSTER outperforms state-of-the-art clustering methods, allowing clustering-based program repair tools to fix around 13% more IPAs in a shorter amount of time.
- The INVAASTCLUSTER framework is publicly available on GitHub at <https://github.com/pmorvalho/InvAAST-Cluster>.

The structure of the remainder of this paper is as follows. First, Section 2 illustrates the strengths of using invariants for program representation. Section 3 presents important concepts used throughout this manuscript and describes how to gather and represent sets of program invariants. Section 4 discusses several program representations, including a new invariant-based program representation. Section 5 discusses the implementation of INVAASTCLUSTER. Section 7 presents the experimental evaluation that supports our claim that invariant-based program representations are beneficial to cluster programming assignments semantically. Finally, Section 8 presents the related work, and the paper concludes in Section 9.

2. Motivation

Current program representations for repairing students' programming assignments leverage certain program features, such as code syntax [20], abstract syntax tree [64], control flow [19], and data flow [49], to encode each program into a vector representation. However, all of these features have some weaknesses when we want to cluster programs based on their semantics.

Example 1. Consider the following two programs written in C, that compute the sum of all the natural numbers from 1 to a given number n , i.e., $\sum_{i=1}^n i$.

```

1  int n, sum = 0, i;
2  scanf("%d", &n);
3  i = 0;
4  while(i < n) {
5      i++;
6      sum = sum + i;
7  }
8  printf("%d\n", sum);

```

```

1  int j, n, s = 0;
2  scanf("%d", &n);
3
4  for(j = n; j > 0; j--)
5  {
6      s = j + s;
7  }
8  printf("%d\n", s);

```

Observe that the program on the left, in Example 1, uses a while-loop that iterates over the natural numbers from 0 to n . The program on the right uses a for-loop that iterates from n to 0 in decreasing order. However, both programs are semantically equivalent since both have the same result. Nevertheless, if we build a program representation using the programs' syntax or abstract syntax trees, both programs will have very different representations. In terms of syntax, the names of the used variables (e.g. i , j , s , sum) and structures (e.g. `while`, `for`) are different. Additionally, in terms of data flow and dynamic equivalence, both programs are also different since, for example, the values assigned

to the variable i go from 0 to n in the first program while in the other the variable j is assigned the same values but in decreasing order.

Consider that the variable n is always assigned to a natural number, $n > 0$. If a dynamic invariant detector (e.g. DAIKON [13]) is used, the following set of invariants is observed:

- In the first program, at each iteration of the while-loop: $n > 0; sum \geq 0; 0 \leq i \leq n$.
- In the second program, at each iteration of the for-loop: $n > 0; s \geq 0; 0 \leq j \leq n$.

Therefore, after renaming some variables ($sum \rightarrow s; i \rightarrow j$), these two sets of invariants would be considered equivalent. Hence, using sets of invariants allows finding semantically equivalent programs that can differ in their syntax and/or data flow.

Hence, this paper aims to improve the semantic representation of programming exercises using their sets of invariants. These invariants are dynamically detected by DAIKON [13], at the beginning and at the end of each scope, over several program executions using a predefined set of test cases for each programming assignment. In addition to the set of invariants, which provides semantic information about a program, we also leverage the information of a program's structure to its anonymized abstract syntax tree (AAST), i.e., an AST after removing all the variables' names.

CLARA. Furthermore, in this paper, we compare our clustering approach against CLARA's clustering method since, to the best of our knowledge, CLARA [19] stands as the sole publicly accessible state-of-the-art clustering-based repair tool for repairing IPAS. CLARA leverages correct solutions from past years' submissions for a programming assignment to suggest potential semantic repairs for an erroneous program submitted by a new student.

CLARA's Clustering Approach CLARA assigns two programs to the same cluster if they share identical control flow structures and possess a bijective mapping between their variables [10]. However, if there is any deviation in the control flow graphs or a lack of bijective relation between variables, CLARA returns a 'structural mismatch error', resulting in these programs not being clustered together. For each cluster generated, CLARA maintains a json file containing all relevant information about the programs, including expressions and variables.

Revisiting the programs illustrated in Example 1, CLARA does not detect them as matching. Furthermore, when clustering both programs using CLARA's method, they are assigned to separate clusters.

CLARA's Repairing Process. To repair an incorrect program, CLARA receives either a single or multiple correct programs, which can be representative of clusters produced by CLARA itself. In such cases, CLARA also necessitates access to all pertinent information regarding the programs contained within each cluster, stored in individual json files. CLARA proceeds to generate a series of repairs for each cluster autonomously, wherein a repair entails adjusting a program to align an expression from the incorrect submission with an expression from a program within the cluster. It is important to emphasize that CLARA's repair suggestions are confined to modifying program expressions and do not involve altering control flow. If CLARA encounters an incorrect program whose control flow does not precisely match any of the correct programs provided, it flags a "Structural Mismatch" error, indicating an inability to repair the program. Moreover, CLARA is unable to use the programs outlined in Example 1 for repairs, as it does not identify them as matching.

In this paper, our focus is not on enhancing CLARA's repair procedure. Instead, we aim to compare our clustering approach against CLARA's clustering method, and evaluate CLARA's repair performance when utilizing its own clusters versus INVAASTCLUSTER's clusters of programs.

3. Syntax Trees and Invariants

This section provides some background on syntax trees and program invariants that will be used throughout this paper.

3.1. Definitions

This section provides some definitions that will be used throughout this paper.



Figure 2: A small example of an AST and an AAST for the variable declaration, `int i`. An integer variable with identifier `i`.

Definition 3.1 (Context-free Grammar (CFG)). A *context-free grammar* \mathcal{G} is a 4-tuple (V, Σ, R, S) , where V is the set of non-terminals symbols, Σ is the set of terminal symbols, R is the set of rules and S is the start symbol. A CFG describes all the strings permitted in a certain formal language [22].

Definition 3.2 (Domain-Specific Language (DSL)). A Domain-specific Language (DSL) is a tuple $(\mathcal{G}, \text{Ops})$, where \mathcal{G} is a context-free grammar ($\mathcal{G} = (V, \Sigma, R, S)$) and Ops is the semantics of DSL operators. The CFG \mathcal{G} has the rules to generate all the programs in the DSL. The semantics of DSL operators is necessary to analyze conflicts and make deductions.

Definition 3.3 (Abstract Syntax Tree (AST)). An *abstract syntax tree* (AST) is a syntax tree in which each node represents an operation, and the children of the node represent the arguments of the operation for a given programming language described by a Context-free Grammar [22]. An AST depicts a program's grammatical structure [4].

Figure 2a presents a small example of the AST representation for the variable declaration `int i`.

Definition 3.4 (Anonymized Abstract Syntax Tree (AAST)). An *anonymized abstract syntax tree* (AAST) is an AST in which nodes that contain identifiers are anonymized, i.e., a node's identifier (name of a function or variable) is replaced by a special token (`ID`).

Figure 2b shows the AAST representation for the same declaration presented previously, `int i`.

Definition 3.5 (Program Invariant). A program invariant is a logical condition that remains true throughout the execution of a program, regardless of its state or inputs. It serves as a guarantee or assertion about the behavior of the program at specific points or during certain operations [11]. Formally, a program invariant can be defined as a predicate P over the program state variables such that for all program states s reachable during the execution of the program, $P(s)$ holds true.

Consider again, the first program in Example 1, $n > 0$, $sum \geq 0$, and $0 \leq i \leq n$ are program invariants that are always satisfied during the execution of the while-loop.

Definition 3.6 (Bag of Words (BoW)). A *Bag of Words* (BoW) representation [21] is a vector representation where a tokenized sentence is represented as a bag of its words in a vector. The vector representation contains information on the number of times each token in the language appears in the sentence. Note that this model does not take into consideration the language's grammar and even word order. The tokenization step divides a string into n -grams, which are sub-sequences of the original string of n items.

The following example presents a small illustration of a vector representation of a phrase using a BoW model.

Example 2. Let B be a bag of words model computed using the following sentences: $\{ 'aa', 'e i', 'a e i o u', 'o i' \}$. Given the phrase $p = 'a i a u'$, the vector representation of p is, $B(p) = [0.5, 0.0, 0.25, 0.0, 0.25]$. The size of $B(p)$ is 5 since 5 is the size of the vocabulary of B . For each entry s of $B(p)$, $B(p)[s]$ corresponds to the percentage of p that is equal to s . For example, the symbol a appears twice in a four-symbol phrase. Hence $B(p)[a] = 0.5$.

3.2. Program Invariants

Program invariants are conditions that must always be true at a given point during a program's execution. Dynamically generated program invariants are *likely invariants* observed during several program executions for a given program. The dynamically generated set of program invariants provides information about a program's behavior, i.e., its semantics. If two programs share the same program invariants, they are likely semantically equivalent. Hence, an invariant-based representation of programs should allow us to find out which student submissions in a given programming assignment have the same or similar behavior.

In order to compare two sets of program invariants, a relation between the variables in both sets is required. We propose to rename all the variables in a program based on the variables' type and usage. All the variables are renamed the first time they are assigned to some value in a program. The variable's new name is a concatenation between its type and a counter for how many variables have already been renamed in the program. With this technique of variable renaming, two programs' sets of invariants can be easily compared. This method is very simple and fragile, although IPAS are usually relatively small and simple imperative programs. Therefore, this naive approach should work for IPAS. Example 3 shows two programs whose variables were renamed using this variable renaming method.

Example 3. Consider again the programs presented in Example 1, it is important to note that all variables are renamed based on their usage. Specifically, each variable is renamed the first time it is assigned a value in the program. For instance, in the first program, n is renamed to int_2 since it is the second variable to be used, in the `scanf` (line 2). After renaming all variables based on their usage, the following mapping of variables for the first program is obtained: $\{sum \rightarrow int_0; n \rightarrow int_1; i \rightarrow int_2\}$. Applying the same procedure to the second program yields the mapping $\{s \rightarrow int_0; n \rightarrow int_1; j \rightarrow int_2\}$. Thus, the two programs after renaming are as follows:

<pre> 1 int int1, int0 = 0, int2; 2 scanf("%d", &int1); 3 int2 = 0; 4 while(int2 < int1){ 5 int2++; 6 int0 = int0 + int2; 7 } 8 printf("%d\n",int0); </pre>	<pre> 1 int int2,int1,int0 = 0; 2 scanf("%d", &int1); 3 4 for(int2 = int1; int2 >= 0; int2--) 5 { 6 int0 = int2 + int0; 7 } 8 printf("%d\n", int0); </pre>
--	--

Hence, the set of invariants of both cycles (for and while) is the same: $\{int_1 > 0; int_0 \geq 0; 0 \leq int_2 \leq int_1\}$.

In this work, we use DAIKON [13] to compute dynamically-generated likely invariants observed across multiple program executions for each student submission, employing a predefined set of input-output tests for each programming assignment. DAIKON's default format for program invariants combines Java and mathematical logic, aiming to convey meaning concisely to programmers.

To adapt DAIKON for small imperative C programs, we initially apply a method for variable renaming to all student submissions. Next, we inject empty functions into each scope and pass the variables of the respective scope as parameters. A scope is defined as each block of statements without branching, and in cases of nested scopes, we include all variables available in the parent scopes as parameters as well. We adhere to conventional methods of modeling control flows, such as, computing invariants before a loop, before a loop guard, inside a loop guard, inside the loop, and after the loop.

Finally, DAIKON is executed using all input tests for each programming assignment. The dynamically-generated invariants produced by DAIKON are stored for each program's structure or scope (e.g., if statements, loops, blocks). We do not specifically ask DAIKON to generate any type of invariant. The only type of invariants we turned off is the "OneOf" invariants (e.g., "x is OneOf {1,2}") that may cause overfitting to the test suite. Example 4 presents the two programs, previously presented in Example 3, after being injected with empty functions in each program scope.

Example 4. For Daikon [13] to work on small imperative programs, we have to call empty functions at the beginning of each scope and pass all the visible variables, in that scope, as parameters.

Consider again the programs presented in Example 3, after renaming all the variables based on their usage. These programs, after being injected with empty functions, would look like the following programs:


```

1  scope_1();
2  int int1, int0 = 0, int2;
3  scanf("%d", &int1);
4  int2 = 0;
5  while(int2 < int1,
6  while_1(int0, int1, int2)){
7      scope_2(int0, int1, int2);
8      int2++;
9      int0 = int0 + int2;
10 }
11 scope_3(int0, int1, int2);
12 printf("%d\n", int0);

```

```

1  scope_1();
2  int int2, int1, int0 = 0;
3  scanf("%d", &int1);
4
5  for(int2 = int1; int2 >= 0,
6  for_1(int0, int1, int2); int2--)
7  {
8      scope_2(int0, int1, int2);
9      int0 = int2 + int0;
10 }
11 scope_3(int0, int1, int2);
12 printf("%d\n", int0);

```

4. Program Representations

As the primary objective of this work is to cluster programs based on semantic and syntactic features, each program is represented as a feature vector. In particular, we propose to use a *bag of words (BoW)* model (see Definition 3.6). Using BoW models, we generate vector representations for each student submission based on several features. These features may include the Abstract Syntax Tree (AST), set of invariants, or even the program code. It is also possible to combine several of these features. Next, all the vector representations used in this work are described.

4.1. Syntax Vectors

The syntax vector program representation is the simplest to compute since it is based solely on the program syntax (code). In the interest of comparing the syntax of programs independently of the variables' names, first, all the programming solutions are renamed using the method described in Section 3.2. Next, all the student submissions are tokenized, and a vocabulary with all the available tokens is obtained. Then, vectors for each student submission are created, where the i^{th} entry is the number of times the i^{th} word of the vocabulary appears in the program. Finally, the numbers of occurrences in these vectors are normalized.

4.2. Anonymized Abstract Syntax Tree Vectors

An alternative representation is to compute a bag of words using the strings of the abstract syntax trees (ASTs) of all submissions for a given programming assignment. This representation has already been used in program clustering [64, 25]. However, we represent each AST as a string and remove all names of variables and functions, keeping only their respective types in the AST. Thus, for each submission, we have an *anonymized abstract syntax tree* (AAST) (see Definition 3.4). With these AASTs, we keep the information about a program's structure, ignoring the name of its variables. The information about a program's structure is kept since an AAST contains all the non-terminal symbols of the language's grammar. Next, a vocabulary is built with the tokens present in all submissions, and a normalized vector representation for each AAST is computed.

4.3. Invariant Vectors

Another approach is to use an invariant-based vector representation. In this case, we apply the bag of words model to the set of invariants of the programs. We gather all program invariants as described in Section 3.2. Previous work on the use of invariants to detect semantic similarity between possible patches to a program [6] showed that using string distance measure between invariant sets had similar results and was more efficient than computing the logical similarity between their corresponding sets of program invariants. Therefore, we represent our invariants in the form of strings. However, instead of using a string distance measure between invariant sets (e.g., Levenshtein edit distance [30]), we create a bag of words model with those sets of invariants.

4.4. Combination of Program Features

Finally, observe that these vector representations (*Syntax*, *AAST*, *Invariants*) can be combined, thus taking advantage of using several types of features. For example, we use a bag of words in our work using the program's AST and the sets of invariants. In this case, first, we build two BoW representations independently, one based on AASTs and another one based on invariants. Then, we concatenate, for each submission, the submission's vector representations

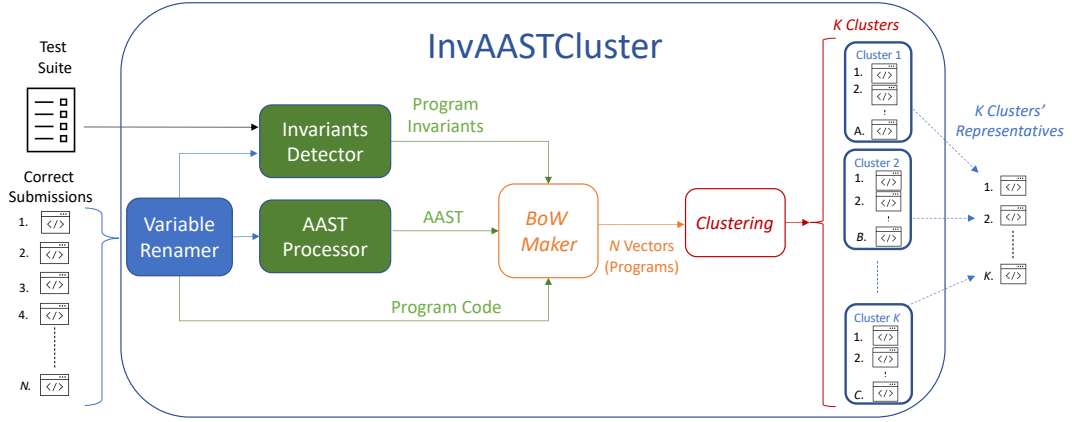


Figure 3: The high-level overview of INVAASTCLUSTER.

using the two BoWs, achieving a vector representation based on the program’s AAST and set of invariants. Both vectors, from the Invariants BoW and the AAST BoW, are normalized before concatenating. Therefore, the invariants and the AASTs have an equal contribution to the AAST + Invariants BoW. The program syntax was not included in this last representation since the BoW based on syntax has a large vocabulary that generates vectors that are too sparse.

5. Implementation

This section presents the implementation of our program clustering technique. We implemented the proposed approach in the tool INVAASTCLUSTER (**I**nvariants and **AAST** Program **C**lustering). INVAASTCLUSTER is publicly available on GitHub at <https://github.com/pmorvalho/InvAASTCluster>. Figure 3 shows the overall architecture of INVAASTCLUSTER. Given a set of N correct submissions and a test suite, INVAASTCLUSTER computes K clusters of programs ($N \geq K$) and returns the set of K clusters’ representatives, i.e., the set of correct programs that are closest to the center of each one of the K clusters. INVAASTCLUSTER is divided into six main modules: variable renamer, invariants detector, AASTs processor, bag of words (BoW) maker, clustering procedure, and the selection of each cluster representative.

Variable Renamer. In this module, INVAASTCLUSTER renames all variables of each one of the N given correct submissions. All variables are renamed based on their usage in each program, as explained in Section 3.2. INVAASTCLUSTER uses `pycparser`¹ to find all variables in a program. Then, when a variable is first used in the program (e.g. assignment) that variable receives a new name considering the variable’s type.

Invariants Detector. INVAASTCLUSTER uses DAIKON [13] to compute dynamically-generated invariants for a given test suite (see Section 3.2). After all the variables have been renamed, this module produces a set of invariants for each program’s scope using the provided test suite. All these sets of invariants are then sent to the BoW maker module.

AAST Processor. In this step, INVAASTCLUSTER also uses `pycparser` to compute a program’s abstract syntax tree (AST). Additionally, INVAASTCLUSTER removes all the variables’ and functions’ identifiers from the AST to transform the program’s AST into an anonymized abstract syntax tree (AAST), conserving only the program’s structure.

Bag of Words (BoW) Maker. This module receives three sets as input: (1) the set of correct program submissions with all their variables renamed from the Variable Renamer module; (2) all the program’s AASTs from the AAST processor, and (3) the set of the programs’ dynamically-generated invariants. The *BoW Maker* computes the bag of words (BoW) model (see Definition 3.6) that is going to be used to generate vector representations for each program.

¹<https://github.com/eliben/pycparser>

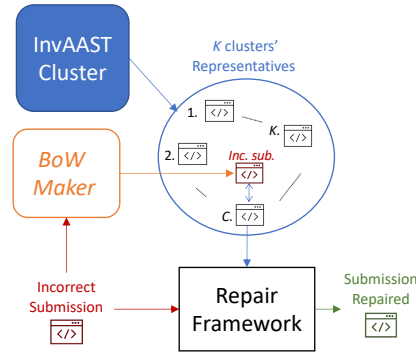


Figure 4: Finding the closest correct program, i.e., the closest correct program representative to the incorrect submission vector representation. This approach passes only one program to the repair tool instead of K programs.

Depending on this module’s parameterization, the BoW maker can compute four different bags of words: (1) based on the programs’ code (syntax), (2) based on the programs’ AASTs (structure), (3) using the set of programs’ invariants (semantics) and (4) joining the programs’ AASTs and their sets of invariants (structure + semantics). To compute these BoW models, INVAASTCLUSTER uses scikit-learn package, `feature_extraction`².

INVAASTCLUSTER tokenizes the input strings into tokens of size n (n -grams) to build a vocabulary with all the submissions’ information, i.e., invariants, syntax, or AASTs. In our case, we define $n = 3$ (3-grams) for this parameter of the BoW maker. Afterward, once a vocabulary has been collected, INVAASTCLUSTER computes a vector representation for each program by counting the number of times each token appears in the program’s information string (invariants, syntax, or AASTs) normalizing the vector by the length of the BoW’s vocabulary.

Clustering. The main goal of INVAASTCLUSTER is to reduce the vast number of correct submissions, N , into a significantly smaller number of program representatives, K , to help program repair frameworks to become more scalable (if $N \gg K$). Therefore, INVAASTCLUSTER accepts as parameter the number of desired clusters K , which is by default 10% of N . The BoW maker module passes the set of vector representations for each one of the N correct submissions to the clustering procedure. Then, INVAASTCLUSTER uses the *KMeans* algorithm to cluster these submissions into K different clusters. The *KMeans* algorithm receives as a parameter the number of clusters it should return (K). The *KMeans* algorithm divides the set of observations, in our case, students’ programs, into K clusters, where each program is assigned to the cluster with the nearest mean [59]. INVAASTCLUSTER uses *KMeans* but other clustering algorithms can be applied. INVAASTCLUSTER provides users with the option to select from various clustering algorithms offered in scikit-learn³: Affinity Propagation, MeanShift, MiniBatch *KMeans*, Agglomerative Clustering, Ward, Spectral Clustering, DBSCAN, OPTICS, BIRCH, and Gaussian Mixture.

Clusters’ representatives selection. In this last module, INVAASTCLUSTER chooses a program representative for each cluster. For each one of the K clusters, INVAASTCLUSTER computes the program closest to the center of the cluster using the Euclidean distance. Afterward, INVAASTCLUSTER returns a set of K clusters’ representatives.

Easy upgradability. INVAASTCLUSTER was designed with modularity in mind. On that account, one can easily remove, add or modify any module of INVAASTCLUSTER. For example, one can use other models instead of the bag of words model, only needing to replace that specific procedure.

Several repair tools [23, 33, 2] are implementation-based program repair tools, i.e., they receive a single correct program to act as a reference implementation for repairing any given incorrect program. Hence, these frameworks are not designed to take advantage of a vast number of semantically different correct submissions. These frameworks can be run in parallel. However, they typically do not have a procedure to choose which among several possible repairs is the best (minimal). These tools compute the set of repairs based on the reference implementation. INVAASTCLUSTER

²sklearn.feature_extraction.text.TfidfVectorizer

³<https://scikit-learn.org/stable/modules/clustering.html>

Table 1

Description of our dataset of IPAs.

Labs	#IPAs	#Correct Submissions	#Incorrect Submissions	#IPAs (Clara)	#Correct Submissions (Clara)
Lab02	10	789	118	10	738
Lab03	7	363	35	5	244
Lab04	8	468	43	5	159
Total	25	1620	196	20	1141

cannot be used on these frameworks since its output is a set of K clusters' representatives. Consequently, in order to allow these non-clustering-based frameworks to take advantage of INVAASTCLUSTER, we developed an additional module that finds the closest correct program representative to an incorrect program. Our motivation is that with the closest correct submission from previous years, these tools can provide the student with a minimal set of repairs.

Closest correct program finder. The overall idea of this module is presented in Figure 4. Given a student's incorrect submission, INVAASTCLUSTER finds which of the K clusters' representatives, returned by INVAASTCLUSTER's selection module, is the closest program to the incorrect submission. This is done by identifying the smallest Euclidean distance between the vector representation of each one of the clusters' representatives and the incorrect submission. Hence, we can identify one correct program that is most likely the reference implementation to use for repairing a specific student's program. In the example of Figure 4, INVAASTCLUSTER would return only the program *C* to the repair framework since it is the closest program to the incorrect submission.

6. Introductory Programming Assignments (IPAs) Datasets

C-PACK-IPAS. In order to evaluate the program representations described in Section 4, we gathered a benchmark of students' programs developed during an introductory programming course in C language. These programs were collected over three distinct practical classes at Instituto Superior Técnico for 25 different IPAs. Each lab class focuses on a different topic of the C programming language. Lab02 deals with integers and input-output operations. Secondly, Lab03 focuses on loops and chars. Lastly, in Lab04, the students learn to use vectors and strings. The textual description of each programming assignment can be found in Appendix A.

This dataset of introductory programming exercises, C-PACK-IPAS, is publicly available at <https://github.com/pmorvalho/C-Pack-IPAs> and can be used by other IPAs repair/clustering frameworks. In this git repository, the interested reader can find all the information about the description and the input-output tests used to evaluate each IPA [42]. Furthermore, there is also a reference implementation for each IPA in the public git repository that can be used by program repair frameworks that only accept a single reference implementation to repair incorrect programs. Moreover, C-PACK-IPAS [44] has also proven successful in evaluating various works across program analysis [47, 48], fault localization [45], program repair [46] and program transformation [43].

Since this work focuses only on program semantics, only submissions that compile without any errors were selected. The set of submissions was split into two sets: correct submissions and incorrect submissions. The students' submissions that satisfied a set of input-output test cases for each IPA were considered correct and selected as benchmark instances. The submissions that failed at least one input-output test were considered incorrect.

Table 1 presents the number of submissions gathered. For 25 different programming exercises, this dataset contains 1620 different correct and 196 incorrect submissions. CLARA's clustering method does not support all the features present in the correct submissions collected. Hence, as shown in Table 1, after removing the set of exercises and correct programs that CLARA does not support, we achieved a final set of 1141 correct submissions for 20 IPAs.

ITSP. The ITSP dataset has been used by other automated program analysis tools [67, 2, 43]. This dataset is also a collection of C programs although it is well balanced, i.e., the number of correct submissions is closer to the number of incorrect submissions in this dataset. Table 2 presents the number of programs in the ITSP dataset after we removed the programs that CLARA and our variable renamer module do not support.

Table 2

Description of ITSP [67] dataset. Correct programs that our approach and CLARA do not support were removed.

ITSP Dataset	#IPAs	#Correct Submissions	#Incorrect Submissions
Lab3	4	45	63
Lab4	6	74	75
Lab5	7	64	62
Lab6	6	19	24
Total	23	202	224

7. Experiments

The experimental results presented in this section aim to support our claims that the proposed novel program representation based on a program’s AAST and its set of program invariants help (1) to efficiently cluster semantically equivalent small imperative programs submitted in IPAs, and (2) to repair faster and significantly more IPAs’ incorrect submissions in current state-of-the-art clustering-based program repair tools, such as CLARA [19].

The goal of our experiments was to answer the following research questions:

RQ1. How does invariant-based program clustering compare against AAST and syntax-based clustering on a set of correct submissions? (Section 7.1)

RQ2. Does CLARA repair more programs using INVAASTCLUSTER’s closest correct submission or its set of KMEANS clusters’ representatives? (Section 7.2)

To answer these research questions, we evaluate INVAASTCLUSTER in two different use cases: (1) clustering IPAs (Section 7.1), and (2) repairing IPAs (Section 7.2). For this evaluation, we have gathered a set of IPAs, previously described in Section 6, developed during an introductory programming university course in C language. Section 7.1 presents the first use case where we perform clustering on the students’ submissions for different IPAs and evaluate its accuracy on different program representations. Afterward, Section 7.2 shows the second use case where we integrate our program representations into a state-of-the-art program repair tool, CLARA [19], to evaluate if our clustering technique is able to outperform CLARA’s clustering method, which is the only current publicly available state-of-the-art clustering method for repairing IPAs. We are going to give our program clusters for each IPA to CLARA and use CLARA’s repairing process. The idea is to evaluate our clustering approach being integrated into a clustering-based program repair tool. Program repair is only one of the several possible applications for our program clustering method.

All of the experiments were conducted on an Intel(R) Xeon(R) Silver computer with 4210R CPUs @ 2.40GHz, using a memory limit of 64GB.

7.1. Use Case 1: Clustering IPAs

A study was performed to evaluate different program representations by applying program clustering to the set of correct programs described in Table 1. The main idea of this experiment was to evaluate if program invariants help identify different IPAs’ submissions.

INVAASTCLUSTER was used to cluster the 1620 correct submissions, to different IPAs, into 25 distinct clusters since our dataset has 25 different programming exercises. Other works [49] that perform program clustering on IPAs perform an equivalent study on clustering submissions for different exercises. The main reason to cluster programs from 25 different exercises is that we know the ground truth label for each program since we know for which specific IPA the students submit their assignments. Otherwise, we would have to manually choose semantically different implementations for the same IPA and assign labels, which might be subjective.

INVAASTCLUSTER, as explained in Section 5, starts by renaming all the variables in the student submissions. Then uses DAIKON [13] to collect the student submissions’ dynamically generated invariants sets as described in Section 3.2. Lastly, it uses the python library, `pycparser`⁴, to compute all the anonymized abstract syntax trees (AAST) (see Section 4.2). Using all these program features, we computed four different bags of words models. One model for each program representation (syntax, AAST, and invariants) and one additional model using a combination of a program’s AAST and its invariants set. The program syntax is not included in this last representation since the bag of words based on program syntax has a large vocabulary that generates vectors that are too sparse.

⁴<https://github.com/eliben/pycparser>

Table 3

The values for the cluster accuracy using four different clustering algorithms on each program representation after ten different runs, each run using a different seed.

Clustering Algorithm	Program Representation	Average	Median	Variance	Standard deviation
KMeans	AAST+Invariants	81.44%	80.65%	0.02%	1.35%
	AAST	73.63%	73.70%	0.03%	1.69%
	Invariants	78.69%	78.73%	0.01%	0.88%
	Syntax	58.05%	58.15%	0.01%	1.22%
MiniBatch KMeans	AAST+Invariants	79.09%	79.63%	0.05%	2.23%
	AAST	72.33%	72.96%	0.10%	3.12%
	Invariants	75.83%	76.23%	0.03%	1.68%
	Syntax	58.46%	58.21%	0.03%	1.81%
Birch	AAST+Invariants	80.10%	80.09%	0.00%	0.51%
	AAST	73.92%	73.55%	0.08%	2.81%
	Invariants	77.99%	77.90%	0.01%	0.86%
	Syntax	59.05%	58.98%	0.01%	0.94%
Gaussian Mixture	AAST+Invariants	78.89%	79.60%	0.05%	2.26%
	AAST	70.97%	71.70%	0.09%	2.92%
	Invariants	75.59%	74.81%	0.07%	2.63%
	Syntax	57.70%	57.35%	0.02%	1.58%

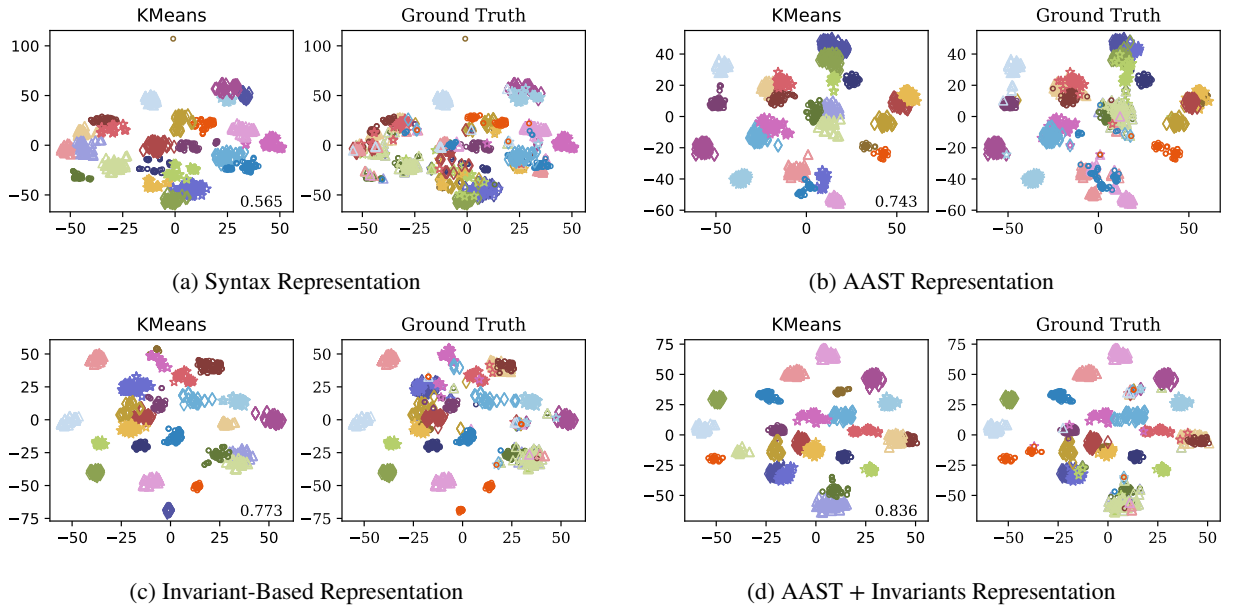


Figure 5: Comparison between the ground truth (on the right) and the clusters and cluster accuracy obtained using the KMEANS algorithm (on the left) for each type of program representation.

The following clustering algorithms available in scikit-learn⁵ were applied to each program representation: KMEANS, MiniBatch KMEANS, BIRCH and Gaussian Mixture. The *KMeans* algorithm divides the set of observations, in our case students' programs, into n clusters where each program is assigned to the cluster with the nearest mean. We used the Euclidean distance as the similarity measurement for KMEANS. MiniBatch KMEANS is similar to KMeans, but instead of using the entire dataset to update the cluster centers at each iteration, this algorithm uses randomly selected subsets. On the other hand, BIRCH is a hierarchical clustering algorithm that builds a tree structure to represent the

⁵<https://scikit-learn.org/stable/modules/clustering.html>

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.84	0.81	0.80	0.81	0.81	0.83	0.81	0.81	0.81	0.83
AAST	0.74	0.75	0.70	0.75	0.73	0.74	0.76	0.72	0.74	0.73
Invariants	0.77	0.80	0.78	0.78	0.80	0.79	0.78	0.79	0.79	0.80
Syntax	0.57	0.58	0.58	0.60	0.58	0.60	0.58	0.59	0.56	0.59

Figure 6: The values for cluster accuracy using the KMEANS algorithm on each program representation after ten different runs, each run using a different seed.

data. It incrementally clusters data points by recursively splitting clusters into subclusters. Finally, a Gaussian Mixture Model is a weighted sum of n component Gaussian densities where n is the number of clusters [52]. Our discussion centers around these clustering algorithms, as they yielded the most favorable outcomes in our experiments.

Since our dataset of IPAS has 25 different programming exercises, the ground truth has 25 different clusters. Each student program is a submission to a specific programming exercise (label) that we know. Consequently, the *cluster accuracy* metric can be used to evaluate the obtained clusters. With this metric, each cluster is assigned the label (exercise) which is most frequent in the cluster. Afterward, the accuracy of this assignment is measured by counting the number of correctly assigned student submissions and dividing by the number of total submissions. This metric is also known as *purity* [57].

Table 3 presents the average, median, variance, and standard deviation values for the cluster accuracy for each clustering algorithm on four different program representations after ten different runs. Each run uses a different seed. Entries highlighted in bold correspond to the highest average/median accuracy values for each different program representation for all clustering algorithms. One can see that the AAST + Invariants representation has the best performance considering all the clustering algorithms. The Invariants BoW has the second highest accuracy, followed by the AAST BoW. Lastly, the Syntax BoW presents the poorest performance of all. From now on, we focus the discussion on the KMEANS results since this clustering algorithm achieved the best results (see Table 3).

The KMEANS clustering algorithm divides the set of observations, in our case, students' programs, into n clusters where each program is assigned to the cluster with the nearest mean [59]. The KMEANS algorithm receives as a parameter the number of clusters it should return, i.e., it always returns 25 different clusters of programs. Figure 6 shows a matrix with the different cluster accuracy values using the KMEANS algorithm on each program representation for ten different seeds. Each entry is highlighted accordingly to its value. The lowest value is highlighted in black, and the highest is highlighted in white. Intermediate values are highlighted in different shades of grey, depending on how far they are from the lowest value. Matrices with values of the cluster accuracy for the other clustering algorithms can be found in Appendix B.1.

Furthermore, Figure 5 presents the results of applying the KMEANS model to each one of the four program representations being analyzed. To present these results graphically, we used a method for visualizing high-dimensional data in a 2-dimension map, called t-SNE [36]. Each subfigure corresponds to a different type of representation. The left side of each subfigure shows the clustering results and the value of the cluster accuracy (right-bottom corner). The right side presents the real clusters of each programming exercise, i.e., the ground truth represented using each program representation. Figure 5a shows the results after clustering all the student submissions using a syntax representation, which resulted in a cluster accuracy of almost 57%. The AAST representation achieved a cluster accuracy of 74.3% as presented in Figure 5b.

Regarding the use of program invariants, Figure 5c and Table 3 support the idea that program invariants improve program clustering since this representation obtained a cluster accuracy of 77.3%. Lastly, Figure 5d presents the representation that uses the combination AASTs and invariants sets, which also shows an improvement compared to the invariant-based representation. This representation outperforms all the other representations with an accuracy of 83.6%. Furthermore, Table 3 shows that this representation based on AAST and invariants achieved the best cluster accuracy for all clustering algorithms. Another advantage of this representation is that it best separates all the students' submissions in different regions of the space, i.e., the majority of the clusters are visibly separated from each other.

Table 4

This table presents the percentage of submissions repaired (success), structural mismatch errors, and timeouts (failure) for each clustering approach. The total number of submissions is 319.

	Clustering Method	%Success	%Failure	
		%Submissions Repaired	%Structural Mismatch	%Timeouts (600s)
1	CLARA	71.79%	7.52%	20.69%
2	KMEANS - Invs	82.45%	11.91%	5.64%
3	KMEANS - Syntax	84.01%	10.97%	5.02%
4	KMEANS - AAST	84.64%	9.72%	5.64%
5	KMEANS - AAST + Invs	84.95%	10.03%	5.02%
6	Closest Program (KMEANS) - AAST + Invs	84.33%	10.66%	5.02%

Other evaluation metrics for the KMEANS algorithm can be found in Appendix B.2 for the interested reader's convenience. Clustering metrics such as the *Rand index*, the *adjusted Rand index*, the *normalized mutual information*, the *adjusted mutual information*, the *Fowlkes–Mallows index*, the *completeness score*, the *homogeneity score*, and the *V measure*.

7.2. Use Case 2: Repairing IPAS

This section presents the results of integrating INVAASTCLUSTER as the clustering approach for CLARA [19], a publicly available state-of-the-art clustering-based program repair tool. Since our set of IPAS, described in Table 1, has a small number of incorrect submissions, only 196, for this evaluation, we have also considered the ITSP dataset [67] described in Section 6. Thus, overall we have a total of 420 incorrect submissions (196 from our dataset plus 224 from the ITSP dataset) and 1343 correct submissions (1141 from our dataset plus 202 from the ITSP dataset) for 43 different IPAS (20 from our dataset plus 23 from the ITSP dataset). To fully evaluate our clustering technique for repairing IPAS, we are going to compare INVAASTCLUSTER's results against CLARA's in terms of: (1) the number of student submissions repaired; (2) the number of clusters produced by each clustering approach for each IPA; and (3) the time spent to repair each incorrect submission.

We would like to point out that in this experiment, we are not trying to improve CLARA's repair process. Instead, we are comparing the performance of CLARA's repair process when using its own or INVAASTCLUSTER's clusters of programs.

Different procedures for program clustering using INVAASTCLUSTER (see Table 4) are evaluated:

- **KMEANS - BoW**: Uses KMEANS and four different bag of words (BoW) based on AAST, syntax, and invariants (lines 2–5 in Table 4);
- **Closest Program (KMEANS) - AAST + Invs**: Uses the closest program (see Section 5) using the AAST + Invs BoW, from a set of clusters' representatives using KMEANS (line 6);

Table 4 presents the overall repair evaluation on 319 incorrect submissions since CLARA's repair algorithm does not support the C implementation of 101 incorrect submissions (24.05% of the instances). Entries in bold correspond to the highest rate of submissions repaired, the lowest percentage of structural mismatch errors, or the lowest rate of timeouts, i.e., executions that did not repair a program using a timeout of 10 minutes (600s).

One can see in line 1 (Table 4) that CLARA, using its own clusters, can only repair 229 (around 72%) of the incorrect submissions and shows the largest percentage of instances that were not repaired due to timeout (20.69%). Secondly, the configuration using INVAASTCLUSTER's KMEANS and the BoW based on AAST and Invariants achieved the highest score, repairing 84.95% of the incorrect submissions. Furthermore, the BoW based only on invariants has the highest percentage of structural mismatch (11.91%), which may be explained by CLARA's inability to use a program with a different control flow in the repair process. Using only invariants on a vector representation helps clustering programs with similar semantics, although it does not take into account the programs' structure (control flow). Hence, a higher rate of structural mismatch is observed.

Since the BoW based on AASTs and program invariants achieved the best results both in the program clustering experiment (see Section 7.1) as well as when repairing submissions (lines 2-5 in Table 4), we opted to use only this

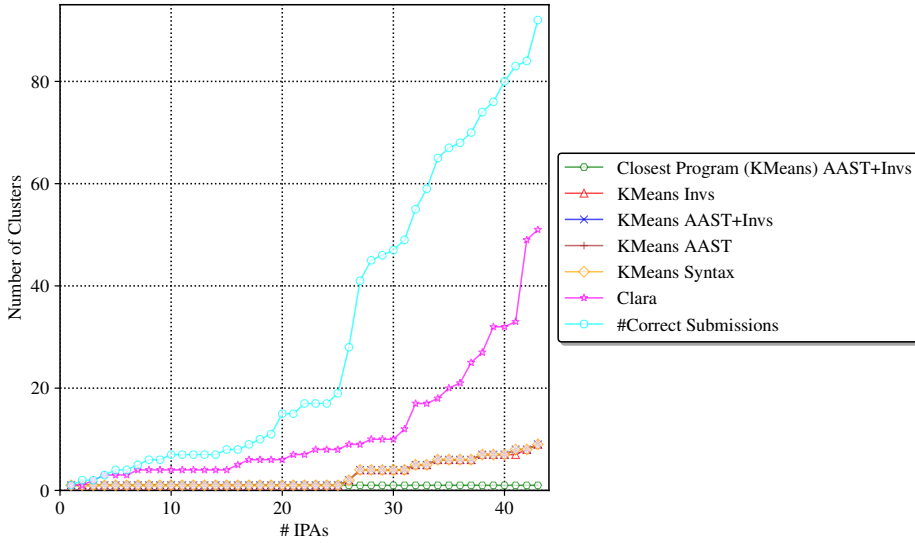


Figure 7: Cactus plot - The number of clusters generated by each clustering technique for 43 different IPAs.

BoW when finding the closest correct program (line 6, Table 4). Regarding the use of just one correct solution to fix an incorrect submission, the *Closest Program* (KMEANS) approach did not achieve better results than using the set of clusters' representatives.

We have also analyzed the closest program technique using all submissions, i.e., use the closest program among all submissions (no clustering step). This approach, the *Closest Program (All Submissions)*, was able to repair 86.5% of the submissions. The number of timeouts in this approach and using KMEANS was similar. Once again, this high rate of repaired programs (86.5%) may be explained by CLARA's strict requirements for both programs, the program being repaired and the correct program used by the repair process, to have the same control flow. Therefore, when INVAASTCLUSTER finds the closest program among all submissions instead of using clusters, INVAASTCLUSTER has a more diverse collection of programs' structures. Consequentially, the *Closest Program (All Submissions)* approach also achieved the lowest score of structural mismatch errors (only 9.4%). Although we would like to draw the reader's attention to the difference between the number of submissions repaired using KMEANS (85%) or using the closest correct program (86.5%), which is less than 2%. Furthermore, the computation to find the closest correct program among all correct submissions can only be done online since it requires the student's incorrect program. On the other hand, the computation of the KMEANS clusters can be done offline since it only requires past students' correct submissions. In this evaluation, this is not a concern since each IPA has at most a hundred correct submissions. However, in a large-scale MOOC with thousands of correct submissions per exercise, the process of finding the closest correct program among all of the submissions may become impractical to compute in a short period of time.

7.2.1. Number of Clusters

Figure 7 illustrates a cactus plot detailing the number of clusters generated by each clustering technique for each of the 43 different IPAs used. One can see that CLARA generates an enormous quantity of clusters, almost half of the correct submissions of each IPA. This large number of clusters is explained by CLARA's strict clustering method, which does not allow two programs to be in the same cluster if there is no exact match between both programs' control flows. Furthermore, even if two programs are semantically equivalent and share the same control flow but one of them uses a for-loop and the other uses a while-loop, then CLARA assigns these two programs to different clusters.

INVAASTCLUSTER produces K clusters, which in this experiment is always set to 10% of the number of correct submissions of each exercise. The technique that uses the closest correct program has only a single cluster which is the closest correct program. This evaluation of the number of clusters used by each approach allows us to observe that CLARA produces a large number of clusters, resulting in a detriment of performance. In contrast, our approach can generate fewer clusters resulting in a more effective repairing process.

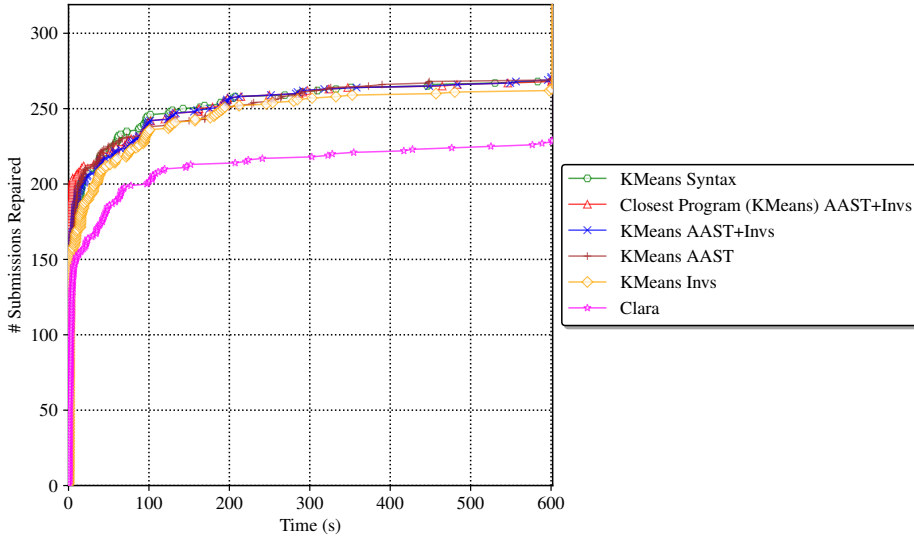


Figure 8: Cactus plot - Time Performance (timeout=600s).

Example 5. The following two programs are correct implementations for the IPA where students are asked to print the maximum value among three given numbers. These programs are clustered together using INVAASTCLUSTER (AASTs + Invariants) because their AASTs are quite similar, and their sets of program invariants are identical. In contrast, CLARA assigns these two programs to different clusters because they have a different number of variables. Therefore, because CLARA is highly strict in comparing programs, it generates an excessive number of program clusters.

```

1  int main(){
2      int n1, n2, n3, max;
3      scanf("%d%d%d", &n1, &n2, &n3);
4      max = n1;
5      if(n2 > max){
6          max = n2;
7      }
8      if(n3 > max){
9          max = n3;
10     }
11
12     printf("%d\n", max);
13     return 0;
14 }
```

```

1  int main(){
2      int res,b,c;
3      scanf("%d%d%d", &res, &b, &c);
4
5      if (b > res){
6          res = b;
7      }
8      if (c > res){
9          res = c;
10     }
11
12     printf("%d\n", res);
13     return 0;
14 }
```

7.2.2. CPU time

Regarding the time performance of each clustering technique, Figure 8 shows a cactus plot that presents the CPU time spent on repairing each program (x-axis) against the number of repaired programs (y-axis) using different clustering techniques. The legend in this plot is not sorted.

One can clearly see a gap between CLARA's time performance and INVAASTCLUSTER's (considering any clustering approach). For example, after 100 seconds CLARA using its own clusters, can only repair around 200 programs while using our clusters, it can repair around 230/250 programs. Furthermore, Figure 9 presents a scatter plot comparing the CPU time spent using CLARA's clusters against the KMEANS AAST + Invariants clusters. Each point in this plot represents a program where the x-value (resp. y-value) is the CPU time spent to repair the program using the KMEANS AAST + Invariants Clusters (resp. CLARA's clusters). If a point is above the diagonal, then it means

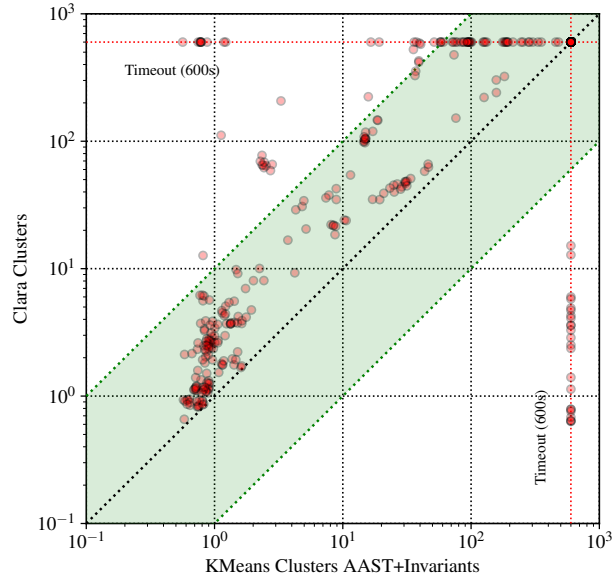


Figure 9: Scatter plot - Time Performance (timeout=600s) - CLARA VS INVAASTCLUSTER (KMEANS w/ AAST + Invariants)

that our clusters outperformed CLARA’s clusters since using our clusters CLARA is able to repair each program above the diagonal faster than using its own clusters. Thus, if we consider the programs repaired by both clustering methods, using our clusters is always faster than using CLARA’s.

There are two main reasons for CLARA’s time performance. Firstly, CLARA generates a significantly larger number of clusters compared to INVAASTCLUSTER (see Figure 7). Consequently, CLARA needs to compute a set of repairs for each cluster’s representative, resulting in more time spent in the repair process due to the larger number of clusters. Secondly, as explained in Section 2, CLARA maintains a json file containing data (e.g., expressions) of all the programs belonging to the cluster. During its repair algorithm, CLARA takes advantage of all this data to repair a given submission, leading to more time spent on the repair process.

The main goal of educational program repair frameworks is to provide real-time feedback to students on how they should repair their submissions. In this evaluation, we used a timeout of 10 minutes. However, a student expects a faster result. Therefore, Figure 10 shows another cactus plot that shows the time performance of the clustering approaches, although with a timeout of 10 seconds. One can see that after 10 seconds, CLARA using its own clusters, can only repair around 150 submissions (47%). On the other hand, using INVAASTCLUSTER, CLARA can repair around 200 submissions (63%). Furthermore, one can also verify that after 2 seconds CLARA can only repair around 75 programs using its own cluster while using our clusters CLARA is able to repair around 150 programs.

7.2.3. Program invariants of incorrect submissions

We have also tried the method *Closest Program* (KMEANS) with INVAASTCLUSTER not taking into account incorrect submissions’ invariants to check if incorrect submissions’ sets of invariants had a negative impact on program representations. First, INVAASTCLUSTER clustered all the correct programs considering their AASTs and their sets of invariants. Secondly, to find the closest correct program to an incorrect submission INVAASTCLUSTER used only the AAST BoW. However, this combined approach of clustering with one BoW (AASTs + invariants) and calculating the programs’ distances with another (AASTs only) was only able to repair 269 submissions (84%). Thus, according to this experiment, incorrect submissions’ sets of program invariants do not cause negative effects on program representations.

To summarize, in Section 7.1, we used INVAASTCLUSTER to cluster different IPAS correct students’ submissions. The obtained results support that this work’s novel program representation based on a program’s AAST and invariants performs better when compared to representations solely based on a program’s code, AST, or set of program invariants.

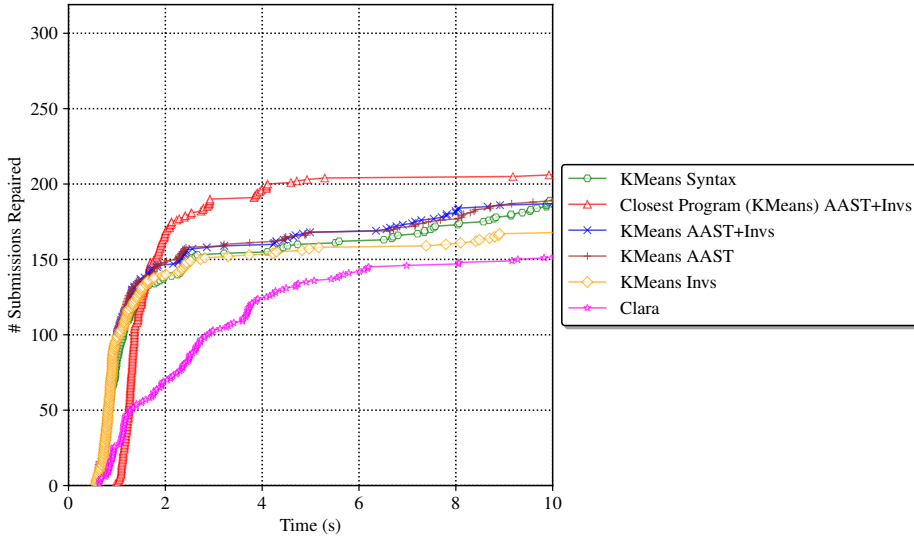


Figure 10: Cactus plot - Time Performance (timeout=10s).

Furthermore, in Section 7.2, we integrated INVAASTCLUSTER into CLARA to evaluate our tool’s performance when integrated into a clustering-based framework for repairing IPAS. This study shows that INVAASTCLUSTER significantly increases the performance of the state-of-the-art clustering technique by allowing CLARA to repair more student submissions and doing so notably faster.

7.3. Threats to Validity

This work relies on DAIKON to compute dynamically-generated likely invariants. Using another tool to detect likely invariants may produce different results. Dynamically-generated likely invariants depend highly on the test suite used for each programming exercise. Therefore, using a different set of input-output tests may produce different sets of program invariants for each student submission.

Furthermore, our evaluation of INVAASTCLUSTER exclusively employed small imperative programs commonly encountered in IPAS. The application of INVAASTCLUSTER to more intricate programs is deferred to future investigations, although we anticipate no major scalability issues with larger datasets or more complex programs when utilizing INVAASTCLUSTER. Additionally, it is worth noting that the number of clusters expands based on the diversity of semantic and syntactic implementations for each IPA. In theory, as the complexity of IPAS increases, we can expect a greater diversity of implementations. Nevertheless, the scalability of INVAASTCLUSTER should remain unaffected, and there is no imposed limit on the maximum number of supported IPAS.

While our evaluation focused solely on C programs, our clustering methodologies are agnostic and can be adapted to other programming languages. This adaptation entails substituting (1) the variable renaming module with one compatible with the target language and (2) INVAASTCLUSTER’s invariant detection module (DAIKON) with another suitable program invariant detector tailored for the desired language.

A different repair tool may produce different results since the repair process may differ. For example, another tool may support C features that CLARA currently does not support or the opposite. Only small imperative programs, usually found in IPAS, were used to evaluate INVAASTCLUSTER. Using INVAASTCLUSTER on more complex programs is left for future work, although we anticipate no scalability issues with larger programs.

INVAASTCLUSTER could provide an intriguing perspective on utilizing AASTs and program invariants to identify potential similarities/copies among students’ submissions. However, implementing this would require some time to adapt the plagiarism detection tool. Additionally, it is worth noting that MOSS is not open-source.

In this paper, we compared INVAASTCLUSTER against CLARA’s clustering method since, to the best of our knowledge, CLARA is currently the only publicly available state-of-the-art clustering-based repair tool for repairing IPAS. However, evaluating INVAASTCLUSTER on other repair frameworks would be valuable. Unfortunately, we

could not find public implementations for the other tools [34, 16, 49, 64] nor the datasets of IPAs used for their evaluation, except for the ITSP dataset [67].

8. Related Work

Program clustering has also been used to find different semantic solutions for a given programming exercise [17, 16, 9]. PACTON [16] clusters programming assignments based on their symbolic analysis. PACTON clusters two submissions together if their path conditions are equivalent. PACTON only takes into consideration a program’s semantics, while INVAASTCLUSTER also considers a program’s structure. OVERCODE [17] lets the user visualize and explore different implementations for the same exercise. CODEBERT [14], CODE2SEQ [5] and other deep learning models [69] build vector representations of programs by training machine learning models using the programs’ code and ASTs. However, unlike INVAASTCLUSTER, these techniques only consider the programs’ syntax, not their semantics.

CODERASSIST [26] is a counter-example guided feedback generation tool that provides feedback on student implementations of *dynamic programming algorithms*. CODERASSIST starts by clustering both correct and incorrect programs based on these dynamic programs’ syntactic features. Afterward, CODERASSIST generates feedback for a buggy program by calling an SMT solver, using a counterexample obtained from an equivalence check against a correct implementation in the same cluster. However, CODERASSIST only works for dynamic programs since the implemented feature extraction procedure searches for syntactic features of dynamic programming algorithms. Thus, this tool cannot be used for clustering our dataset of IPAs. Rocha et al. [54] introduce a framework to improve how teaching assistants provide feedback in introductory programming courses, helping them understand different feedback approaches and resources. This framework offers structured guidance for pedagogical decision-making regarding adaptive feedback.

CODEHOUND [29] is a system that automatically tracks pedagogical code dependencies by using static analysis to detect function introductions and reuse throughout an entire course. It offers instructors assistance in creating new content, collaborating on content refactoring, and estimating future course change costs. Furthermore, ODS [66] is a system for detecting overfitting patches in automatic program repair. ODS utilizes supervised learning with AST level code features and patch correctness labels to automatically learn a probabilistic model, which can then classify new program repair patches. Moreover, CLACER [31] is a neural network model designed to classify compilation errors by proposing categories based on program tokens, aiming to improve localization effectiveness and prediction performance, thereby enhancing the students’ learning process.

Code search techniques [51, 1, 27, 60, 61, 35, 37] are another family of semantic program repair techniques. Code search uses a specification (e.g., input-output tests) to find code in a large repository of code snippets that satisfies that specification. To repair an incorrect program, semantic code search methods do not find the closest correct implementation for the same IPA or use a reference implementation provided by the lecturer. Instead, these methods search for code fragments of several correct programs to repair a given incorrect submission. These methods have no knowledge about the program’s structure and where that code fragment came from. Therefore, there is no guarantee that the set of repairs proposed by these tools is the minimal set required to fix a program.

Solution-driven program repair tools use one reference implementation to repair a given incorrect submission [2, 23, 33]. AUTOGRADER [33] finds potential path differences between the executions of a student’s submission and a reference implementation using symbolic execution. Then, AUTOGRADER provides feedback to students using counterexamples for each path difference found. VERIFIX [2] aligns an incorrect program with the reference solution into an automaton. Then, using that alignment relation and MaxSMT solving, VERIFIX proposes fixes to the incorrect program. TEGCER [3] is an automated feedback tool for novice programmers. This tool uses supervised learning to match *compilation errors* in new code submissions with pre-existing errors submitted by previously enrolled students. TEGCER only works for syntactic errors. However, similarly to INVAASTCLUSTER, TEGCER also performs variable renaming to each program. This tool renames the variables with their generic types using the LLVM [28], a standard static analysis tool.

Regarding *clustering-based program repair tools* [49, 34, 64, 19], CLARA [19] was already described, and differences were highlighted in previous sections. More recently, Contractor and Rivero [10], Chowdhury et al. [7] improved CLARA’s matching algorithm to a new graph matching algorithm that is more relaxed in terms of control flow restrictions. However, this new graph-matching algorithm only works for Python programs. SARFGEN [64], used for C# programs, creates program embeddings based on the programs’ ASTs [25]. Then, given an incorrect program, finds the closest correct submission using those embeddings and tries to repair the program by aligning the variables

in both programs. Thus, unlike INVAASTCLUSTER, SARFGEN only considers a program's structure (AST) and not its semantics during the clustering process.

SEMCLUSTER [49] clusters programs based on their control and data flow features. SEMCLUSTER creates vector representations, *program features vectors* (PFV), for each program using a test suite. This PFV takes into account control flow features as well as data flow features. For each program, SEMCLUSTER counts the number of times each control flow path is used in each test and builds a vector with this data. Afterward, SEMCLUSTER builds another vector containing the data flow features, i.e., the number of occurrences of consecutive values a variable takes during its lifetime. Finally, SEMCLUSTER merges these two feature vectors into a single vector, the PFV. Similar to INVAASTCLUSTER, SEMCLUSTER uses the KMEANS clustering algorithm. However, unlike INVAASTCLUSTER, which uses each program's AAST, SEMCLUSTER does not account for any syntactic features. Furthermore, SEMCLUSTER tries to capture the semantics of a program by counting the number of different values each variable takes. On the other hand, INVAASTCLUSTER considers the program's set of invariants which can be more robust and independent of the test suite used. Some research has been conducted regarding *the use of invariants to promote patch diversity* or to help with patch selection on a search-based program repair [6, 12, 65].

9. Conclusions and Future Work

In the context of introductory programming assignments (IPAS) in university courses or Massive Open Online Courses (MOOCs), it is possible to collect a large number of correct implementations for the proposed IPAS. Hence, when a student submits an incorrect program, one can take advantage of previously correct submissions to automatically suggest repairs that help the student. However, it is not feasible to analyze all possible previous correct submissions to find an appropriate reference implementation for the repair tool. Therefore, clustering is often used to identify similar program implementations. Afterward, the automated repair tool analyses a single reference program from each cluster to find the most suitable correction to the student's incorrect submission.

This work proposes INVAASTCLUSTER, a novel approach for program clustering based on their semantic and syntactic features. INVAASTCLUSTER uses AASTs and invariant-based program representations to distinguish small imperative programs according to their semantics (invariants) and structure (AAST). Results show that the proposed AASTs and invariant-based representation improve upon syntax-based representations when performing program clustering on several correct student submissions for different programming exercises. Additionally, given an incorrect student submission and a set of correct students' submissions, INVAASTCLUSTER can also find the closest correct submission to the faulty program using INVAASTCLUSTER's program vector representations.

Furthermore, INVAASTCLUSTER has also been integrated into a state-of-the-art clustering-based program repair framework to evaluate the proposed clustering techniques for repairing IPAS. This evaluation showed that INVAASTCLUSTER outperforms the current state-of-the-art clustering method used in clustering-based program repair. Using INVAASTCLUSTER, the automated repair tool CLARA can repair significantly more IPAS, around 13%, with a better time performance and with a smaller number of program clusters.

To conclude, INVAASTCLUSTER is a program clustering tool based on programs' invariants and AASTs. INVAASTCLUSTER can be used: (1) to cluster semantically equivalent implementations for programming exercises; (2) by any clustering-based program repair tool; and (3) by any program repair framework that requires a single reference implementation (INVAASTCLUSTER's closest correct program).

As future directions, we propose to evaluate the new program representations described in this paper (i.e., using AASTs and invariants) as program encodings to use on deep learning models on several tasks such as fault localization [32] or program synthesis [58]. As we expand to consider more complex programs, we plan to evaluate INVAASTCLUSTER on clustering-based repair tools focused on repairing industrial software. Finally, INVAASTCLUSTER will be evaluated on other clustering-based program repair frameworks with a more permissive repair algorithm than CLARA. Moreover, it could also be interesting to explore using a neural architecture to learn deep representations of the features built up in our model, as the current Bag of Words (BoW) approach may lead to a significant loss of information.

Acknowledgments

This research was supported by Fundação para a Ciência e Tecnologia (FCT) through grant SFRH/BD/07724/2020 (DOI: 10.54499/2020.07724.BD) and projects UIDB/50021/2020 (DOI: 10.54499/UIDB/50021/2020), PTDC/CCI-COM/2156/2021 (DOI: 10.54499/PTDC/CCI-COM/2156/2021) and 2023.14280.PEX (DOI: 10.54499/2023.14280.PEX) and grant SFRH/BD/07724/2020 (DOI: 10.54499/2020.07724.BD). This work was also supported by the MEYS within the program ERC CZ under the project POSTMAN no. LL1902 and co-funded by the EU under the project *ROBOPROX* (reg. no. CZ.02.01.01/00/22_008/0004590).

References

- [1] Afzal, A., Motwani, M., Stolee, K.T., Brun, Y., Goues, C.L., 2019. Sosrepair: Expressive semantic search for real-world program repair. *IEEE Trans. Software Eng.* 47, 2162–2181. doi:10.1109/TSE.2019.2944914.
- [2] Ahmed, U.Z., Fan, Z., Yi, J., Al-Bataineh, O.I., Roychoudhury, A., 2022. Verifix: Verified repair of programming assignments. *ACM Trans. Softw. Eng. Methodol.* URL: <https://doi.org/10.1145/3510418>, doi:10.1145/3510418.
- [3] Ahmed, U.Z., Sindhgatta, R., Srivastava, N., Karkare, A., 2019. Targeted example generation for compilation errors, in: 34th IEEE/ACM International Conference on Automated Software Engineering, ASE 2019, San Diego, CA, USA, November 11–15, 2019, IEEE. pp. 327–338. URL: <https://doi.org/10.1109/ASE.2019.00039>, doi:10.1109/ASE.2019.00039.
- [4] Aho, A.V., Sethi, R., Ullman, J.D., 1986. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley series in computer science / World student series edition, Addison-Wesley.
- [5] Alon, U., Brody, S., Levy, O., Yahav, E., 2019. code2seq: Generating sequences from structured representations of code, in: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6–9, 2019, OpenReview.net. pp. –. URL: <https://openreview.net/forum?id=H1gKY09tX>.
- [6] Cashin, P., Martinez, C., Weimer, W., Forrest, S., 2019. Understanding automatically-generated patches through symbolic invariant differences, in: 34th IEEE/ACM International Conference on Automated Software Engineering, ASE 2019, San Diego, CA, USA, November 11–15, 2019, IEEE. pp. 411–414. URL: <https://doi.org/10.1109/ASE.2019.00046>.
- [7] Chowdhury, M.T.A., Contractor, M.R., Rivero, C.R., 2024. Flexible control flow graph alignment for delivering data-driven feedback to novice programming learners. *J. Syst. Softw.* 210, 111960. URL: <https://doi.org/10.1016/j.jss.2024.111960>, doi:10.1016/J.JSS.2024.111960.
- [8] Churchill, B.R., Padon, O., Sharma, R., Aiken, A., 2019. Semantic program alignment for equivalence checking, in: McKinley, K.S., Fisher, K. (Eds.), *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019, Phoenix, AZ, USA, June 22–26, 2019*, ACM. pp. 1027–1040.
- [9] Clune, J., Ramamurthy, V., Martins, R., Acar, U.A., 2020. Program equivalence for assisted grading of functional programs. *Proc. ACM Program. Lang.* 4, 171:1–171:29. URL: <https://doi.org/10.1145/3428239>, doi:10.1145/3428239.
- [10] Contractor, M.R., Rivero, C.R., 2022. Improving program matching to automatically repair introductory programs, in: Crossley, S., Popescu, E. (Eds.), *Intelligent Tutoring Systems*, Springer International Publishing, Cham. pp. 323–335.
- [11] Dijkstra, E.W., 1976. *A Discipline of Programming*. Prentice-Hall. URL: <https://www.worldcat.org/oclc/01958445>.
- [12] Ding, Z.Y., Lyu, Y., Timperley, C.S., Goues, C.L., 2019. Leveraging program invariants to promote population diversity in search-based automatic program repair, in: Petke, J., Tan, S.H., Langdon, W.B., Weimer, W. (Eds.), *Proceedings of the 6th International Workshop on Genetic Improvement, GI@ICSE 2019, Montreal, Quebec, Canada, May 28, 2019*, ACM. pp. 2–9. URL: <https://doi.org/10.1109/GI.2019.00011>.
- [13] Ernst, M.D., Perkins, J.H., Guo, P.J., McCamant, S., Pacheco, C., Tschantz, M.S., Xiao, C., 2007. The daikon system for dynamic detection of likely invariants. *Sci. Comput. Program.* 69, 35–45. URL: <https://doi.org/10.1016/j.scico.2007.01.015>.
- [14] Feng, Z., Guo, D., Tang, D., Duan, N., Feng, X., Gong, M., Shou, L., Qin, B., Liu, T., Jiang, D., Zhou, M., 2020. Codebert: A pre-trained model for programming and natural languages, in: Cohn, T., He, Y., Liu, Y. (Eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16–20 November 2020*, Association for Computational Linguistics. pp. 1536–1547.
- [15] Fowlkes, E.B., Mallows, C.L., 1983. A method for comparing two hierarchical clusterings. *Journal of the American statistical association* 78, 553–569.
- [16] Fu, Y., Osei-Owusu, J., Astorga, A., Zhao, Z.N., Zhang, W., Xie, T., 2021. Pacon: a symbolic analysis approach for tactic-oriented clustering of programming submissions, in: Curtsinger, C., Nguyen, T.N. (Eds.), *Proceedings of the 2021 ACM SIGPLAN International Symposium on SPLASH-E, Chicago, IL, USA. October 20, 2021*, ACM. pp. 32–42.
- [17] Glassman, E.L., Scott, J., Singh, R., Guo, P.J., Miller, R.C., 2015. Overcode: Visualizing variation in student solutions to programming problems at scale. *ACM Trans. Comput. Hum. Interact.* 22, 7:1–7:35.
- [18] Gulwani, S., Radicek, I., Zuleger, F., 2014. Feedback generation for performance problems in introductory programming assignments, in: *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, (FSE-22), Hong Kong, China, November 16 - 22, 2014*, ACM. pp. 41–51.
- [19] Gulwani, S., Radicek, I., Zuleger, F., 2018. Automated clustering and program repair for introductory programming assignments, in: *PLDI 2018*, ACM. pp. 465–480.
- [20] Gupta, R., Pal, S., Kanade, A., Shevade, S.K., 2017. Deepfix: Fixing common C language errors by deep learning, in: Singh, S.P., Markovitch, S. (Eds.), *AAAI 2017*, AAAI Press. pp. 1345–1351.
- [21] Harris, Z.S., 1954. Distributional structure. *Word* 10, 146–162.

- [22] Hopcroft, J.E., Motwani, R., Ullman, J.D., 2007. Introduction to automata theory, languages, and computation, 3rd Edition. Pearson international edition, Addison-Wesley.
- [23] Hu, Y., Ahmed, U.Z., Mechtaev, S., Leong, B., Roychoudhury, A., 2019. Re-factoring based program repair applied to programming assignments, in: 34th IEEE/ACM International Conference on Automated Software Engineering, ASE 2019, San Diego, CA, USA, November 11-15, 2019, IEEE. pp. 388–398.
- [24] Hubert, L., Arabie, P., 1985. Comparing partitions. *Journal of classification* 2, 193–218.
- [25] Jiang, L., Misherghi, G., Su, Z., Glondou, S., 2007. DECKARD: scalable and accurate tree-based detection of code clones, in: 29th International Conference on Software Engineering (ICSE 2007), Minneapolis, MN, USA, May 20-26, 2007, IEEE Computer Society. pp. 96–105.
- [26] Kaleeswaran, S., Santhiar, A., Kanade, A., Gulwani, S., 2016. Semi-supervised verified feedback generation, in: Zimmermann, T., Cleland-Huang, J., Su, Z. (Eds.), *Proceedings of the 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering, FSE 2016, Seattle, WA, USA, November 13-18, 2016*, ACM. pp. 739–750. URL: <https://doi.org/10.1145/2950290.2950363>, doi:10.1145/2950290.2950363.
- [27] Ke, Y., Stolee, K.T., Goues, C.L., Brun, Y., 2015. Repairing programs with semantic code search (T), in: Cohen, M.B., Grunke, L., Whalen, M. (Eds.), 30th IEEE/ACM International Conference on Automated Software Engineering, ASE 2015, IEEE Computer Society. pp. 295–306.
- [28] Lattner, C., Adve, V.S., 2004. LLVM: A compilation framework for lifelong program analysis & transformation, in: 2nd IEEE / ACM International Symposium on Code Generation and Optimization (CGO 2004), 20-24 March 2004, San Jose, CA, USA, IEEE Computer Society. pp. 75–88. URL: <https://doi.org/10.1109/CGO.2004.1281665>, doi:10.1109/CGO.2004.1281665.
- [29] Lau, S., Guo, P.J., 2022. Codehound: Helping instructors track pedagogical code dependencies in course materials, in: Henz, M., Lerner, B.S. (Eds.), *Proceedings of the 2022 ACM SIGPLAN International Symposium on SPLASH-E, SPLASH-E 2022, Auckland, New Zealand, 5 December 2022*, ACM. pp. 1–6. URL: <https://doi.org/10.1145/3563767.3568126>, doi:10.1145/3563767.3568126.
- [30] Levenshtein, V.I., 1966. Binary codes capable of correcting deletions, insertions, and reversals, in: *Soviet physics doklady akademii nauk*, Soviet Union. pp. 707–710.
- [31] Li, Z., Sun, F., Wang, H., Ding, Y., Liu, Y., Chen, X., 2021. CLACER: A deep learning-based compilation error classification method for novice students’ programs, in: *IEEE 45th Annual Computers, Software, and Applications Conference, COMPSAC 2021, Madrid, Spain, July 12-16, 2021*, IEEE. pp. 74–83. URL: <https://doi.org/10.1109/COMPSAC51774.2021.00022>, doi:10.1109/COMPSAC51774.2021.00022.
- [32] Li, Z., Wu, S., Liu, Y., Shen, J., Wu, Y., Zhang, Z., Chen, X., 2023. Vsusfl: Variable-suspiciousness-based fault localization for novice programs. *Journal of Systems and Software* 205, 111822.
- [33] Liu, X., Wang, S., Wang, P., Wu, D., 2019. Automatic grading of programming assignments: an approach based on formal semantics, in: Beecham, S., Damian, D.E. (Eds.), *Proceedings of the 41st International Conference on Software Engineering: Software Engineering Education and Training, ICSE (SEET) 2019*, IEEE / ACM. pp. 126–137.
- [34] Lu, Y., Meng, N., Li, W., 2021. FAPR: fast and accurate program repair for introductory programming courses. *CoRR* abs/2107.06550.
- [35] Luan, S., Yang, D., Barnaby, C., Sen, K., Chandra, S., 2019. Aroma: code recommendation via structural code search. *Proc. ACM Program. Lang.* 3, 152:1–152:28.
- [36] Van der Maaten, L., Hinton, G., 2008. Visualizing data using t-sne. *Journal of machine learning research* 9.
- [37] Mathew, G., Stolee, K.T., 2021. Cross-language code search using static and dynamic analyses, in: Spinellis, D., Gousios, G., Chechik, M., Penta, M.D. (Eds.), *ESEC/FSE ’21: 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, Athens, Greece, August 23-28, 2021, ACM. pp. 205–217.
- [38] Mechtaev, S., Yi, J., Roychoudhury, A., 2015. Directfix: Looking for simple program repairs, in: Bertolino, A., Canfora, G., Elbaum, S.G. (Eds.), 37th IEEE/ACM International Conference on Software Engineering, ICSE 2015, IEEE Computer Society. pp. 448–458.
- [39] Mechtaev, S., Yi, J., Roychoudhury, A., 2016. Angelix: scalable multiline program patch synthesis via symbolic analysis, in: Dillon, L.K., Visser, W., Williams, L.A. (Eds.), *ICSE 2016*, ACM. pp. 691–701.
- [40] Milner, R., 1971. An algebraic definition of simulation between programs, in: Cooper, D.C. (Ed.), *Proceedings of the 2nd International Joint Conference on Artificial Intelligence*. London, UK, September 1-3, 1971, William Kaufmann. pp. 481–489.
- [41] Nguyen, H.D.T., Qi, D., Roychoudhury, A., Chandra, S., 2013. Semfix: program repair via semantic analysis, in: Notkin, D., Cheng, B.H.C., Pohl, K. (Eds.), 35th International Conference on Software Engineering, ICSE ’13, IEEE Computer Society. pp. 772–781.
- [42] Orvalho, P., Janota, M., Manquinho, V., 2022a. C-Pack of IPAs: A C90 Program Benchmark of Introductory Programming Assignments. *CoRR* abs/2206.08768. URL: <https://doi.org/10.48550/arXiv.2206.08768>, doi:10.48550/arXiv.2206.08768, arXiv:2206.08768.
- [43] Orvalho, P., Janota, M., Manquinho, V., 2022b. MultiIPAs: Applying Program Transformations to Introductory Programming Assignments for Data Augmentation, in: *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2022*, ACM, Singapore. pp. 1657–1661. doi:10.1145/3540250.3558931.
- [44] Orvalho, P., Janota, M., Manquinho, V., 2024a. C-Pack of IPAs: A C90 Program Benchmark of Introductory Programming Assignments, in: *2024 IEEE/ACM International Workshop on Automated Program Repair (APR)*, ACM. pp. 14–21. URL: <https://doi.org/10.1145/3643788.3648010>, doi:10.1145/3643788.3648010.
- [45] Orvalho, P., Janota, M., Manquinho, V., 2024b. CFaults: Model-Based Diagnosis for Fault Localization in C Programs with Multiple Test Cases, in: *Formal Methods - 26th International Symposium, FM 2024, Milan, Italy, 2024*, *Proceedings*, Springer, Cham. pp. 463–481. doi:10.1007/978-3-031-71162-6_24.
- [46] Orvalho, P., Janota, M., Manquinho, V.M., 2025. Counterexample guided program repair using zero-shot learning and maxsat-based fault localization, in: Walsh, T., Shah, J., Kolter, Z. (Eds.), *AAAI-25*, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA, AAAI Press. pp. 649–657. URL: <https://doi.org/10.1609/aaai.v39i1.32046>, doi:10.1609/AAAI.V39I1.32046.
- [47] Orvalho, P., Piepenbrock, J., Janota, M., Manquinho, V., 2022c. Project proposal: Learning variable mappings to repair programs. 7th Conference on Artificial Intelligence and Theorem Proving, AITP 2022.

- [48] Orvalho, P., Piepenbrock, J., Janota, M., Manquinho, V.M., 2023. Graph neural networks for mapping variables between programs, in: ECAI 2023 - 26th European Conference on Artificial Intelligence, IOS Press, Poland. pp. 1811–1818. URL: <https://doi.org/10.3233/FAIA230468>.
- [49] Perry, D.M., Kim, D., Samanta, R., Zhang, X., 2019. Semcluster: clustering of imperative programming assignments based on quantitative semantic features, in: McKinley, K.S., Fisher, K. (Eds.), Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019, Phoenix, AZ, USA, June 22–26, 2019, ACM. pp. 860–873.
- [50] Rand, W.M., 1971. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical association* 66, 846–850.
- [51] Reiss, S.P., 2009. Semantics-based code search, in: 31st International Conference on Software Engineering, ICSE 2009, May 16–24, 2009, Vancouver, Canada, Proceedings, IEEE. pp. 243–253. doi:10.1109/ICSE.2009.5070525.
- [52] Reynolds, D.A., 2009. Gaussian mixture models. *Encyclopedia of biometrics* 741, 659–663.
- [53] Rice, H.G., 1953. Classes of recursively enumerable sets and their decision problems. *Transactions of the American Mathematical Society* 74, 358–366.
- [54] Rocha, H.J.B., de Barros Costa, E., Tedesco, P.C.d.A.R., 2023. Helping to provide adaptive feedback to novice programmers: a framework to assist the teachers, in: 2023 18th Iberian Conference on Information Systems and Technologies (CISTI), IEEE. pp. 1–6.
- [55] Rosenberg, A., Hirschberg, J., 2007. V-measure: A conditional entropy-based external cluster evaluation measure, in: Proceedings of the 2007 joint conference on empirical methods in natural language processing and computational natural language learning (EMNLP-CoNLL), pp. 410–420.
- [56] Schleimer, S., Wilkerson, D.S., Aiken, A., 2003. Winnowing: Local algorithms for document fingerprinting, in: Halevy, A.Y., Ives, Z.G., Doan, A. (Eds.), Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data, San Diego, California, USA, June 9–12, 2003, ACM. pp. 76–85.
- [57] Schütze, H., Manning, C.D., Raghavan, P., 2008. Introduction to information retrieval. volume 39. Cambridge University Press Cambridge.
- [58] da Silva, P., 2019. SQUARES : A SQL Synthesizer Using Query Reverse Engineering. Master’s thesis. Instituto Superior Técnico, Lisboa, Portugal.
- [59] Steinhaus, H., 1956. Sur la division des corps matériels en parties. *Bull. Acad. Polon. Sci* 1, 801.
- [60] Stolee, K.T., Elbaum, S.G., 2012. Toward semantic search via SMT solver, in: Tracz, W., Robillard, M.P., Bultan, T. (Eds.), 20th ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE-20), SIGSOFT/FSE’12, Cary, NC, USA - November 11 - 16, 2012, ACM. p. 25. doi:10.1145/2393596.2393625.
- [61] Stolee, K.T., Elbaum, S.G., Dwyer, M.B., 2016. Code search with input/output queries: Generalizing, ranking, and assessment. *J. Syst. Softw.* 116, 35–48. doi:10.1016/j.jss.2015.04.081.
- [62] Strehl, A., Ghosh, J., 2002. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. *Journal of machine learning research* 3, 583–617.
- [63] Vinh, N.X., Epps, J., Bailey, J., 2009. Information theoretic measures for clusterings comparison: is a correction for chance necessary?, in: Proceedings of the 26th annual international conference on machine learning, pp. 1073–1080.
- [64] Wang, K., Singh, R., Su, Z., 2018. Search, align, and repair: data-driven feedback generation for introductory programming exercises, in: PLDI 2018, ACM. pp. 481–495.
- [65] Yang, B., Yang, J., 2020. Exploring the differences between plausible and correct patches at fine-grained level, in: 2020 IEEE 2nd International Workshop on Intelligent Bug Fixing (IBF), IEEE. pp. 1–8.
- [66] Ye, H., Gu, J., Martinez, M., Durieux, T., Monperrus, M., 2022. Automated classification of overfitting patches with statically extracted code features. *IEEE Trans. Software Eng.* 48, 2920–2938. URL: <https://doi.org/10.1109/TSE.2021.3071750>, doi:10.1109/TSE.2021.3071750.
- [67] Yi, J., Ahmed, U.Z., Karkare, A., Tan, S.H., Roychoudhury, A., 2017. A feasibility study of using automated program repair for introductory programming assignments, in: Bodden, E., Schäfer, W., van Deursen, A., Zisman, A. (Eds.), ESEC/FSE 2017, ACM. pp. 740–751.
- [68] Zimmerman, K., Rupakheti, C.R., 2015. An automated framework for recommending program elements to novices (N), in: Cohen, M.B., Grunski, L., Whalen, M. (Eds.), 30th IEEE/ACM International Conference on Automated Software Engineering, ASE 2015, Lincoln, NE, USA, November 9–13, 2015, IEEE Computer Society. pp. 283–288.
- [69] Zügner, D., Kirschstein, T., Catasta, M., Leskovec, J., Günnemann, S., 2021. Language-agnostic representation learning of source code from structure and context, in: 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3–7, 2021, OpenReview.net. pp. —.

A. Description of C-PACK-IPAS

The set of IPAS corresponds to three different lab classes of the introductory programming course to the C programming language at Instituto Superior Técnico, Universidade de Lisboa. Each lab class focuses on a different topic of the C programming language. Lab02 is described in Section A.1 deals with integers and input-output operations. Section A.2 presents Lab03, which focuses on loops and chars. Lastly, Section A.3 describes Lab04, where the students learn to use vectors and strings. The textual description of each programming assignment can be found in the next sections. The input/output tests used to evaluate semantically the set of students’ submissions can be found at the public GitHub <https://github.com/pmorvalho/C-Pack-IPAs> of C-PACK-IPAS [44]. Moreover, there is also a reference implementation for each IPA in the public git repository that can be used by program repair frameworks that only accept a single reference implementation to repair incorrect programs.

A.1. Lab02 - Integers and IO operations.

In Lab02, the students learn how to program with integers, floats, IO operations (mainly `printf` and `scanf`), conditionals (if-statements), and simple loops (for and while-loops).

IPA #1: *Lab02 - Ex01*. Write a program that determines and prints the largest of three integers given by the user.

IPA #2: *Lab02 - Ex02*. Write a program that reads two integers 'N, M' and prints the smallest of them in the first row and the largest in the second.

IPA #3: *Lab02 - Ex03*. Write a program that reads two positive integers 'N, M' and prints "yes" if 'M' is a divisor of 'N', otherwise prints "no".

IPA #4: *Lab02 - Ex04*. Write a program that reads three integers and prints them in order on the same line. The smallest number must appear first.

IPA #5: *Lab02 - Ex05*. Write a program that reads a positive integer 'N' and prints the numbers '1..N', one per line.

IPA #6: *Lab02 - Ex06*. Write a program that determines the largest and smallest number of 'N' real numbers given by the user. Consider that 'N' is a value requested from the user. The result must be printed with the command `'printf("min: %f, max: %f\n", min, max)'`. *Hint: initialize the largest and smallest to the first read value.*

IPA #7: *Lab02 - Ex07*. Write a program that asks the user for a positive integer 'N' and prints the number of divisors of 'N'. Remember that prime numbers have 2 divisors.

IPA #8: *Lab02 - Ex08*. Write a program that calculates and prints the average of 'N' real numbers given by the user. The program should first ask the user for an integer 'N', representing the number of numbers to be entered. The real numbers must be represented by float type. The result must be printed with the command `'printf("%.2f", avg)'`.

IPA #9: *Lab02 - Ex09*. Write a program that asks the user for a value 'N' corresponding to a certain period of time in seconds. The program should output this period of time in the format 'HH:MM:SS'.

Hint: use the operator that calculates the remainder of division ('%').

IPA #10: *Lab02 - Ex10*. Write a program that asks the user for a positive value 'N'. The output should present the number of digits that make up 'N' (on the first line), as well as the sum of the digits of 'N' (on the second line). For example, the number 12345 has 5 digits, and the sum of these digits is 15.

A.2. Lab03 - Loops and Chars.

In this lab, the students learn how to program with loops, nested loops, auxiliary functions, and chars.

IPA #11: *Lab03 - Ex01*. Write a program that draws a square of numbers like the following using the function `'void square(int N)'`. The value of 'N', given by the user, must be greater than or equal to 2. The tab (character ' ') must be used as the separator. The square shown is the example for 'N = 5'.

```
1 2 3 4 5
2 3 4 5 6
3 4 5 6 7
4 5 6 7 8
5 6 7 8 9
```

IPA #12: *Lab03 - Ex02*. Write a program that draws a pyramid of numbers using the `'void pyramid(int N)'` function. The value of 'N', given by the user, must be greater than or equal to 2. The space (character ' ') must be used as the separator. The pyramid shown is the example for 'N = 5'.

```
1
1 2 1
1 2 3 2 1
```

```

1 2 3 4 3 2 1
1 2 3 4 5 4 3 2 1

```

IPA #13: *Lab03 - Ex03*. Write a program that draws a cross on diagonals using the ‘void cross(int N);’ function. The asterisk (‘*’) character must be used to draw the cross; hyphen (‘-’) character must be used as the separator. The crosses shown are the examples for ‘N = 3’ and ‘N=8’.

```

  * - *
  - * -
  * - *

* - - - - *
- * - - - * -
- - * - - * -
- - - * * - -
- - - * * - -
- - * - - * -
- * - - - * -
* - - - - *

```

IPA #14: *Lab03 - Ex04*. Write a program that reads a sequence of numbers separated by spaces and newlines, and print the same string, but the numbers in the output should not contain 0 at the beginning, eg ‘007’ should print ‘7’. The exception is the number 0, which should be printed as 0. The string in the input ends with ‘EOF’.

Warning: Number values may be greater than the maximum value of type ‘int’ or any primitive type in C.

Hint: the ‘int getchar()’ function can be used to read a character.

IPA #15: *Lab03 - Ex05*. Write a program that reads a sequence of messages and prints them out, one per line. Each message is delimited by quotation marks (character ‘”’). The message can contain an "escape sequence" - the character loses special meaning if it is preceded by the character ‘\’ (backslash). For example, the input “afoõbar” matches the message ‘a"foõbar"’. So the backslash allows you to include quotes in the message just like the backslash itself.

Hint: use the concept of state as we did in the word counter in the theoretical class.

IPA #16: *Lab03 - Ex06*. Write a program that reads a positive integer from the input (such as a sequence of characters up to 100 chars) and that decides whether the number read is divisible by 9. If the number is divisible by 9, the program should print the message ‘yes’, and should print ‘no’ otherwise.

Warning: Number values can be greater than the maximum value of type ‘int’ or any primitive type in C.

Hint: A number is divisible by 9 if and only if the sum of its digits is divisible by 9.

For example, the sum of the digits of the number 729 is 18, so it is divisible by 9. The fact can be seen from the following equation: $7 \times 100 + 2 \times 10 + 9 = (7 \times 99 + 7) + (2 \times 9 + 2) + 9$.

IPA #17: *Lab03 - Ex07*. Write a program that takes a sequence of numbers and operators (‘+’, ‘-’) representing an arithmetic expression and returns the result of that arithmetic expression. The string in the input ends with ‘ñ’. You can assume that every two numbers are always separated by ‘space, operator, space’, i.e., ‘op’, for either of the 2 operators above. Example: Input ‘70 + 22 - 3’ should return ‘89’.

Hint: You should start by converting a sequence of digits (characters) to an integer.

A.3. Lab04 - Vectors and Strings.

In this lab, the students learn how to program with integers arrays, and strings.

IPA #18: *Lab04 - Ex01*. Write a program that asks the user for a positive integer ‘n < VECMAX’, where ‘VECMAX=100’. Then read ‘n’ positive integers. In the end, the program should write a graphical representation of the values read as follows. The graph shown is the example for ‘n = 3’ and values ‘1 3 4’.

```

  *
  ***
  ****

```

IPA #19: *Lab04 - Ex02*. Write a program that asks the user for a positive integer ‘ $n < \text{VECMAX}$ ’, where ‘ $\text{VECMAX}=100$ ’. Then read ‘ n ’ positive integers. In the end, the program should write a graphical representation of the values read as follows. The graph shown is the example for ‘ $n = 3$ ’ and values ‘1 3 4’.

```
***
**
**
*
```

IPA #20: *Lab04 - Ex03*. Write a program that asks the user for a positive integer ‘ $n < \text{VECMAX}$ ’, where ‘ $\text{VECMAX}=100$ ’. Then read ‘ n ’ positive integers. In the end, the program should write a graphical representation of the values read as follows. The graph shown is the example for ‘ $n = 3$ ’ and values ‘1 3 4’.

```
*
**
**
***
```

Consider that in the following IPAs, all strings have a maximum of ‘ $\text{MAX} = 80$ ’ characters (including the end-of-string character).

IPA #21: *Lab04 - Ex04*. Write a program that reads a word from the terminal and checks whether the word is a palindrome or not. A word is a palindrome if it is spelled the same way from left to right and vice versa (eg "AMA" is a palindrome). If the word is a palindrome, the program should print the value ‘yes’, and ‘no’ if not.

Hint: You can use ‘scanf("%s", s)’ to read a word. Note that the string ‘s’ does not ask for ‘&’ in ‘scanf’.

IPA #22: *Lab04 - Ex05*. Write a program that reads characters from the keyboard, character by character until it finds the character ‘ñ’ or EOF and writes the line read to the terminal. Implement the ‘int leLinha(char s[])’ function, which reads the line into the string ‘s’ and returns the number of characters read. *Hint: After solving this exercise, try using the ‘fgets’ command.*

IPA #23: *Lab04 - Ex06*. Write a program that reads a line from the terminal (use the function from the previous exercise) and writes the same text to the terminal but with the lowercase letters replaced by the respective uppercase letters. Implement the ‘void uppercase(char s[])’ function. *Note: Remember that the string ‘s’ is changed by the ‘uppercase’ function.*

IPA #24: *Lab04 - Ex07*. Write a program that reads a line and a character and writes to the terminal the same line where all occurrences of the character were removed. Implement the ‘void eraseCharacter(char s[], char c)’ function that erases the character ‘c’ from the string ‘s’.

IPA #25: *Lab04 - Ex08*. Write a program that reads two integers in decimal representation and prints the larger of those two numbers. You can assume that the two numbers have the same number of digits and a maximum of 100 characters.

Note: The numbers may be too large to be stored in a ‘long long’ variable, for example ‘99888888888888888887’ and ‘99888888888888888888’.

B. Use Case #1: Clustering IPAs

B.1. Clustering Accuracy

Figure 11 shows a matrix with the different values of the cluster accuracy using the MINIBATCH KMEANS algorithm on each program representation using ten different seeds. Each entry is highlighted accordingly to its value. The lowest value is highlighted in black, and the highest is highlighted in white. Intermediate values are highlighted in different shades of grey, depending on how far they are from the lowest value.

Secondly, Figure 12 shows a matrix with the different values of the cluster accuracy using the BIRCH algorithm on each program representation using ten different seeds.

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.81	0.79	0.82	0.76	0.81	0.79	0.79	0.80	0.74	0.81
AAST	0.74	0.72	0.75	0.71	0.71	0.75	0.76	0.75	0.66	0.70
Invariants	0.77	0.77	0.74	0.73	0.74	0.76	0.77	0.76	0.79	0.75
Syntax	0.58	0.58	0.57	0.59	0.61	0.61	0.58	0.60	0.59	0.55

Figure 11: The values for cluster accuracy using the **MINIBATCH KMEANS** algorithm on each program representation after ten different runs, each run using a different seed.

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.80	0.81	0.79	0.80	0.80	0.80	0.80	0.80	0.80	0.81
AAST	0.79	0.72	0.77	0.74	0.72	0.73	0.76	0.75	0.70	0.71
Invariants	0.78	0.79	0.77	0.77	0.78	0.78	0.77	0.77	0.79	0.80
Syntax	0.58	0.59	0.59	0.59	0.59	0.60	0.60	0.60	0.57	0.59

Figure 12: The values for cluster accuracy using the **BIRCH** algorithm on each program representation after ten different runs, each run using a different seed.

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.74	0.80	0.76	0.81	0.80	0.78	0.80	0.79	0.79	0.81
AAST	0.72	0.65	0.67	0.74	0.72	0.72	0.71	0.72	0.71	0.74
Invariants	0.75	0.78	0.75	0.79	0.74	0.72	0.73	0.76	0.74	0.80
Syntax	0.58	0.56	0.56	0.60	0.56	0.58	0.60	0.58	0.57	0.57

Figure 13: The values for cluster accuracy using the **GAUSSIAN MIXTURE** algorithm on each program representation after ten different runs, each run using a different seed.

Lastly, Figure 13 shows a matrix with the different values of the cluster accuracy using the **GAUSSIAN MIXTURE** algorithm on each program representation for ten different seeds.

B.2. Other Evaluation Metrics

In this section, we present other clustering evaluation metrics for the **KMEANS** algorithm, such as: the *Rand index*, the *adjusted Rand index*, the *normalized mutual information*, the *adjusted mutual information*, the *Fowlkes–Mallows index*, the *completeness score*, the *homogeneity score*, and the *V measure*.

Rand Index. The *Rand index* measures the similarity of the two assignments, ignoring permutations [50]. The Rand index is given by the following equation 1:

$$RI = \frac{TP + TN}{TP + FP + FN + TN} \quad (1)$$

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98
AAST	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.97
Invariants	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98
Syntax	0.96	0.96	0.96	0.96	0.96	0.96	0.96	0.96	0.96	0.96

Figure 14: The values for the Rand index using the KMEANS algorithm on each program representation after ten different runs, each run using a different seed.

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.79	0.76	0.76	0.77	0.77	0.79	0.77	0.77	0.77	0.78
AAST	0.68	0.68	0.64	0.68	0.67	0.67	0.68	0.67	0.67	0.66
Invariants	0.73	0.74	0.73	0.75	0.75	0.75	0.73	0.74	0.74	0.75
Syntax	0.49	0.49	0.50	0.51	0.51	0.53	0.50	0.50	0.47	0.50

Figure 15: The values for the adjusted Rand index using the KMEANS algorithm on each program representation after ten different runs, each run using a different seed.

Figure 14 presents the values for the Rand index using the KMEANS algorithm on each program representation after ten different runs.

Adjusted Rand Index. The *adjusted Rand index* is the corrected-for-chance version of the Rand index since the Rand index does not guarantee that random label assignments will get a value close to zero [24]. The adjusted Rand index is given by equation 2, where TP is the number of true positives, TN is the number of true negatives, FP is the number of false positives, and FN is the number of false negatives.

$$ARI = \frac{RI - E[RI]}{\max(RI) - E[RI]} \quad (2)$$

Figure 15 presents the values for the adjusted Rand index using the KMEANS algorithm on each program representation after ten different runs.

Normalized Mutual Information. The *normalized mutual information* of two random variables is a measure of the mutual dependence between the two variables [62]. Figure 16 shows the values for the normalized mutual information using the KMEANS algorithm on each program representation after 10 different runs.

Adjusted Mutual Information. The *adjusted mutual information* corrects the effect of the agreement solely due to chance between clusterings, similar to the way the adjusted rand index corrects the Rand index [63]. Figure 17 presents the values for the adjusted mutual information using the KMEANS algorithm on each program representation after ten different runs.

Fowlkes–Mallows index. The *Fowlkes–Mallows index* measures the similarity between two clusters. A high value for the Fowlkes–Mallows index indicates a great similarity between the clusters and the benchmark classifications [15]. The Fowlkes–Mallows index can be computed using equation 3, where TP is the number of true positives, FP is the

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.88	0.87	0.87	0.87	0.87	0.88	0.87	0.87	0.87	0.88
AAST	0.83	0.83	0.82	0.84	0.83	0.83	0.83	0.84	0.83	0.83
Invariants	0.85	0.85	0.86	0.86	0.86	0.86	0.85	0.85	0.85	0.86
Syntax	0.69	0.70	0.70	0.71	0.70	0.71	0.70	0.71	0.69	0.70

Figure 16: The values for the normalized mutual information using the KMEANS algorithm on each program representation after ten different runs, each run using a different seed.

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.87	0.86	0.86	0.86	0.86	0.87	0.86	0.86	0.86	0.87
AAST	0.82	0.82	0.81	0.83	0.82	0.82	0.82	0.82	0.82	0.81
Invariants	0.84	0.84	0.85	0.85	0.85	0.85	0.83	0.84	0.84	0.85
Syntax	0.67	0.68	0.68	0.69	0.68	0.69	0.68	0.69	0.67	0.68

Figure 17: The values for the adjusted mutual information using the KMEANS algorithm on each program representation after ten different runs, each run using a different seed.

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.80	0.77	0.77	0.78	0.78	0.80	0.78	0.78	0.78	0.79
AAST	0.69	0.69	0.66	0.69	0.68	0.68	0.70	0.69	0.68	0.67
Invariants	0.75	0.75	0.74	0.76	0.76	0.76	0.74	0.76	0.75	0.76
Syntax	0.52	0.51	0.52	0.53	0.53	0.55	0.52	0.52	0.50	0.53

Figure 18: The values for the Fowlkes–Mallows index using the KMEANS algorithm on each program representation after ten different runs, each run using a different seed.

number of false positives, and FN is the number of false negatives. TPR is the true positive rate, also called sensitivity or recall, and PPV is the positive predictive rate, also known as precision.

$$FM = \sqrt{PPV \cdot TPR} = \sqrt{\frac{TP}{TP + FP} \cdot \frac{TP}{TP + FN}} \quad (3)$$

Figure 18 shows the values for the Fowlkes–Mallows index using the KMEANS algorithm on each program representation after ten different runs.

Completeness score. The *completeness score* measure if all members of a given class are assigned to the same cluster [55]. Figure 19 shows the values for the completeness score using the KMEANS algorithm on each program representation after ten different runs.

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.88	0.87	0.87	0.87	0.87	0.88	0.87	0.87	0.87	0.88
AAST	0.84	0.83	0.82	0.84	0.83	0.83	0.83	0.84	0.83	0.83
Invariants	0.86	0.86	0.86	0.86	0.86	0.86	0.85	0.86	0.85	0.86
Syntax	0.70	0.71	0.71	0.72	0.70	0.71	0.70	0.71	0.70	0.71

Figure 19: The values for the completeness score using the KMEANS algorithm on each program representation after ten different runs, each run using a different seed.

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.88	0.87	0.87	0.87	0.87	0.88	0.87	0.87	0.87	0.88
AAST	0.83	0.83	0.82	0.84	0.83	0.83	0.83	0.83	0.83	0.83
Invariants	0.85	0.85	0.85	0.85	0.85	0.86	0.84	0.85	0.85	0.85
Syntax	0.68	0.68	0.69	0.70	0.69	0.70	0.69	0.70	0.68	0.69

Figure 20: The values for the homogeneity score using the KMEANS algorithm on each program representation after ten different runs, each run using a different seed.

	1	2	3	4	5	6	7	8	9	10
AAST+Invariants	0.88	0.87	0.87	0.87	0.87	0.88	0.87	0.87	0.87	0.88
AAST	0.83	0.83	0.82	0.84	0.83	0.83	0.83	0.84	0.83	0.83
Invariants	0.85	0.85	0.86	0.86	0.86	0.86	0.85	0.85	0.85	0.86
Syntax	0.69	0.70	0.70	0.71	0.70	0.71	0.70	0.71	0.69	0.70

Figure 21: The values for the V measure using the KMEANS algorithm on each program representation after ten different runs, each run using a different seed.

Homogeneity score. The *homogeneity score* checks if each cluster contains only members of a single class [55]. Figure 20 presents the values for the homogeneity score using the KMEANS algorithm on each program representation after ten different runs.

V measure. The *V measure* is the harmonic mean between the homogeneity and completeness scores [55]. Figure 21 shows the values for the V measure using the KMEANS algorithm on each program representation after ten different runs.

C. Use Case #2: Repairing IPAs

Table 5 presents the number of clusters each clustering method uses for each IPA. This table was used to generate the cactus plot in Figure 7.

Table 5

The number of clusters generated using each clustering approach for each IPA.

Exercise	#Correct Submissions	Clara Clusters	KMeans AAST	KMeans AAST+Invs	KMeans Invs	KMeans Syntax	Closest Program (KMeans) AAST+Invs
lab02/ex01	92	18	9	9	9	9	1
lab02/ex02	84	6	8	8	8	8	1
lab02/ex03	83	4	8	8	7	8	1
lab02/ex04	76	20	7	7	7	7	1
lab02/ex05	80	10	7	7	7	7	1
lab02/ex06	68	25	6	6	6	6	1
lab02/ex07	67	17	6	6	6	6	1
lab02/ex08	49	21	4	4	4	4	1
lab02/ex09	74	12	7	7	7	7	1
lab02/ex10	65	17	6	6	6	6	1
lab03/ex01	70	51	6	6	6	6	1
lab03/ex02	55	49	5	5	5	5	1
lab03/ex03	45	27	4	4	4	4	1
lab03/ex04	28	8	2	2	2	2	1
lab03/ex06	46	8	4	4	4	4	1
lab04/ex01	59	32	5	5	5	5	1
lab04/ex02	47	32	4	4	4	4	1
lab04/ex03	41	33	4	4	4	4	1
lab04/ex08	8	6	1	1	1	1	1
lab05/ex01	4	3	1	1	1	1	1
itsp/lab3/ex2810	17	9	1	1	1	1	1
itsp/lab3/ex2811	7	3	1	1	1	1	1
itsp/lab3/ex2812	17	7	1	1	1	1	1
itsp/lab3/ex2813	4	4	1	1	1	1	1
itsp/lab4/ex2824	15	5	1	1	1	1	1
itsp/lab4/ex2825	10	4	1	1	1	1	1
itsp/lab4/ex2827	6	6	1	1	1	1	1
itsp/lab4/ex2831	7	4	1	1	1	1	1
itsp/lab4/ex2832	17	7	1	1	1	1	1
itsp/lab4/ex2833	19	9	1	1	1	1	1
itsp/lab5/ex2865	7	4	1	1	1	1	1
itsp/lab5/ex2866	11	10	1	1	1	1	1
itsp/lab5/ex2867	7	4	1	1	1	1	1
itsp/lab5/ex2868	8	6	1	1	1	1	1
itsp/lab5/ex2869	7	4	1	1	1	1	1
itsp/lab5/ex2870	9	8	1	1	1	1	1
itsp/lab5/ex2871	15	10	1	1	1	1	1
itsp/lab6/ex2932	3	3	1	1	1	1	1
itsp/lab6/ex2933	1	1	1	1	1	1	1
itsp/lab6/ex2936	5	4	1	1	1	1	1
itsp/lab6/ex2937	2	2	1	1	1	1	1
itsp/lab6/ex2938	6	4	1	1	1	1	1
itsp/lab6/ex2939	2	1	1	1	1	1	1