# Search, Align, and Repair: Data-Driven Feedback Generation for Introductory Programming Exercises

Ke Wang University of California, Davis kbwang@ucdavis.edu Rishabh Singh Microsoft Research risin@microsoft.com

# Zhendong Su University of California, Davis su@cs.ucdavis.edu

Approach	No Manual Effort	Minimal Repair	Complex Repairs	Data-Driven	Production Deployment
AutoGrader [24]	X	1	X	×	X
CLARA [11]	1	X	1	1	×
QLOSE [3]	X	X	X	×	X
sk_p [22]	1	X	1	1	X
REFAZER [23]	1	X	×	1	×
CoderAssist [16]	X	X	1	~	×
SARFGEN	1	1	1	1	1

**Table 1.** Comparison of SARFGEN against the existing feedback generation approaches.

learning. Most of these techniques model the problem as a *program repair* problem: repairing an incorrect student submission to make it functionally equivalent *w.r.t.* a given specification (*e.g.*, a reference solution or a set of test cases). Table 1 summarizes some of the major techniques in terms of their capabilities and requirements, and compares them with our proposed technique realized in the SARFGEN<sup>1</sup> system.

As summarized in the table, existing systems still face important challenges to be effective and practical. In particular, AutoGrader [24] requires a custom error model per programming problem, which demands manual effort from the instructor and her understanding of the system details. Moreover, its reliance on constraint solving to find repairs makes it expensive and unsuitable in interactive settings at the MOOC scale. Systems like CLARA [11] and sk\_p [22] can generate repairs relatively quickly by using clustering and machine learning techniques on student data, but the generated repairs are often imprecise and not minimal, reducing the quality of the feedback. CLARA's use of Integer Linear Programming (ILP) for variable alignment during repair generation also hinders its scalability. Section 5 provides a detailed survey of related work.

To tackle the weaknesses of existing systems, we introduce "Search, Align, and Repair", a conceptual framework for program repair from a data-driven perspective. First, given an incorrect program, we search for similar reference solutions. Second, we align each statement in the incorrect program with a corresponding statement in the reference solutions to identify discrepancies for suggesting changes. Finally, we pinpoint minimal fixes to patch the incorrect program. Our key observation is that the diverse approaches and errors students make are captured in the abundant previous submissions because of the MOOC scale [5]. Thus, we aim at a fully automated, data-driven approach to generate *instant* (to be

# Abstract

1

2

3

4

5

6

7

8

9

10

29

30

This paper introduces the "Search, Align, and Repair" data-11 driven program repair framework to automate feedback gen-12 eration for introductory programming exercises. Distinct 13 from existing techniques, our goal is to develop an efficient, 14 fully automated, and problem-agnostic technique for large or 15 MOOC-scale introductory programming courses. We lever-16 age the large amount of available student submissions in such 17 settings and develop new algorithms for identifying simi-18 lar programs, aligning correct and incorrect programs, and 19 repairing incorrect programs by finding minimal fixes. We 20 have implemented our technique in the SARFGEN system and 21 evaluated it on thousands of real student attempts from the 22 Microsoft-DEV204.1x edX course and the Microsoft Code-23 Hunt platform. Our results show that SARFGEN can, within 24 two seconds on average, generate concise, useful feedback 25 for 89.7% of the incorrect student submissions. It has been 26 integrated with the Microsoft-DEV204.1X edX class and de-27 ployed for production use. 28

# 1 Introduction

31 The unprecedented growth of technology and computing 32 related jobs in recent years [25] has resulted in a surge in 33 Computer Science enrollments at colleges and universities, 34 and hundreds of thousands of learners worldwide signing 35 up for Massive Open Online Courses (MOOCs). While larger 36 classrooms and MOOCs have made education more accessi-37 ble to a much more diverse and greater audience, several key 38 challenges remain to ensure comparable education quality 39 to that of traditional smaller classroom settings. This paper 40 tackles one such challenge: providing fully automated, per-41 sonalized feedback to students for introductory programming 42 exercises without requiring any instructor effort. Even though 43 introductory programming exercises require relatively small 44 program size, relevant literature [7, 24] has shown that stu-45 dents still struggle and need effective tools to help them, 46 further highlighting the need of automated feedback tech-47 nology. 48

The problem of automated feedback generation for introductory programming courses has seen much recent interest — many systems have been developed using techniques from formal methods, programming languages, and machine

**54** 2017.

55

49

50

51

52

 $^1\underline{S}earch, \underline{A}lign, and \underline{R}epair for \underline{F}eedback \, \underline{GEN}eration$ 

<sup>&</sup>lt;sup>53</sup> PL'17, January 01–03, 2017, New York, NY, USA

116

153

154

155

156

158

159

160

161

162

163

164

165

interactive), *minimal* (to be precise) and *semantic* (to allow
complex repairs) fixes to incorrect student submissions. At
the technical level, we need to address three key challenges:

- **Search:** Given an incorrect student program, how to efficiently identify a set of closely-related candidate programs among all correct solutions?
- Align: How to efficiently and precisely align each selected
   program with the incorrect program to compute a
   correction set that consists of expression- or statement level discrepancies?
- **Repair:** How to quickly identify a minimal set of fixes out of an entire correction set?

124 For the "Search" challenge, we identify syntactically most 125 similar correct programs w.r.t. the incorrect program in a 126 coarse-to-fine manner. First, we perform exact matching 127 on the control flow structures of the incorrect and correct 128 programs, and rank matched programs w.r.t. the tree edit dis-129 tance between their abstract syntax trees (ASTs). Although 130 introductory programming assignments often feature only 131 lightweight small programs, the massive number of submis-132 sions in MOOCs still makes the standard tree edit distance 133 computation too expensive for the setting. Our solution is 134 based on a new tree embedding scheme for programs using 135 numerical embedding vectors with a new distance metric 136 in the Euclidean space. The new program embedding vec-137 tors (called *position-aware characteristic vectors*) efficiently 138 capture the structural information of the program ASTs. Be-139 cause the numerical distance metric is easier and faster to 140 compute, program embedding allows us to scale to a large 141 number of programs.

142 For the "Align" challenge, we propose a usage-based  $\alpha$ -143 conversion to rename variables in a correct program using 144 those from the incorrect program. In particular, we represent 145 each variable with the new embedding vector based on its 146 usage profile in the program, and compute the mapping be-147 tween two sets of variables using the Euclidean distance met-148 ric. In the next phase, we split both programs into sequences 149 of basic blocks and align each pair accordingly. Finally, we 150 construct discrepancies by matching statements only from 151 the aligned basic blocks. 152

For the final "Repair" challenge, we *dynamically* minimize the set of corrections needed to repair the incorrect program from the large set of possible corrections generated by the alignment step. We present some optimizations that gain significant speed-up over the enumerative search.

Our automated program repair technique offers several important benefits:

- Fully Automated: It does not require any manual effort during the complete repair process.
- **Minimal Repair**: It produces a *minimal* set of corrections (*i.e.*, any included correction is relevant and necessary) that can better help students.

- Unrestricted Repair: It supports both simple and complex repair modifications to the incorrect program without changing its control-flow structure.
- **Portable**: Unlike most previous approaches, it is independent of the programming exercise it only needs a set of correct submissions for each exercise.

We have implemented our technique in the SARFGEN system and extensively evaluated it on thousands of programming submissions obtained from the Microsoft-DEV204.1x edX course and the CodeHunt platform [28]. SARFGEN is able to repair 89.7% of the incorrect programs with *minimal* fixes in under two seconds on average. In addition, it has been integrated with the Microsoft-DEV204.1x course and deployed for production use. The feedback collected from online students demonstrates its practicality and usefulness.

This paper makes the following main contributions:

- We propose a high-level data-driven framework search, align and repair to fix programming submissions for introductory programming exercises.
- We present novel instantiations for different framework components. Specifically, we develop novel program embeddings and the associated distance metric to efficiently and precisely identify similar programs and compute program alignment.
- We present an extensive evaluation of SARFGEN on repairing thousands of student submissions on 17 different programming exercises from the Microsoft-DEV204-.1x edx course and the Microsoft CodeHunt platform.

# 2 Overview

This section gives an overview of our approach by introducing its key ideas and high-level steps.

X 0 X 0 X 0 X 0	
0 X 0 X 0 X 0 X	
X 0 X 0 X 0 X 0	
0 X 0 X 0 X 0 X	
X 0 X 0 X 0 X 0	
0 X 0 X 0 X 0 X	
X 0 X 0 X 0 X 0	
0 X 0 X 0 X 0 X	

Figure 1. Desired output for the chessboard exercise.

#### 2.1 Example: The Chessboard printing problem

The Chessboard printing assignment, taken from the edX course of C# programming, requires students to print the pattern of chessboard using "X" and "0" characters to represent the squares as shown in Figure 1.

The goal of this problem is to teach students the concept of conditional and looping constructs. On this problem, students struggled with many issues, such as loop bounds or conditional predicates. In addition, they also had trouble 212

213

214

215

216

217

218

219



bining looping and conditional constructs. For example, one common error we observed among student attempts was that they managed to alternate "X" and "0" for each separate row/column, but (mistakenly) repeated the same pattern for all rows/columns (rather than flipping for consecutive rows/columns).

One key challenge in providing feedback on programming submissions is that a given problem can be solved using many different algorithms. In fact, among all the correct student submissions, we found 337 different control-flow structures, indicating the fairly large solution space one needs to consider. It is worth mentioning that modifying incorrect programs to merely comply with the specification using a reference solution may be inadequate. For example, completely rewriting an incorrect program into a correct program will achieve this goal, however such a repair might lead a student to a completely different direction, and would not help her understand the problems with her own approach. Therefore, it is important to pinpoint minimal modifications in their particular implementation that addresses the root cause of the problem. In addition, efficiency is important especially in an interactive setting like MOOC where students expect instant feedback within few seconds and also regarding the deployment cost for thousands of students.

317

318

319

320

321

322

323

324

325

326

327

328

329

330

274 275

262

263

264

265

266

267

268

269

270

271

272

Given the three different incorrect student submissions shown in Figure 2, SARFGEN generates the feedback depicted in Figure 4 in under two seconds each. During these two seconds, SARFGEN searches over more than two thousand reference solutions, and selects the programs in Figure 3 for each incorrect submission to be compared against. For brevity, we only show the top-1 candidate program in the comparison set. SARFGEN then collects and minimizes the set of changes required to eliminate the errors. Finally, it produces the feedback which consists of the following information (highlighted in bold in Figure 4 for emphasis): 

- The number of changes required to correct the errors in the program.
- The location where each error occurred denoted by the line number.
- The problematic expression/statement in the line.
  - The problematic sub-expression/expression that needs to be corrected.
  - The new value of sub-expression/expression.

SARFGEN is customizable in terms of the level of the feedback an instructor would like to provide to the students. The generated feedback could consist of different combinations of the five kinds of information, which enables a more personalized, adaptive, and dynamic tutoring workflow.

## 2.2 Overview of the Workflow

 We now briefly present an overview of the workflow of the SARFGEN system. The three major components are outlined in Figure 5: (1) *Searcher*, (2) *Aligner*, and (3) *Repairer*.



Figure 5. The overall workflow of the SARFGEN system.

- (1) **Searcher**: The Searcher component takes as input an incorrect program  $P_e$  and searches for the top k closest solutions among all the correct programs  $P_1, P_2, P_3, ...$  and  $P_n$  for the given exercise. The key challenge is to search a large number of programs in an efficient and precise manner. The Searcher component consists of:
  - Program Embedder: The Program Embedder converts P<sub>e</sub> and P<sub>1</sub>, P<sub>2</sub>, P<sub>3</sub>, ...P<sub>n</sub> into numerical vectors. We propose a new scheme of embedding programs that improves the original characteristic vector representation used previously in Jiang et al. [12].
  - **Distance Calculator**: Using program embeddings, the Distance Calculator computes the top *k* closest reference solutions in the Euclidean space. The advantage is such distance computations are much more scalable than the tree edit distance algorithm.

- (2) **Aligner**: After computing a set of top k candidate programs for comparison, the Aligner component aligns each of the candidate programs  $P_1, ..., P_k$  w.r.t.  $P_e$ .
  - α-conversion: We propose a usage-set based α-conversion. Specifically, we profile each variable on its usage and summarize it into one single numeric vector via our new embedding scheme. Next, we minimize the distance between two sets of variables based on the Euclidean distance metric. This novel technique not only enables efficient computation but also achieves high accuracy.
  - Statement Matching: After the  $\alpha$ -conversion, we align basic blocks, and in turn individual statements within each aligned basic blocks.

Based on the alignment, we produce a set of syntactical discrepancies  $C(P_e, P_k)$ . This is an important step as misaligned programs would result in an imprecise set of corrections.

(3) **Repairer**: Given  $C(P_e, P_k)$ , the Repairer component produces a set of fixes  $\mathcal{F}(P_e, P_k)$ . Later it minimizes  $\mathcal{F}(P_e, P_k)$  by removing the syntactic/semantic redundancies that are unrelated to the root cause of the error. We propose our minimization technique based on an optimized dynamic analysis to make the minimal repair computation efficient and scalable.

# 3 The Search, Align, and Repair Algorithm

This section presents our feedback generation algorithm. In particular, it describes the three key functional components: *Search, Align,* and *Repair* to cope with the challenges discussed in Section 1.

# 3.1 Search

To realize our goal of using correct programs to repair an incorrect program, the very first problem we need to solve is to identify a small subset of correct programs among thousands of submissions that are relevant to fixing the incorrect program in an efficient and precise manner. To start with, we perform exact matching between reference solutions and the incorrect program *w.r.t.* their control-flow structures.

**Definition 3.1.** (Control-Flow Structure) Given a program P, its control-flow structure, CF(P), is a syntactic registration of how control statements in P are coordinated. For brevity, we denote CF(P) to simply be a sequence of symbols (*e.g.* For<sub>start</sub>, For<sub>start</sub>, If<sub>start</sub>, If<sub>end</sub>, Else<sub>start</sub>, Else<sub>end</sub>, For<sub>end</sub>, For<sub>end</sub> for the program in Figure 2a).

Given the selected programs with the same control-flow structure, we search for similar programs using a syntactic approach (in contrast to a dynamic approach) for two reasons: (1) less overhead (*i.e.* faster) and (2) more fault-tolerant as runtime dynamic *traces* likely lead to greater differences if students made an error especially on control predicates. However, using the standard tree edit distance [27] as the

497



**Figure 6.** The eight 2-level atomic tree patterns for a label set of  $\{a, \epsilon\}$  on the left. Using 2-level characteristic vector to embed the tree on the right:  $\langle 1, 1, 0, 0, 0, 0, 0, 1 \rangle$ ; and 2-level position-aware characteristic vector:  $\langle \langle 4 \rangle, \langle 3 \rangle, ..., \langle 2 \rangle \rangle$ .

449

450

451

452

453

distance measure does not scale to a large number of programs. We propose a new method of embedding ASTs into
numerical vectors, namely the position-aware characteristic
vectors, with which Euclidean distance can be computed
to represent the syntactic distance. Next, we briefly revisit
Definitions 3.2 and 3.3 proposed in [12], upon which our
new embedding is built.

Given a binary tree, we define a family of atomic tree
patterns to capture structural information of a tree. They are
parametrized by a parameter q, the height of the patterns.

**Definition 3.2.** (*q*-Level Atomic Tree Patterns) A *q*-level atomic pattern is a complete binary tree of height *q*. Given a label set *L*, including the empty label  $\epsilon$ , there are at most  $|L|^{2^{q}-1}$  distinct *q*-level atomic patterns.

**Definition 3.3.** (*q*-Level Characteristic Vector) Given a tree *T*, its *q*-level characteristic vector is  $\langle b_1, ..., b_{|L|^{2^q-1}} \rangle$ , where *b<sub>i</sub>* is the number of occurrences of the *i*-th *q*-level atomic pattern in *T*.

473 Figure 6 depicts an example. Given the label set  $\{a, \epsilon\}$ , the 474 tree on the right can be embedded into  $\langle 1, 1, 0, 0, 0, 0, 0, 1 \rangle$  us-475 ing 2-level characteristic vectors. The benefit of such embed-476 ding schemes is the realization of gauging program similarity 477 using Hamming distance/Euclidean distance metric which is 478 much faster to compute. In order to further enhance the pre-479 cision of the AST embeddings, we introduce a position-aware 480 characteristic vector embedding which incorporates more 481 structural information into the encoding vectors. 482

483 **Definition 3.4.** (*q*-Level Position-Aware Characteristic Vec-484 tor) Given a tree *T*, its *q*-level position-aware characteristic 485 vector is  $\langle \langle b_1^1, ..., b_1^{n_1} \rangle, \langle b_2^1, ..., b_2^{n_2} \rangle, \langle b_{|L|^{2q}-1}^L, ..., b_{|L|^{2q}-1}^{n_{|L|^{2q}-1}} \rangle \rangle$ , 486 487 where  $b_i^j$  is the height of the root node of the *j*-th occurrence 488 of the *i*-th *q*-level atomic pattern in *T*.

489 Essentially we expand each element  $b_i$  defined in the q-490 level characteristic vector into a list of heights for each *i*-th 491 q-level atomic pattern in T ( $b_i = |\langle b_i^1, ..., b_i^{n_i} \rangle|$ ). Using a 2-492 level position-aware characteristic vector, the same tree in 493 Figure 6 can be embedded into  $\langle \langle 4 \rangle, \langle 3 \rangle, ..., \langle 2 \rangle \rangle$ . For brevity, 494 we shorten the q-level position-aware characteristic vector 495 to be  $\langle h_{b_1}, ..., h_{b_{|L|^{2^q-1}}} \rangle$  where  $h_{b_i}$  represents the vector of heights of all *i*-th *q*-level atomic tree patterns. The distance between two *q*-level position-aware characteristic vector is:

$$\sum_{i=1}^{\lfloor 2^{q-1}} ||sort(h_{b_i}, \phi) - sort(h'_{b_i}, \phi)||_2^2$$
(1)

where 
$$\phi = max(|h_{b_i}|, |h'_{b_i}|)$$
 (2)

where *sort*( $h_{b_i}$ ,  $\phi$ ) means sorting  $h_{b_i}$  in descending order followed by padding zeros to the end if  $h_{b_i}$  happens to be the smaller vector of the two. The idea is to normalize  $h_{b_i}$  and  $h'_{b_i}$ prior to the distance calculation.  $||...||_2^2$  denotes the square of the L2 norm. We are given the incorrect program (having mnodes in its AST) and p candidate solutions (assuming each has the same number of nodes *n* in their ASTs to simplify the calculation). Creating the embedding as well as computing the Euclidean distance on the resulting vectors has a worstcase complexity of  $O(m + \rho * n + \rho * (|h_{b_1}|, ..., |h_{b_{|L|^{2q}-1}}|))$ . Because the position-aware characteristic vectors can be computed offline for the correct programs, we can further reduce the time complexity to  $O(m + \rho * (|h_{b_1}|, ..., |h_{b_{|r|}, 2q_{-1}}|))$ . In comparison, the state-of-the art Zhang-Shasha tree edit distance algorithm [31] runs in  $O(\rho * m^2 n^2)$ . Our large-scale evaluation shows that this new program embeddings using position-aware characteristic vectors and the new distance measure not only leads to significantly faster search than tree edit distance on ASTs but also negligible precision loss.

# 3.2 Align

The goal of the align step is to compute Discrepancies (Definition 3.5). The rationale is that after a syntactically similar program is identified, it must be correctly aligned with the incorrect program such that the differences between the aligned statements can suggest potential corrections.

**Definition 3.5.** (Discrepancies) Given the incorrect program  $P_e$  and a correct program  $P_c$ , discrepancies, denoted by  $C(P_e, P_c)$ , is a list of pairs,  $(S_e, S_c)$ , where  $S_e/S_c$  is a non-control statement<sup>2</sup> in  $P_e/P_c$ .

Aligning a reference solution with the incorrect program is a crucial step in generating accurate fixes, and in turn feedback. Figure 7 shows an example, in which the challenges that need to be addressed are: (1) renaming s1 in the correct solution to s — failing to do so will result in an incorrect let alone precise fix; (2) computing the correct alignment which leads to the minimum fix of changing *char* s = O' on line 3 in Figure 7a to *char* s = X'; and (3) solving the previous two tasks in a timely manner to ensure a good user experience. Our key idea for solving these challenges is to reduce the alignment problem to a distance computation problem,

<sup>&</sup>lt;sup>2</sup>Hereinafter we denote non-control statement to include loop headers, branch conditions, *etc.* 

specifically, of finding an alignment of two programs that minimizes their syntactic distances. We realize this idea in two-steps: variable-usage based  $\alpha$ -conversion and two-level statement matching.



Figure 7. Highlighting the usage of variable s and s1 foraligning the two programs.

*Variable-usage based*  $\alpha$ *-Conversion* The dynamic app-roach of comparing the runtime traces for each variable suffers from the same scalability issues mentioned in the search procedure. Instead, we present a syntactic approach -*variable-usage* based  $\alpha$ -conversion. Our intuition is that how a variable is being used in a program serves as a good indi-cator of its identity. To this end, we represent each variable by profiling its usage in the program. 

**Definition 3.6.** (Usage Set) Given a program *P* and the variable set Vars(P), a usage set of a variable  $v \in Vars(P)$  consists of all the non-control statements featuring v in *P*.

We then collect the usage set for each variable in  $P_e/P_c$ to form  $\mathcal{U}(P_e)/\mathcal{U}(P_c)$ . Now the goal is to find a one-to-one mapping between  $Vars(P_e)$  and  $Vars(P_c)$  which minimizes the distance between  $\mathcal{U}(P_e)$  and  $\mathcal{U}(P_c)$ . Note that if the number of variables between the two programs  $Vars(P_e)$ and  $Vars(P_c)$  are different before alignment, new (existing) variables will be added (deleted) during the alignment step.

$$\alpha\text{-conversion} = \underset{Vars(P_c) \leftrightarrow Vars(P_c)}{\arg \min} \Delta(\mathcal{U}(P_c), \mathcal{U}(P_e))$$
(3)

We can now compute the tree-edit distance between statements in any two usage sets in  $\mathcal{U}(P_e)$  and  $\mathcal{U}(P_c)$ , and then find the matching that adds up to the smallest distance between  $\mathcal{U}(P_e)$  and  $\mathcal{U}(P_c)$ . However, the total number of usage sets in  $\mathcal{U}(P_e)$  and  $\mathcal{U}(P_c)$  (denoted by the level of usage set) and the number of statements in each usage set (denoted  by the level of statement) will lead this approach to a combinatorial explosion that does not scale in practice. Instead, we rely on the new program embeddings to represent each usage set with only one single position-aware characteristic vector. Using  $H_{v_i}/H'_{v_i}$  to denote the vector for usage set of  $v_i \in Vars(P_e)/v'_i \in Vars(P_c)$ , we instantiate Equation 3 into

$$\alpha\text{-conversion} = \underset{v_i \leftrightarrow v'_i}{\operatorname{arg\,min}} \sum_{i=1}^{\omega} ||H_{v_i} - H'_{v_i}||_2 \qquad (4)$$

where 
$$\omega = min(|Vars(P_c)|, |Vars(P_e)|)$$
 (5)

The benefits of this instantiation are: (1) it only focuses on the combination at the level of usage set, and therefore eliminates a vast majority of combinations at the level of statement; and (2) it can quickly compute the Euclidean distance between two usage sets. Note the computed mapping between  $Vars(P_c)$  and  $Vars(P_e)$  in Equation 4 does not necessarily lead to the correct  $\alpha$ -conversion. To deal with the phenomenon that variables may need to be left unmatched if their usages are too dissimilar, we apply the Tukey's method [29] to remove statistical outliers based on the Euclidean distance. After the  $\alpha$ -conversion, we turn  $P_c$  into  $P_{\alpha c}$ .

**Two-Level Statement Matching** We leverage program structure to perform statement matching at two levels. At the first level, we fragment  $P_e$  and  $P_{\alpha c}$  into a collection of basic blocks according to their control-flow structure, and then align each pair in order. This alignment will result in a 1-1 mapping of the basic blocks due to  $CF(P_e) \equiv CF(P_{\alpha c})$ . Next, we match statements/expressions only within aligned basic blocks. In particular, we pick the matching that minimizes the total syntactic distance (tree edit distance) among all pairs of matched statements within each aligned basic blocks. Figure 8 depicts the result of aligning programs in Figure 7.

#### 3.3 Repair

Given  $C(P_e, P_{\alpha c})$ , we generate  $\mathcal{F}(P_e, P_{\alpha c})$  as follows:

- **Insertion Fix**: Given a pair  $(S_e, S_c)$ , if  $S_e = \emptyset$  meaning  $S_c$  is not aligned to any statement in  $P_e$ , an insertion operation will be produced to add  $S_c$ .
- **Deletion Fix**: Given a pair  $(S_e, S_c)$ , if  $S_c = \emptyset$  meaning  $S_e$  is not aligned to any statement in  $P_c$ , a deletion operation will be produced to delete  $S_e$ .
- **Modification Fix:** Given a pair  $(S_e, S_c)$ , if  $S_c \neq \emptyset$  and  $S_e \neq \emptyset$ , a modification operation will be produced consisting the set of standard tree operations to turn the AST of  $S_e$  to that of  $S_c$ .

The computed set of potential fixes  $\mathcal{F}(P_e, P_{\alpha c})$  typically contains a large set of redundancies that are unrelated to the root cause of the error in  $P_e$ . Such redundancies can be classified into two categories:

Variable Initialization		1		
	— int length = 8 bool ch = true	←→	Variable Initialization	bool ch = true
Outer Loop Initialization	inti=0		Outer Loop Initialization	int i = 0
	inter-0	]	outer Loop initialization	1111-0
Outer Loop Condition	i < length	]⊷→	Outer Loop Condition	i < 8
Outer Loop Incremento	r <i>i++</i>	]>	Outer Loop Incrementor	i++
Inner Loop Initialization	int i = 0	<b> </b> ←→→	Inner Loop Initialization	int i = 0
		]		<i>incj</i> = 0
Inner Loop Condition	j < length	<b> </b> ←→	Inner Loop Condition	j < 8
Inner Loop Incrementor	· j++	]↔	Inner Loop Incrementor	j++
If Condition	ch	}	If Condition	ch
If Body	Console.Write("X")	]↔	If Body	Console.Write("X",
Else Body	Console.Write("O")	•>	Else Body	Console.Write("O"
Inner Loop Body	ch = !ch	<b> </b> ←→	Inner Loop Body	ch = !ch

Figure 8. Aligning basic blocks for programs in Figure 2b and 3b after  $\alpha$ -conversion (aligned blocks are connected by arrows); matching statements within each pair of aligned basic blocks (italicized statements denotes matching; statements annotated by +/- denotes insertion/deletion.).

- Syntactic Differences:  $P_e$  and  $P_{\alpha c}$  may use different expressions in method calls, parameters, constants, etc. Figure 9a highlights such syntactic differences.
- Semantic Differences: More importantly,  $\mathcal{F}(P_e, P_{\alpha c})$ may also contain semantically-equivalent differences such as different conditional statements vet logically equivalent if statement; different initialization and conditional expressions that evaluate to the same number of loop iterations; or expressions computing similar values composed of different variables (not variable names). Figure 9b highlights such semantic differences.

The Repair component is responsible for filtering out all 703 those benign differences that are mixed in  $\mathcal{F}(P_e, P_{\alpha c})$ , in 704 other words, picking  $\mathcal{F}_m(P_e, P_{\alpha c})$   $(\mathcal{F}_m(P_e, P_{\alpha c}) \subseteq \mathcal{F}(P_e, P_{\alpha c}))$ 705 that only fix the issues in  $P_e$  while leaving the correct parts 706 unchanged (Definition 3.7). A direct brute-force approach is 707 to exhaustively incorporate each subset of  $\mathcal{F}(P_e, P_{\alpha c})$  into the 708 original incorrect program, and test the resultant program 709 by dynamic execution. This approach yields  $2^{|\mathcal{F}(P_e, P_{\alpha c})|} - 1$ 710 number of trials in total. As the number of operations in 711  $\mathcal{F}(P_e, P_{\alpha c})$  increases, the search space become intractable. 712 Instead, we apply several optimization techniques to make 713 the procedure more efficient and scalable (cf. Section 4). 714

**Definition 3.7.** (Minimality) Given an incorrect program *P* and a set of possible changes *U*, a set of changes  $F \subset U$ to correct *P* is defined to be minimal if there does not exist F' s.t. |F'| < |F| and F' fixes P. Correctly fixing P is w.r.t. a given test suite, *i.e.* the fixed program should pass all tests.

716

717

718

719

720 721

722

723

724

725

726

727

728

729 730

767

768

769

770

#### 3.4 The Repair Algorithm in SARFGEN

We now present the complete repair algorithm in SARFGEN. The system first matches the control flow (Line 3-6); then selects top k candidates for repair (Line 7-12); Line 16-19 denotes the search, align and repair procedure, whereas Line 23 generates the feedback.

Al	gorithm 1: SARFGEN's feedback generation.	
/	* $P_e$ : an incorrect program; $P_s$ : all correct	
	solutions; $\ell$ : feedback level	*/
f	unction FeedbackGeneration( $P_e, P_s, \ell$ )	
	begin	
	/* identify $P_{cs}$ that have the same control	
	flow of $P_e$ among all solutions $P_s$	*/
	$P_{cs} \leftarrow \emptyset$	
	for $P \in P_s$ do	
	if $CF(P) = CF(P_e)$ then	
	$ P_{cs} \leftarrow \{P\} \cup P_{cs} $	
	/* collect $P_{dis}$ that consist of k most	
	similar programs with $P_e$ .	*/
	$P_{dis} \leftarrow \emptyset$	
	for $P \in P_{cs}$ do	
	$  \mathbf{if}   P_{dis}   < k \parallel$	
	$\mathcal{D}(P, P_e) < \mathcal{D}(P_{kth} \in P_{dis}, P_e)$ then	
	$P_{dis} \leftarrow P_{dis} \setminus \{P_{kth}\}$	
	$P_{dis} \leftarrow \{P\} \cup P_{dis}$	
	$n \leftarrow \infty$	
	$\mathcal{F}_m(P_e) \leftarrow \text{null}$ // minimum set of fix for	Pe
	for $P_c \in P_{dis}$ do	c
	$P_{\alpha c} \leftarrow \alpha$ -Conversion $(P_{\alpha}, P_{c})$	
	$C(P_a, P_{ac}) \leftarrow \text{Discrepencies}(P_a, P_{ac})$	
	$\mathcal{F}(P_{\alpha}, P_{\alpha}) \leftarrow \text{Fives}(C(P_{\alpha}, P_{\alpha}))$	
	$\mathcal{F}_{\mathrm{rr}}(P_{\alpha}, P_{\alpha \alpha}) \leftarrow \operatorname{Minimization}(\mathcal{F}(P, P))$	
	$  \mathbf{f}   \mathcal{F} (P P)   < n \text{ than}$	
	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	
	$\int m(re) \leftarrow \int m(re, r\alpha c)$	
	<pre>/* translate fixes into feedback message.</pre>	*/
	$f = \text{Translating}(\mathcal{F}_m(P_e), \ell)$	

715

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

7

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880



(a) Syntactic differences between programs in Figures 2a and 3a. (b) Semantic differences between programs in Figures 2c and 3c.

Figure 9. Syntactic and semantic differences.

#### Implementation and Evaluation 4

We briefly describe some of the implementation details of SARFGEN, and then report the experimental results on benchmark problems. We also conduct an in-depth analysis for measuring the usefulness of the different techniques and framework parameters, and conclude with an empirical comparison with the CLARA tool [11].

### 4.1 Implementation

782

783

784

785

786

787

788

789

790

791

792

793

809

810

794 SARFGEN is implemented in C#. We use the Microsoft Roslyn 795 compiler framework for parsing ASTs and dynamic execu-796 tion. We keep 5 syntactically most similar programs as the 797 reference solutions. For the printing problem from Microsoft-798 DEV204.1X, we convert the console operation to string op-799 eration using the StringBuilder class. As for all the program 800 embeddings, we use 1-level characteristic vectors only due 801 to the suitable dimensionality. In other words, levels greater 802 than one will yield characteristic vectors of excessively high 803 dimensions (i.e. over millions for any realistic programming 804 language such as C/C++/C#) that do not provide good effi-805 ciency/precision tradeoffs. All experiments are conducted 806 on a Dell XPS 8500 with a 3rd Generation Intel Core®<sup>TM</sup> 807 i7-3770 processor and 16GB RAM. 808

#### 4.2 Results

We evaluate SARFGEN on the chessboard printing problem 811 from Microsoft-DEV204.1X as well as 16 out of the 24 prob-812 lems (the other seven are rarely attempted by students) on the 813 CodeHunt education platform (ignoring the submissions that 814 are syntactically incorrect).<sup>3</sup> Table 2 shows the results. Over-815 all, SARFGEN generates feedback based on minimum fixes for 816 4,311 submissions out of 4,806 incorrect programs in total ( $\approx$ 817 818 90%) within two seconds on average. Our evaluation results 819 validate the assumption we made since all the minimal fixes 820 modify up to three lines only. Another interesting finding is 821 that SARFGEN performed better on CodeHunt programming 822

an extreme, we find students writing straight-line programs even in several different ways, and therefore those programs are more diverse and difficult for SARFGEN to precisely find the closest solutions. On the other hand, CodeHunt programs are functional and more structured. Even though there are fewer reference implementations, SARFGEN is capable of repairing more incorrect submissions. 4.3 In-depth Analysis We now present a few in-depth experiments to better under-

submissions than on Microsoft-DEV204.1x's assignment de-

spite the fewer number of correct programs. After further

investigation, we conclude that the most likely cause for

this is the printing nature of Microsoft-DEV204.1x's exer-

cise placing little constraint on the control-flow structure. In

stand the contributions of different techniques for the search, align, and repair phases. First, we investigate the effect of different program embeddings adopted in the search procedure. Then, we investigate the usefulness of our  $\alpha$ -conversion mechanism compared against other variable alignment approaches. For these experiments, we use two key metrics: (i) performance (i.e. how long does SARFGEN take to generate feedback), and (2) capability (i.e. how many incorrect programs are repaired with minimum fixes). To keep our experiments feasible, we focus on three problems that have the most reference solutions (i.e. Printing, MaximumDifference and FibonacciSum).

**Program Embeddings** As the number k of top-k closest programs used for the reference solutions increases, we compare the precision of our program embeddings using position-aware characteristic vector against (1) program embeddings using the characteristic vector and (2) using the original AST representation. We adopt a cut-off point of  $\mathcal{F}_m(P_e, P_{\alpha c})$  to be three. Otherwise, if we let SARFGEN traverse the power set of  $\mathcal{F}(P_e, P_{\alpha c})$ , the capability criterion will be redundant. In all settings, SARFGEN adopts the usage-set based  $\alpha$ -conversion via position-aware characteristic vectors. Also, SARFGEN runs without any optimization techniques. Figure 10 shows the results. The embeddings using

<sup>823</sup> <sup>3</sup>Please refer to https://kbwang.bitbucket.io/misc/pldi18-paper93-supple-824 mental-text.pdf for a brief description of each benchmark problem.

### Data-Driven Feedback Generation for Introductory Programming Exercises PL'17, January 01-03, 2017, New York, NY, USA



Figure 11. Compare different  $\alpha$ -conversion techniques based on the same set of reference solutions.

position-aware characteristic vectors are almost as precise as ASTs and achieves exactly the same accuracy as ASTs when the number of closest solutions equals five. In addition, the position-aware characteristic vector embedding consistently outperforms the characteristic vector embedding by more than 10% in terms of capability. Moreover, although

927

928

929

930

931

932

933

934

935

position-aware characteristic vector uses higher dimensions to embed ASTs, it is generally faster due to the higher accuracy in selected reference solutions, and in turn lowering the computational cost spent on reducing  $\mathcal{F}(P_e, P_{\alpha c})$  to  $\mathcal{F}_m(P_e, P_{\alpha c})$ . When SARFGEN uses fewer reference solutions, the two embeddings do not display noticeable differences 982

983

984

985

986

987

988

989

990

since the speedup offered by the later is insignificant. Both
embedding schemes are significantly faster than the AST
representation.

 $\alpha$ -Conversion We next evaluate the precision and efficiency of the  $\alpha$ -conversion in the align step. Given the same set of candidate solutions (identified by AST representation in the previous step), we compare our usage-based  $\alpha$ -conversion via position-aware characteristic vector against (1) the same usage-based  $\alpha$ -conversion via standard tree-edit operations and (2) a dynamic trace-based  $\alpha$ -conversion via sequence alignment. We adopt the same configurations as in the previ-ous experiment which is to set the cut-off point of  $\mathcal{F}_m(P_e, P_{\alpha c})$ to be three and perform minimization without optimizations. As shown in Figure 11, our usage-set based  $\alpha$ -conversion via position-aware characteristic vectors outperforms that via standard tree-edit operations by more than an order of magnitude at the expense of little precision loss measured by capability. In the meanwhile, dynamic variable-trace based sequence alignment displays the best performance at a con-siderable cost of capability mainly due to its poor tolerance against erroneous variable traces. 

 $\mathcal{P}(P_e, P_{\alpha c})$ . To tackle this challenge, we design and realize 

• *Reachability-based Pruning:* We leverage the reachability of a fix to examine its applicability before it is attempted.<sup>4</sup> In particular, if a fix location is never reached on any execution path *w.r.t.* any of the provided inputs, it is safe to exclude it from the minimization procedure.

• Co-ordinated Changes: Certain corrections are coor-dinated, i.e. they should only be included/excluded simultaneously due to their co-occurrence nature. For example, when operations in  $\mathcal{F}(P_e, P_{\alpha c})$  introduces a new variable ( $v \in Vars(P_{\alpha c}) \land v \notin Vars(P_e)$ ) or delete existing variable  $(v \notin Vars(P_{\alpha c}) \land v \in Vars(P_e))$ , we bundle the operations that deletes/inserts an ex-pression composed of Var together to be co-existing members that can either all be included or excluded in  $\mathcal{F}_m(P_e, P_{\alpha c})$ . The intuition is if a new/existing variable is inserted/deleted, it's illegitimate to separate the cas-cading expressions that are not sensible on its own, *i.e.* expressions composed of a deleted variable will cause compilation to fail. 

• *Reusing Dynamic Executions:* The exhaustive dynamic executions on the power set of  $\mathcal{F}(P_e, P_{\alpha c})$  can be used

to recycle the computations to detect the syntactic/semantic redundancy. For example, while exhausting all subsets of one operation, even if we did not discover  $\mathcal{F}_m(P_e, P_{\alpha c})$ , we can detect the correction sets that are functionally redundant (*i.e.* do not change the semantics of incorrect program), and consequently remove them from the future computations. The more the number of iterations required to find  $\mathcal{F}_m(P_e, P_{\alpha c})$ , the larger the set of redundancies that can be discovered and eliminated resulting in a mutual beneficial relationship.

According to our evaluation, these optimizations are able to gain approximately one order of magnitude speedup.

*Minimization Effectiveness* In Table 3, we show the effectiveness of SARFGEN's repair component by comparing the number of fixes pre- and post-minimization procedure.

Programming Problems	Number of Fixes Before Minimization	Number of Fixes After Minimization
Divisibility	2.1	1.6
ArrayIndexing	1.8	1.2
StringCount	3.5	1.9
Average	3.8	2.1
ParenthesisDepth	8.3	2.8
Reversal	6.5	2.3
LCM	8.2	3.1
MaximumDifference	5.5	2.3
BinaryDigits	6.4	2.1
Filter	7.9	3.5
FibonacciSum	7.6	3.1
K-thLargest	6.8	2.6
SetDifference	10.1	3.6
Combinations	9.7	3.6
MostOnes	10.8	4.2
ArrayMapping	9.5	3.2
Printing	12.7	3.5

**Table 3.** Evaluating SARFGEN's repair component.

#### 4.4 Reliance on Data

We conducted a further experiment to understand the degree to which SARFGEN relies on the correct programming submissions to have a reasonable utility. Initially, we use all the correct programs from all the programming problems, then we gradually down-sample them to observe the effects this may have on SARFGEN's capability and performance. Figure 12a/12b depicts the capability/performance change as the number of correct solutions is reduced from 100% to 1% under the standard configuration. In terms of capability, SARFGEN maintains almost the same power as the number of correct programs drops to half of the total. Even using only 1% of the total correct programs, SARFGEN still manages to produce feedback based on minimal fixes for almost 60% of the incorrect programs in total. The reason for this phenomenon is that the vast majority of students generally

 <sup>&</sup>lt;sup>4</sup>Assuming test cases provided by the instructors are thorough and cover
 all corner cases.

adopt an algorithm that their peers have correctly adopted. 1101 1102 So even though the correct programs are down-sampled, 1103 there still exist some solutions of common patterns that can help a large portion of students who also attempt to solve 1104 1105 the problem in a common way. Consequently, those students will not be affected. On the other hand, SARFGEN will un-1106 1107 derstandably have more difficulties to deal with corner-case 1108 programming submissions. However, because the number of 1109 such programs is small, it generally does not have a severe impact. As for performance, the changes are twofold. On 1110 1111 one hand, with fewer correct solutions, SARFGEN performs less computation. On the other hand, SARFGEN generally 1112 spends more time reducing  $\mathcal{F}(P_e, P_{\alpha c})$  to  $\mathcal{F}_m(P_e, P_{\alpha c})$  since 1113 the reference solutions become more dissimilar due to down-1114 1115 sampling. Since SARFGEN is able to find precise reference 1116 solutions for most of the incorrect programs when the num-1117 ber of reference solutions is not too low, *i.e.* above 1%, the 1118 final outcome is that SARFGEN becomes slightly faster.



**Figure 12.** Capability and performance changes as the number of correct solutions decrease.



Figure 13. Capability and performance changes as the sizeof program grows.

### 1144 4.5 Scalability with Program Size

1129

1130

1142

1143

1155

To understand how our system scales as the size of program 1145 increases, we conducted another experiment in which we 1146 1147 partitioned the incorrect programs into ten groups according to size of the program, *i.e.* the number of lines in code. 1148 1149 In terms of performance (Figure 13a), SARFGEN manages to generate the repairs within ten seconds in almost all cases. 1150 1151 Regarding capability (Figure 13b), SARFGEN generates minimal fixes for close to 90% of incorrect programs when their 1152 1153 size is under 20 lines of code and still over 50% when the size is less than 40 lines of code. 1154

#### 4.6 Comparison against CLARA

In this experiment, we conduct an empirical comparison against CLARA [11] using the same benchmark set. Since CLARA works on C programs, we followed the following procedure. First, we convert our C# programs into C programs that CLARA supports. In fact, we only converted the *Console* operation into *printf* for the Printing problem, as a result we have 488 out of 742 programs from edX that still compile, but only 395 out of 4,311 programs from CodeHunt since they generally contain more complex data structures. In total, we have 883 programs as the benchmark set for both systems to compare. Both systems use exactly the same set of correct programs in different languages. Second, because we experienced issues when invoking the provided clustering API<sup>5</sup> to cluster correct programs, we instead run CLARA on each correct program separately to repair the incorrect programs and select the fixes that are minimal and fastest (prioritize minimality over performance when necessary). On the other hand, SARFGEN is set up with standard configuration following Algorithm 1.

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

The results are shown in Table 4. As the solution set is down-sampled from 100% to 1%, SARFGEN generates consistently more minimal fixes than CLARA (Table 4a). For the results on performance shown in Table 4b, CLARA outperforms SARFGEN marginally on the programs of small size, i.e., fewer than 15 lines of code, whereas CLARA scales significantly worse than SARFGEN when the size of programs grows, *i.e.*, slower by more than one order of magnitude when dealing with program of more than 25 lines. Regarding the performance comparison, in reality, CLARA usually compares an incorrect program with hundreds of reference solutions according to [11] to pick the smaller fixes, therefore the performance measured is a significant under-estimation. Furthermore, CLARA shows better performance in part due to the less work it undertakes since it does not guarantee minimal repairs.

#### 4.7 User Study

The latest version of SARFGEN has been deployed onto the Microsoft-DEV204.1x course website on edX. Here we report the feedback users have submitted.

The goal of this study is to measure the usefulness of SARFGEN after being deployed to integrate with the C# edX course<sup>6</sup>. We study two research questions: 1) Can SARFGEN help the learning efficiency of the students? and 2) Do students consider the feedback generated by SARFGEN to be helpful? To answer the first question we randomly selected 200 hundred students based on the edX Id and evenly separated them into two groups (*i.e.* 100 students per group): pre-deployment (Group A) and post-deployment (Group B).

<sup>&</sup>lt;sup>5</sup>The match command provided in the Example section at https://github.com/iradicek/clara produces unsound result <sup>6</sup>http://mslexcodegrader.azurewebsites.net/

1228

1229

1230

1231

1245

1251

1252

1265

Percentage of	Number of Minimum	Number of Minimur
Solutions	Fixes (Sarfgen)	Fixes (CLARA)
100%	688	573
80%	661	540
60%	647	492
40%	603	426
20%	539	349
1%	466	271
	(a) Comparison on capa	ability.
LOC (Total # of	(a) Comparison on cap	ability. Average Time Take
LOC (Total # of Programs)	(a) Comparison on cap Average Time Taken (SARFGEN)	ability. Average Time Take (CLARA)
LOC (Total # of Programs) 0-5 (122)	(a) Comparison on cap Average Time Taken (SARFGEN) 1.14s	ability. Average Time Take (CLARA) 0.84s
LOC (Total # of Programs) 0-5 (122) 5-10 (227)	(a) Comparison on capa Average Time Taken (SARFGEN) 1.14s 1.38s	ability. Average Time Take (CLARA) 0.84s 1.71s
LOC (Total # of Programs) 0-5 (122) 5-10 (227) 10-15 (269)	(a) Comparison on capa Average Time Taken (SARFGEN) 1.14s 1.38s 2.69s	Average Time Take (CLARA) 0.84s 1.71s 2.25s
LOC (Total # of Programs) 0-5 (122) 5-10 (227) 10-15 (269) 15-20 (137)	(a) Comparison on capa Average Time Taken (SARFGEN) 1.14s 1.38s 2.69s 3.02s	Average Time Take (CLARA) 0.84s 1.71s 2.25s 9.71s
LOC (Total # of Programs) 0-5 (122) 5-10 (227) 10-15 (269) 15-20 (137) 20-25 (91)	(a) Comparison on capa Average Time Taken (SARFGEN) 1.14s 1.38s 2.69s 3.02s 3.81s	Average Time Take (CLARA) 0.84s 1.71s 2.25s 9.71s 20.84s

Table 4. SARFGEN vs. CLARA on the same dataset.

Then we compare the interaction history across the two 1232 groups of students using two metrics: i) the number of sub-1233 1234 mission iterations and ii) the length of time until they sub-1235 mitted a correct program. Results show that 96% of students 1236 in Group B (vs. 37% of students in Group A) completed their 1237 assignment within next two attempts. In addition 97% of post-deployment students (vs. 40% of pre-deployment stu-1238 dents) finished within 30 minutes. This preliminary evidence 1239 supports that SARFGEN provides useful feedback. We also 1240 1241 added a couple of qualitative metrics in terms of a student 1242 rating and comments about the provided feedback. The latest overall user rating on a 1-5 scale is 3.74, showing overall 1243 1244 positive feedback.

#### **Related Work** 1246 5

1247 This section describes several strands of related work from 1248 the areas of automated feedback generation, automated pro-1249 gram repair, fault localization and automated debugging. 1250

## 5.1 Automated Feedback Generation

Recent years have seen the emergence of automated feedback 1253 generation for programming assignments as a new, active 1254 research topic. We briefly review the recent techniques. 1255

AutoGrader [24] proposed a program synthesis based au-1256 tomated feedback generation for programming exercises. 1257 The idea is to take a reference solution and an error model 1258 consisting of potential corrections to errors student might 1259 make, and search for the minimum number of corrections 1260 using a SAT-based program synthesis technique. In con-1261 trast, SARFGEN advances the technology in the following 1262 aspects: (1) SARFGEN completely eliminates the manual ef-1263 fort involved in the process of constructing the error model; 1264

and (2) SARFGEN can perform more complex program repairs such as adding, deleting, swapping statements, etc.

CLARA [11] is arguably the most similar work to ours. Their approach is to cluster the correct programs and select a canonical program from each cluster to form the reference solution set. Given an incorrect student solution, CLARA runs a trace-based repair procedure w.r.t. each program in the solution set, and then selects the fix consisting of the minimum changes. Despite the seeming similarity, SARFGEN is fundamentally different from CLARA. At a conceptual level, CLARA assumes for every incorrect student program, there is a correct program whose execution traces/internal states only differs because of the presence of the error. Even though the program repairer generally enjoys the luxury of abundant data in this setting, there are a considerable amount of incorrect programs which yield new (partially) correct execution traces. Since the trace-based repair procedure does not distinguish a benign difference from a fix, it will introduce semantic redundancies which likely will have a negative impact on student's learning experience. As we have presented in our evaluation, CLARA scales poorly with increasing program size, and does not generate minimal repairs on the benchmark programs.

*sk\_p* [22] was recently proposed to use deep learning techniques for program repair. Inspired by the skipgram model, a popular model used in natural language processing [20, 21], sk p treats a program as a collection of code fragments, consisting of a pair of statements with a hole in the middle, and learns to generate the statement based on the local context. Replacing the original statement with the generated statement, one can infer the generated statement contains the fix if the resulting program is correct. However, sk\_p suffers from low capability results, as the system only perform syntactic analysis. Another issue with the deep learning based approaches is low reusability. Significant efforts are needed to retrain new models to be applied across new problems.

QLOSE [3] is another recent work for automatically repairing incorrect solutions to programming assignments. The major contribution of this work is the idea of measuring the program distance not only syntactically but also semantically, *i.e.*, preserving program behavior regardless of syntactic changes. One way to achieve this is by monitoring runtime execution. However, the repair changes to an incorrect program is based on a pre-defined template corresponding to a linear combination of constants and all program variables in scope at the program location. As we have shown, more complex modifications are necessary for real-world benchmarks.

**REFAZER** [23] is another approach applicable in the domain of repairing program assignments. The idea is to learn a syntactic transformation pattern from examples of statement/expression instances before and after the modification.

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

Despite the impressive results, this approach also suffers
from similar issues as QLOSE, *i.e.*, there are many incorrect
programs that require changes more complex than simple
syntactic changes.

1325 *CoderAssist* [16] presents a new methodology for gener-1326 ating verified feedback for student programming exercises. 1327 The approach is to first cluster the student submissions ac-1328 cording to their solution strategies and ask the instructor to 1329 identify a correct submission in each cluster (or add one if 1330 none exists). In the next phase, each submission in a clus-1331 ter is verified against the instructor-validated submission 1332 in the same cluster. Despite the benefit of generating veri-1333 fied feedback, there are several weaknesses. As mentioned, 1334 CoderAssist requires manual effort from the instructor. More 1335 importantly, the quality of the generated feedback relies on 1336 how similar the provided solution is to the incorrect sub-1337 missions in the same cluster. In contrast, SARFGEN searches 1338 through all possible solutions automatically and uses those 1339 that it considers to be the most similar to repair the incorrect 1340 program. In addition, CoderAssist targets dynamic program-1341 ming assignments only. Its utility and scalability would need 1342 require further validation on other problems.

1343 There is also work that addresses other learning aspects in 1344 the MOOC setting. For example, Gulwani et al. [10] proposed 1345 an approach to help students write more efficient algorithms 1346 to solve a problem. Its goal is to teach students about the 1347 performance aspects of a computing algorithm other than its 1348 functional correctness. However, this approach only works 1349 with correct student submissions, i.e., it cannot repair incor-1350 rect programs. Kim et al. [17] is another interesting piece of 1351 work focusing on explaining the root cause of a bug in stu-1352 dents' programs by comparing their execution traces. This 1353 approach works by first matching the assignment statement 1354 symbolically and then propagating to match predicates by 1355 aligning the control dependencies of the matched assign-1356 ment statements. The key difference is that our work can 1357 automatically repair the student's code while Kim et al. [17] 1358 can only illustrate the cause of a bug. 1359

#### 5.2 Automated Program Repair

1361 Gopinath et al. [8] propose a SAT-based approach to gener-1362 ate repairs for buggy programs. The idea is to encode the 1363 specification constraint on the buggy program into a SAT 1364 constraint, whose solutions lead to fixes. Könighofer and 1365 Bloem [18] present an approach based on automated error 1366 localization and correction. They localize faulty components 1367 with model-based diagnosis and then produce corrections 1368 based on SMT reasoning. They only take into account the 1369 right hand side (RHS) of the assignment statements as re-1370 placeable components. Prophet [19] learns a probabilistic, 1371

application-independent model of correct code by generalizing a set of successful human patches. There is also work [13, 26] that models the problem of program repair as a game. The two actors are the environment that provides the inputs and a system that provides correct values for the buggy expressions, so ultimately the specification is satisfied. These approaches use simple corrections (*e.g.*, correcting the RHS sides of expressions) since they aim to repair large programs with arbitrary errors. Another line of approaches use program mutation [4], or genetic programming [1, 6] for automated program repair. The idea is to repeatedly mutate statements ranked by their suspiciousness until the program is fixed. In comparison our approach is more efficient in pinpointing the error and fixes as those mutation-based approaches face extremely large search space of mutants (10<sup>12</sup>). 1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

# 5.3 Automated Debugging and Fault localization

Test cases reduction techniques like Delta Debugging [30] and QuickXplain [15] can complement our approach by ranking the likely fixes prior to dynamic analysis. The hope is to expedite the minimization loop and ultimately speed up performance. A major research direction of fault localization [2, 9] is to compare faulty and successful executions. Jose and Majumdar [14] propose an approach for error localization from a MAX-SAT aspect. However, such approaches suffer from their limited capability in producing fixes.

# 6 Conclusion

We have presented the "Search, Align, and Repair" datadriven framework for generating feedback on introductory programming assignments. It leverages the large number of available student solutions to generate instant, minimal, and semantic fixes to incorrect student submissions without any instructor effort. We introduce a new program representation mechanism using position-aware characteristic vectors that are able to capture rich structural properties of the program AST. These program embeddings allow for efficient algorithms for searching similar correct programs and aligning two programs to compute syntactic discrepancies, which are then used to compute a minimal set of fixes. We have implemented our approach in the SARFGEN system and extensively evaluated it on thousands of real student submissions. Our results show that SARFGEN is effective and improves existing systems w.r.t. automation, capability, and scalability. Since SARFGEN is also language-agnostic, we are actively instantiating the framework to support other languages such as Python. SARFGEN has also been integrated with the Microsoft-DEV204.1X edX course and the early feedback obtained from online students demonstrates its practicality and usefulness.

1372 1373

1374

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

1512

1513

1514

1515

1516

1517

1518

1519

1520

1521

1522

1523

1524

1525

1526

1527

1528

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

#### 1431 References

- [1] Andrea Arcuri. 2008. On the Automation of Fixing Software Bugs. In
   *Companion of the 30th International Conference on Software Engineering*.
   ACM, New York, NY, USA, 1003–1006.
- [2] Thomas Ball, Mayur Naik, and Sriram K. Rajamani. 2003. From Symptom to Cause: Localizing Errors in Counterexample Traces. In Proceedings of the 30th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. ACM, New York, NY, USA, 97–105.
- [3] Loris D'Antoni, Roopsha Samanta, and Rishabh Singh. 2016. Qlose:
   Program repair with quantitative objectives. In *International Conference on Computer Aided Verification*. Springer, 383–401.
- [4] Vidroha Debroy and W Eric Wong. 2010. Using mutation to automatically suggest fixes for faulty programs. In 2010 Third International Conference on Software Testing, Verification and Validation. IEEE, 65– 74.
- [5] Anna Drummond, Yanxin Lu, Swarat Chaudhuri, Christopher Jermaine, Joe Warren, and Scott Rixner. 2014. Learning to grade student programs in a massive open online course. In 2014 IEEE International Conference on Data Mining. IEEE, 785–790.
- [6] Stephanie Forrest, ThanhVu Nguyen, Westley Weimer, and Claire
  Le Goues. 2009. A Genetic Programming Approach to Automated
  Software Repair. In *Proceedings of the 11th Annual Conference on Genetic and Evolutionary Computation*. ACM, New York, NY, USA, 947–954.
- [7] Elena L Glassman, Jeremy Scott, Rishabh Singh, Philip J Guo, and Robert C Miller. 2015. OverCode: Visualizing variation in student solutions to programming problems at scale. ACM Transactions on Computer-Human Interaction (TOCHI) 22, 2 (2015), 7.
- [8] Divya Gopinath, Muhammad Zubair Malik, and Sarfraz Khurshid. 2011.
  Specification-based Program Repair Using SAT. In Proceedings of the 17th International Conference on Tools and Algorithms for the Construction and Analysis of Systems: Part of the Joint European Conferences on Theory and Practice of Software. Springer-Verlag, Berlin, Heidelberg, 173–188.
- [9] Alex Groce, Sagar Chaki, Daniel Kroening, and Ofer Strichman. 2006.
  Error explanation with distance metrics. Tools and Algorithms for the Construction and Analysis of Systems: 10th International Conference, TACAS 2004, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2004, Barcelona, Spain, March 29 - April 2, 2004. Proceedings 8, 3 (2006), 229–247.
- [10] Sumit Gulwani, Ivan Radiček, and Florian Zuleger. 2014. Feedback
   generation for performance problems in introductory programming
   assignments. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*. ACM, 41–51.
- [11] Sumit Gulwani, Ivan Radiček, and Florian Zuleger. 2016. Automated
  [1468 Clustering and Program Repair for Introductory Programming Assignments. arXiv preprint arXiv:1603.03165 (2016).
- [12] Lingxiao Jiang, Ghassan Misherghi, Zhendong Su, and Stephane Glondu. 2007. Deckard: Scalable and accurate tree-based detection of code clones. In *Proceedings of the 29th international conference on Software Engineering*. IEEE Computer Society, 96–105.
- 1473 [13] Barbara Jobstmann, Andreas Griesmayer, and Roderick Bloem. 2005.
  1474 Program Repair as a Game. In *Computer Aided Verification: 17th International Conference, CAV 2005, Edinburgh, Scotland, UK, July 6-10, 2005.* 1476 *Proceedings*, Kousha Etessami and Sriram K. Rajamani (Eds.). Springer 1476 Berlin Heidelberg, Berlin, Heidelberg, 226–238.
- [14] Manu Jose and Rupak Majumdar. 2011. Cause Clue Clauses: Error Localization Using Maximum Satisfiability. In Proceedings of the 32Nd ACM SIGPLAN Conference on Programming Language Design and Implementation. ACM, New York, NY, USA, 437–446.
- [15] Ulrich Junker. 2004. QUICKXPLAIN: Preferred Explanations and Relaxations for Over-constrained Problems. In *Proceedings of the 19th National Conference on Artifical Intelligence*. AAAI Press, 167–172.
- [16] Shalini Kaleeswaran, Anirudh Santhiar, Aditya Kanade, and Sumit
   Gulwani. 2016. Semi-supervised Verified Feedback Generation. In

Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering. ACM, New York, NY, USA, 739–750.

- [17] Dohyeong Kim, Yonghwi Kwon, Peng Liu, I. Luk Kim, David Mitchel Perry, Xiangyu Zhang, and Gustavo Rodriguez-Rivera. 2016. Apex: Automatic Programming Assignment Error Explanation. In Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications. ACM, New York, NY, USA, 311–327.
- [18] Robert Könighofer and Roderick Bloem. 2011. Automated error localization and correction for imperative programs. In *Proceedings of the International Conference on Formal Methods in Computer-Aided Design*. IEEE, 91–100.
- [19] Fan Long and Martin Rinard. 2016. Automatic Patch Generation by Learning Correct Code. In Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. ACM, New York, NY, USA, 298–312.
- [20] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In Advances in neural information processing systems. 3111–3119.
- [21] Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global Vectors for Word Representation.. In *EMNLP*, Vol. 14. 1532–1543.
- [22] Yewen Pu, Karthik Narasimhan, Armando Solar-Lezama, and Regina Barzilay. 2016. Sk\_P: A Neural Program Corrector for MOOCs. In Companion Proceedings of the 2016 ACM SIGPLAN International Conference on Systems, Programming, Languages and Applications: Software for Humanity. ACM, New York, NY, USA, 39–40.
- [23] Reudismam Rolim, Gustavo Soares, Loris D'Antoni, Oleksandr Polozov, Sumit Gulwani, Rohit Gheyi, Ryo Suzuki, and Björn Hartmann. 2017. Learning syntactic program transformations from examples. In Proceedings of the 39th International Conference on Software Engineering. IEEE Press, 404–415.
- [24] Rishabh Singh, Sumit Gulwani, and Armando Solar-Lezama. 2013. Automated Feedback Generation for Introductory Programming Assignments. In Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation. ACM, New York, NY, USA, 15–26.
- [25] Taylor Soper. 2014. Analysis: the exploding demand for computer science education, and why America needs to keep up. Geekwire. (2014). http://www.geekwire.com/2014/ analysis-examining-computer-science-education-explosion
- [26] Stefan Staber, Barbara Jobstmann, and Roderick Bloem. 2005. Finding and Fixing Faults. In Correct Hardware Design and Verification Methods: 13th IFIP WG 10.5 Advanced Research Working Conference, CHARME 2005, Saarbrücken, Germany, October 3-6, 2005. Proceedings, Dominique Borrione and Wolfgang Paul (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 35–49.
- [27] Kuo-Chung Tai. 1979. The tree-to-tree correction problem. J. ACM 26, 3 (1979), 422–433.
- [28] Nikolai Tillmann, Jonathan de Halleux, Tao Xie, and Judith Bishop. 2014. Code hunt: gamifying teaching and learning of computer science at scale. In *Learning@Scale*. 221–222.
- [29] John W Tukey. 1977. Exploratory data analysis. Addison-Wesley Series in Behavioral Science: Quantitative Methods, Reading, Mass.: Addison-Wesley, 1977 (1977).
- [30] Andreas Zeller and Ralf Hildebrandt. 2002. Simplifying and isolating failure-inducing input. *IEEE Transactions on Software Engineering* 28, 2 (2002), 183–200.
- [31] Kaizhong Zhang and Dennis Shasha. 1989. Simple fast algorithms for the editing distance between trees and related problems. *SIAM journal* on computing 18, 6 (1989), 1245–1262.

1485