

Context-Aware and Data-Driven Feedback Generation for Programming Assignments

Dowon Song
Korea University
Republic of Korea
dowon_song@korea.ac.kr

Woosuk Lee
Hanyang University
Republic of Korea
woosuk@hanyang.ac.kr

Hakjoo Oh*
Korea University
Republic of Korea
hakjoo_oh@korea.ac.kr

ABSTRACT

Recently, various techniques have been proposed to automatically provide personalized feedback on programming exercises. The cutting edge of which is the data-driven approaches that leverage a corpus of existing correct programs and repair incorrect submissions by using similar reference programs in the corpus. However, current data-driven techniques work under the strong assumption that the corpus contains a solution program that is close enough to the incorrect submission. In this paper, we present CAFE, a new data-driven approach for feedback generation that overcomes this limitation. Unlike existing approaches, CAFE uses a novel context-aware repair algorithm that can generate feedback even if the incorrect program differs significantly from the reference solutions. We implemented CAFE for OCaml and evaluated it with 4,211 real student programs. The results show that CAFE is able to repair 83% of incorrect submissions, far outperforming existing approaches.

CCS CONCEPTS

• **Software and its engineering** → **Automatic programming.**

KEYWORDS

Program Repair, Program Synthesis

ACM Reference Format:

Dowon Song, Woosuk Lee, and Hakjoo Oh. 2021. Context-Aware and Data-Driven Feedback Generation for Programming Assignments. In *Proceedings of the 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '21)*, August 23–28, 2021, Athens, Greece. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3468264.3468598>

1 INTRODUCTION

In recent years, there has been a surge of interest in automatic feedback generation for programming assignments [1, 5, 12, 19, 20, 22, 33, 35, 36, 39, 40, 43]. As the demand for programming education grows, it is becoming increasingly difficult for an instructor to provide personalized feedback to a large number of students. Simply

providing an instructor’s solution as feedback is unsatisfactory, as students’s attempts typically diverge from the reference solution. The goal of automatic feedback generation technology is to help students to understand what they did wrong and how to fix it without manual effort of instructors.

Data-Driven Feedback Generation. Among prior techniques, data-driven approaches [12, 17, 34, 42, 43] are arguably the current state-of-the-art. The main idea of these techniques is to leverage a corpus of existing correct programs, and repair an incorrect program by using similar reference solutions in the corpus. In contrast to approaches that require intervention of instructors [5, 19, 39], data-driven techniques are fully automatic and yet show impressive performance in repairing introductory programming exercises.

However, existing data-driven techniques have a significant shortcoming. That is, they rely on a strong assumption that the corpus contains a solution program that is *close* to the incorrect program. For example, two notable techniques, CLARA [12] and SARFGEN [43], assume a solution exists that is equivalent to the incorrect program modulo control flows. This assumption, however, does not hold always [22], especially when providing feedback beyond introductory-level exercises. In this case, constructing a corpus with the close-program assumption becomes a challenge.

Our Approach. In this paper, we present CAFE, a new data-driven feedback generation technique that overcomes the above limitation. Unlike existing approaches, CAFE can generate feedback even when the incorrect submission is substantially different from reference solutions.

The keystone of CAFE is its context-aware, function-level repair algorithm. CAFE primarily targets sizable programming exercises, where students are freely allowed to define and use their own helper functions. To repair such a program, CAFE does not seek to find a solution program that matches the submission in its entirety; instead, it leverages multiple, partially-matching references. More specifically, CAFE works at the function level, aiming to separately repair each function in the incorrect program by (1) finding a matching function from the corpus, (2) computing their difference, and (3) extracting a patch from the difference. A main challenge with this approach is how to find the matching function that is useful for repair. Our key idea to solve this problem is to infer and compare the *original intent* of the functions by analyzing their calling contexts in the respective programs, which robustly identifies useful references even when functions have different syntax and semantics.

We evaluated CAFE in a real classroom setting. The original motivation of this work was to develop a feedback generation system for our own programming course, where we use OCaml and newcomers to functional programming often have a hard time. Thus, we

*Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
ESEC/FSE '21, August 23–28, 2021, Athens, Greece
© 2021 Association for Computing Machinery.
ACM ISBN 978-1-4503-8562-6/21/08...\$15.00
<https://doi.org/10.1145/3468264.3468598>

implemented CAFE for OCaml and evaluated it with 664 incorrect and 3,547 correct student programs collected from the course over the past few years. In total, CAFE successfully repaired 83% (548/664) of incorrect submissions, vastly outperforming FixML [22], a recent feedback generation technique for OCaml, whose fix rate was 35% (234/664). We also confirmed that existing data-driven approaches are ineffective for our dataset; replacing our context-aware approach by the matching algorithm of SARFGEN [43] decreased the fix rate from 83% to 59%. Finally, we conducted a user study, which shows CAFE is actually helpful for students.

Contributions. We summarize our contributions below:

- We present CAFE, a new context-aware and data-driven feedback generation technique for programming assignments.
- We evaluate CAFE in a realistic setting and make the tool and benchmarks publicly available.¹

2 OVERVIEW

2.1 Motivating Example

Let us consider a programming exercise asking students to write a function `diff: aexp * string -> aexp`, which takes an arithmetic expression (`aexp`) and a variable name (`string`), and performs symbolic differentiation. Arithmetic expressions are defined in OCaml datatype as follows:

```
type aexp = Const of int | Var of string | Power of (string * int)
          | Sum of aexp list | Times of aexp list
```

An arithmetic expression is either constant integer (`Const`), variable (`Var`), exponentiation (`Power`), sum of arithmetic expressions (`Sum`), or product of arithmetic expressions (`Times`). The function `diff` should produce an expression that results from differentiating the given expression with respect to the given variable. For example, `diff (Sum [Power ("x", 2); Const 1], "x")`, i.e., differentiating $x^2 + 1$ w.r.t. x , outputs `Times [Const 2; Var "x"]` denoting $2 * x$.

Figure 1 shows an incorrect program written by a student in our class, which implements `diff` with three helper functions: `timediff`, `sumdiff`, and `differ`. The functions `timediff` and `sumdiff` are intended to compute the derivatives of the product and sum of `aexp` lists, respectively. The function `differ` performs actual differentiation using `timediff` and `sumdiff`, and handles other base cases. Note that the program erroneously handles the case of multiplication (`Times`). For example, `diff (Times [Var "x"; Var "y"], "x")` produces `Times [Const 1; Const 0]`, while the expected answer is `Var "y"`.

Despite its simple manifestation, fixing the bug correctly and providing right feedback is nontrivial even for instructors. For correct repair, we need to change three places. First, the student implemented `timediff` based on a wrong product rule, $(f \cdot g)' = f' \cdot g'$, and therefore the body of `timediff` needs to be rewritten based on the proper rule, i.e., $(f \cdot g)' = f' \cdot g + f \cdot g'$. Second, we need to rewrite `sumdiff` because it depends on the incorrect definition of `timediff`. Finally, the last line of `differ` (line 18) should be changed in accordance with the correct product rule.

Given the buggy program, test cases, and a corpus of 218 solution programs, CAFE repaired the program as shown in Figure 1 in 5 sec. To use the proper product rule, it replaced lines 4–6 of `timediff` by line 7 and line 18 of `differ` by line 19. Then, it modified the body of `sumdiff` to reflect the change. Note that the generated repair is not only correct but also instructive; indeed, it is identical to what we would manually provide to the student. For example, CAFE shows that the redundancy between `sumdiff` and `differ` can be effectively eliminated by making a recursive call to `diff`.

Compared to existing data-driven techniques [12, 34, 42, 43], the most distinguishing feature of CAFE is its ability to generate feedback by collectively using multiple, dissimilar reference solutions. Each of `timediff`, `sumdiff`, and `differ` in Figure 1 was fixed using different solutions; CAFE repaired `timediff` using `gettimes` in Figure 2(a), `sumdiff` using `diff_sum` in Figure 2(b), and `differ` using `diffh` in Figure 2(c). All of these reference programs are substantially different from the program in Figure 1, as none of the 218 solutions in our corpus had a matching control-flow structure, which implies that existing data-driven techniques [12, 34, 42, 43] would fail to generate the desired feedback.

2.2 How CAFE Works

Now we discuss the high-level ideas of our approach on the student attempts to apply: “Given a list l of integers and a target operation $o \in \{\text{ADD}, \text{SUB}\}$, write a function that increments (resp. decrements) each element of l by 1 if the operation is addition, i.e., $o = \text{ADD}$ (resp. subtraction, i.e., $o = \text{SUB}$)”. Figure 3a–3c show three student submissions of the programming assignment: P is incorrect, while P_1 and P_2 are functionally correct. We differently name the top-level functions that the three submissions implement to distinguish them (P : `apply`, P_1 : `apply1`, and P_2 : `apply2`).

Our goal is to automatically generate modifications that make the incorrect submission P correct as a guided feedback by using the existing correct student solutions P_1 and P_2 . The program P is functionally wrong in decrementing list elements. We can correct P by (i) modifying the expression `h+1` at line 7 to be `h-1` and (ii) changing `(dec_all tl)` at line 13 to `(hd-1)::(dec_all tl)`.

Context-Aware Matching. The first step is to find a matching relation between functions in the student submissions. We say that two functions f and g match, written $f \sim g$ if the functions are invoked and call other functions under *compatible contexts*. Here, we mean contexts by conditions over execution paths when function calls occur. We say two path conditions *compatible* if there exists an input that exercises both of the two execution paths.

Based on this notion, CAFE finds the following matching relation: `inc_all` \sim `add_list`, `dec_all` \sim `sub_list`, and `apply` \sim `apply1`. We consider two types of function contexts: (a) incoming contexts describing under what path conditions the function is invoked by other functions, and (b) outgoing contexts describing under what conditions the function invokes other functions. For example, the incoming contexts of `dec_all`, `add_list`, and `sub_list` (denoted $\delta_{\text{dec_all}}$, $\delta_{\text{add_list}}$, and $\delta_{\text{sub_list}}$, respectively) are as follows:

$$\begin{aligned} \delta_{\text{dec_all}} &\equiv i = (l, o) \wedge l = (\text{hd} :: \text{tl}) \wedge (o = \text{SUB}) \\ \delta_{\text{add_list}} &\equiv i = (l_1, o_1) \wedge l_1 = (\text{hd}_1 :: \text{tl}_1) \wedge (o_1 = \text{ADD}) \\ \delta_{\text{sub_list}} &\equiv i = (l_2, o_2) \wedge l_2 = (\text{hd}_2 :: \text{tl}_2) \wedge (o_2 = \text{SUB}) \end{aligned}$$

¹<https://github.com/kupl/LearnML>

```

1 let rec diff (e,x) =
2   let rec timediff (t1st,x) =
3     match t1st with [] -> [] | hd::t1 ->
4     (-) (match hd with Const c -> (Const c)::timediff(t1,x)
5     (-) | Var v -> if v=x then (Const 1)::timediff(t1,x) else (Const 0)::timediff(t1,x)
6     (-) | Power (v,c) -> if v=x then (Times[Const c;Power(v,c-1)]::timediff(t1,x) else (Const 0)::timediff(t1,x))
7     (+) if t1=[] then [diff (hd, x)] else [Times ([diff (hd, x)]@t1)] @ [Times ([hd]@[Sum (timediff (t1, x))])]
8     in let rec sumdiff (slst,x) =
9       match slst with [] -> [] | hd::t1 ->
10      (-) (match hd with Const c -> (Const 0)::(sumdiff (t1,x))
11      (-) | Var v -> if v=x then (Const 1)::(sumdiff (t1,x)) else (Const 0)::(sumdiff (t1,x))
12      (-) | Power (v,c) -> if v=x then (Times [Const c;Power(v,c-1)]::(sumdiff (t1,x)) else (Const 0)::(sumdiff (t1,x))
13      (-) | Sum lst -> (Sum (sumdiff (lst,x)))::(sumdiff (t1,x)) | Times lst -> (Times (timediff (lst,x)))::(sumdiff (t1,x))
14      (+) (diff (hd, x)):(sumdiff (t1, x))
15      in let rec differ (e, x) =
16        match e with Const c -> Const 0 | Var v -> if v=x then Const 1 else Const 0
17        | Power (v,c) -> if v=x then Times [Const c;Power(v,c-1)] else Const 0 | Sum lst -> Sum (sumdiff (lst,x))
18        (-) | Times lst -> Times (timediff (lst,x))
19        (+) | Times lst -> Sum (timediff (lst,x))
20      in differ (e, x)

```

Figure 1: A real incorrect student submission and the feedback generated by CAFE

<pre> 1 let rec check (l, str) = ... 2 let rec diff (aexp, str) = ... 3 and gettimes (l, str) = 4 match l with [] -> [] hd::t1 -> 5 if t1=[] then [diff (hd, str)] ... 6 and getsum (l, str) = ... </pre> <p>(a) Program used to fix timediff</p>	<pre> 1 let rec diff (e, x) = ... 2 and diff_sum (lst, key) = 3 match lst with hd::t1 -> ... 4 and diff_time (lst, key, x, y) = 5 let rec f (lst, key, p, q) = ... in 6 if x > y then [] else ... </pre> <p>(b) Program used to fix sumdiff</p>	<pre> 1 let rec diff (aexp, str) = 2 let rec diffh (aexp, str) = ... 3 Sum l -> Sum (sum (l, str)) 4 and times (aexp, str, acc) = ... 5 and sum (aexp, str) = ... 6 in diffh (aexp, str) </pre> <p>(c) Program used to fix differ</p>
---	---	--

Figure 2: Three different solution programs chosen by CAFE to repair the program in Figure 1

<pre> 1 let rec inc_all l = 2 match l with [] -> [] 3 h::t -> (h+1)::(inc_all t) 4 5 let rec dec_all l = 6 match l with [] -> [] 7 h::t -> (h-1)::(dec_all t) 8 9 let apply ((l, o) as i) = 10 match l with [] -> [] 11 hd::t1 -> match o with 12 ADD -> (hd+1)::(inc_all t1) 13 SUB -> (dec_all t1) </pre> <p>(a) Incorrect program (P)</p>	<pre> 1 let rec add_list l1 = 2 match l1 with [] -> [] 3 h1::t1 -> (h1+1)::(add_list t1) 4 5 6 let rec apply1 ((l1, o1) as i) = 7 match l1 with [] -> [] hd1::t1 -> 8 let i' = (t1, o1) in match o1 with 9 ADD -> (hd1+1)::(add_list t1) 10 SUB -> (hd1-1)::(apply1 i') </pre> <p>(b) Correct program (P_1)</p>	<pre> 1 let id x = x 2 3 4 5 let rec sub_list l2 = 6 match l2 with [] -> [] 7 h2::t2 -> (h2-1)::(sub_list t2) 8 9 let rec apply2 ((l2, o2) as i) = 10 match l2 with [] -> id [] 11 hd2::t2 -> match o2 with 12 ADD -> (hd2+1)::(apply2 (t2, o2)) 13 SUB -> (hd2-1)::(sub_list t2) </pre> <p>(c) Correct program (P_2)</p>
---	---	--

Figure 3: A running example to illustrate how CAFE works

where i is called *input variable* representing the input simultaneously provided to the top-level functions. Over these contexts, CAFE concludes $\text{dec_all} \sim \text{sub_list}$ because the formula $\delta_{\text{dec_all}} \wedge \delta_{\text{sub_list}}$ is satisfiable meaning that there exists a value for i that leads to invoking both functions. On the other hand, CAFE concludes $\text{dec_all} \not\sim \text{add_list}$ because the formula $\delta_{\text{dec_all}} \wedge \delta_{\text{add_list}}$ is unsatisfiable ($\because \text{SUB} \neq \text{ADD}$) meaning that there is no value for i that results in invoking both functions. After a similar process,

CAFE concludes $\text{inc_all} \sim \text{add_list}$. Also, based on outgoing contexts, CAFE concludes $\text{apply} \sim \text{apply1}$.

Note that context-aware matching differs crucially from conventional syntactic or semantic matching used in prior data-driven techniques [12, 34, 42, 43]. For example, existing approaches would match dec_all with add_list because they are equivalent in syntax and semantics. However, this matching is undoubtedly useless;

we need to find a function useful for repair, rather than merely finding similar one.

Extracting Repair Templates. Next, CAFE learns fixes for each function in a buggy submission from its corresponding function in a correct submission. We support insertion/deletion of conditional branches, modifying subexpressions, and adding new function definitions. Especially, changing control flows and defining new functions are beyond the capability of the state-of-the-art feedback generation systems.

We learn such a fix by extracting a *repair template* from the correct submissions and *instantiating* it to be fit into the buggy submission. CAFE learns a repair template by syntactically differencing matched functions. For example, CAFE identifies a discrepancy between the subexpression $h+1$ in `decl_all` and h_2-1 in `sub_list`, and derives the template

$$\text{Modify}(7, h+1 \rightarrow \square_{\text{int}-1}) \quad (1)$$

which means one way to fix P is to replace the expression $h+1$ at line 7 by an expression of the form $\square_{\text{int}-1}$ where \square_{int} can be filled by some integer variable. Similarly, CAFE produces the following template by differencing `apply` and `apply1`:

$$\text{Modify}(13, \text{dec_all } t1 \rightarrow (\square_{\text{int}} - 1) :: (\square_{\text{int list} \times \tau} \square_{\text{int list} \times \tau}) \rightarrow_{\text{int list}}) \quad (2)$$

where τ denotes the algebraic data type for ADD or SUB.

Instantiating Templates. Next CAFE instantiates the templates by filling the holes with proper variables to obtain concrete fixes. Each hole with a type is filled with a variable available at the location with the same type. From the template (1), `Modify(7, $h+1 \rightarrow h-1$)` is generated because h is the only available integer variable at line 7, which is correct. From the template (2), the first hole can be filled with `hd` because it is the only available integer variable at line 13, and the following partially completed template is obtained.

$$\text{Modify}(13, \text{dec_all } t1 \rightarrow (\text{hd} - 1) :: (\square_{\text{int list} \times \tau} \square_{\text{int list} \times \tau}) \rightarrow_{\text{int list}})$$

Unfortunately, there is no variable of type $\text{int list} \times \tau \rightarrow \text{int list}$ (`apply` is unavailable as it is not recursive). In such a case, we just enumerate all function identifiers as candidates for the hole and obtain the partially completed templates:

$$\begin{aligned} &\text{Modify}(13, \text{dec_all } t1 \rightarrow (\text{hd} - 1) :: (\text{dec_all } \square_{\text{int list}})) \\ &\text{Modify}(13, \text{dec_all } t1 \rightarrow (\text{hd} - 1) :: (\text{inc_all } \square_{\text{int list}})) \end{aligned}$$

Note that the annotated type of the last hole has been changed in accordance with the type of `dec_all` and `inc_all`. The last hole is filled with `t1`, which is an available variable of type `int list`. Finally, we obtain the following set \mathcal{A} :

$$\left\{ \begin{array}{l} \text{Modify}(7, h+1 \rightarrow h-1), \\ \text{Modify}(13, \text{dec_all } t1 \rightarrow (\text{hd} - 1) :: (\text{dec_all } t1)), \\ \text{Modify}(13, \text{dec_all } t1 \rightarrow (\text{hd} - 1) :: (\text{inc_all } t1)) \end{array} \right\}$$

Finding a Patch. By trying each subset of \mathcal{A} , CAFE finds that applying the first two fixes in sequence corrects the buggy submission. Note that CAFE produced a patch that preserves the original intent of the program as much as possible by using the existing helper function `dec_all` instead of simply removing it and making `apply` recursive.

3 PROBLEM DEFINITION

In this section, we define our problem of data-driven feedback generation for programming assignments. We first define a program model that captures key aspects of Meta Language (ML)-like languages and introduce notations that allow us to formalize our algorithm in the next section.

Language. To formalize our approach, we consider an idealized functional language similar to the core of ML, with the additional property that we *label* all expressions. Our target language features algebraic data types and recursive functions. A program is an expression defined as follows:

$$\begin{array}{lll} e \in \text{Exp} & (\text{Expressions}) & x \in \text{Vid} & (\text{Variables}) \\ f \in \text{Fld} & (\text{Functions}) & \text{Id} = \text{Vid} \cup \text{Fld} & (\text{Identifiers}) \\ \ell \in \text{Label} & (\text{Labels}) & \tau \in \text{Type} & (\text{Types}) \end{array}$$

$$\begin{aligned} e & ::= n \mid x \mid \lambda x. e \mid e_1 \oplus e_2 \mid e_1 e_2 \mid \kappa(e_1, \dots, e_{a(\kappa)}) \\ & \mid \kappa^{-i}(e) \mid \text{let } x = e_1 \text{ in } e_2 \mid \text{let rec } f(x) = e_1 \text{ in } e_2 \\ & \mid \text{match } e \text{ with } \overline{p_i} \rightarrow e_i^k \\ p & ::= \kappa(x_1, \dots, x_n) \mid _ \quad \tau ::= \text{int} \mid T \mid \tau_1 \rightarrow \tau_2 \end{aligned}$$

We assume each expression is associated with a unique label. Expression e associated with a label $\ell \in \mathcal{L}$ is denoted by e^ℓ . For the sake of better readability, we will often elide ℓ when the label is not necessary for discussion. In addition, when we determine equality of two expressions, we do not consider their labels and only check if they are syntactically equivalent.

The syntax of expressions is standard: application is written $e_1 e_2$, κ ranges over data type constructors, $a(\kappa)$ denotes the arity of κ , κ^{-i} denotes a destructor which extracts the i -th subcomponent of a constructor κ , and let bindings for variables and recursive functions are allowed. For conciseness, we assume that all functions take a single argument and are not mutually recursive (our implementation in Section 5 is not limited by these restrictions though). We use ML-style pattern match expressions in which each pattern p binds subcomponents of a constructor κ , or the underscore ($_$) called the wildcard pattern. We use $\overline{p_i} \rightarrow e_i^k$ to denote $p_1 \rightarrow e_1 \mid \dots \mid p_k \rightarrow e_k$. Types include the integer type `int`, user-defined algebraic data types T , and function types $\tau_1 \rightarrow \tau_2$.

We will use some notations regarding expressions throughout the remaining sections. We use \rightarrow^* to denote the standard multi-step call-by-value operational relation. To denote the set of all subexpressions of expression e , we will use $\text{Sub}(e)$. The size of expression e will be denoted by $|e|$. Identifiers of functions defined in an expression e will be denoted by $\text{functions}(e)$ (i.e., $\text{functions}(e) = \{f \in \text{Fld} \mid \text{let rec } f(x) = e_1 \text{ in } e_2 \in \text{Sub}(e)\}$). Identifiers used in an expression e will be denoted by $\text{vars}(e)$.

Setting. We assume that each student submission $P \in \text{Exp}$ has no type error, and is in the following form: `let rec $f(x) = e$ in $f x_i$` , where f is a top-level function that the student is asked to implement, and x_i is a special variable which we call *input variable*. In Fig. 3a – 3c, `apply`, `apply1`, and `apply2` are the top-level functions, and the variable `i` can be regarded as the input variable. For convenience, we also assume that all labels and identifiers in the pile of entire submissions are unique with exception of the input variable that all the submissions have in common. This assumption enables

to use the following global functions: (1) $\text{body} : \text{Fld} \rightarrow \text{Exp}$ that returns the body of a function, (2) $\text{param} : \text{Fld} \rightarrow \text{Fld}$ that returns the parameter variable of a function, and (3) $\text{type} : (\text{Label} \cup \text{Id}) \rightarrow \text{Type}$ that returns the type of an expression with a given label or a variable.

Problem. Assuming some (possibly infinite) set Val of values, a set of test cases $\mathcal{T} \subseteq \text{Val} \times \text{Val}$ is used to determine the correctness of each submission. The submission P is correct (denoted $\text{correct}(P, \mathcal{T})$) iff $\forall (i, o) \in \mathcal{T}. (\lambda x_i. P) i \longrightarrow^* o$. Otherwise, the submission is buggy.

Our problem is defined as follows: given a buggy submission $P_b \in \text{Exp}$, a corpus of correct submissions $\mathcal{P}_c \subseteq \text{Exp}$, and a set of test cases \mathcal{T} , derive a correct submission $P \notin \mathcal{P}_c$ from P_b with minimal changes (the notion of the minimality will be detailed in Section 4.2.3).

4 ALGORITHM

In this section, we describe our data-driven feedback generation algorithm. Section 4.1 formalizes calling contexts and how to find a matching relation between functions. Section 4.2 describes the repair algorithm that extracts repair templates from reference functions and uses them to correct a buggy submission.

4.1 Context-Aware Matching

We formally define the notion of context-aware matching. From each function call in all the given submissions, we collect calling contexts. A context $\delta \in \text{Ctx}$ is a path condition on the input variable under which a function is invoked. Formally, path conditions are defined below:

$$\delta ::= \text{true} \mid \text{false} \mid e = e \mid \neg\delta \mid \delta \wedge \delta.$$

Calling contexts are defined as follows:

Definition 4.1 (Calling context). A calling context is a triple $\langle f, \delta, g \rangle \in \text{Fld} \times \text{Ctx} \times \text{Fld}$ where f is a function, g is another function called in the body of f , and δ is a path condition under which the function call happens.

We first perform a *path-sensitive 0-CFA* on all the submissions to obtain calling contexts. Like the standard 0-CFA [32], the analysis information at any given expression is the set of possible evaluation results of the expression. Here, evaluation results are expressions in which the input variable is a free variable. We add path sensitivity by making our analysis track information separately for different execution paths. The analysis computes a dataflow state $\sigma \in (\text{Id} \cup \text{Label}) \times \text{Ctx} \rightarrow \mathcal{P}(\text{Exp}) \cup \{\top\}$.

Figure 4 depicts a subset of the constraint generation rules; the full set can be found in the supplementary material. The judgement $\delta \vdash \llbracket e \rrbracket^\ell \hookrightarrow C$ can be read as “the analysis of expression e with label ℓ generates set constraints C over dataflow state σ under a current path condition δ ”. While solving the constraints via a least-fix point computation, special constraints of kinds **fn**, **cn**, and **pat** are interpreted by the constraint solver to generate additional concrete constraints by referring to an intermediate analysis result. For example, from a constraint $\text{fn}_\delta \ell_1 : \ell_2 \implies \ell$, for every function $\lambda x. e_0^{\ell_0}$ that the analysis (eventually) concludes the expression labeled ℓ_1 may evaluate to, additional constraints are generated to

capture value flow from the actual argument expression ℓ_2 to formal function argument x , and from the function result to the calling expression ℓ . To enforce termination, we use a standard widening operator that transforms each collected expression whose size is greater than a threshold into \top .

Note that the analysis for collecting calling contexts does not affect the correctness of the overall algorithm but just determines the effectiveness of matching.

We derive calling contexts from a result of the path-sensitive 0-CFA as follows: given submissions $P_1^{\ell_1}, \dots, P_m^{\ell_m}$, we collect set constraints C_i for each submission such that $\text{true} \vdash \llbracket P_i \rrbracket^{\ell_i} \hookrightarrow C_i$, and obtain the least solution σ_i . We collect a set of calling contexts $\Delta = \Delta_1 \cup \dots \cup \Delta_m$ where each Δ_i is

$$\bigcup_{\substack{f \in \text{functions}(P_i) \\ \delta \in \text{Ctx}}} \{ \langle f, \delta, g \rangle \mid (e_1^{\ell_1} e_2) \in \text{Sub}(\text{body}(f)), g \in \sigma_i(\ell_1, \delta) \}.$$

We also conjoin the analysis results from the submissions and obtain $\sigma = \bigsqcup_{1 \leq i \leq m} \sigma_i$, which will also be used in Section 4.2.

Example 4.2. Consider the invocation to `sub_list` in the body of the function `apply2` in Figure 3c. We will show how a calling context representing this function invocation is derived. From the parameter definition binding variables ℓ_2 and ℓ_2 in `apply2` where the initial path condition is `true`, we first generate the following constraints over a dataflow state σ .

$$\text{pair}^{-1}(\ell_2) \in \sigma(\ell_2, \text{true}) \quad (3)$$

$$\text{pair}^{-2}(\ell_2) \in \sigma(\ell_2, \text{true}) \quad (4)$$

In the outer pattern matching match ℓ_2 with $\text{cons}(\text{hd}_2, \text{tl}_2) \rightarrow \dots$, from (3), we generate the following constraint making ℓ_2 under a path condition: $\sigma(\ell_2, \text{true}) \subseteq \sigma(\ell_2, \underbrace{\text{pair}^{-1}(\ell_2) = \text{cons}(\text{hd}_2, \text{tl}_2)}_{\delta})$.

Under the current path condition δ , in the inner pattern matching match ℓ_2 with $\text{SUB} \rightarrow \dots \text{sub_list}^\ell \text{tl}_2$ where we assume `sub_list` is associated with label ℓ , the constraint $\text{sub_list} \in \sigma(\ell, \delta \wedge \text{pair}^{-2}(\ell_2) = \text{SUB})$ is generated from (4). After computing a least solution satisfying these constraints, we obtain a calling context $\langle \text{apply2}, \delta \wedge \text{pair}^{-2}(\ell_2) = \text{SUB}, \text{sub_list} \rangle$.

Now we are ready to measure similarity between arbitrary two functions using the calling contexts. Given two functions f and g , we compute a distance between the two functions as follows:

$$\text{dist}(f, g) = w_1 \times |CC_{\text{in}}^{f,g}|^{-1} + w_2 \times |CC_{\text{out}}^{f,g}|^{-1}$$

where $w_{\{1,2\}}$ are coefficients that can be adjusted via statistical learning. $CC_{\text{in}}^{f,g}$ (resp. $CC_{\text{out}}^{f,g}$) is called incoming (resp. outgoing) compatible calling context and defined as follows:

$$CC_{\text{in}}^{f,g} = \{ (\delta, \delta') \mid \delta, \delta' \in \text{Ctx}, \langle _, \delta, f \rangle, \langle _, \delta', g \rangle \in \Delta, \text{SAT}(\delta \wedge \delta') \}$$

$$CC_{\text{out}}^{f,g} = \{ (\delta, \delta') \mid \delta, \delta' \in \text{Ctx}, \langle f, \delta, _ \rangle, \langle g, \delta', _ \rangle \in \Delta, \text{SAT}(\delta \wedge \delta') \}$$

If both f and g do not have any callers (resp. callees), we do not consider the term involving $CC_{\text{in}}^{f,g}$ (resp. $CC_{\text{out}}^{f,g}$). Note that the more compatible pairs of calling contexts two functions have, the shorter distance they have between. Based on this notion of distance, we are equipped with the following matching function that takes a

$$\begin{array}{c}
\frac{}{\delta \vdash \llbracket n \rrbracket^\ell \hookrightarrow n \in \sigma(\ell, \delta)} \quad \frac{}{\delta \vdash \llbracket x \rrbracket^\ell \hookrightarrow \begin{cases} x \in \sigma(x, \delta) \cap \sigma(\ell, \sigma) & (x \in \{x_i\} \cup Fld) \\ \sigma(x, \delta) \subseteq \sigma(\ell, \sigma) & (\text{o.w.}) \end{cases}} \\
\frac{\delta \vdash \llbracket e_1 \rrbracket^{\ell_1} \hookrightarrow C_1 \quad \delta \vdash \llbracket e_2 \rrbracket^{\ell_2} \hookrightarrow C_2}{\delta \vdash \llbracket e_1 e_2 \rrbracket^\ell \hookrightarrow C_1 \cup C_2 \cup \text{fn } \ell_1 : \ell_2 \implies \ell} \quad \frac{\lambda x. e_0^{\ell_0} \in \sigma(\ell_1, \delta)}{\text{fn}_\delta \ell_1 : \ell_2 \implies \ell \hookrightarrow \{\sigma(\ell_2, \delta) \subseteq \sigma(x, \delta), \sigma(\ell_0, \delta) \subseteq \sigma(\ell, \delta)\}} \\
\frac{}{\delta \vdash \llbracket e \rrbracket^{\ell'} \hookrightarrow C} \quad \frac{e \in \sigma(\ell', \delta)}{\text{cn}_\delta^{\kappa^{-i}} \ell' \implies \ell \hookrightarrow \kappa^{-i}(e) \in \sigma(\ell, \delta)} \quad \frac{\delta \vdash \llbracket e \rrbracket^{\ell_0} \hookrightarrow C_0}{\delta \vdash \llbracket \text{match } e^{\ell_0} \text{ with } \overline{p_i} \rightarrow \overline{e_i^k} \rrbracket^\ell \hookrightarrow C_0 \cup \text{pat}_\delta \ell_0 : \overline{p_i} \rightarrow \overline{e_i^k} \implies \ell} \\
\frac{\delta \wedge r_1 \vdash e_1 \hookrightarrow C_1 \quad \delta \wedge r_2 \vdash e_2 \hookrightarrow C_2 \quad \dots \quad \delta \wedge r_k \hookrightarrow C_k}{\text{pat}_\delta \ell_0 : \overline{p_i} \rightarrow \overline{e_i^k} \implies \ell \hookrightarrow \bigcup_{1 \leq i \leq k} C_k \cup C_{FV} \cup C_\sqcup \cup \bigcup_{1 \leq i \leq k, p_i = \kappa(x_1, \dots, x_n)} \{\kappa^{-j}(e_0) \in \sigma(x_j, \delta) \mid 1 \leq j \leq n\}} \\
\begin{array}{l}
e_0 \in \sigma(\ell_0, \delta) \quad r_i = \begin{cases} e_0 = p_i & (i = 1) \\ e_0 = p_i \wedge \bigwedge_{1 \leq j < i} e_0 \neq p_j & (\text{o.w.}) \end{cases} \\
C_{FV} = \bigcup_{1 \leq i \leq k} \{\sigma(x, \delta) \subseteq \sigma(x, \delta \wedge r_i) \mid x \in FV(e_i)\} \\
C_\sqcup = \bigcup_{1 \leq i \leq k} \{\sigma(\ell_i, \delta \wedge r_i) \subseteq \sigma(\ell, \delta)\}
\end{array}
\end{array}$$

Figure 4: Constraint generation rules for our path-sensitive 0-CFA (selected). The special constraints of kinds `fn`, `cn`, and `pat` are interpreted by the constraint solver to generate additional concrete constraints.

function in a buggy submission and returns a function in a correct submission to be referred for correction:

$$M = \lambda P_b. \{f \mapsto \underset{\substack{P_c \in \mathcal{P}_c \\ g \in \text{functions}(P_c) \\ \text{type}(f) = \text{type}(g)}}}{\text{argmin}} \text{dist}(f, g) \mid f \in \text{functions}(P_b)\}.$$

Note that we only consider functions of the same type.

Example 4.3. Suppose we want to measure the distance between `dec_all` and `sub_list`, and the distance between `dec_all` and `add_list` in Figure 3. After the 0-CFA analysis, we obtain the following calling contexts.

$$\begin{array}{l}
\langle \text{apply}, \underbrace{\text{pair}^{-1}(i) = \text{cons}(\text{hd}, \text{tl}) \wedge \text{pair}^{-2}(i) = \text{SUB}, \text{dec_all}}_{\delta_1} \rangle \\
\langle \text{dec_all}, \underbrace{\delta_1 \wedge \text{cons}^{-2}(\text{pair}^{-1}(i)) = \text{cons}(h, t), \text{dec_all}}_{\delta_2} \rangle \\
\langle \text{apply1}, \underbrace{\text{pair}^{-1}(i) = \text{cons}(\text{hd}_1, \text{tl}_1) \wedge \text{pair}^{-2}(i) = \text{ADD}, \text{add_list}}_{\delta_3} \rangle \\
\langle \text{add_list}, \underbrace{\delta_3 \wedge \text{cons}^{-2}(\text{pair}^{-1}(i)) = \text{cons}(h_1, t_1), \text{add_list}}_{\delta_4} \rangle \\
\langle \text{apply2}, \underbrace{\text{pair}^{-1}(i) = \text{cons}(\text{hd}_2, \text{tl}_2) \wedge \text{pair}^{-2}(i) = \text{SUB}, \text{sub_list}}_{\delta_5} \rangle \\
\langle \text{sub_list}, \underbrace{\delta_5 \wedge \text{cons}^{-2}(\text{pair}^{-1}(i)) = \text{cons}(h_2, t_2), \text{sub_list}}_{\delta_6} \rangle \\
\dots
\end{array}$$

With $w_1 = 1$ and $w_2 = 2$ that we are using in our implementation, $\text{dist}(\text{dec_all}, \text{sub_list}) = |\{(\delta_1, \delta_5), (\delta_2, \delta_6)\}^{-1} + 2 \cdot |\{(\delta_2, \delta_6)\}^{-1}| = 2.5$

whereas $\text{dist}(\text{dec_all}, \text{add_list}) = \infty$ as the two functions do not share any compatible calling contexts.

In case of tie, we pick the most syntactically similar function. To measure the syntactic similarity, we use the method of embedding ASTs into numerical vectors, which is called the position-aware characteristic vectors proposed by Wang et al. [43]. We compute Euclidean distances between vectors to obtain the syntactic distances.

4.2 Repair Algorithm

In this subsection, we explain how to extract repair templates from correct submissions and instantiate them to generate patches. In particular, our goal is to obtain a sequence of *edit actions* that transform a given buggy submission into a new correct one. This sequence is called *edit script*. We consider the following edit actions:

- `Modify`(ℓ, e) replaces the old subexpression at label ℓ by the new expression e .
- `Insert`($\ell, p \rightarrow e$) adds a new pattern matching case $p \rightarrow e$ into a match expression associated with label ℓ .
- `Delete`($\ell, p \rightarrow e$) removes an existing pattern matching case $p \rightarrow e$ from a match expression associated with label ℓ .
- `Define`(f) adds a new definition of function f into the expression. If we apply this action into an expression e , the resulting expression would be `let rec $f(\text{param}(f)) = \text{body}(f)$ in e .`

4.2.1 Learning Repair Templates. We generate edit scripts by *instantiating* templates (which we call *repair templates*) collected from correct submissions. A repair template is a variant of an edit action where each expression in `Modify` or `Insert` action does not have any variables but just *holes*. Each hole is annotated with a type and plays a role as a placeholder that can be replaced with a variable of the type. The set Exp_\square of expressions with holes is similarly defined as Exp in the followings.

$$\begin{array}{l}
e_\square \in \text{Exp}_\square \\
e_\square ::= \square_\tau \mid n \mid \lambda x. e_\square \mid e_{\square,1} \oplus e_{\square,2} \mid e_{\square,1} e_{\square,2} \\
\quad \mid \kappa(e_{\square,1}, \dots, e_{\square,a(\kappa)}) \mid \kappa^{-i}(e_\square) \\
\quad \mid \text{let } x = e_{\square,1} \text{ in } e_{\square,2} \\
\quad \mid \text{let rec } f(x) = e_{\square,1} \text{ in } e_{\square,2} \\
\quad \mid \text{match } e_\square \text{ with } \overline{p_{\square,i}} \rightarrow \overline{e_{\square,i}^k} \\
p_\square ::= \kappa(\square_{\tau_1}, \dots, \square_{\tau_k}) \mid _
\end{array}$$

An expression with holes can be considered an abstraction of multiple expressions. The abstraction function $\alpha_e : \text{Exp} \rightarrow \text{Exp}_\square$, which we apply to expressions in correct submissions to extract templates, is defined as follows (to avoid unnecessary clutter, we omit simple

inductive cases):

$$\begin{aligned} \alpha_e(n) &= n & \alpha_e(x^\ell) &= \square_{\text{type}(x)}^\ell & \alpha_e(\lambda x.e) &= \lambda x.\alpha_e(e) \\ \dots & & & & & \\ \alpha_e(\text{match } e \text{ with } \overline{p_i \rightarrow e_i^k}) & & & & & \\ &= \text{match } \alpha_e(e) \text{ with } \overline{p_i \rightarrow \alpha_e(e_i)^k} & & & & \\ \alpha_p(\kappa(x_1, \dots, x_{a(\kappa)})) &= \kappa(\alpha_e(x_1), \dots, \alpha_e(x_{a(\kappa)})) & \alpha_p(_) &= _ \end{aligned}$$

When abstracting a variable into a hole, we preserve its label. The label is used in various ways, which will be described later in the next subsection.

Now we describe how to generate repair templates. Given a buggy submission P and the matching function \mathcal{M} , we obtain a set of templates $\mathcal{T} = \mathcal{T}_D \cup \mathcal{T}_M$ where \mathcal{T}_D is a set of templates of kind Define defined as follows:

$$\mathcal{T}_D = \{\text{Define}(g) \mid f \in \text{functions}(P), g \in \text{callees}(\text{body}(\mathcal{M}(f)))\}.$$

In other words, we collect all the auxiliary functions used in reference solutions. The function $\text{callees} : \text{Exp} \rightarrow \mathcal{P}(\text{Fld})$ returns all functions that may be invoked in a given expression. The set \mathcal{T}_M includes templates of kinds Modify, Insert, and Delete defined as follows:

$$\mathcal{T}_M = \bigcup \{T \mid f \in \text{functions}(P), \llbracket \text{body}(f), \text{body}(\mathcal{M}(f)) \rrbracket \rightsquigarrow T\}.$$

The judgement $\llbracket e, e' \rrbracket \rightsquigarrow T$ can be read as “by differencing a buggy expression e and a reference expression e' , we extract a set T of edit action templates that can be potentially used to correct e' ”. Figure 5 depicts a subset of inference rules for extracting templates for a given pairs of expressions. The full set is deferred to the supplementary material.

Example 4.4. Suppose we extract a set T of templates from `sub_list` for correcting `dec_all` in Figure 3. We extract templates T such that $\llbracket e_1, e_2 \rrbracket \rightsquigarrow T$ where e_1 is `body(sub_list)` with labels $\ell_{1..6}$ and e_2 is `body(dec_all)` with labels $\ell'_{1..6}$:

$$\begin{aligned} & \text{match } l^{\ell_1} \text{ with} \\ e_1 = & p_1 \rightarrow \text{empty} \\ & | p_2 \rightarrow \text{cons}((h^{\ell_2} + 1)^{\ell_3}, \text{dec_all}^{\ell_4} \ t^{\ell_5})^{\ell_6} \\ & \text{match } l_2^{\ell'_1} \text{ with} \\ e_2 = & p'_1 \rightarrow \text{empty} \\ & | p'_2 \rightarrow \text{cons}((h_2^{\ell'_2} - 1)^{\ell'_3}, \text{sub_list}^{\ell'_4} \ t_2^{\ell'_5})^{\ell'_6} \end{aligned}$$

and $p_1 = p'_1 = \text{empty}$, $p_2 = \text{cons}(h, t)$, and $p'_2 = \text{cons}(h_2, t_2)$. By the inference rule for differencing two matching expressions in Figure 5, we first compare l^{ℓ_1} and $l_2^{\ell'_1}$ and derive a template $\text{Modify}(\ell_1, \square_{\text{int list}})$. Since both $\{\alpha_p(p_1), \alpha_p(p_2)\}$ and $\{\alpha_p(p'_1), \alpha_p(p'_2)\}$ are $\{\text{empty}, \text{cons}(\square_{\text{int}}, \square_{\text{int list}})\}$ and the expressions matched for p_1 and p'_1 are the same, we compare the matched expressions for p_2 and p'_2 labeled ℓ_6 and ℓ'_6 respectively. By the rule for differencing two constructors, we additionally derive $\text{Modify}(\ell_3, \alpha_e(h_2 - 1)) = \text{Modify}(\ell_3, \square_{\text{int}}^{\ell'_2} - 1)$, $\text{Modify}(\ell_4, \square_{\text{int list} \rightarrow \text{int list}})$, and $\text{Modify}(\ell_5, \square_{\text{int list}})$.

4.2.2 Generating Edit Scripts. We instantiate the collected templates into edit actions using the following concretization function

$\gamma_e^P : \text{Exp}_\square \rightarrow \mathcal{P}(\text{Exp})$ parametrized by a buggy submission P .

$$\begin{aligned} \gamma_e^P(n) &= \{n\} & \gamma_e^P(\lambda x.e_\square) &= \{\lambda x.e \mid e \in \gamma_e^P(e_\square)\} \\ \dots & & & \\ \gamma_e^P(\text{match } e_{\square,0} \text{ with } \overline{p_i \rightarrow e_{\square,i}^k}) & & & \\ &= \{\text{match } e_0 \text{ with } \overline{p_i \rightarrow e_i^k} \mid e_i \in \gamma_e^P(e_{\square,i})\} \end{aligned}$$

Most importantly,

$$\gamma_e^P(\square_\tau^\ell) = \begin{cases} \{x \in V \mid \text{type}(x) = \tau\} & (\exists x \in V. \text{type}(x) = \tau) \\ V & (\text{otherwise}) \end{cases}$$

where $V = \text{vars}(P) \cup \{x \in \text{Fld} \mid \delta \in \text{Ctx}. x \in \sigma(\ell, \delta)\}$ is the set of candidate variables for the given hole \square_τ^ℓ . Note that the label ℓ always originates from a correct submission. We consider not only variables in P but also function identifiers reachable at ℓ as candidates for the hole. This is because we may copy function definitions in correct submissions (via Define actions) into P and let P invoke the newly added functions. In case of no candidate variable of type τ , we just enumerate all the variables in V .

Over a buggy submission P using the concretization function, we may obtain the following set \mathcal{A} of edit actions.

$$\begin{aligned} \mathcal{A} = & \{\text{Modify}(\ell, e) \mid \text{Modify}(\ell, e_\square) \in \mathcal{T}_M, e \in \gamma_e^P(e_\square)\} \\ & \cup \{\text{Insert}(\ell, p \rightarrow e) \mid \text{Insert}(\ell, p \rightarrow e_\square) \in \mathcal{T}_M, e \in \gamma_e^P(e_\square)\} \\ & \cup \{\text{Delete}(\ell, p \rightarrow e) \mid \text{Delete}(\ell, p \rightarrow e) \in \mathcal{T}_M\} \\ & \cup \mathcal{T}_D. \end{aligned}$$

However, this method is not scalable in practice as the number of concretized edit actions is potentially exponential to the number of holes in the template (despite the type-based pruning). To reduce the number of candidates for the holes, we use the following improved concretization function $\tilde{\gamma}_e^P$, which is similarly defined as γ_e^P except for the following case:

$$\tilde{\gamma}_e^P(\square_\tau^\ell) = \begin{cases} V_\tau \mid_\ell & (V_\tau \mid_\ell \neq \emptyset) \\ \gamma_e^P(\square_\tau^\ell) & \text{otherwise} \end{cases}$$

where $V_\tau \mid_\ell$ is the set of variables of type τ in V that may take the same value reachable at label ℓ . Formally,

$$V_\tau \mid_\ell = \{x \in V \mid \text{type}(x) = \tau, \exists \delta, \delta'. \sigma(x, \delta) \cap \sigma(\ell, \delta') \neq \emptyset\}.$$

This heuristic is inspired by the variable-usage based α -conversion of SARFGEN [43] but we analyze the usage more accurately using the result of our path-sensitive 0-CFA.

Example 4.5. Recall the template $\text{Modify}(\ell_3, \square_{\text{int}}^{\ell'_2} - 1)$ derived in Example 4.4. When concretizing $\square_{\text{int}}^{\ell'_2}$, we only consider the variable h as a candidate for the hole because

$$\begin{aligned} \sigma(h, \delta_1) &\ni \text{cons}^{-1}(\text{cons}^{-2}(\text{pair}^{-1}(i))) \\ \sigma(\ell'_2, \delta_5) &\ni \text{cons}^{-1}(\text{cons}^{-2}(\text{pair}^{-1}(i))) \end{aligned}$$

where δ_1 and δ_5 are the path conditions defined in Example 4.3.

Changing Annotated Types during Instantiation. For ease of presentation, we have presented our instantiation method as if types associated with holes could never change after they were determined. In the actual implementation, we change the types during the course of instantiation. Whenever a hole is filled with a variable, we perform type inference to change types of the other holes accordingly. An example case is the instantiation of the template (2) described in Section 2.2.

$$\begin{array}{c}
\frac{}{\llbracket n_1, n_2 \rrbracket \rightsquigarrow \emptyset} \quad n_1 = n_2 \quad \frac{}{\llbracket x_1^{\ell_1}, x_2^{\ell_2} \rrbracket \rightsquigarrow \{\text{Modify}(\ell_1, \square_{\text{type}(\ell_2)}^{\ell_2})\}} \quad \frac{\llbracket e_1, e_2 \rrbracket \rightsquigarrow T}{(\llbracket \lambda x_1. e_1, (\lambda x_2. e_2) \rrbracket \rightsquigarrow T)} \quad \frac{\llbracket e_1, e'_1 \rrbracket \rightsquigarrow T_1 \quad \llbracket e_2, e'_2 \rrbracket \rightsquigarrow T_2}{\llbracket (e_1 e_2), (e'_1 e'_2) \rrbracket \rightsquigarrow T_1 \cup T_2} \\
\\
\frac{\llbracket e_0, e'_0 \rrbracket \rightsquigarrow T_0 \quad T_M = \cup \{T_j \mid 1 \leq j \leq k, \llbracket e_j, e'_j \rrbracket \rightsquigarrow T_j, \alpha_p(p_j) = \alpha_p(p'_j)\} \quad T_I = \{\text{Insert}(\ell, p_j \rightarrow \alpha_e(e_j)) \mid 1 \leq j \leq m, \alpha_p(p_j) \in \{\alpha_p(p'_i)\}_{i=1}^n \setminus \{\alpha_p(p_i)\}_{i=1}^m\} \quad T_D = \{\text{Delete}(\ell, p_j \rightarrow e_j) \mid 1 \leq j \leq n, \alpha_p(p_j) \in \{\alpha_p(p_i)\}_{i=1}^n \setminus \{\alpha_p(p'_i)\}_{i=1}^m\}}{\llbracket (\text{match } e_0 \text{ with } \overline{p_i \rightarrow e_i^n}, (\text{match } e'_0 \text{ with } \overline{p'_i \rightarrow e'_i^m})) \rrbracket \rightsquigarrow T_0 \cup T_M \cup T_I \cup T_D} \\
\\
\frac{\llbracket e_1, e'_1 \rrbracket \rightsquigarrow T_1 \quad \dots \quad \llbracket e_k, e'_k \rrbracket \rightsquigarrow T_n}{\llbracket \kappa(e_1, \dots, e_n), \kappa(e'_1, \dots, e'_n) \rrbracket \rightsquigarrow \bigcup_{1 \leq i \leq n} T_i} \quad \frac{}{\llbracket e_1^{\ell_1}, e_2^{\ell_2} \rrbracket \rightsquigarrow \{\text{Modify}(\ell_1, \alpha_e(e_2))\}} \quad (e_1 \text{ and } e_2 \text{ are of different kinds.})
\end{array}$$

Figure 5: Inference rules for extracting edit action templates for given two expressions (selected).

Algorithm 1 The CAFE Algorithm

Input: A buggy submission P_b , a set of correct submissions \mathcal{P}_c , and a set of test cases \mathcal{T}

Output: A program P_c satisfying all the test cases in \mathcal{T}

```

1:  $\mathcal{A} \leftarrow \emptyset$  ▷ Set of edit actions
2:  $\mathcal{P} \leftarrow \mathcal{P}_c \cup \{P_b\}$ 
3:  $\sigma \leftarrow$  result of the path-sensitive OCFA on  $\mathcal{P}$ 
4:  $\Delta \leftarrow$  all calling contexts derivable from  $\sigma$ 
5: Derive  $\mathcal{M}$  from  $\Delta$  ▷  $\mathcal{M} : Fld \rightarrow Fld$ 
6:  $\mathcal{T} \leftarrow \text{ExtractTemplates}(\mathcal{M}, \Delta, P_b, \mathcal{P}_c)$ 
7: for  $T \in \mathcal{T}$  do
8:    $\mathcal{A} \leftarrow \mathcal{A} \cup \text{InstantiateTemplate}(T, P_b, \sigma)$ 
9:  $n \leftarrow 1$ 
10: repeat ▷  $E$ : edit script comprising  $n$  edit actions
11:   for each permutation  $E$  of  $n$  elements of  $\mathcal{A}$  do
12:      $P \leftarrow$  apply  $E$  into  $P_b$ 
13:     if  $\text{correct}(P, \mathcal{T})$  then
14:       return  $P$ 
15:      $n \leftarrow n + 1$ 
16: until  $n \leq |\mathcal{A}|$ 

```

4.2.3 Overall Algorithm. Putting all together, Algo. 1 depicts the CAFE algorithm. We first perform the path-sensitive OCFA on all the submissions and obtain the result σ (line 3). Then, we derive calling contexts from the analysis result (line 4). From the calling contexts, we obtain the matching function \mathcal{M} that maps each function in the buggy submission P_b to a function in a correct submission that is most likely to be useful for repair (line 5). Using the matching function, we collect repair templates (line 6). By instantiating the templates, we obtain a set of edit actions that can be applicable to P_b (lines 7 – 8). The main loop (lines 10 – 16) applies each possible sequence of the edit actions into P_b in turn. The variable n denotes the number of edit actions that can be used, which is initialized to be 1 (line 9). We apply each edit script into P_b (line 12) and check if the submission has been fixed based on the given test cases (line 13). If we have fixed the submission, we return it as a final result (line 14). Otherwise, we increase n by 1 (line 15) and repeat the main loop.

The algorithm enumerates edit scripts in increasing size, guaranteeing to find a *minimal* edit script in the following sense.

Definition 4.6 (Minimality). Given an incorrect submission P_b and a set of possible edit actions \mathcal{A} , an edit script E comprising the edit actions in \mathcal{A} to correct P_b is minimal if there does not exist an edit script E' such that $|E'| < |E|$ and E' fixes P_b .

4.2.4 Optimizations. When generating edit actions of kind Define that add new function definitions into a target buggy submission, we avoid functions that incur a long subsequent call chain to prevent CAFE from generating huge patches (currently, we only consider immediate callees of a reference function). Additionally, when generating edit scripts by permuting edit actions, we avoid generating duplicated edit scripts that lead to the same effect by not respecting orders between edit actions targetting different labels.

4.3 Discussion

Limitation of Context-Aware Matching. Our context-aware matching may produce imprecise results for functions called with trivial contexts. For example, suppose that there are two functions called under empty contexts (i.e., no path conditions are accumulated in their callers), and they should not be matched. In this case, context-aware matching would match them because their contexts are equivalent as empty contexts. We mitigated this shortcoming by applying the idea of context tunneling [18] and updating contexts at call-sites only when they are non-empty. For example, suppose we measure the similarity between two functions that have empty incoming contexts. In this case, we can apply context tunneling, so the callee inherits the callers' incoming contexts rather than producing empty contexts. We can further mitigate the issue of trivial contexts by falling back to existing syntactic matching: when the contexts are trivial (even after context tunneling), we can use the syntactic matching of SARFGEN rather than our context-aware matching.

Applicability to Other Languages. Although we formalized our approach for a functional language, the ideas of CAFE may work for other languages (e.g., Python) as well. First, the core idea can be applicable to any programs that consist of multiple functions, regardless of whether the language is functional or imperative. Second, the ideas of extracting repair templates from syntactic discrepancies and instantiating obtained templates can also be easily adapted for other languages whose syntax is defined inductively.

One feature of OCaml that CAFE particularly assumes is that it is a statically typed language. CAFE uses static type information to

prune out unnecessary function/variable mappings (i.e., excluding type-inconsistent pairs). This type-based optimization can be readily used for imperative languages with static typing such as C and Java. For dynamic languages like Python, CAFE is still applicable simply without the optimization (which might degrade the performance of CAFE slightly) or with extra type information computed by a static analysis.

5 EVALUATION

We evaluate CAFE to answer the following research questions:

- **Performance of CAFE:** How effectively can CAFE repair incorrect programs? How does it compare to the state-of-the-art for OCaml [22]? (Section 5.1)
- **Comparison with Prior Data-Driven Approaches:** How does our approach compare to the existing data-driven approaches (Section 5.2)
- **Helpfulness:** How helpful is CAFE for students? Is the generated feedback useful for students? (Section 5.3)

We implemented CAFE in about 7,000 lines of OCaml code. Although we formalized our approach for an ideal language, our implementation can handle all student programs in our class without any modification. We used the Z3 SMT solver for checking compatibility of path conditions. All experiments were conducted on an iMac with Intel i5 CPU and 16GB memory.

5.1 Performance of CAFE

Setting. We collected 4,211 OCaml programs from 10 exercises used in our class over the last few years, which include most of benchmarks used in our prior work [22].² To distinguish correct and incorrect programs, we used 10–33 test cases per exercise. These test cases have been carefully designed to detect various types of errors over the few years. All programs are compilable with no syntax or type errors. The description of the benchmark programs is given in Table 1. We classify the programming exercises into three levels, i.e., introductory (#1–#4), intermediate (#5–#7), and advanced (#8–#10), based on the code size and the ratio of incorrect to correct programs. Although code sizes look rather small, our benchmark set includes a number of notable programs (e.g., with 7 helper functions). Some examples programs are in the supplementary material. We compared CAFE with FixML [22], the state-of-the-art feedback generation tool for OCaml programs. FixML repairs incorrect programs using search-based program synthesis. Because FixML requires test cases and a solution program, we gave FixML an instructor-provided solution for each exercise and the same set of test cases as ours. We set time budget to 60 seconds per program for both CAFE and FixML.

Result. Table 1 shows that CAFE is far more effective than FixML in repairing student submissions. In total, CAFE successfully fixed 83% (548/664) of the buggy submissions, while the fix rate of FixML was 35% (234/664). Note that CAFE consistently achieves high fix rates for intermediate (82%, 95/116) and advanced problems (79%, 351/443), while FixML does not perform well for intermediate (47%,

54/116) and advanced problems (25%, 109/443). One key contributor to the high fix rate of CAFE was its capability of modifying multiple expressions with diverse repair strategies (e.g., insertion and deletion of branches, introduction of new functions). By contrast, FixML is limited to fixing single-location bugs.

We manually validated the correctness of patches, and Table 1 reports correct patches only. Since both CAFE and FixML use test cases as correctness specifications, they may produce *test-suite-overfitted patches* that satisfy given test cases but still contain errors. Originally, FixML generated 264 patches, among which 30 were overfitted to test cases. Because those patches are incorrect feedback, we only include 234 correct patches in Table 1. On the other hand, CAFE produced no incorrect patches. This was because CAFE leverages common templates extracted from solutions.

More qualitative analysis on the result is as follows. Notably, CAFE successfully fixed complex programs such as one with 7 functions and the call depth of 4. The generated patches were also non-trivial. CAFE modified 31% expressions of the original programs on average. Furthermore, we found that 25% (135/548) of the fixed errors are patched by repairing at least two functions simultaneously. When we investigated the statistics of edit actions used in patches, the distribution of each four templates (Modify, Insert, Delete, Define) was 79%, 8%, 2%, and 10%, respectively.

The performance of CAFE was not very sensitive to the amount of available data. For example, when we used 50% of the correct programs as a corpus, the fix rate remained almost the same: 82% (542/664, averaged over five random trials). When we used 10% of the correct programs, the fix rate decreased to 78% (517/664).

Limitation. We identified representative cases that CAFE can fail to produce patches. Timeout due to the large search space was the most common reason. For example, a buggy submission required CAFE to modify all (eight) if-then-else expressions in the program; CAFE is unlikely to generate such a large fix. Also, CAFE sometimes failed due to the lack of proper patch candidates caused by matching failure. In our experiments, potential imprecise matching discussed in Section 4.3 rarely happened (after applying context tunneling). For example, in problem 10, we observed that only five failures (9.8%, 5/51) were due to imprecise matching.

5.2 Comparison with Prior Techniques

We could not compare CAFE directly with existing data-driven tools as they target different languages [12, 34, 43] and rely on language-specific features [34]. However, to see how much CAFE advances the existing techniques, we implemented two variants, called Prog and Func, of CAFE. Prog and Func are identical to CAFE except that

- Prog uses the program-level matching of SARFGEN [43],
- Func uses it at the function level.

From the corpus of correct programs, Prog selects a program that is most similar to the given incorrect program, where we compute the similarity using the technique of SARFGEN [43], i.e., position-aware characteristic vector embedding. Thus, the performance gap between Prog and CAFE hints at how CAFE performs compared to SARFGEN. Func applies the matching algorithm of SARFGEN at the function level and therefore it partially enjoys the benefit of our approach (i.e., using multiple solutions via function-level matching).

²We excluded four exercises from [22] because they do not contain sufficient number of programs, or they require program-specific testing drivers that make the specification of problem unclear.

Table 1: Performance comparison of CAFE and FIXML. “#Wrong” and “#Correct” report the numbers of incorrect and correct submissions for each problem, respectively. “#Func”: the average, minimum, and maximum numbers of functions in buggy submissions. “LOC”: the average, minimum, and maximum lines of code of buggy submissions. “Time”: average patch-generation time (in sec). “#Fix (Rate)”: #correct patches generated by each tool and the patch rate.

No	Problem Description	#Wrong	#Correct	#Func avg(min-max)	LOC avg(min-max)	FIXML		CAFE	
						Time	#Fix (Rate)	Time	#Fix (Rate)
1	Finding a maximum element in a list	45	171	1.6 (1-3)	5 (1-9)	0.3	40 (89%)	0.0	45 (100%)
2	Checking membership in a binary tree	19	117	1.1 (1-3)	9 (5-13)	3.1	12 (63%)	0.0	19 (100%)
3	Mirroring a binary tree	9	88	1.0 (1-1)	6 (3-9)	0.1	7 (78%)	0.0	9 (100%)
4	Computing $\sum_{i=j}^k f(i)$ for j, k , and f	32	704	1.1 (1-2)	4 (2-10)	1.6	12 (38%)	2.1	29 (91%)
5	Composing functions	49	454	1.5 (1-3)	5 (2-11)	11.8	26 (53%)	1.1	42 (86%)
6	Removing redundant elements in a list	32	125	2.3 (1-5)	12 (5-24)	4.1	10 (31%)	3.0	23 (72%)
7	Arithmetic of user-defined natural numbers	35	412	2.2 (1-5)	13 (7-23)	23.3	18 (51%)	1.0	30 (86%)
8	Evaluating a propositional formula	111	597	2.1 (1-8)	29 (13-64)	1.5	44 (40%)	0.5	78 (70%)
9	Checking the validity of a lambda term	141	661	2.7 (1-7)	20 (6-47)	1.5	23 (16%)	1.8	133 (94%)
10	Differentiating an algebraic expression	191	218	2.1 (1-9)	29 (7-114)	1.1	42 (22%)	3.0	140 (73%)
Total / Average		664	3,547	2.0 (1-9)	20 (1-114)	3.9	234 (35%)	1.6	548 (83%)

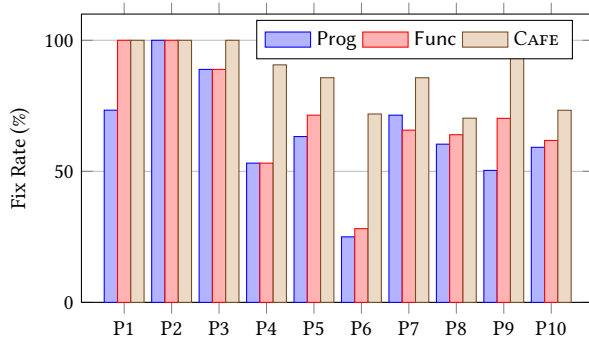


Figure 6: Comparison with existing data-driven techniques.

Thus, comparing Func and CAFE shows the sole impact of using our context-aware matching against the matching of SARFGEN.

The result shows that existing data-driven techniques are unlikely to be effective for our dataset. Figure 6 compares the fix rates of Prog, Func, and CAFE on the same dataset as Figure 1. On average, Prog achieved a fix rate of 59%. Simply extending the existing technique at the function level (Func) did not improve its performance significantly (67%). This is because, as illustrated in Section 2.2, simply aiming to find syntactically or semantically similar functions is unlikely to find a useful reference. CAFE was slightly superior to the existing syntactic approach in terms of efficiency and patch quality. When we compared the CAFE to the variant of SARFGEN (Prog), CAFE was faster than Prog (1.6 seconds vs. 1.9 seconds), which implies that context-aware matching can be as efficient as the embedding method of SARFGEN. In addition, we found that CAFE modified less expressions than Prog (31% and 33% respectively).

5.3 User Study

We recruited 16 undergraduate students from our course. The students were asked to solve the 10 programming exercises in Table 1. Then, we graded their submissions and provided feedback for erroneous ones using CAFE. Finally, students answered the survey

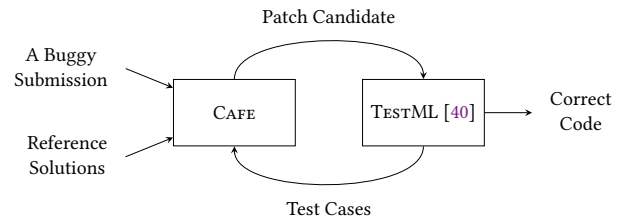


Figure 7: Enhanced CAFE with automatic test generator

questions about the generated feedback: (Q1) Is the feedback correct and easy to understand? (Q2) Does the provided feedback help you understand the mistake? (Q3) Do you think CAFE can be actually useful in our class? For each question, students were asked to choose between 1 (strongly disagree) and 5 (strongly agree). Also, we asked student to leave additional textual comments for each question. For Q1, Q2, and Q3, 14, 13, and all participants agreed (with scores 4 or 5), respectively. The comments the participants left include: "CAFE generates clever patches while keeping my code's structure", "Unlike the TA's solution code, the personalized feedback is easier to understand because it is derived from mine", and "It will help because it not only tells me the error, but also teaches me how to fix it".

5.4 CAFE with Automatic Test Case Generation

A limitation of CAFE is the reliance on manually-provided test cases, which hinders its use in real deployment: the quality of feedback depends on the quality of test cases but coming up with high-quality test cases requires massive human effort. In Table 1, we used manual test cases carefully refined over the few years, but such test cases are not always available. Below, we check if this limitation can be alleviated with the aid of automatic test generation techniques.

As shown in Figure 7, we built an enhanced version of CAFE in combination with TESTML [40]. TESTML is a recent counter-example generation tool for OCaml programming exercises, which takes two programs and tries to generate a test case on which

Table 2: CAFE with automatic test generation.

No	CAFE w/o TESTML			CAFE with TESTML		
	#Tests	#Fix	Time	#Tests (min-max)	#Fix	Time
8	18	78	0.5	4 (1-14)	79	124.5
9	33	133	1.8	3 (1-8)	133	123.6
10	28	140	3.0	5 (1-13)	141	124.8

the two programs behave differently. The enhanced CAFE in Figure 7 combines CAFE and TESTML in a loop. Note that it no longer requires manual test cases for patch validation. Instead, it uses TESTML as a correctness oracle. When TESTML fails to generate a counter-example for the submission and a solution (randomly chosen from the corpus), we regard the current patch candidate as a correct repair. Otherwise, TESTML augments the set of test cases, which is initially empty, with the generated counter-example, and CAFE is re-run with new test cases. The loop repeats until one of two components fails. We set the time budget for CAFE and TESTML to 60 and 120 seconds, respectively.

Table 2 compares CAFE with and without TESTML (results for Problems 8–10 only due to the lack of space). For CAFE without TESTML, #Tests reports the number of manual test cases for each exercise. For CAFE with TESTML, #Tests shows the average, smallest, and largest number of test cases generated by TESTML during the process in Figure 7. The results show that the enhanced system with TESTML reproduces the results in Table 1. Indeed, CAFE with TESTML generated three new patches and we confirmed they are correct. This was because CAFE now uses a smaller number of test cases and spends less time in patch validation. Despite the overhead, we found that combining CAFE and TESTML increases the usability significantly by reducing the instructor’s burden of crafting test cases and validating generated patches.

6 RELATED WORK

Automatic feedback generation has received an increasing amount of attention over the last years [1, 2, 5, 7, 10–14, 19, 20, 22, 23, 27, 31, 33, 35, 36, 39, 40, 43]. We discuss closely-related work below.

Our work builds upon but represents a significant departure from prior data-driven feedback generation techniques [12, 17, 34, 42, 43]. CLARA [12] is a clustering-based method that uses control flows and dynamic traces to find a correct solution. Similarly, SARFGEN [43] represents ASTs as vectors to find similar solution programs. These syntactic approaches could be improved using semantic features [34, 42]; however, those features are designed with imperative languages in mind and not readily applicable to CAFE. These data-driven approaches assume that there exists a *close enough* correct program in the corpus and CAFE aims to address this limitation by leveraging multiple, partially-matching programs.

Refactory [17] resolves the strict control flow matching problem of existing works [12, 43] by generating semantic-preserving references through refactoring. Note, however, that *Refactory* still relies on the control flow structures of refactored reference programs. On the other hand, our key observation in this paper is that conventional syntactic/semantic matching is not suitable for finding useful references when target programs consist of multiple

sub-functions. In this setting, “similar” reference functions are unlikely to be useful for repair (Section 2.2). To address this issue, we present a new method, context-aware matching, which matches functions by analyzing how they are used in programs (rather than comparing syntactic similarity).

RITE [37] and FixML [22] are state-of-the-art techniques for repairing student programs written in OCaml. Unlike ours, however, RITE can fix type errors only. FixML [22] uses program synthesis to repair general errors specified by test cases. However, FixML cannot fix multi-location errors, which consequently leads to a low fix rate. In this paper, we showed that CAFE outperforms FixML.

AutoGrader [39] and CoderAssist [19] are approaches that require manual effort. AutoGrader is a model-based technique to generate feedback by using constraint-based program synthesis. sk_p [35] addresses the limitation of AutoGrader by using a seq2seq neural network but is limited to small Python programs. CoderAssist generates verified feedback for introductory programming assignments but requires instructor-validated submissions.

Our work belongs to program repair techniques that use test cases as correctness criteria. Automatic program repair has a large volume of prior work [9, 30], which broadly classified into techniques for particular error types [4, 8, 16, 21, 26, 38, 41, 44, 47] and general-purpose techniques [3, 24, 25, 28, 29, 45, 46]. In particular, our work is similar to [28] in that both techniques use reference programs but our goal is to provide feedback on student programs.

Our context-aware matching is also related to existing program embedding techniques [6, 15]. For example, FUNC2VEC [6] is used to identify similar functions based on call relationships of functions. By contrast, our technique finds similar functions based on function contexts.

7 CONCLUSION

We presented a new technique that advances the existing data-driven approaches for automatically generating feedback on programming assignments. Unlike prior approaches, which works under the assumption that close enough reference programs exist in the corpus, CAFE can repair an incorrect submission by using multiple, partially-matching reference programs. To achieve this, we presented a new, context-aware repair algorithm. Evaluation results with real student submissions show that CAFE has a high fix rate and produces quality feedback actually useful for students.

ACKNOWLEDGMENTS

This work was supported by Institute of Information & communications Technology Planning & Evaluation(IITP) grant funded by the Korea government(MSIT) (No.2020-0-01337,(SW STAR LAB) Research on Highly-Practical Automated Software Repair) and Samsung Research Funding & Incubation Center of Samsung Electronics under Project Number SRFC-IT1701-51. This research was partly supported by the MSIT(Ministry of Science and ICT), Korea, under the ICT Creative Consilience program(IITP-2021-2020-0-01819) supervised by the IITP(Institute for Information & communications Technology Planning & Evaluation) and the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (2020R1C1C1014518, 2021R1A5A1021944).

REFERENCES

- [1] U. Z. Ahmed, P. Kumar, A. Karkare, P. Kar, and S. Gulwani. 2018. Compilation Error Repair: For the Student Programs, From the Student Programs. In *2018 IEEE/ACM 40th International Conference on Software Engineering: Software Engineering Education and Training (ICSE-SEET)*. 78–87.
- [2] Umair Z. Ahmed, Nisheeth Srivastava, Renuka Sindhgatta, and Amey Karkare. 2020. Characterizing the Pedagogical Benefits of Adaptive Feedback for Compilation Errors by Novice Programmers. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering: Software Engineering Education and Training (Seoul, South Korea) (ICSE-SEET '20)*. Association for Computing Machinery, New York, NY, USA, 139–150. <https://doi.org/10.1145/3377814.3381703>
- [3] Johannes Bader, Andrew Scott, Michael Pradel, and Satish Chandra. 2019. Getafix: Learning to Fix Bugs Automatically. *Proc. ACM Program. Lang.* 3, OOPSLA, Article 159 (Oct. 2019), 27 pages. <https://doi.org/10.1145/3360585>
- [4] Xi Cheng, Min Zhou, Xiaoyu Song, Ming Gu, and Jianguang Sun. 2017. IntPTI: Automatic Integer Error Repair with Proper-type Inference. In *Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering (Urbana-Champaign, IL, USA) (ASE 2017)*. IEEE Press, Piscataway, NJ, USA, 996–1001. <http://dl.acm.org/citation.cfm?id=3155562.3155693>
- [5] Loris D'Antoni, Roopsha Samanta, and Rishabh Singh. 2016. Qclose: Program Repair with Quantitative Objectives. In *27th International Conference on Computer Aided Verification (CAV 2016)* (27th international conference on computer aided verification (cav 2016) ed.).
- [6] Daniel DeFreez, Aditya V. Thakur, and Cindy Rubio-González. 2018. Path-Based Function Embedding and Its Application to Error-Handling Specification Mining. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (Lake Buena Vista, FL, USA) (ESEC/FSE 2018)*. Association for Computing Machinery, New York, NY, USA, 423–433. <https://doi.org/10.1145/3236024.3236059>
- [7] A. Drummond, Y. Lu, S. Chaudhuri, C. Jermaine, J. Warren, and S. Rixner. 2014. Learning to Grade Student Programs in a Massive Open Online Course. In *2014 IEEE International Conference on Data Mining*. 785–790. <https://doi.org/10.1109/ICDM.2014.142>
- [8] Qing Gao, Yingfei Xiong, Yaqing Mi, Lu Zhang, Weikun Yang, Zhaoping Zhou, Bing Xie, and Hong Mei. 2015. Safe Memory-leak Fixing for C Programs. In *Proceedings of the 37th International Conference on Software Engineering - Volume 1 (Florence, Italy) (ICSE '15)*. IEEE Press, Piscataway, NJ, USA, 459–470. <http://dl.acm.org/citation.cfm?id=2818754.2818812>
- [9] L. Gazzola, D. Micucci, and L. Mariani. 2019. Automatic Software Repair: A Survey. *IEEE Transactions on Software Engineering* 45, 1 (Jan 2019), 34–67. <https://doi.org/10.1109/TSE.2017.2755013>
- [10] 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 Trans. Comput.-Hum. Interact.* 22, 2, Article 7 (March 2015), 35 pages. <https://doi.org/10.1145/2699751>
- [11] 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 (Hong Kong, China) (FSE 2014)*. Association for Computing Machinery, New York, NY, USA, 41–51. <https://doi.org/10.1145/2635868.2635912>
- [12] Sumit Gulwani, Ivan Radiček, and Florian Zuleger. 2018. Automated Clustering and Program Repair for Introductory Programming Assignments. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation (Philadelphia, PA, USA) (PLDI 2018)*. Association for Computing Machinery, New York, NY, USA, 465–480. <https://doi.org/10.1145/3192366.3192387>
- [13] Rahul Gupta, Aditya Kanade, and Shirish Shevade. 2019. Neural Attribution for Semantic Bug-Localization in Student Programs. In *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.). Curran Associates, Inc., 11884–11894.
- [14] Andrew Head, Elena Glassman, Gustavo Soares, Ryo Suzuki, Lucas Figueredo, Loris D'Antoni, and Björn Hartmann. 2017. Writing Reusable Code Feedback at Scale with Mixed-Initiative Program Synthesis. In *Proceedings of the Fourth (2017) ACM Conference on Learning @ Scale (Cambridge, Massachusetts, USA) (L@S '17)*. Association for Computing Machinery, New York, NY, USA, 89–98. <https://doi.org/10.1145/3051457.3051467>
- [15] Jordan Henkel, Shuwendu K. Lahiri, Ben Liblit, and Thomas Reps. 2018. Code Vectors: Understanding Programs through Embedded Abstracted Symbolic Traces. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (Lake Buena Vista, FL, USA) (ESEC/FSE 2018)*. Association for Computing Machinery, New York, NY, USA, 163–174. <https://doi.org/10.1145/3236024.3236085>
- [16] Seongjoon Hong, Junhee Lee, Jeongsoo Lee, and Hakjoo Oh. 2020. SAVER: Scalable, Precise, and Safe Memory-Error Repair. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering (Seoul, South Korea) (ICSE '20)*. Association for Computing Machinery, New York, NY, USA, 271–283. <https://doi.org/10.1145/3377811.3380323>
- [17] Yang Hu, Umair Z. Ahmed, Sergey Mehtaev, Ben Leong, and Abhik Roychoudhury. 2019. Re-Factoring Based Program Repair Applied to Programming Assignments. In *Proceedings of the 34th IEEE/ACM International Conference on Automated Software Engineering (San Diego, California) (ASE '19)*. IEEE Press, 388–398. <https://doi.org/10.1109/ASE.2019.00044>
- [18] Minseok Jeon, Seungho Jeong, and Hakjoo Oh. 2018. Precise and Scalable Points-to Analysis via Data-Driven Context Tunneling. *Proc. ACM Program. Lang.* 2, OOPSLA, Article 140 (Oct. 2018), 29 pages. <https://doi.org/10.1145/3276510>
- [19] 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 (Seattle, WA, USA) (FSE 2016)*. Association for Computing Machinery, New York, NY, USA, 739–750. <https://doi.org/10.1145/2950290.2950363>
- [20] Dohyeon Kim, Yonghui 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 (Amsterdam, Netherlands) (OOPSLA 2016)*. Association for Computing Machinery, New York, NY, USA, 311–327. <https://doi.org/10.1145/2983990.2984031>
- [21] Junhee Lee, Seongjoon Hong, and Hakjoo Oh. 2018. MemFix: Static Analysis-based Repair of Memory Deallocation Errors for C. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (Lake Buena Vista, FL, USA) (ESEC/FSE 2018)*. ACM, New York, NY, USA, 95–106. <https://doi.org/10.1145/3236024.3236079>
- [22] Junho Lee, Dowon Song, Sunbeom So, and Hakjoo Oh. 2018. Automatic Diagnosis and Correction of Logical Errors for Functional Programming Assignments. *Proc. ACM Program. Lang.* 2, OOPSLA, Article 158 (Oct. 2018), 30 pages. <https://doi.org/10.1145/3276528>
- [23] X. Liu, S. Wang, P. Wang, and D. Wu. 2019. Automatic Grading of Programming Assignments: An Approach Based on Formal Semantics. In *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering Education and Training (ICSE-SEET)*. 126–137. <https://doi.org/10.1109/ICSE-SEET.2019.00022>
- [24] Fan Long and Martin Rinard. 2015. Staged Program Repair with Condition Synthesis. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering (Bergamo, Italy) (ESEC/FSE 2015)*. ACM, New York, NY, USA, 166–178. <https://doi.org/10.1145/2786805.2786811>
- [25] 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 (St. Petersburg, FL, USA) (POPL '16)*. ACM, New York, NY, USA, 298–312. <https://doi.org/10.1145/2837614.2837617>
- [26] A. Marginean, J. Bader, S. Chandra, M. Harman, Y. Jia, K. Mao, A. Mols, and A. Scott. 2019. SapFix: Automated End-to-end Repair at Scale. In *Proceedings of the 41st International Conference on Software Engineering: Software Engineering in Practice (Montreal, Quebec, Canada)*. IEEE Press, Piscataway, NJ, USA, 269–278. <https://doi.org/10.1109/ICSE-SEIP.2019.00039>
- [27] V. J. Marin, T. Pereira, S. Sridharan, and C. R. Rivero. 2017. Automated Personalized Feedback in Introductory Java Programming MOOCs. In *2017 IEEE 33rd International Conference on Data Engineering (ICDE)*. 1259–1270. <https://doi.org/10.1109/ICDE.2017.169>
- [28] Sergey Mehtaev, Manh-Dung Nguyen, Yannic Noller, Lars Grunske, and Abhik Roychoudhury. 2018. Semantic Program Repair Using a Reference Implementation. In *Proceedings of the 40th International Conference on Software Engineering (Gothenburg, Sweden) (ICSE '18)*. Association for Computing Machinery, New York, NY, USA, 129–139. <https://doi.org/10.1145/3180155.3180247>
- [29] Sergey Mehtaev, Jooyong Yi, and Abhik Roychoudhury. 2016. Angelix: Scalable Multiline Program Patch Synthesis via Symbolic Analysis. In *Proceedings of the 38th International Conference on Software Engineering (Austin, Texas) (ICSE '16)*. ACM, New York, NY, USA, 691–701. <https://doi.org/10.1145/2884781.2884807>
- [30] Martin Monperrus. 2018. Automatic Software Repair: A Bibliography. *ACM Comput. Surv.* 51, 1, Article 17 (Jan. 2018), 24 pages. <https://doi.org/10.1145/3105906>
- [31] Andy Nguyen, Christopher Piech, Jonathan Huang, and Leonidas Guibas. 2014. Codewebs: Scalable Homework Search for Massive Open Online Programming Courses. In *Proceedings of the 23rd International Conference on World Wide Web (Seoul, Korea) (WWW '14)*. Association for Computing Machinery, New York, NY, USA, 491–502. <https://doi.org/10.1145/2566486.2568023>
- [32] Flemming Nielson, Hanne R. Nielson, and Chris Hankin. 1999. *Principles of Program Analysis*. Springer-Verlag, Berlin, Heidelberg.
- [33] Sagar Parihar, Ziyaan Dadachanji, Praveen Kumar Singh, Rajdeep Das, Amey Karkare, and Arnab Bhattacharya. 2017. Automatic Grading and Feedback Using Program Repair for Introductory Programming Courses. In *Proceedings of the 2017 ACM Conference on Innovation and Technology in Computer Science Education (Bologna, Italy) (ITICSE '17)*. Association for Computing Machinery, New York, NY, USA, 92–97. <https://doi.org/10.1145/3059009.3059026>

- [34] David M. Perry, Dohyeong Kim, Roopsha Samanta, and Xiangyu Zhang. 2019. SemCluster: Clustering of Imperative Programming Assignments Based on Quantitative Semantic Features. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Phoenix, AZ, USA) (PLDI 2019). Association for Computing Machinery, New York, NY, USA, 860–873. <https://doi.org/10.1145/3314221.3314629>
- [35] Yewen Pu, Karthik Narasimhan, Armando Solar-Lezama, and Regina Barzilay. 2016. *SKP*: 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* (Amsterdam, Netherlands) (SPLASH Companion 2016). Association for Computing Machinery, New York, NY, USA, 39–40. <https://doi.org/10.1145/2984043.2989222>
- [36] 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* (Buenos Aires, Argentina) (ICSE '17). IEEE Press, 404–415. <https://doi.org/10.1109/ICSE.2017.44>
- [37] Georgios Sakkas, Madeline Endres, Benjamin Cosman, Westley Weimer, and Ranjit Jhala. 2020. Type Error Feedback via Analytic Program Repair. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation* (London, UK) (PLDI 2020). Association for Computing Machinery, New York, NY, USA, 16–30. <https://doi.org/10.1145/3385412.3386005>
- [38] Alex Shaw, Dusten Doggett, and Munawar Hafiz. 2014. Automatically Fixing C Buffer Overflows Using Program Transformations. In *Proceedings of the 2014 44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks* (DSN '14). IEEE Computer Society, Washington, DC, USA, 124–135. <https://doi.org/10.1109/DSN.2014.25>
- [39] 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* (Seattle, Washington, USA) (PLDI '13). Association for Computing Machinery, New York, NY, USA, 15–26. <https://doi.org/10.1145/2491956.2462195>
- [40] Dowon Song, Myungho Lee, and Hakjoo Oh. 2019. Automatic and Scalable Detection of Logical Errors in Functional Programming Assignments. *Proc. ACM Program. Lang.* 3, OOPSLA, Article 188 (Oct. 2019), 30 pages. <https://doi.org/10.1145/3360614>
- [41] Rijnaard van Tonder and Claire Le Goues. 2018. Static Automated Program Repair for Heap Properties. In *Proceedings of the 40th International Conference on Software Engineering* (Gothenburg, Sweden) (ICSE '18). ACM, New York, NY, USA, 151–162. <https://doi.org/10.1145/3180155.3180250>
- [42] Ke Wang, Rishabh Singh, and Zhendong Su. 2018. Dynamic Neural Program Embeddings for Program Repair. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net. <https://openreview.net/forum?id=BjuWrGW0Z>
- [43] Ke Wang, Rishabh Singh, and Zhendong Su. 2018. Search, Align, and Repair: Data-Driven Feedback Generation for Introductory Programming Exercises. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Philadelphia, PA, USA) (PLDI 2018). Association for Computing Machinery, New York, NY, USA, 481–495. <https://doi.org/10.1145/3192366.3192384>
- [44] Westley Weimer. 2006. Patches as better bug reports. In *Proceedings of the 5th international conference on Generative programming and component engineering*. ACM, 181–190.
- [45] Westley Weimer, ThanhVu Nguyen, Claire Le Goues, and Stephanie Forrest. 2009. Automatically Finding Patches Using Genetic Programming. In *Proceedings of the 31st International Conference on Software Engineering* (ICSE '09). IEEE Computer Society, Washington, DC, USA, 364–374. <https://doi.org/10.1109/ICSE.2009.5070536>
- [46] Yingfei Xiong, Jie Wang, Runfa Yan, Jiachen Zhang, Shi Han, Gang Huang, and Lu Zhang. 2017. Precise Condition Synthesis for Program Repair. In *Proceedings of the 39th International Conference on Software Engineering* (Buenos Aires, Argentina) (ICSE '17). IEEE Press, Piscataway, NJ, USA, 416–426. <https://doi.org/10.1109/ICSE.2017.45>
- [47] Xuezheng Xu, Yulei Sui, Hua Yan, and Jingling Xue. 2019. VFix: Value-flow-guided Precise Program Repair for Null Pointer Dereferences. In *Proceedings of the 41st International Conference on Software Engineering* (Montreal, Quebec, Canada) (ICSE '19). IEEE Press, Piscataway, NJ, USA, 512–523. <https://doi.org/10.1109/ICSE.2019.00063>