

Technical Report (PRL-2020-01)

Context-Aware and Data-Driven Program Repair

Dowon Song
dowon_song@korea.ac.kr
Korea University

Hakjoo Oh
hakjoo_oh@korea.ac.kr
Korea University

1 Problem Definition

In this section, we define our program-repair problem. We first define a program model that captures key aspects of ML-like languages and introduce notations that allow us to formalize our algorithm in the next section.

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:

$e \in Exp$	(Expressions)	$x \in Vid$	(Variables)
$f \in Fld$	(Functions)	$Id = Vid \cup Fld$	(Identifiers)
$\ell \in Label$	(Labels)	$\tau \in Type$	(Types)

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

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. 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 \longrightarrow^* 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., $\text{functions}(e) = \{f \in Fld \mid$

$\text{let rec } f(x) = e_1 \text{ in } e_2 \in \text{Sub}(e)\}$). Identifiers used in an expression e will be denoted by $\text{vars}(e)$.

We assume that each program $P \in Exp$ has no type error, and is in the following form:

$$\text{let rec } f(x) = e \text{ in } f x_i$$

where f is a top-level function, and x_i is a special variable which we call *input variable*. For convenience, we also assume that all labels and identifiers in the pile of entire programs are unique with exception of the input variable that all the programs have in common. This assumption enables to use the following global functions: (1) $\text{body} : Fld \rightarrow Exp$ that returns the body of a function, (2) $\text{param} : Fld \rightarrow Vid$ that returns the parameter variable of a function, and (3) $\text{type} : (Label \cup Id) \rightarrow Type$ that returns the type of an expression with a given label or a variable.

Assuming some (possibly infinite) set Val of values, a set of test cases $\mathcal{T} \subseteq Val \times Val$ is used to determine the correctness of each program. The program 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 program is buggy.

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

2 Repair Algorithm

In this section, we describe our repair algorithm. Section ?? formalizes calling contexts and how to find a matching relation between functions. Section ?? describes the repair algorithm that extracts repair templates from reference functions and uses them to correct a buggy program.

2.1 Context-Aware Matching

We formally define the notion of context-aware matching. From each function call in all the given programs, we collect calling contexts. A context $\delta \in Ctx$ 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 2.1 (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 programs to obtain calling contexts. Like the standard 0-CFA, 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

$$\sigma \in (\text{Id} \cup \mathcal{L}) \times \text{Ctx} \rightarrow \mathcal{P}(\text{Exp}) \cup \{\top\}.$$

Figure ?? depicts the constraint generation rules. 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, four kinds of special constraints of the forms

$$\begin{array}{ll} \text{fn}_\delta \ell_1 : \ell_2 \implies \ell & \text{cn}_\delta^\kappa \ell_1, \dots, \ell_{a(\kappa)} \implies \ell \\ \text{cn}_\delta^{\kappa-i} \ell' \implies \ell & \text{pat}_\delta \ell_0 : p_i \rightarrow e_i^{\ell_i} \implies \ell \end{array}$$

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 programs $P_1^{\ell_1}, \dots, P_m^{\ell_m}$, we collect set constraints C_i for each program 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 for each program and obtain $\sigma = \bigsqcup_{1 \leq i \leq m} \sigma_i$, which will also be used in Section ??.

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:

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

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

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

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

Note that we only consider functions of the same type.

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. We compute Euclidean distances between vectors to obtain the syntactic distances.

2.2 Context-Aware Repair

In this subsection, we explain how to extract repair templates from correct programs and instantiate them to generate patches. In particular, our goal is to obtain a sequence of *edit actions* that transform a given buggy program 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 $\text{let rec } f(\text{param}(f)) = \text{body}(f) \text{ in } e$.

2.2.1 Learning Repair Templates

We generate edit scripts by *instantiating* templates (which we call *repair templates*) collected from correct programs. 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

$$\begin{array}{c}
\frac{}{\delta \vdash \llbracket n \rrbracket^\ell \hookrightarrow n \in \sigma(\ell, \delta)} \quad \frac{\llbracket e \rrbracket^{\ell'} \hookrightarrow C}{\delta \vdash \llbracket \lambda x. e' \rrbracket^\ell \hookrightarrow \lambda x. e \in \sigma(\ell, \delta) \cup C} \quad \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 \mathbf{fn} \ell_1 : \ell_2 \implies \ell} \\
\\
\frac{x \in \{x_i\} \cup \mathbf{Fld}}{\delta \vdash \llbracket x \rrbracket^\ell \hookrightarrow x \in \sigma(x, \delta) \cap \sigma(\ell, \sigma)} \quad \frac{x \notin \{x_i\} \cup \mathbf{Fld}}{\delta \vdash \llbracket x \rrbracket^\ell \hookrightarrow \sigma(x, \delta) \subseteq \sigma(\ell, \sigma)} \quad \frac{\delta \vdash \llbracket e \rrbracket^{\ell'} \hookrightarrow C}{\delta \vdash \llbracket \kappa^{-i}(e') \rrbracket^\ell \hookrightarrow C \cup \mathbf{cn}_\delta^{\kappa^{-i}} \ell' \implies \ell} \\
\\
\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^{\ell_1} \oplus e_2^{\ell_2} \rrbracket^\ell \hookrightarrow \sigma(\ell_1, \delta) \oplus \sigma(\ell_2, \delta) \subseteq \sigma(\ell, \delta) \cup C_1 \cup C_2} \quad \frac{\delta \vdash \llbracket e_1 \rrbracket^{\ell_1} \hookrightarrow C_1 \quad \dots \quad \delta \vdash \llbracket e_n \rrbracket^{\ell_n} \hookrightarrow C_n}{\delta \vdash \llbracket \kappa(e_1^{\ell_1}, \dots, e_n^{\ell_n}) \rrbracket^\ell \hookrightarrow \bigcup_{1 \leq i \leq n} C_i \cup \mathbf{cn}_\delta^{\kappa} \ell_1, \dots, \ell_n \implies \ell} \\
\\
\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 e_i^k} \rrbracket^\ell \hookrightarrow C_0 \cup \mathbf{pat}_\delta \ell_0 : \overline{p_i \rightarrow e_i^k} \implies \ell} \quad \frac{e_1 \in \sigma(\ell_1, \delta) \quad \dots \quad e_n \in \sigma(\ell_n, \delta)}{\mathbf{cn}_\delta^{\kappa} \ell_{1..n} \implies \ell \hookrightarrow \kappa(e_1, \dots, e_n) \in \sigma(\ell, \delta)} \quad \frac{e \in \sigma(\ell', \delta)}{\mathbf{cn}_\delta^{\kappa^{-i}} \ell' \implies \ell \hookrightarrow \kappa^{-i}(e) \in \sigma(\ell, \delta)} \\
\\
\frac{\lambda x. e_0^{\ell_0} \in \sigma(\ell_1, \delta)}{\mathbf{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)\}} \quad \frac{f \in \sigma(\ell_1, \delta) \quad x = \text{param}(f) \quad e_0^{\ell_0} = \text{body}(f)}{\mathbf{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{\begin{array}{c} \delta \wedge r_1 \vdash e_1 \hookrightarrow C_1 \\ \delta \wedge r_2 \vdash e_2 \hookrightarrow C_2 \\ \dots \\ \delta \wedge r_k \hookrightarrow C_k \end{array} \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} \quad \begin{array}{c} 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}}{\mathbf{pat}_\delta \ell_0 : \overline{p_i \rightarrow 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\}}
\end{array}$$

Figure 1. Constraint generation rules for our path-sensitive 0-CFA (some simple cases are omitted for brevity). The special constraints of the kinds **fn**, **cn**, and **pat** are interpreted by the constraint solver to generate additional concrete constraints. $E_1 \oplus_\uparrow E_2$ denotes $\{e_1 \oplus e_2 \mid e_1 \in E_1, e_2 \in E_2\}$.

variable of the type. The set Exp_\square of expressions with holes is similarly defined as Exp in the followings.

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

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

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

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 program 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 programs. The function callees : $Exp \rightarrow \mathcal{P}(\mathbf{Fld})$ 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 ?? depicts inference rules for extracting templates for a given pairs of expressions.

2.2.2 Generating Edit Scripts

We instantiate the collected templates into edit actions using the following concretization function $\gamma_e^P : Exp_\square \rightarrow \mathcal{P}(Exp)$ parametrized by a buggy program P .

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

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 $V = \text{vars}(P) \cup \{x \in \mathbf{Fld} \mid \delta \in \text{Ctx}. x \in \sigma(\ell, \delta)\}$ is the set of candidate variables for the given hole \square_τ^ℓ . Note that the label ℓ always originates from a correct program. We consider not only variables in P but also function identifiers

$$\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 \rrbracket, (\lambda x_2. e_2)) \rightsquigarrow T} \quad \frac{\llbracket e_1, e_2 \rrbracket \rightsquigarrow T}{\llbracket \kappa^{-i}(e_1), \kappa^{-i}(e_2) \rrbracket \rightsquigarrow T} \\
\frac{\llbracket e_1, e'_1 \rrbracket \rightsquigarrow T_1 \quad \llbracket e_2, e'_2 \rrbracket \rightsquigarrow T_2}{\llbracket (e_1 \oplus e_2), (e'_1 \oplus e'_2) \rrbracket \rightsquigarrow T_1 \cup T_2} \oplus = \oplus' \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} \quad \frac{\llbracket e_1, e'_1 \rrbracket \rightsquigarrow T_1 \quad \cdots \quad \llbracket e_k, e'_n \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} \\
\frac{\llbracket e, e_2 \rrbracket \rightsquigarrow T}{\llbracket e, (\text{let } x = e_1 \text{ in } e_2) \rrbracket \rightsquigarrow T} \quad \frac{\llbracket e_2, e \rrbracket \rightsquigarrow T}{\llbracket (\text{let } x = e_1 \text{ in } e_2), e \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 (\text{let rec } f_1(x_1) = e_1 \text{ in } e_2), (\text{let rec } f_2(x_2) = e'_1 \text{ in } e'_2) \rrbracket \rightsquigarrow T_1 \cup T_2} \\
\frac{\llbracket e, e_2 \rrbracket \rightsquigarrow T}{\llbracket e, (\text{let rec } f(x) = e_1 \text{ in } e_2) \rrbracket \rightsquigarrow T} \quad \frac{\llbracket e_0, e'_0 \rrbracket \rightsquigarrow T_0 \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\}}{\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, e_2 \rrbracket \rightsquigarrow T}{\llbracket (\text{let rec } f(x) = e_1 \text{ in } e_2), e \rrbracket \rightsquigarrow T} \quad \frac{\llbracket e_2, e \rrbracket \rightsquigarrow T}{\llbracket (\text{let rec } f(x) = e_1 \text{ in } e_2), e \rrbracket \rightsquigarrow T} \quad \frac{\llbracket e_1^{\ell_1}, e_2^{\ell_2} \rrbracket \rightsquigarrow \{\text{Modify}(\ell_1, \alpha_e(e_2))\}}{\text{(the other remaining cases)}}
\end{array}$$

Figure 2. Inference rules for extracting edit action templates for given two expressions.

reachable at ℓ as candidates for the hole. This is because we may copy function definitions in correct programs (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 program 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:

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

where $V|_\ell$ is the set of variables in V that may take the same value reachable at label ℓ . Formally,

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

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.

Algorithm 1 The CAFE Algorithm

Input: A buggy program P_b
Input: A set of correct programs \mathcal{P}_c ,
Input: 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 σ
- 5: Derive \mathcal{M} from Δ ▷ $\mathcal{M} : \text{Fld} \rightarrow \text{Fld}$
- 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** $n \leq |\mathcal{A}|$

2.2.3 Overall Algorithm

Putting all together, Algo. ?? depicts the CAFE algorithm. We first perform the path-sensitive OCFA on all the programs and obtain the result σ (line ??). Then, we derive calling contexts from the analysis result (line ??). From the calling contexts, we obtain the matching function \mathcal{M} that maps each function in the buggy program P_b to a function in a correct program that is most likely to be useful for repair (line ??). Using the matching function, we collect repair templates (line ??). By instantiating the templates, we obtain a set of edit actions that can be applicable to P_b (lines ?? – ??). The main loop (lines ?? – ??) 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 ??). We apply each edit script into P_b (line ??) and check if the program has been fixed based on the given test cases (line ??). If we have fixed the program, we return it as a final result (line ??). Otherwise, we increase n by 1 (line ??) and repeat the main loop.

The algorithm finds a *minimal* edit script that corrects the given buggy program.

Definition 2.2 (Minimality). Given an incorrect program 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 .

2.2.4 Optimizations

When generating edit actions of kind Define that add new function definitions into a target buggy program, 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 targeting different labels.