

Abstract Parsing for Two-staged Languages with Concatenation

Soonho Kong, Wontae Cho, Kwangkeun Yi
{soon,wtchoi,kwang}@ropas.snu.ac.kr
Seoul National University

June 9, 2009

Abstract

This article, based on Doh, Kim, and Schmidt’s “abstract parsing” technique, presents an abstract interpretation for statically checking the syntax of generated code in two-staged programs. Abstract parsing is a static analysis technique for checking the syntax of generated strings. We adopt this technique for two-staged programming languages and formulate it in the abstract interpretation framework. We parameterize our analysis with the abstract domain so that one can choose the abstract domain as long as it satisfies the condition we provide. We also present an instance of the abstract domain, namely an abstract parse stack and its widening with k -cutting.

1 Introduction

1.1 Motivation

For programs that generate and run programs during execution, statically checking the program safety is a challenge. We need to check the safety of generated programs as well as that of the immediate target program. Checking the safety must include checking the programs resulting from evaluating programs.

The semantic safety of such multi-staged programs can be achieved in part by a static type system as reported in [4, 15, 18, 25]. A sound static type system assures that program as data as well as the immediate target program will not have a type error during their executions.

In such a static type system, syntactic errors in the generated code are not an issue. The considered target language is such that primitive code fragments and their compositions are always syntactically correct.

In reality though (as in most web-programming or scripting languages such as PHP, Python, Ruby, Perl, and Javascript), if code is represented as a string and code composition is achieved by string concatenation, syntactically checking the code string value is the foremost issue in the static safety check of such multi-staged programs.

Recently, Doh, Kim, and Schmidt reported a powerful technique called “abstract parsing” [16] that statically analyzes the string values from programs. In abstract parsing, to statically check the generated string is to simulate the parsing actions for possible string values.

In this paper, we report a formalization of abstract parsing in the abstract interpretation framework [10, 11, 12] for two-staged programming languages. Our contribution is to lay a basis to expose the power and, if any, limitation of abstract parsing as static analysis for multi-staged languages.

$$\begin{array}{l|l}
 x = \text{'a'} & X_0 = a \\
 l = \text{'['} & L = [\\
 r = \text{']'} & R =] \\
 x = l . x & X_1 = L.X_0
 \end{array}$$

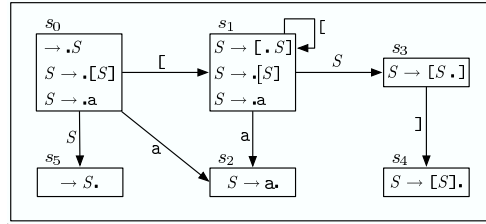
Figure 1: Example program (left) and its data-flow equations (right)

1.2 Abstract Parsing

We here review the abstract parsing idea of [16]. Suppose we want to check that strings generated by the program in Figure 1 conform to the following grammar.

$$S \rightarrow a \mid [S]$$

Abstract parsing derives data-flow equations from the program as in Figure 1. The equation variables are treated as functions that map an input parse state to an output parse stack. They are solved using the goto controller of an LR parser for the grammar, shown in Figure 2.

Figure 2: Goto controller of the LR(0) parser for $S \rightarrow a \mid [S]$. (from [16])

Suppose we want to check that X_1 will accept strings of the target grammar. The analysis starts with $X_1(s_0)$ where s_0 is an initial parse state. To solve

$$X_1(s_0) = (L.X_0)(s_0),$$

the analysis first computes $L(s_0)$. With the state s_0 and the token “[” the goto controller returns $goto(s_0, [) = s_1$. Having $L(s_0) = s_1$, the analysis computes $X_0(s_1)$. After consuming the token “a” and moving to the parse state s_2 , parser reduces with $S \rightarrow a$ and moves the parse state back to s_1 . Then $goto(s_1, S) = s_3$ yields $X_0(s_1) = s_3$. Therefore we have

$$X_1(s_0) = s_3s_1$$

and the analysis concludes that X_1 has a string unacceptable for the grammar because state s_3 is not the accept state.

1.3 Contribution

- We formulate this abstract parsing idea in the abstract interpretation framework for two-staged languages with concatenation. By this formulation we can see what approximations are involved in abstract parsing and what limitations (as a static analysis) to expect from the abstract parsing technique.

Based on the abstract interpretation framework, we present a concise and elegant perspective on the core idea of abstract parsing. In the original work [16], code is abstracted

into the parse stack and the special operator “*” is needed to handle string concatenation. In our formulation, however, we abstract code into a function which maps an input parse stack to an output parse stack. Code concatenation is handled simply by function composition.

- We generalize the abstract-parsing abstract interpretation, as usual, by parameterizing the abstract domain of parse stacks.

This generalization separates the core idea and its implementation of abstract parsing. By choosing an appropriate abstract domain, one can control the analysis precision and cost.

1.4 Organization

Section 2 introduces the syntax and semantics of our target two-staged language with concatenation. Section 3 presents concrete parsing semantics with $LR(k)$ parsing. Section 4 presents abstract parsing semantics and its parameterized framework. Section 5 presents a concrete example of the abstract domain which can be used to instantiate the framework. Section 6 reviews related work and Section 7 concludes.

2 Two-staged Language

We consider a two-staged language with concatenation. The language is an imaginary, first-order language whose only value is code. The language is minimal, so as not to distract our focus on formalizing the abstract parsing method. For example, loops and conditional jumps are without the condition expression, for which abstract interpretation anyway considers all iterations and all branches.

2.1 Syntax and Semantics

A program is an expression e :

$$e \in Exp ::= x \mid \mathbf{let} \ x \ e_1 \ e_2 \mid \mathbf{or} \ e_1 \ e_2 \mid \mathbf{re} \ x \ e_1 \ e_2 \ e_3 \mid 'f$$

An expression can contain code fragments f :

$$f \in Frag ::= x \mid \mathbf{let} \mid \mathbf{or} \mid \mathbf{re} \mid (\mid) \mid f_1.f_2 \mid ,e$$

Operational semantics of the target language is shown in Figure 3. Expression $\mathbf{or} \ e_1 \ e_2$ is for branches. It could be the value of e_1 or the value of e_2 . Expression $\mathbf{re} \ x \ e_1 \ e_2 \ e_3$ is for loops. Variable x has the value of e_1 as its initial value. Loop body e_2 is iterated ≥ 0 times. The result of each iteration e_2 will be bound to x in e_2 for next iteration or in e_3 for the result of the loop. Backquote form $'f$ is for code fragment f . We construct the fragment by using the following tokens: variables, \mathbf{let} , \mathbf{or} , \mathbf{re} , $($, and $)$. Compound fragment $f_1.f_2$ concatenates two code fragments f_1 and f_2 . Comma fragment $,e$ first evaluates e then substitutes its result code value for itself. Note that the meaning of $'f$ and $,e$ is the same as in LISP’s quasi-quotation system.

2.2 Example Program

In our language, it is possible to write a program generating mal-formed code. For instance, the following program generates “a b” (after zero iterations), “or a b” (after one iteration), “or or a b” (after two iterations), and so on.

```
re x 'a '(or . ,x ) '(,x . b)
```

$$\begin{array}{c}
\sigma \in Env = Var \rightarrow Code \\
v \in Code = Token\ sequence \\
e \in Exp \\
f \in Frag
\end{array}$$

$$\boxed{\sigma \vdash^0 e \Rightarrow v}$$

$$\frac{}{\sigma \vdash^0 x \Rightarrow \sigma(x)} \quad \text{(variable)}$$

$$\frac{\sigma \vdash^0 e_1 \Rightarrow v \quad \sigma[x \mapsto v] \vdash^0 e_2 \Rightarrow v'}{\sigma \vdash^0 \mathbf{let}\ x\ e_1\ e_2 \Rightarrow v'} \quad \text{(let binding)}$$

$$\frac{\sigma \vdash^0 e_1 \Rightarrow v}{\sigma \vdash^0 \mathbf{or}\ e_1\ e_2 \Rightarrow v} \quad \frac{\sigma \vdash^0 e_2 \Rightarrow v}{\sigma \vdash^0 \mathbf{or}\ e_1\ e_2 \Rightarrow v} \quad \text{(branch)}$$

$$\frac{\sigma \vdash^0 e_1 \Rightarrow v \quad \sigma[x \mapsto v] \vdash^0 \mathbf{loop}\ x\ e_2\ e_3 \Rightarrow v'}{\sigma \vdash^0 \mathbf{re}\ x\ e_1\ e_2\ e_3 \Rightarrow v'} \quad \text{(loop)}$$

$$\frac{\sigma \vdash^0 e_2 \Rightarrow v \quad \sigma[x \mapsto v] \vdash^0 \mathbf{loop}\ x\ e_2\ e_3 \Rightarrow v'}{\sigma \vdash^0 \mathbf{loop}\ x\ e_2\ e_3 \Rightarrow v'}$$

$$\frac{\sigma \vdash^0 e_3 \Rightarrow v}{\sigma \vdash^0 \mathbf{loop}\ x\ e_2\ e_3 \Rightarrow v}$$

$$\frac{\sigma \vdash^1 f \Rightarrow v}{\sigma \vdash^0 \text{'}\ f \Rightarrow v} \quad \text{(back quote)}$$

$$\boxed{\sigma \vdash^1 f \Rightarrow v}$$

$$\frac{}{\sigma \vdash^1 x \Rightarrow x} \quad \frac{}{\sigma \vdash^1 \mathbf{let} \Rightarrow \mathbf{let}} \quad \text{(token)}$$

$$\frac{}{\sigma \vdash^1 \mathbf{or} \Rightarrow \mathbf{or}} \quad \frac{}{\sigma \vdash^1 \mathbf{re} \Rightarrow \mathbf{re}}$$

$$\frac{}{\sigma \vdash^1 (\Rightarrow (} \quad \frac{}{\sigma \vdash^1) \Rightarrow)}$$

$$\frac{\sigma \vdash^1 f_1 \Rightarrow v_1 \quad \sigma \vdash^1 f_2 \Rightarrow v_2}{\sigma \vdash^1 f_1.f_2 \Rightarrow v_1v_2} \quad \text{(concatenation)}$$

$$\frac{\sigma \vdash^0 e \Rightarrow v}{\sigma \vdash^1 ,e \Rightarrow v} \quad \text{(comma)}$$

Figure 3: Operational semantics of the target language.

Only “or a b” is correct and the rest of them have a syntax error.

However the following program generates “a” (after zero iterations), “or a b” (after one iteration), “or (or a b) b” (after two iterations), and so on,

```
re x 'a '(or . ,x . b) x
```

and all of them are syntactically correct.

2.3 Collecting Semantics and its Abstraction Plan

The collecting semantics of the language is defined as follows. This semantics is the natural set extension for the sets of environments. The *fix* operator is the usual least fixpoint operator to capture all the iteration results from loops.

$$\begin{aligned}
& \text{Code} = \text{Token sequence} \\
& \sigma \in \text{Env} = \text{Var} \rightarrow \text{Code} \\
& \llbracket e \rrbracket^0 \in 2^{\text{Env}} \rightarrow 2^{\text{Code}} \\
& \llbracket f \rrbracket^1 \in 2^{\text{Env}} \rightarrow 2^{\text{Code}} \\
& \llbracket x \rrbracket^0 \Sigma = \{ \sigma(x) \mid \sigma \in \Sigma \} \\
& \llbracket \text{let } x \ e_1 \ e_2 \rrbracket^0 \Sigma = \bigcup_{\sigma \in \Sigma} \bigcup_{c \in \llbracket e_1 \rrbracket^0 \{ \sigma \}} \llbracket e_2 \rrbracket^0 \{ \sigma[x \mapsto c] \} \\
& \llbracket \text{or } e_1 \ e_2 \rrbracket^0 \Sigma = \llbracket e_1 \rrbracket^0 \Sigma \cup \llbracket e_2 \rrbracket^0 \Sigma \\
& \llbracket \text{re } x \ e_1 \ e_2 \ e_3 \rrbracket^0 \Sigma = \bigcup_{\sigma \in \Sigma} \llbracket e_3 \rrbracket^0 \{ \sigma[x \mapsto c] \mid c \in \\
& \quad \text{fix } \lambda C. \llbracket e_1 \rrbracket^0 \{ \sigma \} \cup \llbracket e_2 \rrbracket^0 \{ \sigma[x \mapsto c'] \mid c' \in C \} \} \\
& \llbracket ' f \rrbracket^0 \Sigma = \llbracket f \rrbracket^1 \Sigma \\
& \llbracket x \rrbracket^1 \Sigma = \{ x \} \\
& \llbracket \text{let} \rrbracket^1 \Sigma = \{ \text{let} \} \\
& \llbracket \text{or} \rrbracket^1 \Sigma = \{ \text{or} \} \\
& \llbracket \text{re} \rrbracket^1 \Sigma = \{ \text{re} \} \\
& \llbracket (\rrbracket^1 \Sigma = \{ (\} \\
& \llbracket) \rrbracket^1 \Sigma = \{) \} \\
& \llbracket f_1 . f_2 \rrbracket^1 \Sigma = \bigcup_{\sigma \in \Sigma} \{ xy \mid x \in \llbracket f_1 \rrbracket^1 \{ \sigma \} \wedge y \in \llbracket f_2 \rrbracket^1 \{ \sigma \} \} \\
& \llbracket , e \rrbracket^1 \Sigma = \llbracket e \rrbracket^0 \Sigma
\end{aligned}$$

From the collecting semantics above, we derive a series of abstract semantics. From the collecting semantics' semantic domain

$$2^{\text{Var} \rightarrow \text{Code}} \rightarrow 2^{\text{Code}},$$

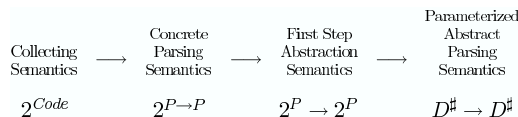
the powerset environment $2^{\text{Var} \rightarrow \text{Code}}$ is abstracted into $\text{Var} \rightarrow 2^{\text{Code}}$, i.e., the semantic domain becomes

$$(\text{Var} \rightarrow 2^{\text{Code}}) \rightarrow 2^{\text{Code}}.$$

Now the abstraction of 2^{Code} becomes the essential part of the abstract interpretation design. Before we abstract 2^{Code} , we formulate a code fragment as a function that maps a parse stack

to a parse stack. We call this formulation “concrete parsing” (Section 3). That is, 2^{Code} becomes $2^{P \rightarrow P}$ (where P is the set of parse stacks). Then we abstract $2^{P \rightarrow P}$ into $2^P \rightarrow 2^P$ (Section 4). Lastly, we present an abstract-parsing abstract interpretation that parameterizes an abstract domain $D^\#$ of 2^P .

In summary, this series of abstraction steps for the value domain in the semantics is:



3 Concrete Parsing

3.1 Analyze-and-parse Strategy

We take the analyze-*and*-parse strategy in abstract parsing [16] into our semantics. The semantics simulates the parsing operations. It is compared to the analyze-*then*-parse strategy which analyzes the program, calculates approximated set of code, then parses them.

Analyze-*and*-parse strategy is more efficient than analyze-*then*-parse strategy as reported in [16]. In analyze-*then*-parse strategy, code is abstracted into a grammar. Then it checks whether the abstracted grammar is included in the reference grammar or not. However, grammar inclusion check is more expensive than parsing. In addition, analyze-*and*-parse directly computes parsing information without approximating the code into the grammar.

We formulate the analyze-*and*-parse strategy in our semantics. The parsing domain is constructed as an abstract domain where code is abstracted into parsing information. We abstract the parsing domain into an abstract parsing domain to control the precision and cost of analysis and to make sure the analysis terminates.

Since our semantics uses an LR(k) parser as a component, it is essential to review its key concepts.

3.2 LR Parsing

The LR(k) parsing technique [1] is an efficient way to determine whether the string conforms to the given grammar or not. An LR parser is a state machine which consists of a parse stack, an action table, and a goto table. The set of parse states $\Sigma = \{s_1, s_2, \dots, s_n\}$ is defined by parser generator from the given grammar. Parse stack $p \in P = \Sigma^+$ is a sequence of parse states which it has been in. Two special parse stacks p_{init} and p_{acc} are defined. Parsing starts with the initial parse stack p_{init} . Successful parsing should stop at the accept parse stack p_{acc} . Otherwise it indicates that the parsed string does not conform to the given grammar. String representation “ $s_{top} \dots s_{bot}$ ” denotes a parse stack whose top state is s_{top} and bottom state is s_{bot} . The action table decides which operation (shift/reduce) to perform from the current state and current token. The goto table determines the state to push after we pop states in the reduce operation.

The process of parsing is a composition of the atomic function $parse_action : P \rightarrow Token \rightarrow P$ which is described in Algorithm 1. It returns the parse stack from the given parse stack p and input token t .

The parsing process $parse : P \rightarrow Token\ sequence \rightarrow P$ is a composition of the $parse_action$.

$$\begin{aligned} parse(p, t_1 \dots t_n) \\ = parse_action(\dots (parse_action(p, t_1)), \dots, t_n) \end{aligned}$$

A parser gets the input code c and returns the parse stack $parse(p_{init}, c)$.

Algorithm 1 *parse_action* algorithm

```

1: procedure parse_action( $p, t$ )
2:    $s_{top} \leftarrow$  the state on top of stack  $p$ 
3:   if  $ACTION[s_{top}, t] = \text{shift } s$  then
4:     push  $s$  onto the stack  $p$ 
5:     return  $p$ 
6:   else if  $ACTION[s_{top}, t] = \text{reduce } A \rightarrow \beta$  then
7:     pop  $|\beta|$  symbol off the stack  $p$ 
8:      $s_{top} \leftarrow$  the state on top of stack  $p$ 
9:     push  $GOTO[s_{top}, A]$  onto the stack  $p$ 
10:    return parse_action( $p, t$ )
11:  end if
12: end procedure

```

3.3 Concrete Parsing Domain : Value $V_P = 2^{P \rightarrow P}$

We define the concrete parsing domain. It is “concrete” in that we use the same parse stack defined in LR(k) parsing.

Because LR(k) parsing computes a parse stack, it is tempting to choose the parse stack $parse(p_{init}, c)$ as an abstraction of the code c . However this setting causes a problem when we handle the concatenation $x.y$. Let $p_x = parse(p_{init}, x)$ and $p_y = parse(p_{init}, y)$ be the resulting parse stacks for the code fragments x and y . The parse stack p for the concatenation $x.y$ is computed as

$$\begin{aligned}
p &= parse(p_{init}, x.y) \\
&= parse(parse(p_{init}, x), y) \\
&= parse(p_x, y).
\end{aligned}$$

However, what we have is $p_y = parse(p_{init}, y)$ not $parse(p_x, y)$. The $parse(p_{init}, y)$ is the parse stack after parsing y from the initial stack not from the p_x stack. We cannot directly compute p from p_x and p_y .

Abstracting the code c into the parse stack transition function $\lambda p.parse(p, c)$ solves the above concatenation problem elegantly. Let f_x and f_y be the parse stack transition functions for code fragments x and y respectively. Then we have

$$\begin{aligned}
f_x &= \lambda p.parse(p, x) \\
f_y &= \lambda p.parse(p, y).
\end{aligned}$$

With the two parse stack transition functions f_x and f_y , we construct the parse stack transition function $f_{x.y}$ as follows.

$$\begin{aligned}
f_{x.y} &= \lambda p.parse(p, x.y) \\
&= \lambda p.parse(parse(p, x), y) \\
&= (\lambda p.parse(p, y)) \circ (\lambda p.parse(p, x)) \\
&= f_y \circ f_x
\end{aligned}$$

3.4 Concrete Parsing Semantics

Using the abstraction from $Code$ to $P \rightarrow P$, the Galois connection $2^{Code} \xrightleftharpoons[\alpha]{\gamma} V_P = 2^{P \rightarrow P}$ is established as follows.

$$\alpha = \lambda C. \{ \lambda p.parse(p, c) \mid c \in C \}$$

$$\gamma = \lambda F. \bigcup_{f \in F} \{c \mid \forall p \in P. parse(p, c) = f(p)\}$$

Concrete parsing semantics is derived from the collecting semantics as follows.

$$\begin{aligned} P &= \text{the set of parse stacks} \\ V_P &= 2^{P \rightarrow P} \\ \sigma \in Env_P &= Var \rightarrow V_P \\ \llbracket e \rrbracket_P^0 &\in Env_P \rightarrow V_P \\ \llbracket f \rrbracket_P^1 &\in Env_P \rightarrow V_P \\ \\ \llbracket x \rrbracket_P^0 \sigma &= \sigma(x) \\ \llbracket \mathbf{let } x \ e_1 \ e_2 \rrbracket_P^0 \sigma &= \llbracket e_2 \rrbracket_P^0 (\sigma[x \mapsto \llbracket e_1 \rrbracket_P^0 \sigma]) \\ \llbracket \mathbf{or } e_1 \ e_2 \rrbracket_P^0 \sigma &= \llbracket e_1 \rrbracket_P^0 \sigma \cup \llbracket e_2 \rrbracket_P^0 \sigma \\ \llbracket \mathbf{re } x \ e_1 \ e_2 \ e_3 \rrbracket_P^0 \sigma &= \llbracket e_3 \rrbracket_P^0 (\sigma[x \mapsto \\ &\quad \mathit{fix } \lambda k. \llbracket e_1 \rrbracket_P^0 \sigma \cup \llbracket e_2 \rrbracket_P^0 (\sigma[x \mapsto k])]) \\ \llbracket \cdot \rrbracket_P^0 \sigma &= \llbracket f \rrbracket_P^1 \sigma \\ \llbracket t \rrbracket_P^1 \sigma &= \{\lambda p. parse_action(p, t)\} \\ \llbracket f_1. f_2 \rrbracket_P^1 \sigma &= \{p_2 \circ p_1 \mid p_1 \in \llbracket f_1 \rrbracket_P^1 \sigma \wedge p_2 \in \llbracket f_2 \rrbracket_P^1 \sigma\} \\ \llbracket \cdot, e \rrbracket_P^1 \sigma &= \llbracket e \rrbracket_P^0 \sigma \end{aligned}$$

4 Abstract Parsing

4.1 First Step Abstraction : Value $V_{\hat{P}} = 2^P \rightarrow 2^P$

First, we abstract the concrete parsing domain $V_P = 2^{P \rightarrow P}$ to $V_{\hat{P}} = 2^P \rightarrow 2^P$ by establishing the Galois connection $V_P \xrightleftharpoons[\alpha]{\gamma} V_{\hat{P}}$ where

$$\begin{aligned} \alpha &= \lambda F. \lambda P. \bigcup_{p \in P} \{f(p) \mid f \in F\} \\ \gamma &= \lambda F. \bigcup \{S \mid \alpha(S) \subseteq f\}. \end{aligned}$$

Then we derive the abstract semantics for this domain as follows.

$$\begin{aligned} \sigma \in Env_{\hat{P}} &= Var \rightarrow V_{\hat{P}} \\ \llbracket e \rrbracket_{\hat{P}}^0 &\in Env_{\hat{P}} \rightarrow V_{\hat{P}} \\ \llbracket f \rrbracket_{\hat{P}}^1 &\in Env_{\hat{P}} \rightarrow V_{\hat{P}} \\ \\ \llbracket x \rrbracket_{\hat{P}}^0 &= \sigma(x) \\ \llbracket \mathbf{let } x \ e_1 \ e_2 \rrbracket_{\hat{P}}^0 \sigma &= \llbracket e_2 \rrbracket_{\hat{P}}^0 (\sigma[x \mapsto \llbracket e_1 \rrbracket_{\hat{P}}^0 \sigma]) \\ \llbracket \mathbf{or } e_1 \ e_2 \rrbracket_{\hat{P}}^0 \sigma &= \llbracket e_1 \rrbracket_{\hat{P}}^0 \sigma \cup \llbracket e_2 \rrbracket_{\hat{P}}^0 \sigma \\ \llbracket \mathbf{re } x \ e_1 \ e_2 \ e_3 \rrbracket_{\hat{P}}^0 \sigma &= \llbracket e_3 \rrbracket_{\hat{P}}^0 (\sigma[x \mapsto \\ &\quad \mathit{fix } \lambda k. \llbracket e_1 \rrbracket_{\hat{P}}^0 \sigma \cup \llbracket e_2 \rrbracket_{\hat{P}}^0 (\sigma[x \mapsto k])]) \\ \llbracket \cdot \rrbracket_{\hat{P}}^0 \sigma &= \llbracket f \rrbracket_{\hat{P}}^1 \sigma \end{aligned}$$

$$\begin{aligned}
\llbracket t \rrbracket_{\hat{P}}^1 \sigma &= \lambda P. \text{Parse_action}(P, t) \\
\llbracket f_1.f_2 \rrbracket_{\hat{P}}^1 \sigma &= \llbracket f_2 \rrbracket_{\hat{P}}^1 \sigma \circ \llbracket f_1 \rrbracket_{\hat{P}}^1 \sigma \\
\llbracket \cdot, e \rrbracket_{\hat{P}}^1 \sigma &= \llbracket e \rrbracket_{\hat{P}}^0 \sigma
\end{aligned}$$

where $\text{Parse_action} : 2^P \rightarrow \text{Token} \rightarrow 2^P$ is the natural set extension of parse_action :

$$\text{Parse_action} = \lambda P. \lambda t. \{ \text{parse_action}(p, t) \mid p \in P \}.$$

The abstract semantic function $\llbracket \cdot \rrbracket_{\hat{P}}^0$ is used to check whether generated code conforms to the grammar. For the given program e , we compute

$$S = \llbracket e \rrbracket_{\hat{P}}^0 \sigma_0 \{ p_{\text{init}} \} : 2^P$$

where $\sigma_0 \in \text{Env}_{\hat{P}}$ is an empty environment. Then we compare S with $\{ p_{\text{acc}} \}$. If they are equal, we conclude that the generated code in the given program conforms to the grammar. Otherwise, the analysis concludes that the program may generate syntactically incorrect code.

4.2 Parameterized Framework

We generalize the analysis by parameterizing abstract domain. Instead of abstracting 2^P (the powerset domain of parse stacks) into a particular domain, we provide conditions which the abstract domain for 2^P should satisfy. Then we define the semantic function on the abstract parsing domain.

Definition 1 (Abstract Parsing Domain). $V^\sharp = D^\sharp \rightarrow D^\sharp$ is an abstract parsing domain if an abstract domain D^\sharp satisfies the following conditions [10, 11, 12].

1. $\langle D^\sharp, \sqsubseteq, \sqcup, \perp_{D^\sharp} \rangle$ is a CPO (Complete Partial Order).
2. Powerset domain of concrete parse stack 2^P and its abstract domain D^\sharp are Galois connected via $\alpha_{2^P \rightarrow D^\sharp}$ and $\gamma_{D^\sharp \rightarrow 2^P}$.
3. $\text{Parse_action}^\sharp : D^\sharp \rightarrow \text{Token} \rightarrow D^\sharp$ is a sound approximation of Parse_action . That is, we have

$$\begin{aligned}
\forall P \in 2^P. \forall t \in \text{Token}. \\
\alpha_{2^P \rightarrow D^\sharp}(\text{Parse_action}(P, t)) \sqsubseteq \text{Parse_action}^\sharp(\alpha_{2^P \rightarrow D^\sharp}(P), t)
\end{aligned}$$

The partial order \sqsubseteq_{V^\sharp} and join operator \sqcup_{V^\sharp} are defined pointwisely.

Definition 2. Semantic function $\llbracket \cdot \rrbracket_{D^\sharp}$ on the abstract parsing domain V^\sharp is defined as follows.

$$\begin{aligned}
\sigma \in \text{Env}_{D^\sharp} &= \text{Var} \rightarrow V^\sharp \\
\llbracket e \rrbracket_{D^\sharp}^0 &\in \text{Env}_{D^\sharp} \rightarrow V^\sharp \\
\llbracket f \rrbracket_{D^\sharp}^1 &\in \text{Env}_{D^\sharp} \rightarrow V^\sharp \\
\llbracket x \rrbracket_{D^\sharp}^0 \sigma &= \sigma(x) \\
\llbracket \text{let } x \ e_1 \ e_2 \rrbracket_{D^\sharp}^0 \sigma &= \llbracket e_2 \rrbracket_{D^\sharp}^0 (\sigma[x \mapsto \llbracket e_1 \rrbracket_{D^\sharp}^0 \sigma]) \\
\llbracket \text{or } e_1 \ e_2 \rrbracket_{D^\sharp}^0 \sigma &= \llbracket e_1 \rrbracket_{D^\sharp}^0 \sigma \sqcup \llbracket e_2 \rrbracket_{D^\sharp}^0 \sigma \\
\llbracket \text{re } x \ e_1 \ e_2 \ e_3 \rrbracket_{D^\sharp}^0 \sigma &= \llbracket e_3 \rrbracket_{D^\sharp}^0 (\sigma[x \mapsto \\
&\quad \text{fix } \lambda k. \llbracket e_1 \rrbracket_{D^\sharp}^0 \sigma \sqcup \llbracket e_2 \rrbracket_{D^\sharp}^0 (\sigma[x \mapsto k])])
\end{aligned}$$

$$\begin{aligned}
\llbracket f \rrbracket_{D^\sharp}^0 \sigma &= \llbracket f \rrbracket_{D^\sharp}^1 \sigma \\
\llbracket t \rrbracket_{D^\sharp}^1 \sigma &= \lambda D. \text{Parse_action}^\sharp(D, t) \\
\llbracket f_1.f_2 \rrbracket_{D^\sharp}^1 \sigma &= \llbracket f_2 \rrbracket_{D^\sharp}^1 \sigma \circ \llbracket f_1 \rrbracket_{D^\sharp}^1 \sigma \\
\llbracket ., e \rrbracket_{D^\sharp}^1 \sigma &= \llbracket e \rrbracket_{D^\sharp}^0 \sigma
\end{aligned}$$

Theorem 1 shows that $\llbracket \cdot \rrbracket_{D^\sharp}$ is a sound approximation of $\llbracket \cdot \rrbracket_{\hat{P}}$.

Theorem 1. *Semantic function $\llbracket \cdot \rrbracket_{D^\sharp}$ on the abstract parsing domain V^\sharp is a sound approximation of $\llbracket \cdot \rrbracket_{\hat{P}}$. That is, we have*

$$\begin{aligned}
\forall e \in \text{Exp}. \forall \sigma \in \text{Env}_{\hat{P}}. \\
\alpha_{V_{\hat{P}} \rightarrow V^\sharp}(\llbracket e \rrbracket_{\hat{P}} \sigma) \sqsubseteq \llbracket e \rrbracket_{D^\sharp}(\alpha_{\text{Env}_{\hat{P}} \rightarrow \text{Env}_{D^\sharp}}(\sigma))
\end{aligned}$$

where

$$\begin{aligned}
\alpha_{V_{\hat{P}} \rightarrow V^\sharp} &= \lambda F. \lambda D. \alpha_{2^P \rightarrow D^\sharp}(F(\gamma_{D^\sharp \rightarrow 2^P}(D))) \\
\alpha_{\text{Env}_{\hat{P}} \rightarrow \text{Env}_{D^\sharp}} &= \lambda \sigma. \lambda x. \alpha_{V_{\hat{P}} \rightarrow V^\sharp}(\sigma(x)).
\end{aligned}$$

Proof. By structural induction on e with the conditions that the abstract parsing domain D^\sharp should satisfy. \square

5 Instantiation : Powerset Domain of Abstract Parse Stack \hat{P} with k -cutting

In this section, we introduce an instance of an abstract domain D^\sharp , which is \hat{D} , a powerset domain of abstract parse stack \hat{P} and its widening with k -cutting. Thus the abstract parsing domain is

$$\hat{V} = \hat{D} \rightarrow \hat{D}.$$

First, we define abstract parse stack \hat{P} which is an abstraction of concrete parse stack P . The abstract domain \hat{D} is constructed with abstract parse stack \hat{P} . By establishing the Galois connection between 2^P and \hat{D} and defining $\widehat{\text{Parse_action}}$, we show that \hat{D} is an instance of D^\sharp . Finally, a widening operator is defined using k -cutting to guarantee the termination of analysis.

5.1 Abstract Parse Stack \hat{P}

Cut parse stack \bar{P} is introduced to represent a parse stack which has been cut and only maintains top n states. Special state ‘-’ $\notin \Sigma$ at the bottom of cut parse stack \bar{p} indicates that it has been cut.

$$\bar{P} = \{p \cdot - \mid p \in \Sigma^*\}$$

Abstract parse stack \hat{P} is defined as a union of concrete parse stack P and cut parse stack \bar{P} .

$$\hat{P} = P \cup \bar{P}$$

Given an abstract parse stack ρ , the function $\text{prefix} : \hat{P} \rightarrow \Sigma^*$ yields the longest prefix of ρ which does not contain special state ‘-’.

$$\text{prefix}(\rho) = \begin{cases} \rho & \text{if } \rho \in P \\ s_1 \dots s_n & \text{if } \rho = s_1 \dots s_n - \\ \epsilon \text{ (empty string)} & \rho = - \end{cases}$$

Then we define partial order \sqsubseteq on \hat{P} as follows.

$$\rho_1 \sqsubseteq_{\hat{P}} \rho_2 \stackrel{\text{def}}{=} \text{prefix}(\rho_1) \text{ starts with } \text{prefix}(\rho_2)$$

5.2 Abstract Domain : \hat{D}

Abstract domain \hat{D} with its partial order \sqsubseteq and join \sqcup is defined as follows.

$$\begin{aligned}\hat{D} &= \{norm(\hat{d}) \mid \hat{d} \in 2^{\hat{P}}\} \\ \hat{d}_1 \sqsubseteq \hat{d}_2 &\stackrel{\text{def}}{=} \forall \rho_1 \in \hat{d}_1. \exists \rho_2 \in \hat{d}_2. \rho_1 \sqsubseteq_{\hat{P}} \rho_2 \\ \hat{d}_1 \sqcup \hat{d}_2 &\stackrel{\text{def}}{=} norm(\hat{d}_1 \cup \hat{d}_2)\end{aligned}$$

Normalization function $norm$ is defined by

$$norm(\hat{d}) = \{\rho \in \hat{d} \mid \forall \rho' \in \hat{d}. \rho \not\sqsubseteq_{\hat{P}} \rho'\}$$

to ensure that elements in $norm(\hat{d})$ are the maximal elements of \hat{d} . It is necessary to eliminate non-maximal elements for binary relation \sqsubseteq to be anti-symmetric and to be a partial order. For instance, $\{-, s_1 s_0-\}$ contains non-maximal element $\{s_1 s_0-\}$ since we have $s_1 s_0- \sqsubseteq -$. If we allow $\{s_1 s_0-, -\}$ in \hat{D} without normalizing it into $\{-\}$, we have

$$\begin{aligned}(\{s_1 s_0-, -\} \sqsubseteq \{-\}) \wedge (\{-\} \sqsubseteq \{s_1 s_0-, -\}) \\ \not\Rightarrow \{s_1 s_0-, -\} = \{-\}.\end{aligned}$$

Then partial order \sqsubseteq on \hat{D} becomes preorder since \sqsubseteq is not anti-symmetric anymore.

5.3 Abstract Domain \hat{D} as an Instance of D^\sharp

Abstract domain \hat{D} is an instance of D^\sharp . To verify this, we need to show that \hat{D} satisfies the conditions in Definition 1.

1. $\langle \hat{D}, \sqsubseteq, \sqcup, \phi \rangle$ is a CPO by definition of \hat{D} .
2. To establish a Galois connection between 2^P and \hat{D} , we first define the function $expand : \hat{P} \rightarrow 2^P$ as follows.

$$expand(\rho) = \begin{cases} \{\rho\} & \text{if } \rho \in P \\ \{prefix(\rho) \cdot p \mid p \in P\} & \text{if } \rho \in \bar{P} \end{cases}$$

It expands an abstract parse stack ρ to all the concrete parse stacks which ρ can represent. For instance, $expand(s_1-) = \{s_1 s_0, s_1 s_2, \dots\}$. Note that we have $\forall p \in expand(\rho). p \sqsubseteq_{\hat{P}} \rho$ and therefore we get $expand(\rho) \sqsubseteq_{\hat{D}} \{\rho\}$.

The function $Expand : \hat{D} \rightarrow 2^P$ is also defined as the natural set extension of $expand$ function by

$$Expand(\hat{d}) = \bigcup_{\rho \in \hat{d}} expand(\rho).$$

From the property of $expand$, it is clear that $Expand(\hat{d}) \sqsubseteq_{\hat{D}} \hat{d}$.

Using the $Expand$ function, we define the Galois connection $2^P \xrightarrow[\alpha]{\gamma} \hat{D}$ where

$$\begin{aligned}\alpha &= id \\ \gamma &= \lambda \hat{d}. Expand(\hat{d})\end{aligned}$$

Proof. To prove α and γ constitute a Galois connection, it is sufficient to show that the following properties hold.

- (a) Trivially, $\alpha = id$ is monotone.
- (b) Monotonicity of γ is immediate from the monotonicity of $expand$.
- (c) We show that $\forall \hat{d} \in \hat{D}. \alpha \circ \gamma(\hat{d}) \sqsubseteq \hat{d}$.

$$\begin{aligned} \alpha \circ \gamma(\hat{d}) &= \gamma(\hat{d}) \\ &= Expand(\hat{d}) \\ &\sqsubseteq \hat{d} \end{aligned}$$

- (d) We show that $\forall P \in 2^P. P \subseteq \gamma \circ \alpha(P)$.

$$\begin{aligned} \gamma \circ \alpha(P) &= \gamma(P) \\ &= Expand(P) \\ &= \bigcup_{p \in P} expand(p) \\ &= \bigcup_{p \in P} p \\ &= P \end{aligned}$$

□

- 3. We first define $parse_action : \hat{P} \rightarrow Token \rightarrow \hat{P}$ as in Algorithm 2. It is a modified version of the $parse_action$ algorithm. The only modifications are adding lines 2 – 4 and lines 12 – 14 to handle $\bar{p} \in \bar{P}$ because the action and goto tables do not have an entry for the special state ‘-’. Note that $\forall p \in P. \forall t \in Token. parse_action(p, t) = parse_action(p, t)$.

Algorithm 2 $parse_action$ algorithm

```

1: procedure  $parse\_action(\rho, t)$ 
2:   if  $\rho = -$  then
3:     return  $\rho$ 
4:   end if
5:    $s_{top} \leftarrow$  the state on top of stack  $\rho$ 
6:   if  $ACTION[s_{top}, t] = \text{shift } s$  then
7:     push  $s$  onto the stack  $\rho$ 
8:     return  $\rho$ 
9:   else if  $ACTION[s_{top}, t] = \text{reduce } A \rightarrow \beta$  then
10:    pop  $|\beta|$  symbol off the stack  $\rho$ 
11:     $s_{top} \leftarrow$  the state on top of stack  $\rho$ 
12:    if  $s_{top} = -$  then
13:      return  $\rho$ 
14:    end if
15:    push  $GOTO[s_{top}, A]$  onto the stack  $\rho$ 
16:    return  $parse\_action(\rho, t)$ 
17:   end if
18: end procedure

```

$Parse_action : \hat{D} \rightarrow Token \rightarrow \hat{D}$ is defined as the natural set extension of $parse_action$ as follows.

$$Parse_action = \lambda \hat{d}. \lambda t. \{ parse_action(\rho, t) \mid \rho \in \hat{d} \}$$

$Parse_action$ is a sound approximation of $Parse_action$.

Proof. For all $P \in 2^P$ and $t \in \text{Token}$, we have

$$\begin{aligned} \alpha_{2^P \rightarrow \hat{D}}(\text{Parse_action}(P, t)) &= \text{Parse_action}(P, t) \\ &= \{\text{parse_action}(p, t) \mid p \in P\} \\ &= \{\widehat{\text{parse_action}}(p, t) \mid p \in P\} \\ &= \widehat{\text{Parse_action}}(P, t) \\ &= \widehat{\text{Parse_action}}(\alpha_{2^P \rightarrow \hat{D}}(P), t). \end{aligned}$$

□

5.4 Widening Operator on \hat{D}

The termination of the analysis parameterized with \hat{D} is not guaranteed because \hat{D} has infinite height. Instead of limiting the height of the domain, we use the widening method to achieve termination.

We define an operator $\nabla_{\hat{D}} : \hat{D} \times \hat{D} \rightarrow \hat{D}$ such that

$$A \nabla_{\hat{D}} B = \{\text{norm}(\text{cut}_k(\rho)) \mid \rho \in A \cup B\}$$

where cut_k is defined by

$$\text{cut}_k(\rho) = \begin{cases} \rho & \text{if } |\rho| \leq k \\ s_1 \dots s_{k-1}^- & \text{if } \rho = s_1 \dots s_{k-1} s_k \dots s_n. \end{cases}$$

Theorem 2 shows that $\nabla_{\hat{D}}$ is a widening operator on \hat{D} .

Theorem 2. $\nabla_{\hat{D}} : \hat{D} \times \hat{D} \rightarrow \hat{D}$ is a widening operator.

Proof. We have to check if $\nabla_{\hat{D}}$ operator satisfies the widening conditions [10, 11, 12]:

- (i) $\forall x, y \in \hat{D}. (x \sqsubseteq x \nabla_{\hat{D}} y) \wedge (y \sqsubseteq x \nabla_{\hat{D}} y)$.
- (ii) for all increasing chains $x_0 \sqsubseteq x_1 \dots$, the increasing chain defined by $y_0 = x_0, \dots, y_{i+1} = y_i \nabla_{\hat{D}} x_{i+1}, \dots$ is not strictly increasing.

To prove (i), we observe that $\rho \sqsubseteq_{\hat{P}} \text{cut}_k(\rho)$ by definition of $\sqsubseteq_{\hat{P}}$ and cut_k . Using this we get

$$\begin{aligned} x &= \{\rho \mid \rho \in x\} \\ &\sqsubseteq \{\text{cut}_k(p) \mid p \in x\} \\ &\sqsubseteq \{\text{cut}_k(p) \mid p \in x \cup y\} \\ &= x \nabla_{\hat{D}} y. \end{aligned}$$

Proof for $y \sqsubseteq x \nabla_{\hat{D}} y$ is analogous.

To prove (ii), we observe that the range of $\nabla_{\hat{D}}$ operator is the finite set $\hat{P}' = \{\rho \in \hat{P} \mid |\rho| \leq k\}$. For all $i \geq 1$, we have $y_i \in \hat{P}'$. Therefore the increasing chain y_0, \dots, y_n, \dots is not strictly increasing. □

Using the widening operator $\nabla_{\hat{D}}$, we define a widening operator for $\hat{V} = \hat{D} \rightarrow \hat{D}$, such that

$$f \nabla_{\hat{V}} g = \lambda \hat{d}. \begin{cases} f(\hat{d}) \nabla_{\hat{D}} g(\hat{d}) & \text{if } \forall \rho \in \hat{d}. |\rho| \leq l \\ \{-\} & \text{otherwise.} \end{cases}$$

where l is a constant to bound the set of meaningful entries finitely fixed in the resulting function. For those finite meaningful entries, we widen their images using the widening $\nabla_{\hat{D}}$. This $\nabla_{\hat{V}}$ operator is a widening operator and the analysis using this widening always terminates.

The analysis using the widening gives better precision than the one using the domain with finite height, for instance, $\hat{D}' = \{\hat{d} \in \hat{D} \mid |\hat{d}| \leq k\}$. The widening approach only restricts the length of the parse stack at the loop head. In the loop body, it allows arbitrary length of parse stacks. Let's consider the following program.

```
re x 'y
  '(let . z . ,x
      . or . z . z)
x
```

Using widening approach, it is sufficient to have cut threshold $k = 2$ to analyze the program without precision loss. In the loop body, the length of parse stack increases from one to seven. At the end of the loop body, however, the parser reduces the parse stack and its length becomes two. Since the maximum length of the parse stack at the loop head is two and the widening only occurs at the loop head, cut threshold $k = 2$ is enough.

However if we restrict the domain and use \hat{D} with $k = 2$, we will lose precision while analyzing the loop body. A parse stack with special state '-' at the bottom will be introduced. To be as precise as the one using the widening approach, we need to increase k to seven.

6 Related Work

In this section, we discuss several areas of related work: multi-staged languages, string analysis, string verification, and code generation.

Multi-staged Languages According to Sheard [24], there are three representations for code in multi-staged languages – quasi-quote, algebraic data type, and strings.

Quasi-quote and algebraic data type representation preserve internal structure of the code and syntactic safety can be guaranteed.

Regarding semantic safety check, various static type systems for multi-staged languages are reported [2, 4, 5, 6, 13, 14, 15, 21, 22, 23, 25]. The static type system has been matured [18] to support almost all of the Lisp-like multi-staged programming practices.

However, using a string representation loses track of internal structure of the code and opens the syntactic safety problem. Even if string representation is strongly discouraged in the research community [24], it is used in most web-programming or scripting languages such as PHP, Python, Ruby, Perl, and Javascript.

String Analysis Before Doh, Kim, and Schmidt's abstract parsing technique [16], string analysis works [8, 19, 7] were all "analyze-then-parse" techniques.

Christensen et al. [8] presents an approach which abstracts the set of generated strings in a program into a regular grammar and performs a grammar inclusion check between the regular grammar and the reference grammar, which is a context-free grammar (CFG). Precision loss occurs when the generated strings are abstracted into the regular grammar.

Minamide [19] takes a slightly different approach. The context-free reference grammar is abstracted into a regular grammar by restricting the nesting depth of generated strings. The set of generated strings is abstracted into a context-free grammar. Then a grammar inclusion check is performed between the regular grammar and the context-free grammar. This approach is practical for HTML/XML document analysis because their nesting depth is usually restricted. However, it is not applicable in general.

Choi et al. [7] presents abstract interpretation based string analysis which uses a heuristic widening method and overcomes the difficulties of handling heap variables and context sensitivity. However, using a regular grammar as an abstraction results in precision loss.

String Verification String analysis techniques are used as a basis in string verification. Christodorescu et al. [9] and Wassermann et al. [26] present string verifiers based on the Christensen et al.’s string analysis [8]. Wassermann and Su [27] and Minamide and Tozawa [20] present string verifiers based on the Minamide’s string analysis [19].

Code Generation There are many researches to achieve syntax safety of generated code other than multi-staged languages area.

Repleo [3] shows a template engine which generates syntax safe code. Having the grammar for the template, it statically checks the template and its sub-templates to detect possible syntactic errors in the generated code. However, it still requires evaluating the template with the model to guarantee the generated code is syntax safe.

Engler et al. presents ‘C [17], an extension of ANSI C which allows dynamic code generation. They employ a version of quasi-quote system and type annotation to achieve a type safety in the generated code. This setting does not allow the generated code to have a syntactic error.

7 Conclusion

We have presented a static analysis technique for checking the syntax of generated code in two-staged languages with concatenation. Formulating abstract parsing in the abstract interpretation framework, we derive abstract semantics which is composed by atomic parse stack transition functions defined in LR(k) parser. This formulation not only gives us a more concise and elegant explanation of the original idea but also decouples the core idea of abstract parsing from its implementation. The provided framework allows us to choose the abstract domain of abstract parsing and control the precision and cost of analysis.

References

- [1] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley, 1986.
- [2] D. Ancona and E. Moggi. A fresh calculus for name management. In *Proceedings of the International Conference on Generative Programming and Component Engineering*, October 2004.
- [3] Jeroen Arnoldus, Jeanot Bijpost, and Mark van den Brand. Repleo: a syntax-safe template engine. In *Proceedings of the International Conference on Generative Programming and Component Engineering*, pages 25–32, New York, NY, USA, 2007. ACM.
- [4] Cristiano Calcagno, Eugenio Moggi, and Tim Sheard. Closed types for a safe imperative MetaML. *Journal of Functional Programming*, 13(3), 2003.
- [5] Cristiano Calcagno, Eugenio Moggi, and Walid Taha. ML-like inference for classifiers. In *Proceedings of the European Symposium on Programming*, pages 79–93. Springer, 2004.
- [6] Chiyen Chen and Hongwei Xi. Meta-programming through typeful code representation. In *ACM International Conference on Functional Programming*, pages 275–286. ACM, August 2003.
- [7] Tae-Hyoung Choi, Oukseh Lee, Hyunha Kim, and Kyung-Goo Doh. A practical string analyzer by the widening approach. In *Proceedings of the Asian Symposium on Programming Languages and Systems*, volume 4729 of *Lecture Notes in Computer Science*, pages 374–388, Sydney, Australia, November 2006. Springer-Verlag.

- [8] Aske Simon Christensen, Anders Mller, and Michael I. Schwartzbach. Precise analysis of string expressions. In *Proceedings of the Static Analysis Symposium*, pages 1–18. Springer-Verlag, 2003.
- [9] Mihai Christodorescu, Nicholas Kidd, and Wen-Han Goh. String analysis for x86 binaries. In *Proceedings of the ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering*, pages 88–95, New York, NY, USA, 2005. ACM.
- [10] P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Proceedings of The ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 238–252, Los Angeles, California, 1977. ACM Press, New York, NY.
- [11] P. Cousot and R. Cousot. Systematic design of program analysis frameworks. In *Proceedings of The ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 269–282, San Antonio, Texas, 1979. ACM Press, New York, NY.
- [12] P. Cousot and R. Cousot. Abstract interpretation frameworks. *Journal of Logic and Computation*, 2(4):511–547, 1992.
- [13] Rowan Davies. A temporal-logic approach to binding-time analysis. In *Proceedings of the Symposium on Logic in Computer Science*, pages 184–195. IEEE Computer Society Press, 1996.
- [14] Rowan Davies and Frank Pfenning. A modal analysis of staged computation. In *Proceedings of The ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 258–270. ACM, 1996.
- [15] Rowan Davies and Frank Pfenning. A modal analysis of staged computation. *Journal of the ACM*, 48(3):555–604, 2001.
- [16] Kyung-Goo Doh, Hyunha Kim, and David Schmidt. Abstract parsing: static analysis of dynamically generated string output using lr-parsing technology. In *Proceeding of the International Static Analysis Symposium*, 2009. Available from <http://santos.cis.ksu.edu/schmidt/dohsas09.pdf>.
- [17] Dawson R. Engler, Wilson C. Hsieh, and M. Frans Kaashoek. C: a language for high-level, efficient, and machine-independent dynamic code generation. In *Proceedings of The ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 131–144, New York, NY, USA, 1996. ACM.
- [18] Ik-Soon Kim, Kwangkeun Yi, and Cristiano Calcagno. A polymorphic modal type system for lisp-like multi-staged languages. In *Proceedings of The ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 257–269, 2006.
- [19] Yasuhiko Minamide. Static approximation of dynamically generated web pages. In *Proceedings of the International Conference on World Wide Web*, pages 432–441, New York, NY, USA, 2005. ACM.
- [20] Yasuhiko Minamide and Akihiko Tozawa. Xml validation for context-free grammars. In *Proceedings of the Asian Symposium on Programming Languages and Systems*, pages 357–373. Springer-Verlag, 2006.
- [21] Aleksandar Nanevski. Meta-programming with names and necessity. In *ACM International Conference on Functional Programming*, pages 206–217. ACM, October 2002.
- [22] Aleksandar Nanevski and Frank Pfenning. Staged computation with names and necessity. *Journal of Functional Programming*, 15(6):893–939, 2005.

- [23] Morten Rhiger. First-class open and closed code fragments. In *Proceedings of the Symposium on Trends in Functional Programming*, September 2005.
- [24] Tim Sheard. Accomplishments and research challenges in meta-programming. In *Proceedings of the International Workshop on Semantics, Applications, and Implementation of Program Generation*, pages 2–44, London, UK, 2001. Springer-Verlag.
- [25] Walid Taha and Michael Florentin Nielsen. Environment classifiers. In *Proceedings of The ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM, 2003.
- [26] Gary Wassermann, Carl Gould, Zhendong Su, and Premkumar Devanbu. Static checking of dynamically generated queries in database applications. *ACM Transactions on Software Engineering and Methodology*, 16(4):14, 2007.
- [27] Gary Wassermann and Zhendong Su. Sound and precise analysis of web applications for injection vulnerabilities. In *Proceedings of the SIGPLAN Conference on Programming Language Design and Implementation*, pages 32–41, New York, NY, USA, 2007. ACM.