

A Type System for String Sanitation Implemented Inside a Python

Nathan Fulton Cyrus Omar Jonathan Aldrich

School of Computer Science
Carnegie Mellon University
{nathanfu, comar, aldrich}@cs.cmu.edu

Abstract

ABSTRACT HERE

1. Introduction

In web applications and other settings, incorrect input sanitation often causes security vulnerabilities. In fact, ...OWASP says it's important... . Mention CVE stats. For this reason, modern web frameworks and libraries use various techniques to ensure proper sanitation of arbitrary user input. On the rare occasions when these methods are unavailable or insufficient, developers (hopefully) create ad hoc sanitation algorithms. In most cases, sanitation algorithms – ad hoc or otherwise – are ultimately implemented using the language's regular expression capabilities. Therefore, a system capable of statically checking properties about operations performed using regular expressions will be expressive enough to capture real-world implementations of sanitation algorithms.

Input sanitation is the problem of ensuring that an arbitrary string is coerced into a safe form before potentially unsafe use. For example, preventing SQL injection attacks requires ensuring that any string coming from user input to a query does not contain unescaped SQL. The point at which arbitrary user input is concatenated into a SQL query is called a use site. Although we believe a general approach extends to a wider class of problems (e.g. sanitation algorithms for preventing XSS attacks might be definable using regular tree languages), these generalizations are beyond the scope of the present paper.

This paper presents a language extension, definable in the Ace programming language, for ensuring that input sanitation algorithms are implemented correctly with respect to use site specifications. If use site specifications are sufficient,

then type checking ensures the absence of vulnerabilities and other bugs which arise from improper sanitation. Our extension focuses on sanitation algorithms which aim to prevent command-injection style attacks (e.g. SQL injection or RFI).

1.1 Related Work and Alternative Approaches

The input sanitation problem is well-understood. There exist a large number of techniques and technologies, proposed by both practitioners and researchers, for preventing injection-style attacks. In this section, we defend the novelty and significance of our approach with respect to the state of the art in practice and in research.

Unlike frameworks provided by languages such as Haskell and Ruby, our type system provides a *static* guarantee that input is always properly sanitized before use. We achieve this by defining a typing relation which captures idiomatic sanitation algorithms. Type safety in our system relies upon several closure and decidability results about regular languages.

Libraries and frameworks available in functional programming language communities often make claims about security and sometimes even sophistically mention sophisticated type systems as evidence of freedom from injection-style attacks. However, even the specification that input is always sanitized properly before use is not actually captured by the type system or anywhere else; in fact, the full specification of these algorithms is rarely characterized by anything more specific than the type $\text{String} \rightarrow \text{String}$, which is a class of safe input sanitation algorithms only in the most degenerate case.

A number of research languages provide actually static guarantees that a program is free of input sanitation vulnerabilities. Most rely on some form of information flow. TODO-nrf give citations and examples. Our extension to Ace differs from these systems in the ways following:

- Our system is a light-weight solution to a single class of sanitation vulnerabilities (e.g. we do not address Cross-Site Scripting). We present our system not as a comprehensive solution to the web security problem, but rather

as evidence that composable, light-weight and simple analyses can address security problems.

- Our system is defined as a library in terms of an extensible type system, as opposed to a stand-alone language. This is important for two reasons. First, our system is one of the first large examples written in Ace, and serves as an extended case study. Second, although our approach requires developers to adopt a new language, it does not require developers to adopt a language specifically for web development.
- Ace is implemented in Python and shares its grammar. Since Python is a popular programming language among web developers, the barrier between our research and adopted technologies is far lower than is e.g. Ur/Web's.
- In general, our is a composable, light-weight and natural solution to a single problem – rather than a comprehensive solution to the web security problem. Our goal is to demonstrate that extensible type systems can capture static properties of common idioms, and thereby ensure safety without introducing much additional complexity.

Finally, incorporating regular expressions into the type system is not novel. The XDuce system [?] typechecks XML schemas using regular expressions. We differ from this and related work in at least two ways. First, our system is defined within an extensible type system; this first difference introduced an interesting class of design problems (see §???). TODO-nrf this is the second time this is used as a warrant. Factor out this argument and place it at the top??? Second, our focus on security required novel design decisions; for instance, the filter elimination form is unique to Ace.

In conclusion, our system is novel in at least two ways:

- The safety guarantees provided by libraries and frameworks in popular languages are not as (statically) justified as is often belived (or even claimed).
- Our extension is the first major demonstration of how an extensible type system may be used to provide light-weight, composable security analyses based upon idiomatic code.

1.2 Outline

An outline of this paper follows:

TODO-nrf real outline!

- In §2, we define the type system's static and dynamic semantics.
- Section 3 recalls some classical results about regular expressions and presents meta-theory for our system, including the main soundness theorem for λ_{CS} .
- Finally, §4 discusses our implementation of λ_{CS} as a type system extension within the Ace programming language.

r	$::= \epsilon \mid . \mid a \mid r \cdot r \mid r + r \mid r^*$	Regex ($a \in \Sigma$).
ψ	$::=$ $\text{string_in}[r]$ $ $	Source language types Regular expression types
S	$::=$ $\text{si_str}[S]$ $ \text{siconcat}(S, S)$ $ \text{subst}[r](S, S)$ $ \text{coerce}[r](S)$	Source language expressions String concatenation
θ	$::=$ string $ \text{rx}[]$	Target language types
P	$::=$ $\text{tstr}[P]$ $ \text{siconcat}(P, P)$ $ \text{tsubst}(r, P, P)$ $ \text{tcheck}(r, P)$	Internal expressions

Figure 1. Syntax for the string sanitation fragment of λ_{CS}

2. A Type System for String Sanitation

The λ_{CS} language is characterized by a type of strings indexed by regular expressions, together with operations on such strings which correspond to common input sanitation patterns.

This section presents the grammar and semantics of λ_{CS} . The semantics are defined in terms of an internal language with at least strings and a regex filter function. These constraints are captured by the internal term valuations. The internal language does not necessarily need a regex filter function because any dynamic conversion is easily definable using a combination of filters and safe casting.

$$\boxed{\Gamma \vdash S : \psi} \quad \Gamma ::= \emptyset \mid \Gamma, x : \psi$$

$$\begin{array}{c}
\text{S-T-INCLUDE} \\
\frac{s \in \mathcal{L}\{r\}}{\Gamma \vdash \text{si_str}[s] : \text{string_in}[r]} \\
\\
\text{S-T-CONCAT} \\
\frac{\Gamma \vdash S_1 : \text{string_in}[r_1] \quad \Gamma \vdash S_2 : \text{string_in}[r_2]}{\Gamma \vdash \text{siconcat}(S_1, S_2) : \text{string_in}[r_1 \cdot r_2]} \\
\\
\text{S-T-SUBST} \\
\frac{\Gamma \vdash S_1 : \text{string_in}[r_1] \quad \Gamma \vdash S_2 : \text{string_in}[r_2] \quad \text{lsubst}(r, r_1, r_2) = r'}{\Gamma \vdash \text{subst}[r](S_1, S_2) : \text{string_in}[r']} \\
\\
\text{S-T-COERCE} \\
\frac{\Gamma \vdash S : \text{string_in}[r]}{\Gamma \vdash \text{coerce}[r'](s) : \text{string_in}[r']}
\end{array}$$

Figure 2. Typing rules for source language terms (S). The metavariable s ranges over string literals.

$$\boxed{S \Downarrow S}$$

$$\begin{array}{c}
\text{S-E-INCLUDE} \\
\frac{}{\text{si_str}[s] \Downarrow \text{si_str}[s]} \\
\\
\text{S-E-CONCAT} \\
\frac{S_1 \Downarrow \text{si_str}[s_1] \quad S_2 \Downarrow \text{si_str}[s_2]}{\text{siconcat}(S_1, S_2) \Downarrow \text{si_str}[s_1 \cdot s_2]} \\
\\
\text{S-E-SUBST} \\
\frac{S_1 \Downarrow \text{si_str}[s_1] \quad S_2 \Downarrow \text{si_str}[s_2] \quad \text{dosubst}(r, s_1, s_2) = s}{\text{subst}[r](S_1, S_2) \Downarrow s} \\
\\
\text{S-E-COERCE-OK} \\
\frac{S \Downarrow \text{si_str}[s] \quad s \in \mathcal{L}\{r\}}{\text{coerce}[r](S) \Downarrow \text{si_str}[s]} \\
\\
\text{S-E-COERCE-NOTOK} \\
\frac{\text{coerce}[r](S) \Downarrow \text{si_str}[s] \quad s \notin \mathcal{L}\{r\}}{S \text{ err}}
\end{array}$$

Figure 3. Big step semantics for the string fragment of the source language. There are also unmentioned error propagation rules for S-E-Concat and S-E-Subst.

Note that we never insert an internal check where the type of the string implies that a check must succeed. Furthermore, expensive calculations (such as language inclusion or dsubst , all occur at compile time. Since sanitation problems are generally expressible with smaller, simple expressions, we believe that the compile-time overhead is not significant enough to prohibit our target usecase. Informal experimentation with our implementation seems to support this assertion.

$$\boxed{\Omega \vdash P : \theta} \quad \Omega ::= \emptyset \mid \Omega, x : \theta$$

$$\begin{array}{c}
\text{P-T-STRING} \\
\frac{}{\Omega \vdash \text{tstr}[s] : \text{string}} \\
\\
\text{P-T-CONCAT} \\
\frac{\Omega \vdash P_1 : \text{string} \quad \Omega \vdash P_2 : \text{string}}{\Omega \vdash \text{tconcat}(P_1, P_2) : \text{string}} \\
\\
\text{P-T-SUBST} \\
\frac{\Omega \vdash P_1 : \text{string} \quad \Omega \vdash P_2 : \text{string}}{\Omega \vdash \text{tsubst}(\text{rx}[r], P_1, P_2) : \text{string}} \\
\\
\text{P-T-CHECK} \\
\frac{\Omega \vdash P : \text{string}}{\Omega \vdash \text{tcheck}(\text{rx}[r], s) : \text{string}}
\end{array}$$

Figure 4. Typing rules for target language terms (P). The metavariable s ranges over string literals.

$$\boxed{P \Downarrow P}$$

$$\begin{array}{c}
\text{P-E-STRIN} \\
\frac{}{\text{rx}[r] \Downarrow \text{rx}[r]} \\
\\
\text{P-E-CONCAT} \\
\frac{P_1 \Downarrow \text{tstr}[s_1] \quad P_2 \Downarrow \text{tstr}[s_2]}{\text{tconcat}(P_1, P_2) \Downarrow \text{tstr}[s_1 \cdot s_2]} \\
\\
\text{P-E-SUBST} \\
\frac{P_1 \Downarrow \text{rx}[r] \quad P_2 \Downarrow \text{tstr}[s_2] \quad P_3 \Downarrow \text{tstr}[s_3] \quad \text{dosubst}(r, s_2, s_3) = s}{\text{tsubst}(r, P_2, P_3) \Downarrow \text{tstr}[s]} \\
\\
\text{P-E-CHECK-OK} \\
\frac{P \Downarrow \text{si_str}[s] \quad s \in \mathcal{L}\{r\}}{\text{tcheck}(r, P) \Downarrow \text{tstr}[s]} \\
\\
\text{P-E-CHECK-NOTOK} \\
\frac{\text{tcheck}(r, P) \Downarrow \text{tstr}[s] \quad s \notin \mathcal{L}\{r\}}{P \text{ err}}
\end{array}$$

Figure 5. Big step semantics for the target language. There are also unimplemented error propagation rules for P-E-Concat and P-E-Subst.

Due to Ei-Replace, all of our meta-theory depends upon a correct implementation of both regular expression replacement and inclusion checking in the internal language. We believe this assumption is okay for two reasons. First, our system is still superior to the status quo, which relies upon the correctness of these libraries *in addition to* correct application logic. Second, regular expression libraries are generally well-tested and there exist verified implementations.

$$\boxed{\llbracket S \rrbracket = P}$$

TR-STRING

$$\frac{}{\llbracket \text{si_str}[s] \rrbracket = \text{tstr}[s]}$$

TR-CONCAT

$$\frac{\llbracket S_1 \rrbracket = P_1 \quad \llbracket S_2 \rrbracket = P_2}{\llbracket \text{siconcat}(S_1, S_2) \rrbracket = \text{tconcat}(P_1, P_2)}$$

TR-SUBST

$$\frac{\llbracket S_1 \rrbracket = P_1 \quad \llbracket S_2 \rrbracket = P_2}{\llbracket \text{subst}[r](S_1, S_2) \rrbracket = \text{tsubst}(\text{rx}[r], P_1, P_2)}$$

TR-COERCE-OK

$$\frac{S : \text{si_str}[r] \quad \mathcal{L}\{r'\} \subseteq \mathcal{L}\{r\}}{\llbracket \text{coerce}[r'](S) \rrbracket = \text{tstr}[s]}$$

TR-COERCE-NOTOK

$$\frac{\llbracket S \rrbracket = P \quad S : \text{si_str}[r] \quad \mathcal{L}\{r'\} \not\subseteq \mathcal{L}\{r\}}{\llbracket \text{coerce}[r'](S) \rrbracket = \text{tcheck}(\text{rx}[r'], P)}$$

Figure 6. Translation from source terms (S) to target terms (P). The translation is type-directed in the Tr-Coerce cases.

2.1 Properties of Regular Languages

The soundness theorems in this section depend upon some well-known properties of regular expressions. Throughout this section, suppose regular expression r and strings s are defined over a single, fixed, finite alphabet Σ .

Theorem 1. *The regular expressions are closed under complements and concatenation.*

Proof. See [?]. \square

Definition 2 (subst). The function $[r/s_1]s_2$ produces a string in which all substrings of s_1 matching r are replaced with s_2 .

Definition 3 (lsubst). The function $\text{lsubst}(r, r_1, r_2)$ produces a regular expression in which and sublanguage r'_1 of r_1 satisfying the condition $\mathcal{L}\{r'_1\} \subseteq r$ is replaced with r_2 .

We now present some properties about substitution.

Theorem 4. *String substitution ($[r/s_1]s_2$) is decidable and defined for all expressions and strings over similar alpha.*

Theorem 5. *Language substitution ($\text{lsubst}(r, r_1, r_2)$) is decidable and defined for all expressions.*

The following relationship between string substitution and language substitution is at the heart of our main correctness result.

Theorem 6. *For all r , if $s_1 \in r_1$ and $s_2 \in r_2$ then $[r/s_1]s_2 \subseteq \text{lsubst}(r, r_1, r_2)$.*

Lemma 7. Coercion Lemma. *Suppose that R and L are regular expressions, and that $s \in R$ is a finite string. Let $s' := \text{coerce}(R, L, s)$ with all maximal substrings recognized by L replaced with ϵ . Then s' is recognized by $(R \setminus L) + \emptyset$ and the construction of $R \setminus L$ is decidable.*

Proof. Let F, G be FAs corresponding to R and L , and let G' be G with its final states inverted (so that G' is the complement of L). Define an FA H as a DFA corresponding to the NFA found by combining F and G' such that H accepts only if R and L' accept or if s is empty (this construction may result in an exponential blowup in state size.) Clearly, H corresponds to $R \setminus L + \emptyset$. Thus, the construction of $R \setminus L + \emptyset$ is decidable.

If $R \subset L$, $s' = \emptyset$. If $L \subset R$, either $s' = \emptyset$, or $s' \in R$ and $s' \notin L$. If R and L are not subsets of one another, then it may be the case that L recognizes part of R . Consider L as the union of two languages, one which is a subset of R and one which is disjoint. The subset language is considered above and the disjoint language is inconsequential. \square

2.2 Safety of the Source and Target Languages

We now prove safety for the string fragment of the source and target languages.

Theorem 8. *Let S be a term in the source language. If $\Gamma \vdash S : \psi$ then $S \Downarrow S'$ and $S' : \psi$.*

Theorem 9. *Let P be a term in the target language. If $\Gamma \vdash P : \theta$ then $P \Downarrow P'$ and $P' : \theta$.*

2.3 Translation Correctness

We now present the main correctness result.

Theorem 10. *If $S : \text{si_str}[r]$ then there exists a T such that $\llbracket S \rrbracket = T$ and either:*

- (a) $T \Downarrow \text{tstr}[s]$ and $S \Downarrow \text{si_str}[s]$, and $s \in \text{langr}$.
- (b) $T \text{ err}$ and $S \text{ err}$.

Proof. By induction on the derivation of $\llbracket S \rrbracket = T$. TODO-nrf add case magic.

- Tr-string. By assumption $\llbracket S \rrbracket = \text{tstr}[s]$. The (a) property follows. By S-E-Str, $S \Downarrow \text{strings}$ and by T-E-Str, $\text{tstr}[s] \Downarrow \text{tstr}[s]$. By source language type safety, it follows that $s \in \mathcal{L}\{r\}$.
- Tr-Concat. By inversion of the typing relation for the source language, $S_1 : \text{si_str}[r_1]$ and $S_2 : \text{string_in}[S_2]$. By induction on the left premise, $T_1 \Downarrow \text{tstr}[s_1]$, $S_1 \Downarrow \text{strings}_1$ and $s_1 \in \mathcal{L}\{r_1\}$ or else $T_1 \text{ err}$. Similarly for T_2 and S_2 . In the error cases, the propagations rules corresponding to S-E-Concat and T-E-Concat imply that $\text{tconcat}(T_1, T_2) \text{ err}$ and $\text{siconcat}(S_1, S_2) \text{ err}$. In the non-error case, $\text{siconcat}(S_1, S_2) \Downarrow \text{si_str}[s_1 \cdot s_2]$ by S-E-Concat and $\text{tconcat}(T_1, T_2) \Downarrow \text{tstr}[s_1 \cdot s_2]$ by T-E-Concat. Finally, $s_1 \in \mathcal{L}\{r_1\}$ and $s_2 \in \mathcal{L}\{r_2\}$ implies that $s_1 \cdot s_2 \in \mathcal{L}\{r\}$ by ??.

tycon STRINGIN of R { (1)
 iana { $\lambda \text{opidx}:\text{String}.\lambda \text{tyidx}:\text{R}.\lambda \text{a}:\text{list}[\text{Arg}]$. (2)
 arity0 a (check opidx tyidx $\triangleright \triangleleft(\text{opidx})$) (3)
 esyn { $\lambda \text{opidx}:1 + (\text{R} + \text{R}).\lambda \text{a}:\text{list}[\text{Arg}]$. (4)
 case opidx of inl(_) \Rightarrow (5)
 arity2 a $\lambda \text{a1}:\text{Arg}.\lambda \text{a2}:\text{Arg}$. (6)
 rsyn a1 $\lambda \text{r1}:\text{R}.\lambda \text{i1}:\text{ITm}$. (7)
 rsyn a2 $\lambda \text{r2}:\text{R}.\lambda \text{i2}:\text{ITm}$. (8)
 (STRINGIN($\text{rseq}(\text{r1}; \text{r2})$), (9)
 $\triangleright (\text{iconcat}(\triangleleft(\text{i1}); \triangleleft(\text{i2})))$) (10)
 | inr(opidx') $\Rightarrow d$ } (11)

Figure 8. Definition of ...

- T-subst. Suppose $\text{subst}[r](S_1, S_2) : \text{si_str}[s]$ for some x . By inverison, $S_1 : \text{si_str}[r_1]$ and $S_2 : \text{si_str}[r_2]$. Note that $s : [r/S_1]S_2 : \text{string_in}[\text{tsubst}(r, r_1, r_2)]$ by S-T-Subst. By induction, $\llbracket S_1 \rrbracket = T_1$ and $\llbracket S_2 \rrbracket = T_2$. Either (a) or (b) holds for each.
 If (a) holds for both, the $S_1 \Downarrow \text{si_str}[S_1]$ and $S_2 \Downarrow \text{si_str}[S_2]$ for $s_1 \in \mathcal{L}\{r_1\}$ and $s_2 \in \mathcal{L}\{r_2\}$. Therefore, $\text{subst}[r](S_1, S_2) \Downarrow \text{si_str}[\text{dosubst}(r, s_1, s_2, \cdot)]$. The same fact can be shown in an analogous manner for T_1 and T_2 . Finally, $s_1 \in \mathcal{L}\{r_1\}$ and $s_2 \in \mathcal{L}\{r_2\}$, which implies that $\text{dosubst}(r, s_1, s_2) \in \mathcal{L}\{\text{subst}[r](s_1, s_2)\}$ by ??.
 If (b) holds for either the left or right premise, then $\text{subst}[r](S_1, S_2)$ err by the propagation rules corresponding to S-T-subst, and $\text{tsubst}(S_1, S_2, \text{err})$ by the propagation rules corresponding to T-T-Subst.

□

2.4 Type Safety Proof

2.5 New Theorems

Theorem 11 (Progress for Internal Language).

Theorem 12 (Preservation for Internal Language).

Theorem 13 (Representational Consistency). *If $\vdash e : \tau \Downarrow i$ then $\vdash i : \sigma$ and $\sigma \equiv \text{rep}(\tau)$.*

Theorem 14 (Main Soundness Result). *If $\vdash e : \text{si_str}[r] \Downarrow i$ then either $i \Downarrow \text{err}$ or if $i \Downarrow v$ for v then $v \in r$.*

3. Implementation Inside a Python

3.1 Background: Ace

TODO: Make a TR out of the OOPSLA submission.

concat := inl[1, R + R](b)
replace := $\lambda \text{r}:\text{R}.\text{inr}[1, R + R](\text{inl}[R, R](\text{r}))$
coerce := $\lambda \text{r}:\text{R}.\text{inr}[1, R + R](\text{inr}[R, R](\text{r}))$

Figure 9. Definitions

3.2 Explicit Conversions

3.3 Adding Subtyping to Ace

3.4 Theory

4. Related Work

5. Discussion

τ		$\lfloor \tau \rfloor_{\lambda_{\text{Ace}}}$	$\lfloor \tau \rfloor_{\text{ace}}$
$\text{si_str}[r]$		$\text{STRINGIN}[r]$	$\text{string_in}["r"]$
e		$\lfloor e \rfloor_{\lambda_{\text{Ace}}}$	$\lfloor e \rfloor_{\text{ace}}$
$\text{str}[r, s]$	synthetic position	$\text{intro}[\text{str}[s]]() : \text{STRINGIN}[\text{rx}[r]]$	$"s" \ (\text{string_in}["r"])$
	analytic position	$\text{intro}[\text{str}[s]]()$	$"s"$
$\text{concat}(e_1; e_2)$		$\lfloor e_1 \rfloor_{\lambda_{\text{Ace}}} \cdot \text{elim}[\text{concat}](\lfloor e_2 \rfloor_{\lambda_{\text{Ace}}})$	$\lfloor e_1 \rfloor_{\text{ace}} + \lfloor e_2 \rfloor_{\text{ace}}$
$\text{replace}[r](e_1; e_2)$		$\lfloor e_1 \rfloor_{\lambda_{\text{Ace}}} \cdot \text{elim}[\text{replace rx}[r]](\lfloor e_2 \rfloor_{\lambda_{\text{Ace}}})$	$\lfloor e_1 \rfloor_{\text{ace}}.\text{replace}("r", \lfloor e_2 \rfloor_{\text{ace}})$
$\text{coerce}[r](e)$		$\lfloor e \rfloor_{\lambda_{\text{Ace}}} \cdot \text{elim}[\text{coerce rx}[r]]()$	$\lfloor e \rfloor_{\text{ace}}.\text{coerce}("r")$

Figure 7. Translation of λ_{CS} to λ_{Ace} and Ace .