

A Type System for String Sanitation Implemented Inside a Python

ABSTRACT

Security-oriented programming languages do not enjoy wide adoption. Conversely, library and framework-based approaches to security are pervasive. Unfortunately, libraries and frameworks are only effective solutions if they are implemented and used properly. In this paper, we propose that extensible programming languages provide a promising mechanism for increasing the amount of trust that can be placed in libraries and frameworks. To demonstrate this approach, we present a conservative extension to the simply-typed lambda calculus for checking the correctness of input sanitation algorithms. We also demonstrate how this language’s security guarantees are preserved under translation into an underlying language containing a regular expression library. We believe this approach – complementing existing techniques with light-weight, composable type-based analyses – constitutes a promising compromise between usability, potential for industry adoption, and theoretically grounded safety guarantees.

1. INTRODUCTION

Improper input sanitation is a leading cause of security vulnerabilities in web applications [OWASP]. Command injection attacks exploit improper input sanitation by inserting malicious code into an otherwise benign command. Modern web frameworks, libraries, and database abstraction layers attempt to ensure proper sanitation of user input. When these methods are unavailable or insufficient, developers implement custom sanitation techniques. In both cases, sanitation algorithms are implemented using the language’s regular expression capabilities and usually *replace* potentially unsafe strings with equivalent escaped strings.

In this paper, we present a type system for implementing and statically checking input sanitation techniques. Our solution suggests a more general approach to the integration of security concerns into programming language design. This approach is characterized by *composable* type system extensions which *complement* existing and well-understood

solutions with compile-time checks.

To demonstrate this approach, we present a simply typed lambda calculus with *constrained strings*; that is, a set of string types parameterized by regular expressions. If $s : \text{stringin}[r]$, then s is a string matching the language r . Additionally, we include an operation $\text{rreplace}[r](s_1, s_2)$ which corresponds to the replace mechanism available in most regular expression libraries; that is, any substring of s_1 matching r is replaced with s_2 . The type of this expression is the computed, and is likely “smaller” or more constrained than the type of s_1 . Libraries, frameworks or functions which construct and execute commands containing input can specify a safe subset $\text{stringin}[r_{\text{spec}}]$ of strings, and input sanitation algorithms can construct such a string using rreplace or, optionally, by coercion (in which case a runtime check is inserted). We also show how this system is translated into a host language containing a regular expression library such that the safety guarantee of the extended language is preserved.

Summarily, we present a simple type system extension which ensures the absence of input sanitation vulnerabilities by statically checking input sanitation algorithms which use an underlying regular expression library. This approach is *composable* in the sense that it is a conservative extension. This approach is also *complementary* to existing input sanitation techniques which use string replacement for input sanitation.

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 explain how our approach to the input sanitation problem differs from each of these approaches. More important than these differences, however, is our more general assertion that language extensibility is a promising approach toward consideration of security goals in programming language design.

Unlike *frameworks and libraries* provided by languages such as Haskell and Ruby, our type system provides a *static* guarantee that input is always properly sanitized before use. Doing so requires reasoning about the operations on regular languages corresponding to standard operations on strings; we are unaware of any production system which contains this form of reasoning. Therefore, even where frameworks and libraries provide a viable interface or wrapper around

input sanitation, our approach is complementary because it ensures the correctness of the framework or library itself. Furthermore, our approach is more general than database abstraction layers because our mechanism is applicable to all forms of command injection (e.g. shell injection or remote file inclusion).

A number of research languages provide static guarantees that a program is free of input sanitation vulnerabilities [Jif][Ur/Web]. Unlike this work, our solution to the input sanitation problem has a very low barrier to adoption; for instance, our implementation conservatively extends Python – a popular language among web developers. We also believe our general approach is better-positioned for security, where continuously evolving threats might require frequent addition of new analyses; in these cases, the composability and generality of our approach is a substantial advantage.

We are also unaware of any extensible programming languages which emphasize applications to security concerns (TRUE?).

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; second, and more importantly, we have demonstrated that regular expression types are applicable to the web security domain.

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 lightweight, composable security analyses based upon idiomatic code.

2. A TYPE SYSTEM FOR STRING SANITATION

The λ_S 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, typing rules and operational semantics for λ_S as well as an underlying language λ_P .

The system λ_S is the simply typed lambda calculus extended with *regular expression types*, which are string types ensuring a string belongs to a specified language. For instance, $S : \text{stringin}[r]$ reads “ s is a string matching r ”. the system includes an operation for replacing all instances of a pattern r in a string s_1 with another string s_2 . Input sanitation algorithms – as implemented by developers or within popular libraries and frameworks – are often implemented in terms of this replace operation.

The language λ_P is a simple functional language extended with a minimal regular expression library. Any general pur-

pose programming language could stand in for λ_P ; for instance, SML or Python. In an implementation, our correctness results are modulo the underlying language’s correct implementation of regular expression matching (see P-E-Replace).

Finally, we define a translation from our type system λ_S into λ_P which preserves the safety guarantee codified in the string types of λ_S . Because our extension is conservative and its guarantees are preserved under translation to an underlying language, it is highly composable with other analyses.

Unfortunately, we are unable to present full proofs in this paper due to space constraints; however, proofs are given in full in the accompanying technical report?.

$$r ::= \epsilon \mid . \mid a \mid r \cdot r \mid r + r \mid r^* \quad a \in \Sigma$$

Figure 1: Regular expressions over the alphabet Σ .

$$\begin{array}{ll} \psi ::= \psi \rightarrow \psi & \text{source types} \\ \quad \mid \text{stringin}[r] & \\ \\ S ::= \lambda x.e & \text{source terms} \\ \quad \mid ee & \\ \quad \mid \text{rstr}[s] & s \in \Sigma^* \\ \quad \mid \text{rconcat}(S, S) & \\ \quad \mid \text{rreplace}[r](S, S) & \\ \quad \mid \text{rcoerce}[r](S) & \end{array}$$

Figure 2: Syntax for the string sanitation fragment of our source language, λ_S .

$$\begin{array}{ll} \theta ::= \theta \rightarrow \theta & \text{target types} \\ \quad \mid \text{string} & \\ \quad \mid \text{regex} & \\ \\ P ::= \lambda x.e & \text{target terms} \\ \quad \mid ee & \\ \quad \mid \text{str}[s] & \\ \quad \mid \text{rx}[r] & \\ \quad \mid \text{concat}(P, P) & \\ \quad \mid \text{preplace}(P, P, P) & \\ \quad \mid \text{check}(P, P) & \end{array}$$

Figure 3: Syntax for the fragment of our target language, λ_P , containing strings and statically constructed regular expressions.

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

$\frac{}{\llbracket \text{rstr}[s] \rrbracket = \text{str}[s]} \text{Tr-STRING}$	$\frac{\llbracket S_1 \rrbracket = P_1 \quad \llbracket S_2 \rrbracket = P_2}{\llbracket \text{rconcat}(S_1, S_2) \rrbracket = \text{concat}(P_1, P_2)} \text{Tr-CONCAT}$	$\frac{\llbracket S_1 \rrbracket = P_1 \quad \llbracket S_2 \rrbracket = P_2}{\llbracket \text{rreplace}[r](S_1, S_2) \rrbracket = \text{replace}(\text{rx}[r], P_1, P_2)} \text{Tr-SUBST}$
$\frac{\text{Tr-COERCE-OK} \quad S : \text{rstr}[r] \quad \mathcal{L}\{r'\} \subseteq \mathcal{L}\{r\}}{\llbracket \text{rcoerce}[r'](S) \rrbracket = \text{str}[s]} \text{Tr-COERCE-OK}$	$\frac{\text{Tr-COERCE-NOTOK} \quad \llbracket S \rrbracket = P \quad S : \text{rstr}[r] \quad \mathcal{L}\{r'\} \not\subseteq \mathcal{L}\{r\}}{\llbracket \text{rcoerce}[r'](S) \rrbracket = \text{check}(\text{rx}[r'], P)} \text{Tr-COERCE-NOTOK}$	

Figure 8: Translation from source terms (S) to target terms (P). The translation is type-directed in the Tr-Coerce cases.

$\boxed{\Psi \vdash S : \psi} \quad \Psi ::= \emptyset \mid \Psi, x : \psi$ $\frac{\text{S-T-STRINGIN-I} \quad s \in \mathcal{L}\{r\}}{\Psi \vdash \text{rstr}[s] : \text{stringin}[r]} \text{S-T-STRINGIN-I}$ $\frac{\text{S-T-CONCAT} \quad \Psi \vdash S_1 : \text{stringin}[r_1] \quad \Psi \vdash S_2 : \text{stringin}[r_2]}{\Psi \vdash \text{rconcat}(S_1, S_2) : \text{stringin}[r_1 \cdot r_2]} \text{S-T-CONCAT}$ $\frac{\text{S-T-REPLACE} \quad \Psi \vdash S_1 : \text{stringin}[r_1] \quad \Psi \vdash S_2 : \text{stringin}[r_2] \quad \text{lreplace}(r, r_1, r_2) = r'}{\Psi \vdash \text{rreplace}[r](S_1, S_2) : \text{stringin}[r']} \text{S-T-REPLACE}$ $\frac{\text{S-T-COERCE} \quad \Psi \vdash S : \text{stringin}[r']}{\Psi \vdash \text{rcoerce}[r](S) : \text{stringin}[r]} \text{S-T-COERCE}$	$\boxed{\Theta \vdash P : \theta} \quad \Theta ::= \emptyset \mid \Theta, x : \theta$ $\frac{}{\Theta \vdash \text{str}[s] : \text{string}} \text{P-T-STRING} \quad \frac{}{\Theta \vdash \text{rx}[r] : \text{regex}} \text{P-T-REGEX}$ $\frac{\text{P-T-CONCAT} \quad \Theta \vdash P_1 : \text{string} \quad \Theta \vdash P_2 : \text{string}}{\Theta \vdash \text{concat}(P_1, P_2) : \text{string}} \text{P-T-CONCAT}$ $\frac{\text{P-T-REPLACE} \quad \Theta \vdash P_1 : \text{regex} \quad \Theta \vdash P_2 : \text{string} \quad \Theta \vdash P_3 : \text{string}}{\Theta \vdash \text{preplace}(P_1, P_2, P_3) : \text{string}} \text{P-T-REPLACE}$ $\frac{\text{P-T-CHECK} \quad \Theta \vdash P_1 : \text{regex} \quad \Theta \vdash P_2 : \text{string}}{\Theta \vdash \text{check}(P_1, P_2) : \text{string}} \text{P-T-CHECK}$
---	--

Figure 4: Typing rules for our fragment of λ_S . The typing context Ψ is standard.

$\boxed{S \Downarrow S} \quad \boxed{S \text{ err}}$ $\frac{}{\text{rstr}[s] \Downarrow \text{rstr}[s]} \text{S-E-RSTR} \quad \frac{\text{S-E-CONCAT} \quad S_1 \Downarrow \text{rstr}[s_1] \quad S_2 \Downarrow \text{rstr}[s_2]}{\text{rconcat}(S_1, S_2) \Downarrow \text{rstr}[s_1 s_2]} \text{S-E-CONCAT}$ $\frac{\text{S-E-REPLACE} \quad S_1 \Downarrow \text{rstr}[s_1] \quad S_2 \Downarrow \text{rstr}[s_2] \quad \text{lsubst}(r, s_1, s_2) = s}{\text{rreplace}[r](S_1, S_2) \Downarrow \text{rstr}[s]} \text{S-E-REPLACE}$ $\frac{\text{S-E-COERCE-OK} \quad S \Downarrow \text{rstr}[s] \quad s \in \mathcal{L}\{r\}}{\text{rcoerce}[r](S) \Downarrow \text{rstr}[s]} \text{S-E-COERCE-OK} \quad \frac{\text{S-E-COERCE-ERR} \quad S \Downarrow \text{rstr}[s] \quad s \notin \mathcal{L}\{r\}}{\text{rcoerce}[r](S) \text{ err}} \text{S-E-COERCE-ERR}$	$\boxed{P \Downarrow P} \quad \boxed{P \text{ err}}$ $\frac{}{\text{str}[s] \Downarrow \text{str}[s]} \text{P-E-STR} \quad \frac{}{\text{rx}[r] \Downarrow \text{rx}[r]} \text{P-E-RX} \quad \frac{\text{P-E-CONCAT} \quad P_1 \Downarrow \text{str}[s_1] \quad P_2 \Downarrow \text{str}[s_2]}{\text{concat}(P_1, P_2) \Downarrow \text{str}[s_1 s_2]} \text{P-E-CONCAT}$ $\frac{\text{P-E-REPLACE} \quad P_1 \Downarrow \text{rx}[r] \quad P_2 \Downarrow \text{str}[s_2] \quad P_3 \Downarrow \text{str}[s_3] \quad \text{lsubst}(r, s_2, s_3) = s}{\text{preplace}(P_1, P_2, P_3) \Downarrow \text{str}[s]} \text{P-E-REPLACE}$ $\frac{\text{P-E-CHECK-OK} \quad P_1 \Downarrow \text{rx}[r] \quad P_2 \Downarrow \text{str}[s] \quad s \in \mathcal{L}\{r\}}{\text{check}(P_1, P_2) \Downarrow \text{str}[s]} \text{P-E-CHECK-OK} \quad \frac{\text{P-E-CHECK-ERR} \quad P_1 \Downarrow \text{rx}[r] \quad P_2 \Downarrow \text{str}[s] \quad s \notin \mathcal{L}\{r\}}{\text{check}(P_1, P_2) \text{ err}} \text{P-E-CHECK-ERR}$
--	---

Figure 5: Big step semantics for our fragment of λ_S . Error propagation rules are omitted.

Figure 6: Typing rules for our fragment of λ_P . The typing context Θ is standard.

$\boxed{P \Downarrow P} \quad \boxed{P \text{ err}}$ $\frac{}{\text{str}[s] \Downarrow \text{str}[s]} \text{P-E-STR} \quad \frac{}{\text{rx}[r] \Downarrow \text{rx}[r]} \text{P-E-RX} \quad \frac{\text{P-E-CONCAT} \quad P_1 \Downarrow \text{str}[s_1] \quad P_2 \Downarrow \text{str}[s_2]}{\text{concat}(P_1, P_2) \Downarrow \text{str}[s_1 s_2]} \text{P-E-CONCAT}$ $\frac{\text{P-E-REPLACE} \quad P_1 \Downarrow \text{rx}[r] \quad P_2 \Downarrow \text{str}[s_2] \quad P_3 \Downarrow \text{str}[s_3] \quad \text{lsubst}(r, s_2, s_3) = s}{\text{preplace}(P_1, P_2, P_3) \Downarrow \text{str}[s]} \text{P-E-REPLACE}$ $\frac{\text{P-E-CHECK-OK} \quad P_1 \Downarrow \text{rx}[r] \quad P_2 \Downarrow \text{str}[s] \quad s \in \mathcal{L}\{r\}}{\text{check}(P_1, P_2) \Downarrow \text{str}[s]} \text{P-E-CHECK-OK} \quad \frac{\text{P-E-CHECK-ERR} \quad P_1 \Downarrow \text{rx}[r] \quad P_2 \Downarrow \text{str}[s] \quad s \notin \mathcal{L}\{r\}}{\text{check}(P_1, P_2) \text{ err}} \text{P-E-CHECK-ERR}$	$\boxed{P \Downarrow P} \quad \boxed{P \text{ err}}$ $\frac{}{\text{str}[s] \Downarrow \text{str}[s]} \text{P-E-STR} \quad \frac{}{\text{rx}[r] \Downarrow \text{rx}[r]} \text{P-E-RX} \quad \frac{\text{P-E-CONCAT} \quad P_1 \Downarrow \text{str}[s_1] \quad P_2 \Downarrow \text{str}[s_2]}{\text{concat}(P_1, P_2) \Downarrow \text{str}[s_1 s_2]} \text{P-E-CONCAT}$ $\frac{\text{P-E-REPLACE} \quad P_1 \Downarrow \text{rx}[r] \quad P_2 \Downarrow \text{str}[s_2] \quad P_3 \Downarrow \text{str}[s_3] \quad \text{lsubst}(r, s_2, s_3) = s}{\text{preplace}(P_1, P_2, P_3) \Downarrow \text{str}[s]} \text{P-E-REPLACE}$ $\frac{\text{P-E-CHECK-OK} \quad P_1 \Downarrow \text{rx}[r] \quad P_2 \Downarrow \text{str}[s] \quad s \in \mathcal{L}\{r\}}{\text{check}(P_1, P_2) \Downarrow \text{str}[s]} \text{P-E-CHECK-OK} \quad \frac{\text{P-E-CHECK-ERR} \quad P_1 \Downarrow \text{rx}[r] \quad P_2 \Downarrow \text{str}[s] \quad s \notin \mathcal{L}\{r\}}{\text{check}(P_1, P_2) \text{ err}} \text{P-E-CHECK-ERR}$
---	---

Figure 7: Big step semantics for our fragment of λ_P . Error propagation rules are omitted.

2.1 Properties of Regular Languages

Our type safety proof for language S relies on a relationship between string substitution and language substitution given in lemma 5. We also rely upon several other properties of regular languages. Throughout this section, we fix an alphabet Σ over which strings s and regular expressions r are defined. throughout the paper, $\mathcal{L}\{r\}$ refers to the language recognized by the expression r . This distinction between the expression and its language – typically elided in the literature – makes our definition and proofs about systems S and P more readable.

Lemma 1. *Properties of Regular Languages and Expressions.* The following are properties of regular expressions which are necessary for our proofs: If $s_1 \in \mathcal{L}\{r_1\}$ and $s_2 \in \mathcal{L}\{r_2\}$ then $s_1 s_2 \in \mathcal{L}\{r_1 r_2\}$. For all strings s and expressions r , either $s \in \mathcal{L}\{r\}$ or $s \notin \mathcal{L}\{r\}$.

Definition 2 (lsubst). The function $\text{lsubst}(r, s_1, s_2)$ produces a string in which all substrings of s_1 matching r are replaced with s_2 .

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

Lemma 4. *Closure and Totality of Replacement.* If r, r_1 and r_2 are regular expressions, then $\text{lreplace}(r, r_1, r_2)$ is also a regular expression.

Proof. By induction on r and closure properties of regular expressions. \square

Lemma 5. *Substitution Correspondence.* If $s_1 \in \mathcal{L}\{r_1\}$ and $s_2 \in \mathcal{L}\{r_2\}$ then $\text{lsubst}(r, s_1, s_2) \in \mathcal{L}\{\text{lreplace}(r, r_1, r_2)\}$.

Proof. The proof proceeds by structural induction on r . The base case ($r = \alpha, \alpha \in \Sigma$) is trivial. The choice and sequential cases require lemma 1, and prove a more general property that $\text{lsubst}(r, s_1, s_2) \in \mathcal{L}\{r'\}$ for some r' equivalent to $\text{lreplace}(r, r_1, r_2)$. The closure case follows from the induction hypothesis applied to any n -unwinding. \square

2.2 Safety of the Source and Target Languages

Lemma 6. *If $\Psi \vdash S : \text{stringin}[r]$ then r is a well-formed regular expression.*

Proof. The only non-trivial case is S-T-Replace, which follows from lemma 4. \square

Lemma 7. *If $\Theta \vdash P : \text{regex}$ then $P \Downarrow \text{rx}[r]$ such that r is a well-formed regular expression.*

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

Theorem 8. *Safety for the String Fragment of P .* Let S be a term in the source language. If $\Psi \vdash S : \text{stringin}[r]$ then $S \Downarrow \text{rstr}[s]$ and $\text{rstr}[s] : \text{stringin}[r]$, or else $S \text{ err}$.

Proof. By induction on the derivation of $\Psi \vdash S : \psi$. The interesting case is S-T-Replace, which requires lemma 5. The Coercion lemma additionally requires the second property of lemma 1. \square

In addition to safety, we proof a correctness result for λ_S which will be preserved under translation.

Theorem 9. *Input Sanitation Correctness of λ_S .* If $\Psi \vdash S : \text{stringin}[r]$ then either (a) $S \Downarrow \text{rstr}[s]$ for $s \in \mathcal{L}\{r\}$ or (b) $S \text{ err}$.

Proof. We may actually want to present this proof. \square

Finally, λ_P is also type-safe. We omit the proof due to space constraints, and because the proof is trivial.

Theorem 10. *Let P be a term in the target language. If $\Theta \vdash P : \theta$ then $P \Downarrow P'$ and $\Theta \vdash P' : \theta$.*

2.3 Translation Correctness

Theorem 11. *Translation Correctness* If $\Psi \vdash S : \text{stringin}[r]$ then there exists a P such that $\llbracket S \rrbracket = P$ and either: (a) $P \Downarrow \text{str}[s]$ and $S \Downarrow \text{rstr}[s]$, or (b) $P \text{ err}$ and $S \text{ err}$.

Proof. The proof proceeds by induction on the typing relation for S and an appropriate choice of P ; in each case, the choice is obvious. The subcases (a) proceed by inversion and appeals to our type safety theorems as well as the induction hypothesis. The subcases (b) proceed by the standard error propagation rules omitted for space. Throughout the proof, properties from the closure lemma for regular languages are necessary. \square

Theorem 12. *If $\Psi \vdash S : \text{stringin}[r]$ and $\llbracket S \rrbracket = P$ then either: (a) $P \Downarrow \text{str}[s]$ and $S \Downarrow \text{rstr}[s]$, or (b) $P \text{ err}$ and $S \text{ err}$.*

Proof. The theorem follows directly from canonical forms for translations and theorem 11. \square

Finally, our main theorem establishes that input sanitation correctness of λ_S is preserved under the translation into λ_P .

Theorem 13. *Input Sanitation Correctness for Translated Terms* If $\llbracket S \rrbracket = P$ and $\Psi \vdash S : \text{stringin}[r]$ then either $P \text{ err}$ or $P \Downarrow \text{str}[s]$ for $s \in \mathcal{L}\{r\}$.

Proof. By theorem 11, $P \Downarrow \text{str}[s]$ implies that $S \Downarrow \text{rstr}[s]$. By theorem 9, this together with the assumption that S is well-typed implies that $s \in \mathcal{L}\{r\}$. \square

3. CONCLUSION

Composable analyses which complement existing approaches constitute a promising approach toward the integration of security concerns into programming languages. In this paper, we presented a system with both of these properties and defined a security-preserving transformation. Unlike other approaches, our solution complements existing, familiar solutions while providing a strong guarantee that traditional library and framework-based approaches are implemented and utilized correctly.

Papers that needs to be cited in this section:

- Ur/Web OSDI paper
- Jif?
- OWASP
- XDuce and related papers.
- src?
- Ace or Wyvern paper?
- hotsos?
- Haskell extension paper
- Maybe some popular FOSS libraries/frameworks that do input sanitation?