

# Statically Typed String Sanitation Inside a Python

Nathan Fulton

Cyrus Omar

Jonathan Aldrich

Carnegie Mellon University  
Pittsburgh, PA

{nathanfu, comar, aldrich}@cs.cmu.edu

## ABSTRACT

Web applications must ultimately generate strings that contain commands to be consumed by systems like web browsers and database engines. If these strings are constructed from user input that has not been properly sanitized, this can expose costly injection vulnerabilities.

In this paper, we introduce *regular string types*, which classify strings known statically to be in a specified regular language. The language of a string is tracked by the type system through operations like concatenation, substitution and coercion, so regular string types can be used to implement, in essentially a conventional manner, the parts of a web application or application framework that constructs commands. Simple type annotations at key points can be used to statically verify that sanitization has been performed correctly without introducing redundant run-time checks.

We take the position that to be practical, such a type system cannot require that programmers adopt a new programming language. Ideally, such a “special-purpose” type system would be distributed as a library that could safely be used together with other such type systems. We support this by specifying a sound translation to a simple language containing only strings and regular expressions, then discuss implementing the type system together with this translation as a library in *atlang*, an extensible static type system for Python (being developed by the authors).

## 1. INTRODUCTION

Improper input sanitation is a leading cause of security vulnerabilities in web applications [7]. 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  $\text{rreplace}$  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 [1]. 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.

Incorporating regular expressions into the type system is not novel. The XDuce system [6, 5] checks XML documents against schemas using regular expressions. Similarly, XHaskell [8] focuses on XML documents. We differ from this and related work in at least three ways:

- Our system is defined within an extensible type system.
- We demonstrate that regular expression types are applicable to the web security domain, whereas previous work on regular expression types focused on XML schemas.
- Although our static replacement operation is definable in some languages with regular expression types, we are the first to expose this operation and connect the semantics of regular language replacement with the semantics of string substitution via a type safety and compilation correctness argument.

In conclusion, our contribution is a type system, implemented within an extensible type system, for checking the correctness of input sanitation algorithms.

## 2. A TYPE SYSTEM FOR STRING SANITATION

In this section we define a language for statically checked string sanitation ( $\lambda_S$ ). The system has regular expression types `stringin[r]` where  $r$  is a regular expression. Expressions of this type evaluate to string literals in the language described by  $r$ . Operations on expressions of type `stringin[r]` preserve this property.

The premier operation for manipulating strings in  $\lambda_S$  is string substitution, which is a familiar operation to any programmer who has used regular expressions. The replacement operation replaces all instances of a pattern in one string with another string; for instance, `lsubst(a|b, a, c) = c`. In order to compute the type resulting from substitution, we also need to compute the result of replacing on language with another inside a given language. Finally, just for convenience, we provide a coerce operation. The introduction of coercion requires handling of runtime errors.

The underlying language  $\lambda_P$  has only one type for strings. We prove that whenever a term is translated from  $\lambda_S$  to  $\lambda_P$ , correctness is preserved. The only exception is in the case of unsafe casts in  $\lambda_S$ , which are unnecessary but are included to demonstrate that the regex library of  $\lambda_P$  may be used to insert dynamic checks whenever even when developers are not careful about using statically checked operations.

A brief outline of this section follows:

- Page 3 contains a definition of  $\lambda_S, \lambda_P$  and the translation from  $\lambda_S$  to  $\lambda_P$ .
- In §2.1 we state some properties about regular expressions which are needed in our correctness proofs.
- In §2.2 we prove type safety for  $\lambda_P$  as well as both type safety and correctness for  $\lambda_S$ .
- In §2.3 we prove that translation preserves the correctness result about  $\lambda_S$ .

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

**Figure 1: Regular expressions over the alphabet  $\Sigma$ .**

$\psi ::= \psi \rightarrow \psi$  source types  
 $\mid \text{stringin}[r]$

$S ::= \lambda x.e$  source terms  
 $\mid ee$   
 $\mid \text{rstr}[s]$   $s \in \Sigma^*$   
 $\mid \text{rconcat}(S, S)$   
 $\mid \text{rreplace}[r](S, S)$   
 $\mid \text{rcoerce}[r](S)$

**Figure 2: Syntax for the string sanitation fragment of our source language,  $\lambda_S$ .**

$\theta ::= \theta \rightarrow \theta$  target types  
 $\mid \text{string}$   
 $\mid \text{regex}$

$P ::= \lambda x.e$  target terms  
 $\mid ee$   
 $\mid \text{str}[s]$   
 $\mid \text{rx}[r]$   
 $\mid \text{concat}(P, P)$   
 $\mid \text{preplace}(P, P, P)$   
 $\mid \text{check}(P, P)$

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

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

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

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

$$\boxed{\Psi \vdash S : \psi}$$

$$\Psi ::= \emptyset \mid \Psi, x : \psi$$

$$\begin{array}{c}
\text{S-T-STRINGIN-I} \quad \frac{s \in \mathcal{L}\{r\}}{\Psi \vdash \text{rstr}[s] : \text{stringin}[r]} \\
\\
\text{S-T-CONCAT} \quad \frac{\Psi \vdash S_1 : \text{stringin}[r_1] \quad \Psi \vdash S_2 : \text{stringin}[r_2]}{\Psi \vdash \text{rconcat}(S_1, S_2) : \text{stringin}[r_1 \cdot r_2]} \\
\\
\text{S-T-REPLACE} \quad \frac{\Psi \vdash S_1 : \text{stringin}[r_1] \quad \Psi \vdash S_2 : \text{stringin}[r_2] \quad \text{lreplace}(r, r_1, r_2) = r'}{\Psi \vdash \text{rreplace}[r](S_1, S_2) : \text{stringin}[r']} \\
\\
\text{S-T-COERCE} \quad \frac{\Psi \vdash S : \text{stringin}[r']}{\Psi \vdash \text{rcoerce}[r](S) : \text{stringin}[r]}
\end{array}$$

**Figure 4: Typing rules for our fragment of  $\lambda_S$ . The typing context  $\Psi$  is standard.**

$$\boxed{S \Downarrow S} \quad \boxed{S \text{ err}}$$

$$\begin{array}{c}
\text{S-E-RSTR} \quad \frac{}{\text{rstr}[s] \Downarrow \text{rstr}[s]} \quad \text{S-E-CONCAT} \quad \frac{S_1 \Downarrow \text{rstr}[s_1] \quad S_2 \Downarrow \text{rstr}[s_2]}{\text{rconcat}(S_1, S_2) \Downarrow \text{rstr}[s_1 s_2]} \\
\\
\text{S-E-REPLACE} \quad \frac{S_1 \Downarrow \text{rstr}[s_1] \quad S_2 \Downarrow \text{rstr}[s_2] \quad \text{lsubst}(r, s_1, s_2) = s}{\text{rreplace}[r](S_1, S_2) \Downarrow \text{rstr}[s]} \\
\\
\text{S-E-COERCE-OK} \quad \frac{S \Downarrow \text{rstr}[s] \quad s \in \mathcal{L}\{r\}}{\text{rcoerce}[r](S) \Downarrow \text{rstr}[s]} \quad \text{S-E-COERCE-ERR} \quad \frac{S \Downarrow \text{rstr}[s] \quad s \notin \mathcal{L}\{r\}}{\text{rcoerce}[r](S) \text{ err}}
\end{array}$$

**Figure 5: Big step semantics for our fragment of  $\lambda_S$ . Error propagation rules are omitted.**

$$\boxed{\Theta \vdash P : \theta}$$

$$\Theta ::= \emptyset \mid \Theta, x : \theta$$

$$\begin{array}{c}
\text{P-T-STRING} \quad \frac{}{\Theta \vdash \text{str}[s] : \text{string}} \quad \text{P-T-REGEX} \quad \frac{}{\Theta \vdash \text{rx}[r] : \text{regex}} \\
\\
\text{P-T-CONCAT} \quad \frac{\Theta \vdash P_1 : \text{string} \quad \Theta \vdash P_2 : \text{string}}{\Theta \vdash \text{concat}(P_1, P_2) : \text{string}} \\
\\
\text{P-T-REPLACE} \quad \frac{\Theta \vdash P_1 : \text{regex} \quad \Theta \vdash P_2 : \text{string} \quad \Theta \vdash P_3 : \text{string}}{\Theta \vdash \text{preplace}(P_1, P_2, P_3) : \text{string}} \\
\\
\text{P-T-CHECK} \quad \frac{\Theta \vdash P_1 : \text{regex} \quad \Theta \vdash P_2 : \text{string}}{\Theta \vdash \text{check}(P_1, P_2) : \text{string}}
\end{array}$$

**Figure 6: Typing rules for our fragment of  $\lambda_P$ . The typing context  $\Theta$  is standard.**

$$\boxed{P \Downarrow P} \quad \boxed{P \text{ err}}$$

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

**Figure 7: Big step semantics for our fragment of  $\lambda_P$ . 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  $\Sigma$  over which strings  $s$  and regular expressions  $r$  are defined. throughout the paper,  $\mathcal{L}\{r\}$  refers to the language recognized by the regular expression  $r$ . This distinction between the regular 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 regular expressions  $r$ , either  $s \in \mathcal{L}\{r\}$  or  $s \notin \mathcal{L}\{r\}$ . Regular languages are closed under difference, right quotient, reversal, and string homomorphism.

If any of these properties are unfamiliar, the reader may refer to a standard text on the subject [4].

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

**Definition 3 (lreplace).** The relation  $\text{lreplace}(r, r_1, r_2) = r'$  relates  $r, r_1$ , and  $r_2$  to a language  $r'$  containing all strings of  $r_1$  except that any substring  $s_{pre} s_{post} \in \mathcal{L}\{r_1\}$  where  $s \in \mathcal{L}\{r\}$  is replaced by the set of strings  $s_{pre} s_2 s_{post}$  for all  $s_2 \in \mathcal{L}\{r_2\}$  (the prefix and postfix positions may be empty).

**Lemma 4.** *Closure.* If  $\mathcal{L}\{r\}, \mathcal{L}\{r_1\}$  and  $\mathcal{L}\{r_2\}$  are regular expressions, then  $\mathcal{L}\{\text{lreplace}(r, r_1, r_2)\}$  is also a regular language.

*Proof.* The theorem follows from closure under difference, right quotient, reversal and string homomorphism.  $\square$

**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 theorem follows from the definitions of  $\text{lsubst}$  and  $\text{lreplace}$ ; note that language substitutions over-approximate string substitutions.  $\square$

## 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 and Sanitation Correctness 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  $\Psi \vdash \text{rstr}[s] : \text{stringin}[r]$ ; or else  $S \text{ err}$ .

*Proof.* By induction on the typing relation, where (a) case holds by lemma 1 in the S-T-Concat case and lemma 5 in the S-T-Replace case.. The (b) cases hold by unstated, but standard, error propagation rules.  $\square$

In addition to safety, we proof a correctness result for  $\lambda_S$  which ensures that well-typed terms of regular string type are in the language associated with their type.

**Theorem 9.** *Correctness of Input Sanitation for  $\lambda_S$ .* If  $\Psi \vdash S : \text{stringin}[r]$  and  $S \Downarrow \text{rstr}[s]$  then  $s \in \mathcal{L}\{r\}$ .

*Proof.* Follows directly from type safety, canonical forms for  $\lambda_S$ .  $\square$

**Theorem 10.** Let  $P$  be a term in the target language. If  $\Theta \vdash P : \theta$  then  $P \Downarrow P'$  and  $\Theta \vdash P' : \theta$ , or else  $P \text{ err}$ .

## 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.  $\square$

Finally, our main theorem establishes that input sanitation correctness of  $\lambda_S$  is preserved under the translation into  $\lambda_P$ .

**Theorem 12.** *Correctness of Input Sanitation 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. IMPLEMENTATION

An implementation of a similar system as a type system extension in the programming language Ace is discussed in [2] and [3]. The implementation only supports a special case of replacement where unsafe substrings are simply stripped. In practice, most libraries and frameworks escape command sequences instead of stripping characters; therefore, we plan to extend our implementation with support for the replace operation as described in this paper..

## 4. 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.

## 5. REFERENCES

- [1] A. Chlipala. Static checking of dynamically-varying security policies in database-backed applications. In *OSDI'10: Proceedings of the 9th USENIX Symposium on Operating Systems Design and Implementation*, Oct. 2010.
- [2] N. Fulton. Security through extensible type systems. In *Proceedings of the 3rd Annual Conference on Systems, Programming, and Applications: Software for Humanity*, SPLASH '12, pages 107–108, New York, NY, USA, 2012. ACM.
- [3] N. Fulton. A typed lambda calculus for input sanitation. Undergraduate thesis in mathematics, Carthage College, 2013.
- [4] J. E. Hopcroft and J. D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley, 1979.
- [5] H. Hosoya and B. C. Pierce. XDuce: A statically typed XML processing language. *ACM Transactions on Internet Technology*, 3(2):117–148, May 2003.
- [6] H. Hosoya, J. Vouillon, and B. C. Pierce. Regular Expression Types for XML. In *ICFP '00*, 2000.
- [7] OWASP. Open web application security project top 10.
- [8] M. Sulzmann and K. Lu. Xhaskell – adding regular expression types to haskell. In O. Chitil, Z. Horváth, and V. Zsák, editors, *Implementation and Application of Functional Languages*, volume 5083 of *Lecture Notes in Computer Science*, pages 75–92. Springer Berlin Heidelberg, 2008.