

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 containing commands that are executed by external systems like web browsers and database engines. If these strings are constructed from user input that has not been properly sanitized, this can lead to a variety of command injection vulnerabilities.

In this paper, we introduce *regular string types*, which classify strings known statically to be in a specified regular language and support operations like concatenation, substitution and coercion. Regular string types can be used to implement, in essentially a conventional manner, the parts of a web application or application framework that construct such command strings. Straightforward type annotations at key interfaces can be used to statically verify that sanitization has been performed correctly without introducing redundant run-time checks. We specify this type system as a minimal typed lambda calculus, λ_{RS} .

To be practical, adopting a type system like this should not require adopting a new programming language. Instead, we favor extensible type systems: new static type systems like this should be distributed as libraries atop a mechanism that guarantees that they can be safely composed. We support this by 1) specifying a translation from λ_{RS} to a language containing only strings and regular expressions, then, taking Python as such a language, 2) implement the type system together with the translation as a library using `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 [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, sani-

tation 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 `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 `stringin`[r_{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 simple can be 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 purpose programming language could stand in for λ_P ; for in-

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

Figure 1: Regular expressions over the alphabet Σ .

stance, 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 here due to space constraints.

$$\begin{array}{ll} \psi ::= \dots & \text{source types} \\ \quad | \text{stringin}[r] & \\ \\ S ::= \dots & \text{source terms} \\ \quad | \text{rstr}[s] & s \in \Sigma^* \\ \quad | \text{rconcat}(S, S) \\ \quad | \text{rreplace}[r](S, S) \\ \quad | \text{rcoerce}[r](S) \end{array}$$

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

$$\begin{array}{ll} \theta ::= \dots & \text{target types} \\ \quad | \text{string} \\ \quad | \text{regex} \\ \\ P ::= \dots & \text{target terms} \\ \quad | \text{str}[s] \\ \quad | \text{rx}[r] \\ \quad | \text{concat}(P, P) \\ \quad | \text{preplace}(P, P, P) \\ \quad | \text{check}(P, P) \end{array}$$

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

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

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

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

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

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

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

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

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

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

Figure 5: Big step semantics for our fragment of λ_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 λ_P . The typing context Θ is standard.

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

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

Figure 7: Big step semantics for our fragment of 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 well-known 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 $s_1s_2 \in \mathcal{L}\{r_1r_2\}$. For all strings s and expressions r , either $s \in \mathcal{L}\{r\}$ or $s \notin \mathcal{L}\{r\}$. Regular languages are closed under complements and concatenation. The regular expressions correspond bijectively to the regular languages.*

Definition 2 (`lsubst`). The function `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 `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 `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 `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 `lsubst(r, s_1, s_2)` $\in \mathcal{L}\{r'\}$ for some r' isomorphic to `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 $P \Downarrow r[x[r]]$ 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 $S \text{ err}$.*

```

1  @fn
2  def sanitize(s : string_in[r'.*']):
3      return (s.replace(r'"', '&quot;')) # TODO: is this right?
4          .replace(r'<', '&lt;')
5          .replace(r'>', '&gt;')
6
7  @fn
8  def results_query(s : string_in[r'[""]*']):
9      return 'SELECT * FROM users WHERE name=' + s + ''
10
11 @fn
12 def results_div(s : string_in[r'^<>]*']):
13     return '<div>Results for ' + s + '</div>'
14
15 def main(db):
16     input = sanitize(user_input())
17     results = db.execute(results_query(input))
18     return results_div(input) + format(results)

```

Figure 8: Regular string types in atlang, a library that enables static type checking for Python.

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

Theorem 9. *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

We now present the main correctness result.

Theorem 10. *If $S : \text{rstr}[r]$ then there exists a P such that $\llbracket s \rrbracket = P$ and either: (a) $P \Downarrow \text{str}[s]$ and $S \Downarrow \text{rstr}[s]$, and $s \in \text{lang}r$; 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 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

3. IMPLEMENTATION IN ATLANG

www Here, we use the argument annotation syntax introduced in Python 3 (a similar syntax is available for Python 2.6+, not shown).

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.

Papers that needs to be cited in this section:

- Ur/Web OSDI paper
- Jif?
- OWASP

```

1 class string_in(atlang.Type):
2     def __init__(self, rx):
3         rx = rx_normalize(rx)
4         atlang.Type.__init__(idx=rx)
5
6     def ana_Str(self, ctx, node):
7         if not in_lang(node.s, self.idx):
8             raise atlang.TypeError("...", node)
9
10    def trans_Str(self, ctx, node):
11        return astx.copy(node)
12
13    def syn_BinOp_Add(self, ctx, node):
14        left_t = ctx.syn(node.left)
15        right_t = ctx.syn(node.right)
16        if isinstance(left_t, string_in):
17            left_rx = left_t.idx
18            if isinstance(right_t, string_in):
19                right_rx = right_t.idx
20                return string_in[lconcat(left_rx, right_rx)]
21            raise atlang.TypeError("...", node)
22
23    def trans_BinOp_Add(self, ctx, node):
24        return astx.copy(node)
25
26    def syn_Method_replace(self, ctx, node):
27        [rx, exp] = node.args
28        if not isinstance(rx, ast.Str):
29            raise atlang.TypeError("...", node)
30        rx = rx.s
31        exp_t = ctx.syn(exp)
32        if not isinstance(exp_t, string_in):
33            raise atlang.TypeError("...", node)
34        exp_rx = exp_t.idx
35        return string_in[lreplace(self.idx, rx, exp_rx)]
36
37    def trans_Method_replace(self, ctx, node):
38        return astx.quote(
39            """__import__(re); re.sub(%0, %1, %2)""" ,
40            astx.Str(s=node.args[0]),
41            astx.copy(node.func.value),
42            astx.copy(node.args[1]))
43
44    def syn_Method_check(self, ctx, node):
45        [rx] = node.args
46        if not isinstance(rx, ast.Str):
47            raise atlang.TypeError("...", node)
48        return string_in[rx.s]
49
50    def trans_Method_check(self, ctx, node):
51        return astx.quote(
52            """__import__(string_in_helper);
53            string_in_helper.coerce(%0, %1)""" ,
54            astx.Str(s=other_t.idx),
55            astx.copy(node))
56
57    def check_Coerce(self, ctx, node, other_t):
58        # coercions can only be defined between
59        # types with the same type constructor,
60        if rx_sublang(other_t.idx, self.idx):
61            return other_t
62        else: raise atlang.TypeError("...", node)

```

Figure 9: Implementation of the `string_in` type constructor in `atlang`.

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

```

1 output of successful compilation

```

Figure 10: Output of successful compilation.

```

1 output of failed compilation

```

Figure 11: Output of failed compilation.