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 command systems like web browsers and database engines using strings. Strings constructed using user input that has not been properly sanitized can thus cause command injection vulnerabilities. In this paper, we introduce regular string types, which classify strings known statically to be in a specified regular language. These types come equipped with common operations like concatenation, substitution and coercion, so they can be used to implement, in essentially a conventional manner, the parts of a web application or application framework that construct command strings. Simple 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 specialized type system like this should not require the adoption of a new programming language. Instead, we favor extensible type systems: new type system fragments 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

Command injection vulnerabilities are among the most common and severe security vulnerabilities in modern web applications [7]. They arise because web applications, at their boundaries, control external systems by issuing commands represented using strings. For example, web browsers are controlled using HTML and Javascript sent from a server as a string, and database engines execute SQL queries also sent as strings. When these commands must include data derived from user input, care must be taken to ensure that the user cannot construct input that subverts the intended

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

command. For example, a SQL query constructed using string concatenation exposes a SQL injection vulnerability:

```
'SELECT * FROM users WHERE name="' + name + '"'
```

If a malicious user enters the name '"; DROP TABLE users --', the entire database could be erased.

To avoid this problem, the program must sanitize user input. For example, in this case, the developer (or, more often, a framework) might define a function sanitize that prepends double quotes (and existing backslashes) with a backslash, which SQL treats as safely. Guaranteeing that user input has already been sanitized before it is used to construct a command is challenging. Note that this function is not idempotent, so it should only be called once.

We observe that most such sanitization techniques can be understood in terms of regular languages [4]. For example, name should be a string in the language described by the regular expression ([^"\]|(\")|(\")| - a sequence of characters other than quotation marks and backslashes; these can only appear escaped. This concrete pattern syntax can be understood to desugar, in a standard way, to the syntax for regular expressions shown in Figure 1, where $r \cdot r$ is sequencing and r + r is disjunction. We will work with this "core" for simplicity.

In this paper, we present a static type system that tracks the regular language a string belongs to. For example, the output of sanitize will be in the regular language described by the regular expression above. We write this as $\mathcal{L}\{r\}$. We take advantage of closure and decidability properties of regular languages to support a number of useful operations on values of such regular string types. These make it possible to implement sanitation protocols like the one just described in an essentially conventional manner. The result is a system where the fact that a string has been correctly sanitized is manifest in its type. Missing calls to sanitization functions are detected statically, and, importantly, so are incorrectly implemented sanitization functions (i.e. these functions need not be trusted). These guarantees require run-time checks only when going from less precise to more precise types (e.g. at the edges of the system, where user input has not yet been validated).

We will begin in Sec. 2 by specifying this type system minimally, as a conservative extension of the simply typed lambda calculus called λ_{RC} . This allows us to specify the guarantees that the type system provides precisely. We also formally specify a translation from this calculus to a typed calculus with only standard strings and regular expressions, intending it as a guide to language implementors interested in building this feature into their own languages. This also

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

Figure 1: Regular expressions over the alphabet Σ .

```
\begin{array}{lll} \sigma & ::= & \sigma \rightarrow \sigma & \text{source types} \\ & | & \mathsf{stringin}[r] & \\ e & ::= & x & \text{source terms} \\ & | & \lambda x : \sigma . e & \\ & | & e & (e) & \\ & | & \mathsf{rstr}[s] & s \in \Sigma^* \\ & | & \mathsf{rconcat}(e,e) & \\ & | & \mathsf{rreplace}[r](e,e) & \\ & | & \mathsf{rcoerce}[r](e) & \end{array}
```

Figure 2: Syntax of λ_{RS} .

demonstrates that no additional space overhead is required.

Waiting for a language designer to build this feature in is unsatisfying. Moreover, we would also face a "chicken-and-egg problem": justifying its inclusion into a commonly used language requires empirically demonstrating that it is useful, but this is difficult to do if developers have no way to use it in practice. As such, we take the position that a better path forward for the community is to work within a programming language where such type system fragments can be introduced modularly and orthogonally, as libraries.

In Sec. 3, we show how to implement the type system fragment we have specified using atlang, an extensible static type system implemented as a library inside Python. atlang leverages local type inference to control the semantics of literal forms, so reguar string types can be introduced using string literals without any run-time overhead. Coercions that are known to be safe due to a sublanguage relationship are performed implicitly, also without run-time overhead. This results in a usably secure system: working with regular strings differs little from working with standard strings.

We conclude after discussing related work in Sec. 4.

2. REGULAR STRING TYPES, MINIMALLY

In this section, we define a minimal typed lambda calculus with regular string types called λ_{RS} . This will serve as the source language for our translation to a calculus with only standard strings and regular expressions, defined in Sec. ??.

The syntax of λ_{RS} is given in Figure 2, its static semantics are specified in Figure 3 and its evaluation semantics are specified in Figure 4.

The system has regular expression types $\mathsf{stringin}[r]$ whre r is a regular expression. Expressions of this type evaluate to string literals in the language described by r. Operations on expressions of type $\mathsf{stringin}[r]$ preserve this property.

The premier operation for manipulating strings in λ_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, $\mathtt{lsubst}(a|b,a,c)=c$. In order to computer the type resulting from aubstitution, we also need to compute the result of replacing on language with another inside a given language. Finally, just for convienance, we provide a coerce operation. The introduction of coercion requires handling of runtime errors.

The underlying language λ_P has only one type for strings.

```
\Psi ::= \emptyset \mid \Psi, x : \sigma
\Psi \vdash S : \sigma
                                      S-T-Stringin-I
                                                 s \in \mathcal{L}\{r\}
                                      \Psi \vdash \mathsf{rstr}[s] : \mathsf{stringin}[r]
                 S-T-Concat
                  \Psi \vdash S_1 : \mathsf{stringin}[r_1]
                                                              \Psi \vdash S_2 : \mathsf{stringin}[r_2]
                         \Psi \vdash \mathsf{rconcat}(S_1, S_2) : \mathsf{stringin}[r_1 \cdot r_2]
                 S-T-Replace
                  \Psi \vdash S_1 : \mathsf{stringin}[r_1]
                                                              \Psi \vdash S_2 : \mathsf{stringin}[r_2]
                                    lreplace(r, r_1, r_2) = r'
                          \Psi \vdash \mathsf{rreplace}[r](S_1, S_2) : \mathsf{stringin}[r']
                               S-T-coerce
                                         \Psi \vdash S : \mathsf{stringin}[r']
                                \Psi \vdash \mathsf{rcoerce}[r](S) : \mathsf{stringin}[r]
```

Figure 3: Typing rules for our fragment of λ_S . The typing context Ψ is standard.

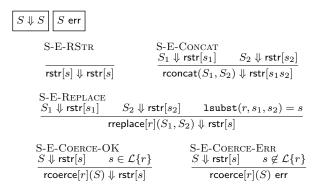


Figure 4: Big step semantics for our fragment of λ_S . Error propagation rules are omitted.

We prove that whenever a term is translated from λ_S to λ_P , correctness is preserved. The only exception is in the case of unsafe casts in λ_S , which are unnecessary but are included to demonstrate that the regex library of λ_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 λ_S, λ_P and the translation from λ_S to λ_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 λ_P as well as both type safety and correctness for λ_S .
- In §2.3 we prove that translation preserves the correctness result about λ_S.

Figure 10: Translation from source terms (S) to target terms (P). The translation is type-directed in the Tr-Coerce cases.

$$\begin{array}{ll} \theta & ::= \theta \rightarrow \theta & \text{target types} \\ \mid & \text{string} \\ \mid & \text{regex} & \end{array}$$

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

Figure 5: Syntax for the fragment of our target language, λ_P , containing strings and statically constructed regular expressions.

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

$$\begin{array}{|c|c|} \hline P \Downarrow P & \hline P \text{ err} \\ \hline \\ \hline P\text{-E-STR} & P\text{-E-Rx} & P\text{-E-Concat} \\ \hline str[s] \Downarrow \text{str}[s] & rx[r] \Downarrow \text{rx}[r] & \hline \\ \hline P\text{-E-Replace} & P_1 \Downarrow \text{rx}[r] \\ \hline \hline P\text{-E-Replace} & P_1 \Downarrow \text{rx}[r] \\ \hline P_2 \Downarrow \text{str}[s_2] & P_3 \Downarrow \text{str}[s_3] & \text{1subst}(r,s_2,s_3) = s \\ \hline \hline & \text{preplace}(P_1,P_2,P_3) \Downarrow \text{str}[s] \\ \hline \hline & P\text{-E-Check-OK} \\ \hline & P_1 \Downarrow \text{rx}[r] & P_2 \Downarrow \text{rstr}[s] & s \in \mathcal{L}\{r\} \\ \hline & \text{check}(P_1,P_2) \Downarrow \text{str}[s] \\ \hline & P\text{-E-Check-Err} \\ \hline & P_1 \Downarrow \text{rx}[r] & P_2 \Downarrow \text{str}[s] & s \not\in \mathcal{L}\{r\} \\ \hline & \text{check}(P_1,P_2) & \text{err} \\ \hline \end{array}$$

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 replies 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 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 $s_1s_2 \in \mathcal{L}\{r_1r_2\}$. 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 (1subst). The replation 1subst $(r, s_1, s_2) = s$ produces a string s in which all substrings of s_1 matching r are replaced with s_2 .

Definition 3 (Ireplace). The relation Ireplace $(r, r_1, r_2) = r'$ relates r, r_1 , and r_2 to a language r' containing all strings of r_1 except that any substring $s_{pre}ss_{post} \in \mathcal{L}\{r_1\}$ where $s \in \mathcal{L}\{r\}$ is replaced by the set of strings $s_{pre}s_2s_{post}$ 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 $\mathcal{L}\{\mathsf{lreplace}(r, r_1, r_2)\}$ is also a regular language.

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

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

Proof. The theorem follows from the definitions of lsubst and lreplace; note that language substitutions over-approximate string substitutions. \Box

2.2 Safety of the Source and Target Languages

Lemma 6. If $\Psi \vdash S$: 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.

Lemma 7. If $\Theta \vdash P$: regex then $P \Downarrow rx[r]$ 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 $\Psi \vdash S$: stringin[r] then $S \Downarrow \mathsf{rstr}[s]$ and $\Psi \vdash \mathsf{rstr}[s]$: stringin[r]; or else S err.

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

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

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. \Box

In addition to safety, we proof a correctness result for λ_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 λ_S . If $\Psi \vdash S$: stringin[r] and $S \Downarrow \mathsf{rstr}[s]$ then $s \in \mathcal{L}\{r\}$.

Proof. Follows directly from type safety, canonical forms for $\lambda \varsigma$.

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 err.

2.3 Translation Correctness

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

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

Theorem 12. Correctness of Input Sanitation for Translated Terms. If $\llbracket S \rrbracket = P$ and $\Psi \vdash S$: stringin[r] then either P err or $P \Downarrow \mathsf{str}[s]$ for $s \in \mathcal{L}\{r\}$.

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

```
1
    class string_in(atlang.Type):
      def __init__(self, rx):
3
        rx = rx_normalize(rx)
4
        atlang.Type.__init__(idx=rx)
5
6
      def ana_Str(self, ctx, node):
        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)
13
      def syn_BinOp_Add(self, ctx, node):
14
        left_t = ctx.syn(node.left)
15
        right_t = ctx.syn(node.right)
        if isinstance(left_t, string_in):
16
17
          left_rx = left_t.idx
18
          if isinstance(right_t, string_in):
19
            right_rx = right_t.idx
            return string_in[lconcat(left_rx, right_rx)]
21
        raise atlang.TypeError("...", node)
23
      def trans_BinOp_Add(self, ctx, node):
24
        return astx.copy(node)
      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)
        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)""",
          astx.Str(s=node.args[0]),
40
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
      def trans Method check(self. ctx. node):
50
51
        return astx.quote(
52
               _import__(string_in_helper);
          string_in_helper.coerce(%0, %1)""",
53
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
50
          types with the same type constructor,
60
        if rx_sublang(other_t.idx, self.idx):
61
          return other_t
        else: raise atlang.TypeError("...", node)
```

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

1 output of successful compilation

Figure 10: Output of successful compilation.

1 output of failed compilation

Figure 11: Output of failed compilation.

3. IMPLEMENTATION AS A LIBRARY IN ATLANG

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

4. RELATED WORK

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 lanuage 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.

5. CONCLUSION

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..

6. 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.

7. 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. HorvÃath, and V. ZsÃşk, editors, Implementation and Application of Functional Languages, volume 5083 of Lecture Notes in Computer Science, pages 75–92. Springer Berlin Heidelberg, 2008.