# Statically Typed String Sanitation Inside a Python

Nathan Fulton   Cyrus Omar   Jonathan Aldrich

## Problem

Web applications must ultimately command systems, such as web browsers and database engines, using strings. Strings derived from improperly sanitized user input are therefore a potential vector for command injection attacks.

## Milieu

To address the largest security threat facing today's web applications:

- **Developers** use libraries and frameworks which ultimately **ground out to operations on strings**.
- **Researchers** propose information flow and taint analyses, which are are often **attack-specific**, can become complicated, and do not generalize to arbitrary validation tasks (e.g., "is this a unix file path?").

## Approach

We introduce **regular expression types** for classifying strings and equip these types with standard operations. Our approach makes it possible to specify and verify correctness of conventional implementations of input sanitation procedures.

## References

N. Fulton, C. Omar, and J. Aldrich. Statically typed string sanitation inside a python. First International Workshop on Privacy and Security in Programming (PSP). ACM, 2014.

## Two Illustrative Excerpts

String **concatenation** is typed using regular expression concatenation:

$$\text{S-T-Concat} \quad \frac{\Psi \vdash e_1 : \text{stringin}[r_1] \qquad \Psi \vdash e_2 : \text{stringin}[r_2]}{\Psi \vdash \text{rconcat}(e_1; e_2) : \text{stringin}[r_1 \cdot r_2]}$$

$$\text{S-E-Concat} \quad \frac{e_1 \Downarrow \text{rstr}[s_1] \qquad e_2 \Downarrow \text{rstr}[s_2]}{\text{rconcat}(e_1; e_2) \Downarrow \text{rstr}[s_1 s_2]}$$

**Substring** operations pattern match on the head of a string. Regular expression derivatives provide a natural approximation (ltl is roughly the derivative of lhd):

$$\text{S-T-Case} \quad \frac{\Psi \vdash e_1 : \text{stringin}[r] \qquad \Psi \vdash e_2 : \sigma \qquad \Psi, x : \text{stringin}[\text{lhd}(r)], y : \text{stringin}[\text{ltl}(r)] \vdash e_3 : \sigma}{\Psi \vdash \text{rstrcase}(e_1; e_2; x, y.e_3) : \sigma}$$

$$\text{S-E-Case-}\epsilon \quad \frac{e_1 \Downarrow \text{rstr}[\epsilon] \qquad e_2 \Downarrow v_2}{\text{rstrcase}(e_1; e_2; x, y.e_3) \Downarrow v_2}$$

$$\text{S-E-Case-Concat} \quad \frac{e_1 \Downarrow \text{rstr}[as] \qquad [\text{rstr}[a], \text{rstr}[s]/x, y]e_3 \Downarrow v_3}{\text{rstrcase}(e_1; e_2; x, y.e_3) \Downarrow v_3}$$

This system ($\lambda_{RS}$) also contains **replacement**, **checked casts** and **dynamically checked coercions**. See the paper for details[1].

## Implementation Example

We are working toward a regular string types library for the extensible programming language Atlang.

```
1  @fn
2  def sanitize(s : stringin[r'.*']):
3      return (s.replace(r'"', '&quot;')
4              .replace(r'<', '&lt;')
5              .replace(r'>', '&gt;'))
6
7  @fn
8  def query(s : stringin[r'[^"]*']):
9      return 'SELECT * FROM users WHERE ' +
10             'name="' + s + '"'
11 @fn
12 def div(s : stringin[r'[^<>]*']):
13     return '<div>Results for '+s+'</div>'
14
15 @fn
16 def main():
17     input = sanitize(user_input())
18     results = db_execute(query(input))
19     return div(input) + format(results)
```

# Carnegie Mellon University