# 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.
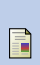
## Milieu

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** and do not generlize to arbitrary validation tasks (e.g. "is this a unix file path?").

## Approach

We introduce **regular expression types** for classifing 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. SPLASH '14. ACM, 2014.

**Carnegie Mellon University**

## 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):

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

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

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

$\lambda_{RS}$ also contains **replacement**, **checked casts** and **dynamically checked coercions**.

## Implementation Example

We are working toward an 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 results_query(s : stringin[r'[^"]*']):
9    return 'SELECT * FROM users WHERE name="' + s + '"'
10
11 @fn
12 def results_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(results_query(input))
19   return results_div(input) + format(results)
```