# Statically Typed String Sanitation Inside a Python

**Nathan Fulton**
Cyrus Omar
Jonathan Aldrich

Carnegie Mellon University

# The Problem

Applications use **strings** to build `SQL` commands

```
sql_exec("SELECT * FROM users WHERE" +
    "username = " + input1 + " AND " +
    "password = " + input2)
```

# The Problem

Applications use **strings** to build `HTML` commands

```
print("You searched for: " + keyword)
```

# The Problem

Applications use **strings** to build `JS`      commands

```
print("<script>" +
    "document.getElementById(" +
    "'" + input + "'"   +
    ")" + "..." +
    "</script>")
```

# **The Problem**

Applications use **strings** to build `shell` commands

```
call("cat " + input)
```

Arbitrary strings are dangerous.

# Existing Solutions

- Web Frameworks

# Existing Solutions

- Web Frameworks
  - may contain bugs

# Existing Solutions

- Web Frameworks
  - may contain bugs
- Prepared Statements

# Existing Solutions

"Drupal is an open source content management platform powering millions of websites… During a code audit of Drupal extensions for a customer **an SQL Injection was found in the way the Drupal core handles prepared statements**. A malicious user can inject arbitrary SQL queries… This leads to a code execution as well."

- Stefan Horst, 6 days ago

# Existing Solutions

- Web Frameworks
  - may contain bugs
- Prepared Statements
  - may contain bugs

# Existing Solutions

- Web Frameworks
  - may contain bugs
- Prepared Statements
  - may contain bugs
- Problem specific parsers

# Existing Solutions

"Three of our Sports API servers had malicious code executed on them… This mutation happened to exactly fit a **command injection bug in a monitoring script** our Sports team was using at that moment to **parse and debug their web logs**."

- Alex Stamos (Yahoo! CISO), two weeks ago

# Existing Solutions

- Web Frameworks
  - may contain bugs
- Prepared Statements
  - may contain bugs
- Problem specific parsers
  - may contain bugs

The Goal: A **general** approach for specifying and verifying input sanitation procedures, **with a minimal trusted core**.

Arbitrary strings are dangerous.

Static reasoning about strings is easy!

# Regular Expression Types

Python, Java, etc:

string

Lambda RS:

string[regex]

# Contributions

- Regular Expression Types corresponding to common string and regex library operations.
- Translation into a language with a bare string type.

Together, these define a **type system extension** which is implemented in the extensible programming language atlang.

# **Typing Rule for String Literals**

If:

● *s* in a string in the language of *r*

Then:

● rstr[s] has type stringin[r].

# Typing Rule for String Literals

$$\frac{s \in \mathcal{L}\{r\}}{\Psi \vdash \mathsf{rstr}[s] : \mathsf{stringin}[r]}$$

# The Security Theorem

If e has type stringin[r], then e evaluates to a string (denoted rstr[s]) such that s ∈ L(r).

```
"""this function will remove quotes."""
def sanitize(s : string): s //TODO

def get_user(u : string):
    sql_exec("select * from users where " +
        "username = '" + u + "'")
```

```
"""this function will remove quotes."""
def sanitize(s : string): s //TODO


def get_user(u : string):
   sql_exec("select * from users where " +
      "username = '" + u + "'")


x = "';DELETE FROM users--"
get_user(sanitize(x))
```

```
"""this function will remove quotes."""
def sanitize(s : string): s //TODO


def get_user(u : string[!']):
   sql_exec("select * from users where " +
      "username = '" + u + "'")


x = "';DELETE FROM users--"
get_user(sanitize(x))
^ type error! L(.*) is not in L(!')
```

```python
"""this function will remove quotes."""
def sanitize(s : string) -> stringin[!']:
  s.replace(r"'", "")


def get_user(u : string[!']):
   sql_exec("select * from users where " +
      "username = '" + u + "'")


x = "';DELETE FROM users--"
get_user(sanitize(x))
^ OK!
```

# Regular Expressions

r ::= a | r·r | r ++ r | r*

# Regular Languages

r ::= a | r·r | r ++ r | r*

L(psp)    = {psp}
L(ps*p)   = {pp, psp, pssp, psssp, ...}
L(a ++ b) = {a, b}

# Regexes as Specs

Often Unstated Specifications:

! '

# Regexes as Specs

Often Unstated Specifications:

```
!'
(a|b|c|...)*
```

# Regexes as Implementations

Often Unstated Specifications:

```
!'

(a|b|c|...)*
```

Implementations:

```
replace(!', "", input)
```

**Unstated Assertion: implementation meets specification.**

# The Core Language (1 / 2)

| Construct | Abstract Syntax | A Python |
|---|---|---|
| **Concat** | `rconcat(e1;e2)` | `e1 + e2` |
| **Substring** | `rstrcase(e1;`<br>`  e2;`<br>`  x,y.e3)` | `if e1 == "":`<br>`  e2`<br>`else:`<br>`  e3(e1[:1], e1[1:])` |
| **Replace** | `rreplace[r](e1; e2)` | `e1.sub(r"r", e2)` |

# The Core Language (2 / 2)

| Concept | Abstract Syntax | A Python |
|---|---|---|
| **Coercion** | `rcoerce[r](e)` | `e` |
| **Checks** | `rcheck[r](e;`<br>`x.e1; e2)` | `if re.search(r"r",e) == None:`<br>`  e2`<br>`else:`<br>`  e1(e)` |

# λ<sub>RS</sub>

**String Concatenation**
   rconcat(e; e)

**Substrings**
   rstrcase(e; e; x,y.e)

**Substitution**
   rreplace[r](e; e)

**Coercions**
   rcoerce[r](e)

**Checked Casts**
   rcheck[r](e; x.e; e)

# String Concatenation

Recall: if e has type stringin[r] then e evaluates to v and v ∈ L(r).

# String Concatenation

Recall: if e has type stringin[r] then e evaluates to v and v $\in$ L(r).

If:

- $e_1$ : stringin[$r_1$]
- $e_2$ : stringin[$r_2$]

then:

- concat($e_1$; $e_2$) : stringin[$r_1 r_2$].

# String Concatenation

Recall: if e has type stringin[r] then e evaluates to v and v $\in$ L(r).

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

# Example Typing Derivation

$$\dfrac{\dfrac{a \in \mathcal{L}\{a*\}}{\Psi \vdash \mathsf{rstr}[a] : \mathsf{stringin}[a*]} \qquad \dfrac{b \in \mathcal{L}\{b*\}}{\Psi \vdash \mathsf{rstr}[b] : \mathsf{stringin}[b*]}}{\Psi \vdash \mathsf{rconcat}(r; \mathsf{rstr}[a])\mathsf{rstr}[b] : \mathsf{stringin}[a * b*]}$$

# Substrings

```
""" S = state code then D.O.B. """
def get_state(s : stringin[(a-z0-9)*]):
    rstrcase(s; ''; x + rstrcase(y; ''; x))
```

# Substrings

```
get_state("WI1956")
```

# Substrings

```
        get_state("WI1956")
                ⇓

rstrcase("WI1956"; ''; x + rstrcase(y; ''; x))
```

# Substrings

```
get_state("WI1956")
          ⇓
rstrcase("WI1956"; ''; x + rstrcase(y; ''; x))
          ⇓
   "W" + rstrcase("I1956"; ''; x)
```

# Substrings

```
get_state("WI1956")
           ⇓
rstrcase("WI1956"; ''; x + rstrcase(y; ''; x))
           ⇓
    "W" + rstrcase("I1956"; ''; x)
           ⇓
       "W" + "I" = "WI"
```

# Substrings

"Get the first n characters of a string s"

# Substrings

"Get the **first** character of a string s"

"Get everything after the first character of s"

# Substrings

"Get the **first** character of a string s"

```
lhead(r)            = lhead(r, ε)
lhead(ε, r')        = ε
lhead(a, r')        = a
lhead(r1·r2, r')    = lhead(r1, r2)
lhead(r1 + r2, r')  = lhead(r1, r') + lhead(r2, r')
lhead(r*, r')       = lhead(r', ε) + lhead(r, ε)
```

# Substrings

"Get the **first** character of a string s"
```
lhead(r)            = lhead(r, ε)
lhead(ε, r')        = ε
lhead(a, r')        = a
lhead(r1·r2, r')    = lhead(r1, r2)
lhead(r1 + r2, r') = lhead(r1, r') + lhead(r2, r')
lhead(r*, r')       = lhead(r', ε) + lhead(r, ε)
```

"Get everything after the first character of s"

$$\delta_a(r) + \delta_b(r) + \delta_c(r) + ...$$

# Substrings

Observation: If s ∈ L((a-z)*(0-9)) then
get_state(rstr[s]) ⇓ rstr[t] such that t ∈ (a-z0-9)*.

# Substrings

Observation: If s $\in$ L((a-z)*(0-9)) then

get_state(rstr[s]) $\Downarrow$ rstr[t] such that t $\in$ (a-z0-9)*.

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

# On the precision of rstrcase

Note that lhead(r)·ltail(r) ≠ r.

# On the precision of rstrcase

Note that lhead(r)·ltail(r) ≠ r.

Example: Choose r = (ab)+(cd), so "ad" ∉ L(r).

Note that:

```
lhead(r) = a + c
ltail(r) = δ (r) + δ (r)
           a        c
         = b + d
```

Therefore, "ad" ∈ L(lhead(r)·ltail(r)).

# String Replacement

$$
\begin{array}{c}
\text{S-E-REPLACE} \\
\dfrac{e_1 \Downarrow \mathsf{rstr}[s_1] \qquad e_2 \Downarrow \mathsf{rstr}[s_2] \qquad \mathbf{subst}(r; s_1; s_2) = s}{\mathsf{rreplace}[r](e_1; e_2) \Downarrow \mathsf{rstr}[s]}
\end{array}
$$

subst(r; s1; s2) reads "substitute s2 for r in s1"

# String Replacement

$$\text{S-T-REPLACE}$$

$$\frac{\Psi \vdash e_1 : \mathsf{stringin}[r_1] \qquad \Psi \vdash e_2 : \mathsf{stringin}[r_2] \qquad \mathsf{lreplace}(r; r_1; r_2) = r'}{\Psi \vdash \mathsf{rreplace}[r](e_1; e_2) : \mathsf{stringin}[r']}$$

# String Replacement

Key Fact: lreplace and subst correspond:

$$subst(r, s1, s2) \text{ is in } lreplace(r, r1, r2)$$

where:
- s1 ∈ r1, and
- s2 ∈ r2.

# String Replacement

subst(r, s1, s2) is in Ireplace(r, r1, r2).

This does **not** entail a definition of Ireplace given a definition of subst.

# Saturation

```
replace("ee", "Kleeene", "e")
```

replace ee in "Kleeene" with e

= "Kleeene"

# Translation

$$\text{Tr-Concat}$$

$$\frac{[\![e_1]\!] = \iota_1 \qquad [\![e_2]\!] = \iota_2}{[\![\mathsf{rconcat}(e_1; e_2)]\!] = \mathsf{concat}(\iota_1; \iota_2)}$$
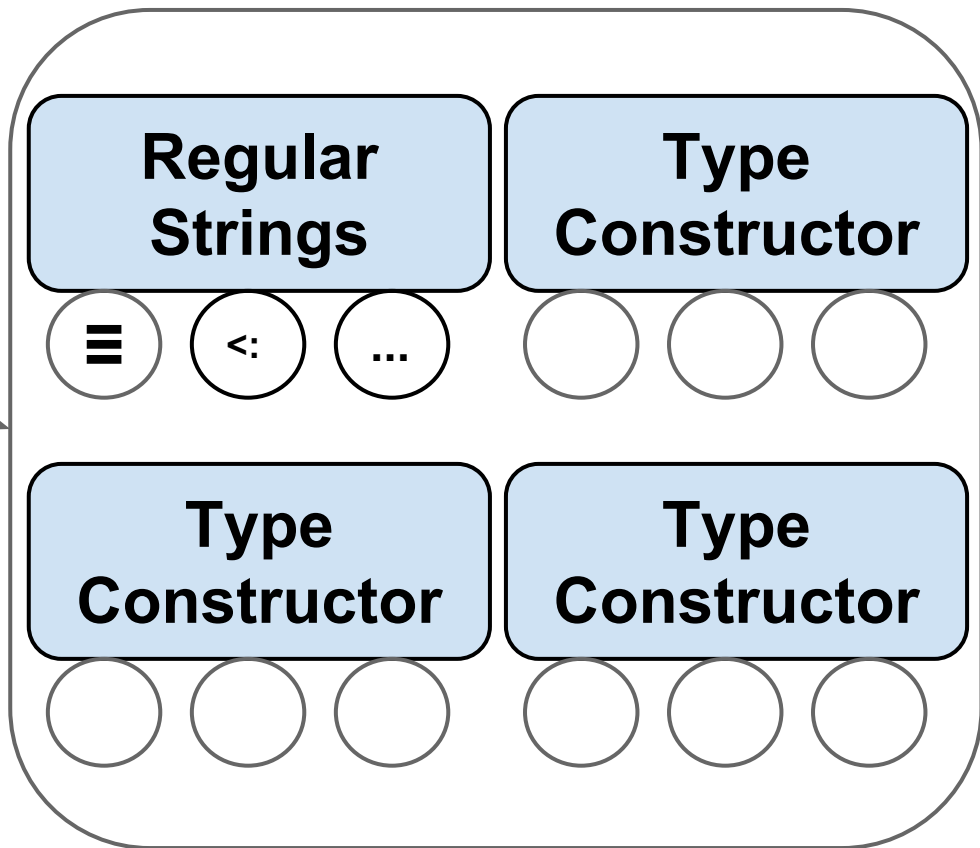
# Translation

Translation defines either an embedding (as a language extension) or, alternatively, an erasure.

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

# Conclusions

Constrained String Types are a *general* approach for specifying and verifying input sanitation procedures.

Unlike other approaches, constrained strings only require a minimal trusted core.

# Future Work

1. Implement a static analysis and verify a realistic query builder.
2. Application of replacement operation to program repair in dynamic logic over trace semantics.
   - replacement on hybrid regular programs.
3. Explore other privacy & security applications of extensible type systems.

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