

# String Solving with Word Equations and Transducers:

*Towards a Logic for Analysing Mutation XSS*

Anthony W. Lin (Yale-NUS), Pablo Barcelo (Univ. of Chile)

# String Solving: A View on the Landscape

# What are String Solvers?

*Solvers for Satisfiability Modulo Theory (SMT) over strings*

**Domain:** the set of all words over  $\Sigma$

**Operations:** concatenation, regex matching, length constraints, replace, replace-all, string transductions, ...

$$s2 = s1.s1 \wedge \text{len}(s2) = \text{len}(s7) \wedge \dots$$

*A different combination of operations gives rise to a different theory over strings!! (Just as for integer domain)*

Many string solvers: CVC, HAMPI, Kaluza, Kudzu, Norn, Pex/Z3, PISA, S3, Saner, Stranger, StrSolve, SUSHI, Z3-str, ...

# Why Develop String Solvers?

- Static analysis of security vulnerabilities in web applications against code injection and XSS
  - *Caused by improper handling of untrusted strings*
- Automatic test case generation for scripting languages
- Path query languages for graph databases

# String Solving: Theory vs. Practice

- Faster heuristics each year
- Much less progress on theory

Which SMT over strings is decidable?

1. Word equations (Makanin'77)

$$s2.a.s3.s4 = s1.s3.s2.b$$

2. Existential theory strings with concat (Buchi&Senger'90)

$$s2 = s1.s1 \wedge s3.s2 \neq s1.s7.s8$$

3. Word equations with regex matching (Schulz'90)

$$s2.a.s3.s4 = s1.s3.s2.b \wedge s1 \in (ab)^*$$

Open Problem: Decidability of Word equations with length constraints

The need to add string  
transductions

# Cross-Site Scripting (XSS)



Created a user profile with name:  
"`<script>  
window.open('http://evilsite.com')  
</script>`"

send a friend request to Dilbert



```
...  
// Obtain friend request from server  
element.innerHTML = friendName;  
...
```



# Sanitising Input Data

- Escape certain characters
- *EVERY* occurrence of `<` should be changed to `&lt;`;
- *EVERY* occurrence of `>` should be changed to `&gt;`;

*A kind of “replace-all” operation*



# Adding Sanitisation

## Google Closure

```
1 var x = goog.string.htmlEscape(friendName);  
2 element.innerHTML = x
```

`<script>...</script>`

will be converted to

`&lt;script&gt;...&lt;script&gt;`

The script won't be executed by Dilbert's browser

# A more tricky example

```
1 var x = goog.string.htmlEscape(name);
2 var y = goog.string.escapeString(x);
3 nameElem.innerHTML = '<a onclick=' +
4   '"viewPerson(\'' + y + '\')">' + x + '</a>';
```

*(Adapted from Kern'14)*

escapeString “backslash-escape” certain metacharacters

‘ is replaced by **&#39;** or **\**

“ is replaced by **&#34;** or **\**”

Q: Is this code vulnerable to XSS?

# Analysis of the code

```
1 var x = goog.string.htmlEscape(name);
2 var y = goog.string.escapeString(x);
3 nameElem.innerHTML = '<a onclick=' +
4   '"viewPerson(\'\' + y + '\')"'>' + x + '</a>';
```

*SWAP*

INPUT 1: name being **Tom & Jerry** gives HTML markup

```
<a onclick="viewPerson('Tom & Jerry')">Tom & Jerry</a>
```

INPUT 2: name being **' );alert(1);//** gives HTML markup

```
<a onclick="viewPerson('&#39;);alert(1);//')">&#39;);alert(1);//</a>
```

innerHTML “mutates” this string to

```
<a onclick="viewPerson('');alert(1);//')">');alert(1);//</a>
```

**XSS!**

# Detecting XSS via a String Solver

**Step 1:** Identify “sink variables” (innerHTML, document.write)

```
1 var x = goog.string.htmlEscape(name);  
2 var y = goog.string.escapeString(x);  
3 nameElem.innerHTML = '<a onclick=' +  
4   '"viewPerson(\'' + y + '\')"'>' + x + '</a>';
```

**Step 2:** Find “attack patterns” from known vulnerabilities (eg, OWASP)

$e1 = /<a\ onclick="viewPerson\(' ( ' | [^']*[^\\]' ) \); [^']*[^\\]' )">.*</a>/$

**Step 3:** Express the program logic in a string logic:

1.  $x = R1(name)$
2.  $y = R2(x)$
3.  $z = w1 . y . w2 . x . w3$
4.  $nameElem.innerHTML = R3(z)$
5.  $nameElem.innerHTML$  matches  $e1$

**Step 4:** Check for satisfiability

# Which String Logic?

1.  $x = R1(\text{name})$
2.  $y = R2(x)$
3.  $z = w1 . y . w2 . x . w3$
4.  $\text{nameElem.innerHTML} = R3(z)$
5.  $\text{nameElem.innerHTML}$  matches  $e1$

concatenation



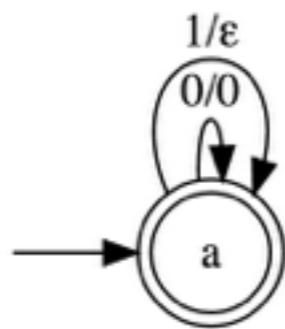
R1, R2, R3 - replace-all kind of operations



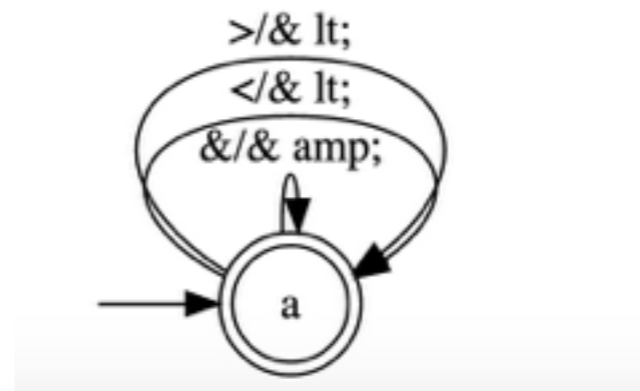
**String transductions!**

# Finite-state I/O Transducers

Just like finite-state automaton, but the transition label is a pair of words:  $v/w$



Erases 1



Replaces some reserved characters by HTML entity names

Relation recognised by  $A$  is

$$\{(v_1 \cdots v_n, w_1 \cdots w_n) : q_0 \xrightarrow{v_1/w_1} \cdots \xrightarrow{v_n/w_n} q_n \in F\}$$

# Modelling sanitisation functions and implicit browser transductions

Lots of works modelling these as FST or extensions  
thereof:

- Saxena et al, *S&P'10*
- D'Antoni&Veanes, *VMCAI'13*
- Hooimejer et al., *USENIX Security'11*
- Veanes et al., *POPL'11*
- ...

Is theory of strings with  
concatenation and FST  
decidable?



# Undecidability

**Proposition** (BFL'13): Checking if the constraint  
 $x = y.z \ \& \ x = R(z)$   
for a transduction  $R$ , is satisfiable is undecidable

**Proposition:** Undecidability still holds when only allowing “erasing” transducers (i.e. replace  $A$  with an empty string)

# The Straight-Line Fragment (SSA Form)

## Inductive Definition:

*(Base)* An empty set  $T$  of conjuncts is in SL

*(Inductive)* If  $S$  is in SL with variables

$$x_1, \dots, x_m$$

then  $S \wedge x_{m+1} = P_{m+1}$  is in SL, where

$$P_{m+1} = R(y) \quad \text{OR} \quad P_{m+1} = y_1 \cdots y_n$$

where the  $y$ 's are variables in  $S$  or new variables

**regex matching:** a boolean combination of

# Decidability of SL

**Theorem:** SATISFIABILITY for the class SL is decidable in exponential space (double-exponential-time)

In fact, EXPSPACE-complete

**Theorem** (Bounded Model Property):  
Every satisfiable constraint in SL has a solution of double-exponential size

**Provides some completeness guarantee of several existing string solvers**

Under a reasonable assumption, we get a single-exponential bound

# Proof idea for decidability (without regex matching)

**Step 1:** Remove concatenation from the formula

$$y = xx \wedge z = R(y)$$

where  $R$  has states  $q_0, \dots, q_n$

$$y_1 = x$$

$$y_2 = x$$

$$\bigvee_{j=0}^n (y_1 = R^{q_0, q_j}(x) \wedge y_2 = R^{q_j, q_n}(x))$$

# Bound on the size of formula without concatenation

“Doubling” Trick

$$y_1 = y_0 y_0$$

$$y_2 = y_1 y_1$$

$$y_3 = y_2 y_2$$

$$L(y_3)$$

Resulting formula uses

$$2^3 + 2^2 + 2^1 + 3 = 17$$

variables

Can use this trick to encode EXPSPACE Turing machines

# Solving the final formula

$$y_1 = x$$

$$y_2 = x$$

$$\bigvee_{j=0}^n (y_1 = R^{q_0, q_j}(x) \wedge y_2 = R^{q_j, q_n}(x))$$

Acyclic (straight-line)

Satisfiability for this kind of formulas is decidable

Post/pre images of regular languages under FST are regular

# Improving the upper bound

The doubling tricks are artificial

Limiting them into a bounded height is reasonable in practice

All the examples we've seen in practice are of height at most 4

**Theorem:** SATISFIABILITY for the restricted SL is decidable in polynomial space (exponential-time)

**Theorem** (Bounded Model Property):  
Every satisfiable constraint in restricted SL has a solution of exponential size

Extending the logic



# Adding integer constraints

Constraints of the form

$$a_1 t_1 + \dots + a_n t_n \leq d$$

where

$a_i$  is a constant integer

$t_i$  is either:

- 1) an integer variable,
- 2)  $|x|$  for some string variable  $x$
- 3)  $|x|_a$  for some string variable  $x$

# Decidability

**Theorem:** SATISFIABILITY for the class SL with integer constraints is decidable in exponential space

In fact, EXPSPACE-complete

**Theorem** (Bounded Model Property):

Every satisfiable constraint in SL with integer constraints has a solution of double-exponential size

# Conclusion and Future Work

- Concatenation and string transductions are both important for XSS applications
- Straight-line fragment of string logic with concatenation and transductions (and even with integer constraints) is decidable
- **Future work 1:** an algorithm for computing a better estimate of the maximum size of solutions
- **Future work 2:** study the extension with symbolic transducers
- **Future work 3:** A more precise model of sanitisation functions and implicit browser transductions as transducers