# Static Syntax Validation for Code Generation with String Templates

## Paper Contribution to SDL 2017

Dorian Weber     Joachim Fischer

Institut für Informatik
Humboldt-Universität zu Berlin

October 9th, 2017

## *Observations*

- core tenet of model-driven engineering is the use of domain models to represent abstract knowledge about a particular application domain
- applications ultimately require a model-to-model or model-to-text transformation

# *Problem*

- performing a model-to-text transformation: what can we say about the generated text?
    1. does its syntax conform to some context-free grammar?
    2. are identifiers declared before being used?
    3. do their data types permit the use in generated operations?
    4. can the execution result in undefined behaviour?
- questions like these may only be answered with an automated proof, not through testing
- we focus on syntax

## HTML in Xtend

```
1    '''
2      <body>
3        <h2>«module.name»</h2>
4        «FOR expr: module.member»
5          <h3>«expr.name»</h3>
6          <dl>
7            «FOR it: expr.member»
8              «IF it.univerval»
9                <dt>«it.name»</dt>
10               <dd>=«it.value»</dt>
11             «ELSEIF it.extensible»
12               <dt />
13               <dd><i>Enumeration is extensible</i></dd>
14             «ENDIF»
15           «ENDFOR»
16         </dl>
17       «ENDFOR»
18     </body>
19   '''
```

## C enumeration in Xtend

```
 1    '''
 2        typedef enum «expr.name» {
 3          «FOR it: expr.member»
 4            «IF it.univerval»
 5              «it.name» = «it.value»,
 6          «ELSEIF it.extensible»
 7              /*
 8               * Enumeration is extensible
 9               */
10          «ENDIF»
11        «ENDFOR»
12        } e_«expr.name»;
13    '''
```

## Research Question (informal)

*Prior to seeing any meta-model instance, under which circumstances can we guarantee that a string template will expand to syntactically correct code?*

*Introduction*
**Context Free Grammars and String Templates**

Formal Prerequisites
*Real-World String Templates*
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

# Context Free Grammar

## Definition (CFG)

A *Context Free Grammar* is defined by the tuple $(V, \Sigma, P, V_0)$ where

- $V$ is a finite set of meta characters,
- $\Sigma$ is a finite set of symbols, disjoint from $V$,
- $P \subseteq V \times (\Sigma \cup V)^*$ is a finite relation,
- $V_0 \in V$ is the start symbol.

We denote the production rule $(S, \alpha) \in P$ as $S \to \alpha$.

Introduction
*Context Free Grammars and String Templates*

*Formal Prerequisites*
*Real-World String Templates*
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

# Context Free Language

## Definition (CFL)

A *Context Free Language* is the set of all strings that can be produced by a CFG through the application of a sequence of production rules to the start symbol via substitution of a meta character by the rule's right-hand-side.

*Introduction*
*Context Free Grammars and String Templates*

Formal Prerequisites
*Real-World String Templates*
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

## Example (Arithmetic Expression CFG)

$G_{\text{Expr}} = (V, \Sigma, P, V_0)$ where $V_0 = E$ and $P$ defined as

$$E \to T \oplus E$$
$$E \to T$$
$$T \to F \odot T$$
$$T \to F$$
$$F \to \langle E \rangle$$
$$F \to t$$

- $\langle t \oplus t \rangle \odot t \oplus t \in L_{\text{Expr}}$

*Introduction*
*Context Free Grammars and String Templates*

Formal Prerequisites
*Real-World String Templates*
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

## String Template System

### Definition (STS)

A *String Template System* can be defined as a tuple $(T, \Sigma, R, T_0)$ where

- $T$ is a finite set of string templates,
- $\Sigma$ is a finite set of symbols, disjoint from $T$,
- $R : T \to (\Sigma \cup T \cup \mathcal{P}(T))^*$ is a function,
- $T_0 \in T$ is a string template.

The symbol $\mathcal{E}$ is used to denote a string template with an empty right-hand-side, i.e. $R(\mathcal{E}) = \varepsilon$.

We denote the mapping $R(A) = \alpha$ as $A \mapsto \alpha$.

*Introduction*
*Context Free Grammars and String Templates*

Formal Prerequisites
*Real-World String Templates*
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

# String Template Language

### Definition (STL)

A *String Template Language* is the set of all strings that can be produced by a STS through recursive substitution of string templates with their mapping, beginning at the start template. For sets of string templates, any member may be expanded.

*Introduction*
**Context Free Grammars and String Templates**

*Formal Prerequisites*
*Real-World String Templates*
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

### Example (Tuple STS)

$S_{\mathrm{Tup}} = (T, \Sigma, R, T_0)$ where $T_0 = S$ and $R$ defined as

$$S \mapsto (F)$$
$$F \mapsto E\,\{C, \mathcal{E}\}$$
$$C \mapsto |F$$
$$E \mapsto \{S, D\}$$
$$D \mapsto p$$

- $\langle p | \langle \langle p | p \rangle | p \rangle \rangle \rangle \in L_{\mathrm{Tup}}$

*Introduction*
*Context Free Grammars and String Templates*

Formal Prerequisites
**Real-World String Templates**
Central Statements
Dynamic Expressions
Discussion and Outlook

# $STS \rightarrow Xtend$

```
1   def S() // S ↦ (F)
2       '''(«F»)'''
3
4   def F() // F ↦ E {C, ℰ}
5       '''«E»«IF c₁»«C»«ENDIF»'''
6
7   def C() // C ↦ |F
8       '''|«F»'''
9
10  def E() // E ↦ {S, D}
11      '''«IF c₂»«S»«ELSE»«D»«ENDIF»'''
12
13  def D() // D ↦ p
14      '''p'''
```

*Introduction*
*Context Free Grammars and String Templates*

*Formal Prerequisites*
***Real-World String Templates***
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

## $Xtend \to STS$

$$S \mapsto \texttt{<body><h2>}D_1\texttt{</h2>}\,\{T_1, \mathcal{E}\}\,\texttt{</body>}$$

$$T_1 \mapsto \texttt{<h3>}D_2\texttt{</h3><dl>}\,\{T_2, \mathcal{E}\}\,\texttt{</dl>}\,\{T_1, \mathcal{E}\}$$

$$T_2 \mapsto \{T_3, T_4, \mathcal{E}\}\,\{T_2, \mathcal{E}\}$$

$$T_3 \mapsto \texttt{<dt>}D_3\texttt{</dt><dd>=}D_4\texttt{</dt>}$$

$$T_4 \mapsto \texttt{<dt/><dd><i>Enumeration is extensible</i></dd>}$$

Introduction
Context Free Grammars and String Templates

Formal Prerequisites
*Real-World String Templates*
Central Statements
Dynamic Expressions
Discussion and Outlook

## Research Question (formal)

*Given a CFG $G = (V, \Sigma, P, V_0)$ describing the target language and a STS $S = (T, \Sigma, R, T_0)$ describing the code generator, can we decide if $L_S \subseteq L_G$?*

Introduction
Context Free Grammars and String Templates

Formal Prerequisites
Real-World String Templates
Central Statements
Dynamic Expressions
Discussion and Outlook

# Lemmas

## Lemma (STS ⊆ CFG)

*Every STS can be expressed as a CFG such that their respective languages are equal.*

## Lemma (STS ⊇ CFG)

*Every CFG can be expressed as a STS such that their respective languages are equal.*

## Corollary (STS = CFG)

*STS and CFG are interchangeable notations for the same set of languages.*

*Introduction*
*Context Free Grammars and String Templates*

*Formal Prerequisites*
*Real-World String Templates*
**Central Statements**
*Dynamic Expressions*
*Discussion and Outlook*

# *Mapping from STS to CFG*

## *Example (Tuple CFG)*

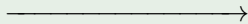$(T, \Sigma, R, T_0) = G_{\mathrm{Tup}} \mapsto S_{\mathrm{Tup}} = (V, \Sigma, P, V_0)$ with $P$ defined as

$$S \mapsto (F) \qquad\qquad\qquad\qquad\qquad S \to (F)$$
$$F \mapsto E\,\{C, \mathcal{E}\} \qquad\qquad\qquad\qquad F \to E$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad F \to EC$$
$$C \mapsto |F \qquad\qquad\longrightarrow\qquad\qquad C \to |F$$
$$E \mapsto \{S, D\} \qquad\qquad\qquad\qquad E \to S$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad E \to D$$
$$D \mapsto p \qquad\qquad\qquad\qquad\qquad D \to p$$

# Mapping from CFG to STS

---

### Example (Arithmetic Expression STS)

$(V, \Sigma, P, V_0) = S_{\text{Expr}} \mapsto G_{\text{Expr}} = (T, \Sigma, R, T_0)$ with $R$ defined as

$$E \mapsto \{A_{E_1}, A_{E_2}\}$$

$E \to T \oplus E \qquad\qquad\qquad A_{E_1} \mapsto T \oplus E$

$E \to T \qquad\qquad\qquad\qquad A_{E_2} \mapsto T$

$$T \mapsto \{A_{T_1}, A_{T_2}\}$$

$T \to F \odot T \qquad \xrightarrow{\hspace{3cm}} \qquad A_{T_1} \mapsto F \odot T$

$T \to F \qquad\qquad\qquad\qquad A_{T_2} \mapsto F$

$$F \mapsto \{A_{F_1}, A_{F_2}\}$$

$F \to \langle E \rangle \qquad\qquad\qquad A_{F_1} \mapsto \langle E \rangle$

$F \to t \qquad\qquad\qquad\qquad A_{F_2} \mapsto t$

*Introduction*
*Context Free Grammars and String Templates*

Formal Prerequisites
Real-World String Templates
**Central Statements**
Dynamic Expressions
Discussion and Outlook

# *Conclusion*

### Theorem

*Given an arbitrary STS $S$ and an arbitrary CFG $G$, the problem $L_S \subseteq L_G$ is undecidable.*

### Corollary

*Given an arbitrary CFG $G$, we can derive an equivalent STS $S$ with $L_S = L_G$. Any subset $S' \subseteq S$ of $S$ fulfills $L_{S'} \subseteq L_G$.*

*Introduction*
*Context Free Grammars and String Templates*

*Formal Prerequisites*
*Real-World String Templates*
*Central Statements*
**Dynamic Expressions**
*Discussion and Outlook*

## *Trojan Horse*

- string templates support control structures and string literals, but no attribute references
- include references to meta-model instances into the formal structure for string templates

*Introduction*
*Context Free Grammars and String Templates*

*Formal Prerequisites*
*Real-World String Templates*
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

# String Template System with Expressions

## Definition (STSE)

A *String Template System with Expressions* can be defined as a tuple $(T, E, \Sigma, F, R, T_0)$ where

- $T$ is a finite set of string templates,
- $E$ is a finite set of expressions, disjoint from $T$,
- $\Sigma$ is a finite set of symbols, disjoint from $T$ and $E$,
- $F : E \to \Sigma^*$ is a function,
- $R : T \to (\Sigma \cup T \cup E \cup \mathcal{P}(T))^*$ is a function,
- $T_0 \in T$ is the expanded string template.

*Introduction*
*Context Free Grammars and String Templates*

*Formal Prerequisites*
*Real-World String Templates*
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

# String Template Language with Expressions

### Definition (STLE)

A *String Template Language with Expressions* is the set of all strings that can be produced by a STSE through recursive substitution of string templates with their mapping and substitution of expressions with their mapping. For sets of string templates, any member may be expanded.

*Introduction*
*Context Free Grammars and String Templates*

*Formal Prerequisites*
*Real-World String Templates*
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

# *Discussion*

- instead of permitting Type-0 languages in dynamic expressions, we have to restrict them to Type-2 or higher
- two potential avenues:
  1. test dynamic expressions before substituting
  2. prove in advance that the expression conforms to the target syntax

*Introduction*
*Context Free Grammars and String Templates*

*Formal Prerequisites*
*Real-World String Templates*
*Central Statements*
*Dynamic Expressions*
*Discussion and Outlook*

## *Outlook*

```
1    grammar HTML {
2        ROOT -> "<html>" BODY "</html>";
3        BODY -> "<body>" (HEAD | ...) "</body>";
4        HEAD -> "<h1>" [^<]* "</h1>" | ...;
5        ...
6    }
7
8    HTML doc = '''
9        <html>
10            <body>«heading("SDL 2017")»</body>
11            ...
12        </html>
13    ''';
14
15    def heading(String<[^<]*> name)
16        '''<h1>«name»</h1>''';
```