

# Types and Analysis for Scripting Languages (Part 3)

Peter Thiemann

Universität Freiburg, Germany

Mini Course, Tallinn, Estland, 25.-27.2.2008

Grammar-Based Analysis of String Expressions

# Grammar-Based Analysis of String Expressions

TLDI 2005

## Motivation: Strings in Inappropriate Contexts

Substitute for “real” datatypes in scripting languages

Examples: Strings contain . . .

- ▶ style descriptors (numbers with units, percentages, etc)
- ▶ HTML and XML
- ▶ XPath expressions
- ▶ SQL statements

in PHP, JavaScript, but also JDBC and many others

# The Problem: Maintenance

- ▶ “strings with semantics”
- ⇒ may lead to runtime errors
- ▶ not prevented by usual type systems
- ▶ hard to test exhaustively

# Objective and NON-Objective

**Objective** Alleviate debugging and maintenance of existing programs

**NON-Objective** Framework for constructing programs

# Related Approach: String Expression Analysis I

Tabuchi, Sumii, Yonezawa.

*Regular expression types for strings in a text processing language.*

TIP 2002.

- ▶ Type and effect system based on simply-typed lambda calculus
- ▶ String type indexed with regular expression
- ▶ Output effect specified with regular expression
- ▶ Regular pattern matching
- ▶ Builds on Hosoya/Pierce regular expression types
- ▶ No type inference

## Related Approach: String Expression Analysis II

Christensen, Møller, and Schwartzbach.

*Precise analysis of string expressions.*

SAS 2003.

- ▶ Java  $\Rightarrow$  extended CFG  $\mathcal{G}$
- ▶  $\mathcal{G} \Rightarrow$  regular grammar  $\mathcal{R}$  with  $L(\mathcal{G}) \subseteq L(\mathcal{R})$   
using Mohri-Nederhof algorithm [2001]
- ▶ regular string assertions in the program can be checked effectively using inclusion of regular languages
- ▶ automatic analysis for full Java language (all string operations)

## This Work's Approach: String Expression Analysis III

- ▶ Based on the **polymorphic** HM (X) typing framework for lambda calculus [Odersky et al, 1999]
- ▶ string operations generate language inclusion constraints (generalizing CFGs with polymorphism)
- ▶ effective check of **context-free assertions** in the program by parsing inclusion constraints
- ▶ a context-free assertion is a sentential form for a context-free reference grammar

# Syntax

Alphabet

$T$

Symbols

$a, b \in T$

Words

$w \in T^*$

Constants

$c \in T^* \cup \{\cdot, \text{if } \}$

Expressions

$e ::= c | x | e(e) | \text{rec } f(x) e |$   
 $\text{let } x = e \text{ in } e$

# Semantics

Values       $v ::= w \mid \text{rec } f(x) e$   
Ev. Contexts     $E ::= [] \mid E(e) \mid e(E) \mid E \cdot e \mid v \cdot E \mid \text{if } E \text{ e e}$

Beta reduction

$$(\text{rec } f(x) e) (v) \longrightarrow e[x \mapsto v, f \mapsto \text{rec } f(x) e]$$

Delta reduction

$$\text{if } (a w) e_1 e_2 \longrightarrow e_1$$

$$\text{if } \varepsilon e_1 e_2 \longrightarrow e_2$$

$$a_1 \dots a_n \cdot b_1 \dots b_m \longrightarrow a_1 \dots a_n b_1 \dots b_m$$

Reduction

$$\frac{e \longrightarrow e'}{E[e] \mapsto E[e']}$$

# Type Language

Types	$\tau ::= \alpha \mid \text{Str}(\varphi) \mid \tau \rightarrow \tau$
Language Variables	$\varphi \in \Phi$
Constrained Types	$\rho ::= C \Rightarrow \tau$
Constraints	$C ::= \text{true} \mid \text{fail} \mid \tau \dot{\leq} \tau \mid r \dot{\subseteq} \varphi \mid C \wedge C$
String Type Indices	$r ::= \varphi \mid \varepsilon \mid a \mid r \cdot r$
Type Schemes	$\sigma ::= \forall \tilde{\varphi} \tilde{\alpha}. \rho$
Type Environments	$\Gamma ::= \emptyset \mid \Gamma(x : \sigma)$

# Typing Rules

Standard HM(X)

$$\frac{(x : \sigma) \in \Gamma}{C, \Gamma \vdash x : \sigma} \quad \frac{C, \Gamma \vdash e : \tau \quad C \Vdash \tau \leq \tau'}{C, \Gamma \vdash e : \tau'}$$

$$\frac{C, \Gamma \vdash e_1 : \tau_2 \rightarrow \tau_1 \quad C, \Gamma \vdash e_2 : \tau_2}{C, \Gamma \vdash e_1(e_2) : \tau_1}$$

$$\frac{C, \Gamma(f : \tau_2 \rightarrow \tau_1)(x : \tau_2) \vdash e : \tau_1}{C, \Gamma \vdash \text{rec } f(x) e : \tau_2 \rightarrow \tau_1}$$

$$\frac{C, \Gamma \vdash e : \sigma \quad C, \Gamma(x : \sigma) \vdash e' : \tau'}{C, \Gamma \vdash \text{let } x = e \text{ in } e' : \tau'}$$

$$\frac{C \wedge D, \Gamma \vdash e : \tau \quad \{\tilde{\alpha}, \tilde{\varphi}\} \cap (fv(C) \cup fv(\Gamma)) = \emptyset}{C, \Gamma \vdash e : \forall \tilde{\alpha} \tilde{\varphi}. D \Rightarrow \tau}$$

$$\frac{C, \Gamma \vdash e : \forall \tilde{\alpha} \tilde{\varphi}. D \Rightarrow \tau \quad S = [\tilde{\alpha} \mapsto \tilde{\tau}, \tilde{\varphi} \mapsto \tilde{\varphi}'] \quad C \Vdash S(D)}{C, \Gamma \vdash e : S(\tau)}$$

# Language Inclusion Constraints

- ▶ arise from the type signatures of the string operations

constant       $w$  :  $\forall \varphi. (w \dot{\subseteq} \varphi) \Rightarrow \text{Str}(\varphi)$

concatenation     $\cdot$  :  $\forall \varphi_1, \varphi_2, \varphi. (\varphi_1 \cdot \varphi_2 \dot{\subseteq} \varphi) \Rightarrow \text{Str}(\varphi_1) \rightarrow \text{Str}(\varphi_2) \rightarrow \text{Str}(\varphi)$

conditional    if :  $\forall \alpha, \varphi. \text{Str}(\varphi) \rightarrow \alpha \rightarrow \alpha \rightarrow \alpha$

- ▶ give rise to subtyping:

$$\frac{\varphi_1 \dot{\subseteq} \varphi_2}{\text{Str}(\varphi_1) \leq \text{Str}(\varphi_2)}$$

# Solution of a Constraint

- ▶ A solution of constraint  $C$  is a pair  $(S, \Theta)$  of
  - ▶ an idempotent substitution  $S$  of type variables by types
  - ▶ a language assignment  $\Theta : \Phi \rightarrow \mathcal{P}(T^*)$so that  $S, \Theta \models C$  holds.
- ▶ The meaning of left hand sides of constraints

$$L_\Theta(\varphi) = \Theta(\varphi)$$

$$L_\Theta(\varepsilon) = \{\varepsilon\}$$

$$L_\Theta(a) = \{a\}$$

$$L_\Theta(r_1 \cdot r_2) = \{w_1 \cdot w_2 \mid w_1 \in L_\Theta(r_1), w_2 \in L_\Theta(r_2)\}$$

# Constraint Satisfaction

$$\frac{\Theta \models S(C)}{S, \Theta \models C}$$

$$\frac{\Theta \models C_1 \quad \Theta \models C_2}{\Theta \models C_1 \wedge C_2}$$

$$\Theta \models \text{true} \quad \Theta \models \alpha \dot{\leq} \alpha$$

$$\frac{\Theta(\varphi_1) \subseteq \Theta(\varphi_2)}{\Theta \models \text{str}(\varphi_1) \dot{\leq} \text{str}(\varphi_2)}$$

$$\frac{\Theta \models \tau'_2 \dot{\leq} \tau_2 \quad \Theta \models \tau_1 \dot{\leq} \tau'_1}{\Theta \models \tau_2 \rightarrow \tau_1 \dot{\leq} \tau'_2 \rightarrow \tau'_1}$$

$$\frac{L_\Theta(r) \subseteq \Theta(\varphi)}{\Theta \models r \dot{\subseteq} \varphi}$$

## Example: Render a Tree in HTML

```
render t = case t of
  Tip num ->
    int2string num
  Fork l r ->
    "<ul><li>" ++ render l ++ "</li><li>" ++
      render r ++ "</li></ul>"
```

Assuming

$$\begin{aligned} C_1 &= (0 \cdot \varphi'' \dot{\subseteq} \varphi' \wedge 1 \cdot \varphi'' \dot{\subseteq} \varphi' \wedge \dots \wedge \varepsilon \dot{\subseteq} \varphi'' \wedge \varphi' \dot{\subseteq} \varphi'') \\ \text{int2string} &: \forall \varphi', \varphi''. (C_1) \Rightarrow \text{Bool} \rightarrow \text{Str}(\varphi') \end{aligned}$$

We find that

$$\begin{aligned} C_2 &= <\ul><\li> \cdot \varphi \cdot </li><\li> \cdot \varphi \cdot </li></ul> \dot{\subseteq} \varphi \wedge \varphi' \dot{\subseteq} \varphi \\ \text{render} &: \forall \varphi, \varphi', \varphi''. (C_2 \wedge C_1) \Rightarrow \text{Tree} \rightarrow \text{Str}(\varphi) \end{aligned}$$

Does render create HTML output?

# Normalizing Constraints

$\text{str}(\varphi) \stackrel{.}{\leq} \text{str}(\varphi')$	$\Leftrightarrow$	$\varphi \stackrel{.}{\subseteq} \varphi'$
$\text{str}(\varphi) \stackrel{.}{\leq} \tau \rightarrow \tau'$	$\Leftrightarrow$	<b>fail</b>
$\text{str}(\varphi) \stackrel{.}{\leq} \alpha$	$\Rightarrow$	$\alpha \stackrel{.}{=} \text{str}(\varphi')$ <b>if <math>\varphi'</math> fresh</b>
$\tau \rightarrow \tau' \stackrel{.}{\leq} \text{str}(\varphi')$	$\Leftrightarrow$	<b>fail</b>
$\tau \rightarrow \tau' \stackrel{.}{\leq} \tau_1 \rightarrow \tau'_1$	$\Leftrightarrow$	$\tau_1 \stackrel{.}{\leq} \tau \wedge \tau' \stackrel{.}{\leq} \tau'_1$
$\tau \rightarrow \tau' \stackrel{.}{\leq} \alpha$	$\Rightarrow$	$\alpha \stackrel{.}{=} \alpha' \rightarrow \alpha''$ <b>if <math>\alpha', \alpha''</math> fresh, <math>\alpha \notin \text{fv}(\tau, \tau')</math></b>
$\alpha \stackrel{.}{\leq} \text{str}(\varphi')$	$\Rightarrow$	$\alpha \stackrel{.}{=} \text{str}(\varphi)$ <b>if <math>\varphi</math> fresh</b>
$\alpha \stackrel{.}{\leq} \tau \rightarrow \tau'$	$\Rightarrow$	$\alpha \stackrel{.}{=} \alpha' \rightarrow \alpha''$ <b>if <math>\alpha', \alpha''</math> fresh, <math>\alpha \notin \text{fv}(\tau, \tau')</math></b>
$\alpha \stackrel{.}{\leq} \alpha$	$\Leftrightarrow$	<b>true</b>
$\alpha \stackrel{.}{\leq} \alpha' \wedge \alpha' \stackrel{.}{\leq} \alpha''$	$\Rightarrow$	$\alpha \stackrel{.}{\leq} \alpha''$
$\alpha \stackrel{.}{\leq} \alpha' \wedge \alpha' \stackrel{.}{\leq} \alpha$	$\Leftrightarrow$	$\alpha \stackrel{.}{=} \alpha'$
$\alpha \stackrel{.}{= \alpha}$	$\Leftrightarrow$	<b>true</b>
$\varphi \stackrel{.}{=} \varphi$	$\Leftrightarrow$	<b>true</b>
$\varphi \stackrel{.}{\subseteq} \varphi$	$\Leftrightarrow$	<b>true</b>
$\varphi \stackrel{.}{\subseteq} \varphi' \wedge \varphi' \stackrel{.}{\subseteq} \varphi$	$\Leftrightarrow$	$\varphi \stackrel{.}{=} \varphi'$

## Normalizing Constraints II

$$\begin{aligned} C @ (C' \wedge & \quad \Rightarrow \quad r_1 rr_2 \dot{\subseteq} \varphi' \quad \text{if } \varphi \notin \text{reach}(C, \varphi) \\ r \dot{\subseteq} \varphi \wedge \\ r_1 \varphi r_2 \dot{\subseteq} \varphi') & \\ C \wedge \varphi \doteq \varphi' & \Leftrightarrow C[\varphi \mapsto \varphi'] \wedge \varphi \doteq \varphi' \\ C \wedge \alpha \doteq \tau & \Leftrightarrow C[\alpha \mapsto \tau] \wedge \alpha \doteq \tau \quad \text{if } \alpha \notin \text{fv}(\tau) \\ C \wedge \alpha \doteq \tau & \Leftrightarrow \text{fail} \quad \text{if } \alpha \in \text{fv}(\tau) \wedge \tau \neq \alpha \\ C \wedge \text{true} & \Leftrightarrow C \\ C \wedge \text{fail} & \Leftrightarrow \text{fail} \end{aligned}$$

where

$$\begin{aligned} \text{reach}(C, \varphi) &= \mu V. F(V) \cup F(\{\varphi\}) \\ F(V) &= \{\varphi' \mid r_1 \varphi' r_2 \dot{\subseteq} \varphi \in C, \varphi \in V\} \end{aligned}$$

# Context-Free Languages

- ▶ A normalized set of inclusion constraints can be read as a context-free grammar  $\mathcal{G}$ .
  - ▶ The language variables serve as nonterminal symbols.
  - ▶ Each constraint  $r \subseteq \varphi$  serves as a production  $\varphi \rightarrow r$ .
- ▶ Normalization unfolds productions up to recursion.

# Context-Free Languages

- ▶ A normalized set of inclusion constraints can be read as a context-free grammar  $\mathcal{G}$ .
  - ▶ The language variables serve as nonterminal symbols.
  - ▶ Each constraint  $r \subseteq \varphi$  serves as a production  $\varphi \rightarrow r$ .
- ▶ Normalization unfolds productions up to recursion.
- ▶ Problem
  - ▶ Given a context-free reference grammar  $\mathcal{G}_0$ ,
  - ▶ How can we check that a sublanguage of  $\mathcal{G}$  is contained in a sublanguage of  $\mathcal{G}_0$ ?

# A Reference Grammar for HTML

$S \rightarrow <\text{html}> H B </\text{html}>$	$C \rightarrow$
$S \rightarrow H B$	$C \rightarrow P$
$H \rightarrow <\text{head}> T </\text{head}>$	$C \rightarrow C C$
$T \rightarrow <\text{title}> P </\text{title}>$	$C \rightarrow <\text{p}> C </\text{p}>$
$B \rightarrow <\text{body}> C </\text{body}>$	$C \rightarrow <\text{b}> C </\text{b}>$
$P \rightarrow$	$C \rightarrow <\text{tt}> C </\text{tt}>$
$P \rightarrow P P$	$C \rightarrow <\text{em}> C </\text{em}>$
$P \rightarrow 0$	$C \rightarrow <\text{i}> C </\text{i}>$
$P \rightarrow 1$	$C \rightarrow <\text{u}> C </\text{u}>$
$\vdots$	$C \rightarrow <\text{s}> C </\text{s}>$
	$C \rightarrow <\text{ol}> L </\text{ol}>$
	$C \rightarrow <\text{ul}> L </\text{ul}>$
	$L \rightarrow <\text{li}> C </\text{li}>$
	$L \rightarrow L L$

Adapted from the paper

# Solving Constraints by Parsing

- ▶ Substitute nonterminals of the reference grammar for the language variables in the constraints (with  $\Theta$ )
- ▶ Check for each constraint  $r \dot{\subseteq} \varphi$  if  $\Theta(\varphi) \stackrel{*}{\Rightarrow} \Theta(r)$
- ▶ If successful, replace constraint by *assignment constraint*  $\Theta$

In the example

$$\begin{array}{ll} 0 \cdot \varphi'' \dot{\subseteq} \varphi' & \rightsquigarrow [\varphi' \mapsto P, \varphi'' \mapsto P] \\ \varepsilon \dot{\subseteq} \varphi'' & \rightsquigarrow [\varphi'' \mapsto C] \vee [\varphi'' \mapsto P] \\ \varphi' \dot{\subseteq} \varphi'' & \rightsquigarrow \varphi' = \varphi'' \vee [\varphi' = P, \varphi'' = C] \\ \varphi' \dot{\subseteq} \varphi & \rightsquigarrow \varphi' = \varphi \vee [\varphi' = P, \varphi = C] \\ <\!\!\text{ul}\!><\!\!\text{li}\!>\cdot \varphi \cdot <\!\!\text{/li}\!><\!\!\text{li}\!>\cdot \varphi \cdot <\!\!\text{/li}\!><\!\!\text{/ul}\!> \dot{\subseteq} \varphi & \rightsquigarrow [\varphi = C] \end{array}$$

- ▶ Result: a disjunction of assignment constraints

# Resolving Assignment Constraints

Simplify disjunction of assignment constraints:

$$\begin{aligned}& [\varphi' \mapsto P, \varphi'' \mapsto P] \bowtie ([\varphi'' \mapsto C] \vee [\varphi'' \mapsto P]) \bowtie (\varphi' = \varphi'' \vee [\varphi' = P, \varphi'' = C]) \\&= [\varphi' \mapsto P, \varphi'' \mapsto P] \bowtie (\varphi' = \varphi'' \vee [\varphi' = P, \varphi'' = C]) \bowtie (\varphi' = \varphi \vee [\varphi' = P, \varphi'' = C]) \\&= [\varphi' \mapsto P, \varphi'' \mapsto P] \bowtie (\varphi' = \varphi \vee [\varphi' = P, \varphi = C]) \bowtie [\varphi = C] \\&= ([\varphi' \mapsto P, \varphi'' \mapsto P, \varphi \mapsto P] \vee [\varphi' \mapsto P, \varphi'' \mapsto P, \varphi \mapsto C]) \bowtie [\varphi = C] \\&= [\varphi' = P, \varphi'' = P, \varphi = C]\end{aligned}$$

Hence the simplified type scheme

$$\text{render} : \forall \varphi. [\varphi = C] \Rightarrow \text{Tree} \rightarrow \text{Str}(\varphi)$$

# Earley's Parser Solves Parsing Constraints

**init**

$[\rightarrow \bullet \gamma, 0] \in E_0.$

**scan**

$[A \rightarrow \alpha \bullet X\beta, j] \in E_i$  and  $X_{i+1} = X \Rightarrow$

$[A \rightarrow \alpha X \bullet \beta, j] \in E_{i+1}$

**pred**

$[A \rightarrow \alpha \bullet B\beta, j] \in E_i$  and  $B \rightarrow \gamma \in P \Rightarrow$

$[B \rightarrow \bullet \gamma, i] \in E_i.$

**red**

$[B \rightarrow \gamma \bullet, j] \in E_i$  and  $[A \rightarrow \alpha \bullet B\beta, k] \in E_j \Rightarrow$

$[A \rightarrow \alpha B \bullet \beta, k] \in E_i.$

- ▶  $X_1 \dots X_n$  input word with  $X_i \in T \cup N$
- ▶  $E_0, \dots, E_n$  sets of items
- ▶ accept, if  $[\rightarrow \gamma \bullet, 0] \in E_n$
- ▶ only difference to standard Earley in **scan**
- ▶ extension: construct assignment while parsing (paper)

# Left/Right Bias

Consider

- ▶ the left-recursive reference grammar  $S \rightarrow a \mid Sb$
- ▶ the right-recursive expression

$a \cdot ((\text{rec } f(x) \text{ if } (x = 0) \varepsilon b \cdot (f(x - 1)))) 10$

# Left/Right Bias

Consider

- ▶ the left-recursive reference grammar  $S \rightarrow a \mid Sb$
- ▶ the right-recursive expression  
 $a \cdot ((\text{rec } f(x) \text{ if } (x = 0) \varepsilon b \cdot (f(x - 1))) \ 10)$
- ▶ has type  $\text{Str}(\varphi)$
- ▶ where  $a \cdot \varphi' \dot{\subseteq} \varphi, \varepsilon \dot{\subseteq} \varphi', b \cdot \varphi' \dot{\subseteq} \varphi', [a \dot{\subseteq} \varphi]$

# Left/Right Bias

Consider

- ▶ the left-recursive reference grammar  $S \rightarrow a \mid Sb$
- ▶ the right-recursive expression  
 $a \cdot ((\text{rec } f(x) \text{ if } (x = 0) \varepsilon b \cdot (f(x - 1))) \ 10)$
- ▶ has type  $\text{Str}(\varphi)$
- ▶ where  $a \cdot \varphi' \dot{\subseteq} \varphi, \varepsilon \dot{\subseteq} \varphi', b \cdot \varphi' \dot{\subseteq} \varphi', [a \dot{\subseteq} \varphi]$
- ▶ We can assign  $[\varphi \mapsto S]$ , but there is no possible assignment for  $\varphi'$

# Left/Right Bias

Consider

- ▶ the left-recursive reference grammar  $S \rightarrow a \mid Sb$
- ▶ the right-recursive expression  
 $a \cdot ((\text{rec } f(x) \text{ if } (x = 0) \varepsilon b \cdot (f(x - 1))) \ 10)$
- ▶ has type  $\text{Str}(\varphi)$
- ▶ where  $a \cdot \varphi' \dot{\subseteq} \varphi, \varepsilon \dot{\subseteq} \varphi', b \cdot \varphi' \dot{\subseteq} \varphi', [a \dot{\subseteq} \varphi]$
- ▶ We can assign  $[\varphi \mapsto S]$ , but there is no possible assignment for  $\varphi'$
- ▶ But for right-recursive grammar  $A \rightarrow aB, B \rightarrow \varepsilon \mid bB$  we find the assignment  $[\varphi \mapsto A, \varphi' \mapsto B]$ .

# Left/Right Bias

Consider

- ▶ the left-recursive reference grammar  $S \rightarrow a \mid Sb$
- ▶ the right-recursive expression  
 $a \cdot ((\text{rec } f(x) \text{ if } (x = 0) \varepsilon b \cdot (f(x - 1))) \ 10)$
- ▶ has type  $\text{Str}(\varphi)$
- ▶ where  $a \cdot \varphi' \dot{\subseteq} \varphi, \varepsilon \dot{\subseteq} \varphi', b \cdot \varphi' \dot{\subseteq} \varphi', [a \dot{\subseteq} \varphi]$
- ▶ We can assign  $[\varphi \mapsto S]$ , but there is no possible assignment for  $\varphi'$
- ▶ But for right-recursive grammar  $A \rightarrow aB, B \rightarrow \varepsilon \mid bB$  we find the assignment  $[\varphi \mapsto A, \varphi' \mapsto B]$ .
- ▶ Moral: do not introduce left or right bias in the grammar

# Conclusions

- + String expression analysis that deals with context-free languages
- + Essentially type inference for HM(Earley)
- + Constraint solver implemented
- + Working and efficient implementation for PHP
- + Can now deal with deconstruction (not shown)
- Special grammar needed
  - ▶ character-based, not token-based
  - ▶ should not be left- or right-biased, ambiguous grammars preferred
- More string operations