

# Type-Based Specialization of XML Transformations

Kazutaka Matsuda\*,  
Zhenjiang Hu\*\*, Masato Takeichi\*

\*The University of Tokyo

\*\*National Institute of Informatics (Japan)

January 19th, 2009 PEPM'09

# XML

- ▶ XML is widely-used.
- ▶ Trans. between XMLs are very important!

```
<members>
  <member>
    <name>Kztk</name>
    <email>kztk@...</email>
    <tel>+81-90...</tel>
  </member>
  ...
</members>
```



```
<html>...
<body>
  <h1>Member List</h1>
  <h2>Kztk</h2>
  <address>kztk@...</address>
  <address>+81-90...</address>
  ...
</body>
</html>
```

Question

Why specialization  
is required?

# A Simple XML Trans.

```
convC(<email>e</email>) = <mailto>e</mailto>  
convC(<phone>t</phone>) = <tel>t</tel>
```

Input XML (Elem.)

Output XML (Elem.)

```
<email>kztk@ipl...</email>
```



```
<mailto>kztk@ipl...</mailto>
```

```
<phone>+81-90-...</phone>
```



```
<tel>+81-90-...</tel>
```

# A Simple XML Trans.

```
convC(<email>(e)) = <mailto>(e)  
convC(<phone>(t)) = <tel>(t)
```

Input XML (Elem.)

Output XML (Elem.)

<email>(kztk@ipl...)



<mailto>(kztk@ipl...)

<phone>(+81-90-...)



<tel>(+81-90-...)

# Type (Schema)

```
convC(<email>(e::String)) = <mailto>(e)
convC(<phone>(t::String)) = <tel>(t)
```

## Src. Type

```
<email>(String)
| <phone>(String)
```

⊔

```
<email>(kztk@ipl...)
```

```
<phone>( +81-90-... )
```



## Tgt. Type

```
<mailto>(String)
| <tel>(String)
| <fax>(String)
```

⊔

```
<mailto>(kztk@ipl...)
```

```
<tel>( +81-90-... )
```

# Type Checking with Subtyping

```
convC(<email>(e::String)) = <mailto>(e)
convC(<phone>(t::String)) = <tel>(t)
```

convC::

```
<email>(String)
| <phone>(String)
```



```
<mailto>(String)
| <tel>(String)
```

Src. Type  $\supseteq$  ?

?  $\supseteq$  Tgt. Type

```
<email>(String)
```

```
<mailto>(String)
| <tel>(String)
| <fax>(String)
```

# Problem

```
convC(<email>(e::String)) = <mailto>(e)
convC(<phone>(t::String)) = <tel>(t)
```

convC::

```
<email>(String)
| <phone>(String)
```



```
<mailto>(String)
| <tel>(String)
```

Src. Type

?

Tgt. Type

?

```
<email>(String)
```

```
<mailto>(String)
```



# Use Trans. Def. for Preciseness

```
convC(<email>(e)) = <mailto>(e)  
convC(<phone>(t)) = <tel>(t)
```

```
data S = <email>(String)  
convC_onS(x::S) = convC(x)
```

Type of "x"

# do Specialization

```
convCs(<email>(e :: String)) = <mailto>(e)
```

Sp. ver. of convC w.r.t. S

```
convC(<email>(e)) = <mailto>(e)  
convC(<phone>(t)) = <tel>(t)
```

```
data S = <email>(String)  
convC_onS(x :: S) = convCs(x)
```

```
convCs :: <email>(String) → <mailto>(String)
```

Src. Type ?

```
<email>(String)
```

? Tgt. Type

```
<mailto>(String)
```

# Short Summary

```
convCs(<email>(e :: String)) = <mailto>(e)
```

*convC* ::  $\begin{matrix} \langle \text{email} \rangle (String) \\ | \\ \langle \text{phone} \rangle (String) \end{matrix} \longrightarrow \begin{matrix} \langle \text{mailto} \rangle (String) \\ | \\ \langle \text{tel} \rangle (String) \end{matrix}$

Specialize

*convCs* ::  $\langle \text{email} \rangle (String) \longrightarrow \langle \text{mailto} \rangle (String)$

Src. Type ?

$\langle \text{email} \rangle (String)$

? Tgt. Type

$\langle \text{mailto} \rangle (String)$

Question

Why specialization  
is required?

Answer

For precise analyses  
of programs.

# Dream

- ▶ We hope to specialize
  - ▶ Types as expressive as XML schemata
    - ▶ Recursively defined types
  - ▶ Transformation on XML trees
    - ▶ Recursively defined transformations

# Difficulty

## ► Infinite specialization

```
ab2cd($) = $  
ab2cd(<a/>.r) = <c/>.ab2cd(r)  
ab2cd(<b/>.r) = <d/>.ab2cd(r)
```

```
data S = <a/>.S.<b/> | $  
ab2cd_onS(x::S) = ab2cd(r)
```

Context-Free  
 $\langle a/\rangle^n \langle b/\rangle^n$

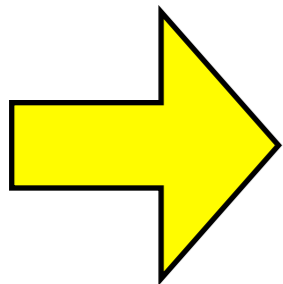
```
ab2cdS(<a/>.r::SB) = <c/>.ab2cdSB(r)  
ab2cdSB(<a/>.r::SBB) = <c/>.ab2cdSBB(r)
```

...

```
data SB = S.<b/>  
data SBB = SB.<b/>
```

...

Infinite!



# 3 important things

1. Types
2. Transformations (functions)
3. Method to specialize

`data S = <email>(String)`  
`convC_onS(x::S) = convC(`

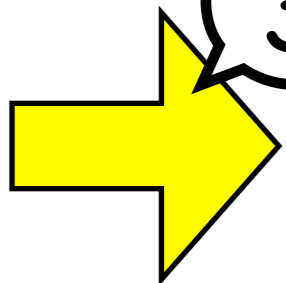
1

`convC(<email>(e::String)) = <mailto>(e)`  
`convC(<tel>(t::String)) = <callto>(t)`

2

3

`convCs(<email>(e::String)) = <mailto>(e)`



# 3 important things

1. Types
2. Transformations (functions)
3. Method to specialize

```
data S = <email>(String)  
convC_onS(x::S) = convC(
```

1

```
convC(<email>(e::String)) = <mailto>(e)  
convC(<tel>(t::String))   = <callto>(t)
```

2

Simple Unfolding

```
convC(<email>(e::String)) = <mailto>(e)
```



# Our Work

## 1. Types

- ▶ **Regular** hedge grammar

## 2. Transformation (functions)

- ▶ Designed to **avoid Context-Freeness**

## 3. Method to specialize

- ▶ Simple unfolding of
  - ▶ Types/Transformations

# (Regular Hedge) Type

- ▶ Regular hedge grammars
  - ▶ A model of schemata [Murata et al. 05]
    - ▶ DTD, XML Schema, Relax NG
  - ▶  $\equiv$  Tree automata [Comon et al. 97]
  - ▶ Types in XDuce/CDuce  
[Hosoya&Pierce 03, Benzaken et al. 03]

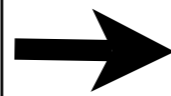
```
data Members = <members>( Member* )  
data Member = <member>( Name . Contact* )  
data Name = <name>( String )  
data Contact = <email>( String ) | <tel>( String )
```

# Transformations

- ▶ Treeless [Wadler 90]
  - ▶ Arguments of funcs. must be variables.
- ▶ + Type-based XML-sequence splitting
  - ▶ A subset of patterns in XDuce

```
data As = ( <a>(String) )*  
data Bs = ( <b>(String) )*  
splitAB(x::As . y::Bs) = <as>(x) . <bs>(y)
```

```
<a>1</a><a>2</a><b>3</b>
```



```
<as><a>1</a><a>2</a></as>  
<bs><b>3</b></bs>
```

```
rev($) = $  
rev(x::Tree . r) = rev(r) . x
```

```
rev(<a/><b/><c/>)
```



```
<c/><b/><a/>
```

# Core Syntax of Lang.

## ► Whole (core) syntax

$program ::= t\_decl\ f\_decl$

$t\_decl ::= \mathbf{data}\ T_1 = t_1; \dots \mathbf{data}\ T_n = t_n;$

$t ::= \$ \mid \sigma(T_1) \cdot T_2$

$f\_decl ::= f_1(p_1) = e_1; \dots f_m(p_m) = e_m;$

$p ::= \$ \mid \sigma(p) \mid p_1 \cdot p_2 \mid x :: T$

$e ::= \$ \mid \sigma(e) \mid e_1 \cdot e_2 \mid x \mid f(x)$

# Specialization Method

Sp.  $rev(x::Tree . r) = rev(r) . x$        $data\ ABS = \$ \mid \langle a/\rangle . \langle b/\rangle . ABS$

## 1. Specialize pattern

$x::Tree . r$   $ABS$   $\rightarrow$   $x::\langle a/\rangle . r::BABS$

where  $data\ BABS = \langle b/\rangle . ABS$

## 2. Generate new rule(s)

$rev_{ABS}(x::\langle a/\rangle . r::BABS) = rev_{BABS}(r) . x$

## 3. Recursively apply specialization

Sp.

$rev(\$) = \$$

$rev(x::Tree . r) = rev(r) . x$

$BABS$

# Termination

Sp.  $rev(x :: Tree . r) = rev(r) . x$

data ABS  
= \$ | <a/> . <b/> . ABS

## 1. Specialize pattern

$x :: Tree . r$  ABS  $\rightarrow$   $x :: <a/> . r :: BABs$

$v$   $<b/> . ABS$

Always Regular

## 2. Generate new rule

$rev_{ABS}(x :: <a/> . r :: BABs) = rev_{BABs}(r) . x$

## 3. Recurs

Sp.

$r$   
 $rev(x .$

#(introduced types) is finite  
 $\Rightarrow$  Always terminating

# Related Work

- ▶ Exact type checking of XML Trans.  
[Maneth et al. 07]
  - ▶ Specialization to handle input XML type
  - ▶ (Affine) Treeless + Acc.
    - ▶ Precisely, (Linear) stay MTT [Fischer 68]
- ▶ Main Difference:  
Type-based XML-sequence splitting

# Related Work

- ▶ Type inference in XDuce/CDuce  
[Hosoya&Pierce 03, Benzaken et al. 03]
- ▶ Related to our pattern specialization

$x :: Tree . r$   $ABs$   $\rightarrow$   $x :: \langle a / \rangle . r :: BABs$

- ▶ Ours: precise (exact) but restricted



# Conclusion

- ▶ Type-based specialization of XML trans.
- ▶ Language:
  - ▶ Regular hedge grammars
  - ▶ Treeless + Type-based XML-seq. splitting
- ▶ Property:
  - ▶ Correct & terminating
- ▶ Applications:
  - (Forward) type inference, inversion, ...

Prototype Implementation:

<http://www.ipl.t.u-tokyo.ac.jp/~kztk/sp/>

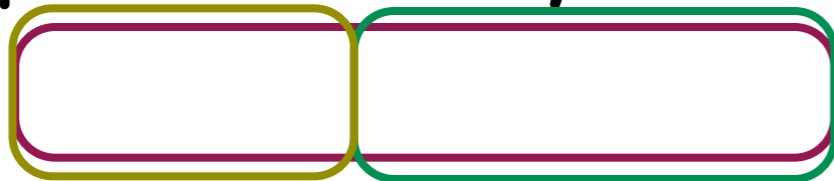
Thank you

# Difficulty of Seq. Decomp.

- ▶ Without XML-seq. decomposition
- ▶ A pattern only decompose the head



- ▶ With XML-seq. decomposition
- ▶ A pattern may extract preceding seq.



- ▶ c.f. Type structure

**data**  $T = \sigma(T_1) \cdot T_2$

# Properties of Sp.

- ▶ Correct

- ▶  $f_S(x) = y \Leftrightarrow x \in S \wedge f(x) = y$

- ▶ Always terminating

- ▶ Thanks to regularity

- ▶ But, slow

- ▶  $\#(\text{introduced types}) = O(2^{2^{(2^{\#\text{types}})^2}})$

- ▶  $\#(\text{generated rules for a rule}) < \#(\text{types})^{\#(\text{vars})}$

- ▶  $\#(\text{output funcs}) <$

- $\#(\text{funcs}) * \#(\text{introduced types})$

# Image of Finiteness

