



Faculty of Science



Coercions and substitutions

Fritz Henglein

DIKU, University of Copenhagen

2nd International Summer School on Advances in Programming Languages
Edinburgh, Scotland

2014-08-22



Recall: Regular expressions as types

- Language of expressions $\text{Reg}_{\mathcal{A}}$:

$$E, F ::= 0 \mid 1 \mid a \mid E|F \mid EF \mid E^*$$

- Type interpretation $\mathcal{T}[E]$:

$$\mathcal{T}[0] = \emptyset$$

$$\mathcal{T}[1] = \{()\}$$

$$\mathcal{T}[a] = \{a\}$$

$$\mathcal{T}[E|F] = \mathcal{T}[E] + \mathcal{T}[F]$$

$$\mathcal{T}[EF] = \mathcal{T}[E] \times \mathcal{T}[F]$$

$$\mathcal{T}[E^*] = \mathcal{T}[E] \text{ list}$$

where

$$S + T = \{\text{inl } v \mid v \in S\} \cup \{\text{inr } w \mid w \in T\}$$

$$S \text{ list} = \{[v_1, \dots, v_n] \mid v_i \in S\}$$



Regular expressions as languages

Language interpretation $\mathcal{L}[E]$:

$$\mathcal{L}[E|F] = \mathcal{L}[E] \cup \mathcal{L}[F]$$

$$\mathcal{L}[EF] = \mathcal{L}[E] \mathcal{L}[F]$$

$$\mathcal{L}[E^*] = \mathcal{L}[E]^*$$

$$S T = \{st \mid s \in S \wedge t \in T\}$$

$$S^* = \{s_1 s_2 \dots s_n \mid s_i \in S\}$$



Values = parse trees = membership proofs

Views of $\mathcal{T}[E]$:

- The values of E read as a *type* built from singletons, sum, product, list
- The parse trees of all the strings matching E
- The certificates (proofs) of membership for strings in $\mathcal{L}[E]$.

Proposition

$$\mathcal{L}[E] = \{\text{flat}(v) \mid v \in \mathcal{T}[E]\}$$



Regular Expression Containment

Definition (Language containment)

$\models E \leq F$ if $\mathcal{L}[[E]] \subseteq \mathcal{L}[[F]]$.

Write $\models E = F$ iff $\models E \leq F$ and $\models F \leq E$.

Definition (Coercion)

$\models f : E \leq F$ if f is function from $\mathcal{T}[[E]]$ to $\mathcal{T}[[F]]$ such that $\text{flat}(f(v)) = \text{flat}(v)$.

Theorem

$\models E \leq F$ if and only if there exists f such that $\models f : E \leq F$.



Example: Language containment by FP

- Conway's denesting rule: $\models (E|F)^* = E^*(FE^*)^*$ for all E, F .
- Proof by functional programming:

Find $f : ('a + 'b) \text{ list} \rightarrow 'a \text{ list} \times ('b \times 'a \text{ list}) \text{ list}$
such that f does not discard, duplicate or reorder its input

$$f([]) = ([], [])$$

$$f(\text{inl } u :: zs) =$$

$$\text{let } (xs, ys) = f(zs) \text{ in } (u :: xs, ys)$$

$$f(\text{inr } v :: zs) =$$

$$\text{let } (xs, ys) = f(zs) \text{ in } ([], (v, xs) :: ys)$$

- f terminates since it is called recursively with smaller sized arguments
- f is string-preserving
- f is polymorphic, so works for all E, F . (Indeed E, F need not even be regular.)
- Therefore: $\models (E|F)^* \leq E^*(FE^*)^*$ for all E, F .
- Reverse direction similar.



Synthesis problem

- Need only *existence* of a? Decision problem, classical automata and formal language theory.
- Need to *construct* a coercion? Synthesis problem. Which coercion, how to construct and represent it?

Challenges:

- Different coercions, extensionally—choose which one?
E.g. $\text{tagL} : a \leq a + a$ or $\text{tagR} : a \leq a + a$
- Different coercions, intensionally—want *efficient* ones.
E.g., $\text{id} : a^* \leq a^*$ better than $\text{map}[\text{id}] : a^* \leq a^*$.
- How to design “sublanguage” (DSL) of well-typed functions that are
 - ▶ guaranteed to be coercions (soundness);
 - ▶ contain at least one coercion for every valid containment (completeness);
 - ▶ can be searched practically efficiently (not too large);
 - ▶ contain efficient coercions (not too small).

Idea: Constructive interpretation of axiomatization of containment



Axiomatization: Weak containment

$$\text{shuffle} : E + (F + G) = (E + F) + G$$

$$\text{retag} : E + F = F + E$$

$$\text{untagL} : 0 + F = F$$

$$\text{untag} : E + E \leq E$$

$$\text{tagL} : E \leq E + F$$

$$\text{assoc} : E \times (F \times G) = (E \times F) \times G$$

$$\text{swap} : E \times 1 = 1 \times E$$

$$\text{proj} : 1 \times E = E$$

$$\text{abortR} : E \times 0 = 0$$

$$\text{abortL} : 0 \times E = 0$$

$$\text{distL} : E \times (F + G) = (E \times F) + (E \times G)$$

$$\text{distR} : (E + F) \times G = (E \times G) + (F \times G)$$

$$\text{wrap} : 1 + E \times E^* = E^*$$

$$\text{id} : E = E$$

$$\frac{c : E \leq E' \quad d : E' \leq E''}{c; d : E \leq E''}$$

$$\frac{c : E \leq E' \quad d : F \leq F'}{c + d : E + F \leq E' + F'}$$

$$\frac{c : E \leq E' \quad d : F \leq F'}{c \times d : E \times F \leq E' \times F'}$$

Double reading:

- Names of axioms and rules for representing derivations as terms.
- Base functions and combinators for building coercions



Computational interpretation of coercions

$$\begin{aligned} \text{shuffle}(\text{inl } v) &= \text{inl } (\text{inl } v) \\ \text{shuffle}(\text{inr } (\text{inl } v)) &= \text{inl } (\text{inr } v) \\ \text{shuffle}(\text{inr } (\text{inr } v)) &= \text{inr } v \\ \text{shuffle}^{-1}(\text{inl } (\text{inl } v)) &= \text{inl } v \\ \text{shuffle}^{-1}(\text{inl } (\text{inr } v)) &= \text{inr } (\text{inl } v) \\ \text{shuffle}^{-1}(\text{inr } v) &= \text{inr } (\text{inr } v) \\ \text{retag}(\text{inl } v) &= \text{inr } v \\ \text{retag}(\text{inr } v) &= \text{inl } v \\ \text{retag}^{-1} &= \text{retag} \\ \text{untagL } (\text{inr } v) &= v \\ \text{untag } (\text{inl } v) &= v \\ \text{untag } (\text{inr } v) &= v \\ \text{tagL } (v) &= \text{inl } v \\ \text{assoc}(v, (w, x)) &= ((v, w), x) \\ \text{assoc}^{-1}((v, w), x) &= (v, (w, x)) \end{aligned}$$



Computational interpretation of coercions (2)

$$\begin{aligned}\text{swap}(v, ()) &= ((), v) \\ \text{swap}^{-1}(() , v) &= (v, ()) \\ \text{proj}(() , w) &= w \\ \text{proj}^{-1}(w) &= ((), w) \\ \text{distL}(v, \text{inl } w) &= \text{inl } (v, w) \\ \text{distL}(v, \text{inr } x) &= \text{inr } (v, x) \\ \text{distL}^{-1}(\text{inl } (v, w)) &= (v, \text{inl } w) \\ \text{distL}^{-1}(\text{inr } (v, x)) &= (v, \text{inr } x) \\ \text{distR}(\text{inl } v, w) &= \text{inl } (v, w) \\ \text{distR}(\text{inr } v, x) &= \text{inr } (v, x) \\ \text{distR}^{-1}(\text{inl } (v, w)) &= (\text{inl } v, w) \\ \text{distR}^{-1}(\text{inr } (v, x)) &= (\text{inr } v, x) \\ \text{wrap}(v) &= \text{fold } v \\ \text{wrap}^{-1}(v) &= \text{fold}^{-1} v\end{aligned}$$



Computational interpretation of coercions (3)

Type constructors as *functors*:

$$\begin{aligned} \text{id}(v) &= v \\ \text{id}^{-1} &= \text{id} \\ (c; d)(v) &= d(c(v)) \\ (c + d)(\text{inl } v) &= \text{inl } (c(v)) \\ (c + d)(\text{inr } w) &= \text{inr } (d(w)) \\ (c \times d)(v, w) &= (c(v), d(w)) \end{aligned}$$



Kozen Axiomatization (1994)

$$\begin{array}{lcl}
 \text{shuffle} & : & E + (F + G) = (E + F) + G \\
 \text{retag} & : & E + F = F + E \\
 \text{untagL} & : & 0 + F = F \\
 \text{untag} & : & E + E \leq E \\
 \text{tagL} & : & E \leq E + F \\
 \text{assoc} & : & E \times (F \times G) = (E \times F) \times G \\
 \text{swap} & : & E \times 1 = 1 \times E \\
 \text{proj} & : & 1 \times E = E \\
 \text{abortR} & : & E \times 0 = 0 \\
 \text{abortL} & : & 0 \times E = 0 \\
 \text{distL} & : & E \times (F + G) = (E \times F) + (E \times G) \\
 \text{distR} & : & (E + F) \times G = (E \times G) + (F \times G) \\
 \text{wrap} & : & 1 + E \times E^* = E^* \\
 \text{id} & : & E = E
 \end{array}$$

$$\frac{c : E \leq E' \quad d : E' \leq E''}{c; d : E \leq E''}$$

$$\frac{c : E \leq E' \quad d : F \leq F'}{c + d : E + F \leq E' + F'}$$

$$\frac{c : E \leq E' \quad d : F \leq F'}{c \times d : E \times F \leq E' \times F'}$$

$$\frac{c : E \times F \leq F}{\text{foldr}[c] : E^* \times F \leq F}$$

$$\frac{c : E \times F \leq E}{\text{foldl}[c] : E \times F^* \leq E}$$



Example proofs

$$\text{cons} = \dots : a \times a^* \leq a^*$$

$$\text{foldl}[\text{foldr}[\text{cons}]] : a^* \times a^{**} \leq a^*$$

Short proof. But: $\text{foldl}[\text{foldr}[\text{cons}]]$ runs in quadratic time!
Does Kozen's axiomatization contain a "faster" proof?



Henglein-Nielsen axiomatization (2011)

- Weak axiomatization + 1 rule:
Safe coinduction (= terminating recursion).
- Contains all proofs of Salomaa, Kozen, Grabmayer.
- Numerous containment proofs that yield the *identity* (noop) on bit-coded parse trees or a *finite state transducer* (streaming linear-time, with constant-sized buffer).
 - ▶ How to efficiently synthesize these? (Future work)



A (Somewhat) Realistic Scenario

Fix malformed CSV data:

12,0,11,2,13,1



12,0;11,2;13,1



RE

$$r = [0-9]^+, [0-9]^+$$

$$t = (r,)^*r$$

$L(t)$ contains all

$$w_0, w_1, \dots, w_{n-1}, w_n$$

where $w_i \in L(r)$.



A PCRE attempt

Perl Compliant Regular Expressions (PCRE): Only one match under Kleene star.

$$s/(r,)* (r) /???/$$


A PCRE attempt

Perl Compliant Regular Expressions (PCRE): Only one match under Kleene star.

$$s / \underbrace{(r,)}_1 * \underbrace{(r)}_2 / ??? /$$



A PCRE attempt

Perl Compliant Regular Expressions (PCRE): Only one match under Kleene star.

$$s / \underbrace{(r,)}_1 * \underbrace{(r)}_2 / ??? /$$

$$\underbrace{w_0, w_1, \dots, w_{n-1}}_{\text{lost!}}, \underbrace{w_{n-1}}_1, \underbrace{w_n}_2$$



Using PCRE

Iteration needs to be hand-coded.

In Python:

```
e = '([0-9]+,[0-9]+)|(,)'
for m in re.finditer(e, text):
    if m.group(1):
        out += m.group(1)
    else:
        out += ";"
```



Substitutions, functorially

$$(r,) * r$$


Substitutions, functorially

$$(r,) * r \sim \text{List}(r \times \{, \}) \times r$$



Substitutions, functorially

$$(r, _) * r \sim \text{List}(r \times \{, \}) \times r$$

Matches are *values*:

$$([(w_1, _), \dots, (w_{n-1}, _)], w_n)$$



Substitutions, functorially

$$(r, _) * r \sim \text{List}(r \times \{, \}) \times r$$



$$(\text{id} \times \text{semic})^* \times \text{id}$$

$$\text{List}(r \times \{; \}) \times r$$

Matches are *values*:

$$([(w_1, _), \dots, (w_{n-1}, _)], w_n)$$



Substitutions, functorially

$$\begin{array}{c}
 (r,) * r \sim \text{List}(r \times \{, \}) \times r \\
 \downarrow (\text{id} \times \text{semic})^* \times \text{id} \\
 \text{List}(r \times \{; \}) \times r
 \end{array}$$

Matches are *values*:

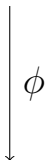
$$([(w_1, \underline{;}), \dots, (w_{n-1}, \underline{;})], w_n)$$



Nested Repetition

$$\text{List}(\text{List}(r \times \{, \}) \times r \times \{\backslash n\})$$


Nested Repetition

$$\text{List}(\text{List}(r \times \{, \}) \times r \times \{\backslash n\})$$

$$\text{List}(\text{List}(r \times \{; \}) \times r \times \{\backslash n\})$$

$$\phi = ((\text{id} \times \text{semic})^* \times \text{id} \times \text{id})^*$$



Using PCRE

Nested iteration needs to be hand-coded.

```
e = '([0-9]+,[0-9]+)|(,)'
for l in text.split('\n'):
    for m in re.finditer(e, l):
        ...
    out += "\n"
```



Summary

- Coercion = Transformation that changes RE but retains string
 - ▶ Proofs of containment induce string representation transformations
- Substitution = Structural transformation that changes strings
 - ▶ Includes *projections*, discarding parts of the input
- Synthesizing *efficient* coercions important (future work)
- Functorial transformations to express substitutions Useful for compact RE-specific representations of strings

