

Inside the Class of REGEX Languages

Markus L. Schmid,
Loughborough University, UK

DLT 2012

Extended Regular Expressions with Backreferences (REGEX)

REGEX are a combination of *regular expressions* and the concept of *homomorphic replacement*.

Extended Regular Expressions with Backreferences (REGEX)

REGEX are a combination of *regular expressions* and the concept of *homomorphic replacement*.

$(a \mid b)^* \quad c b^*$

Extended Regular Expressions with Backreferences (REGEX)

REGEX are a combination of *regular expressions* and the concept of *homomorphic replacement*.

$$(a \mid b)^* \quad c \quad b^*$$
$$\{wcb^n \mid w \in \{a, b\}^*, n \geq 0\}$$

Extended Regular Expressions with Backreferences (REGEX)

REGEX are a combination of *regular expressions* and the concept of *homomorphic replacement*.

$(\color{red}{1} (a \mid b)^* \color{red}{1})_1 c b^*$
 $\{wcb^n \mid w \in \{a, b\}^*, n \geq 0\}$

Extended Regular Expressions with Backreferences (REGEX)

REGEX are a combination of *regular expressions* and the concept of *homomorphic replacement*.

$(\color{red}{1} (a \mid b)^* \color{red}{1})_1 c b^* \backslash 1$
 $\{wcb^n \mid w \in \{a, b\}^*, n \geq 0\}$

Extended Regular Expressions with Backreferences (REGEX)

REGEX are a combination of *regular expressions* and the concept of *homomorphic replacement*.

$(\color{red}{_1} (a \mid b)^* \color{red}{_1} c b^* \backslash \color{red}{1})$
 $\{wcb^n w \mid w \in \{a, b\}^*, n \geq 0\}$

Extended Regular Expressions with Backreferences (REGEX)

REGEX are a combination of *regular expressions* and the concept of *homomorphic replacement*.

$$((a | b)^*)_1 c b^* \backslash 1$$
$$\{ w c b^n w \mid w \in \{a, b\}^*, n \geq 0 \}$$

$$(((a^*)_2 b \backslash 2)_1 c \backslash 1$$
$$\{ a^n b a^n c a^n b a^n \mid n \geq 0 \}$$

Practical Relevance of REGEX

- REGEX are intensely applied in practice...
 - Traditional and Modern grep
 - vi
 - Modern sed
 - GNU Emacs
 - Perl
 - Python
 - Java
 - .Net

Practical Relevance of REGEX

- REGEX are intensely applied in practice...
 - Traditional and Modern grep
 - vi
 - Modern sed
 - GNU Emacs
 - Perl
 - Python
 - Java
 - .Net
- ...even though their membership problem is NP-complete.

An Easy Example

$((a | b)^* c^*)_1 ((b | d))_2 (\backslash 1 | \backslash 2)$

An Easy Example

(**1** (a | b)* c* **1**)₁

(**2** (b | d) **2**)₂

(\b1 | \b2)

An Easy Example

$(\color{red}{1} (a \mid b)^* c^* \color{red}{1})_1$

x_1

$(\color{blue}{2} (b \mid d) \color{blue}{2})_2$

x_2

$(\color{red}{\backslash 1} \mid \color{blue}{\backslash 2})$

$(\color{red}{x_1} \mid \color{blue}{x_2})$

An Easy Example

$(\color{red}{1} (a | b)^* c^* \color{red}{1})$

$(\color{blue}{2} (b | d) \color{blue}{2})$

$(\color{red}{\backslash 1} | \color{blue}{\backslash 2})$

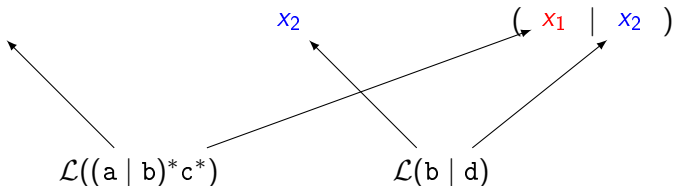
x_1

x_2

$(\color{red}{x_1} | \color{blue}{x_2})$

$\mathcal{L}((a | b)^* c^*)$

$\mathcal{L}(b | d)$



A More Complex Example

(1 (2 (3 (a | b)*)3 c\3)2 d (\2 c)*)1 e \1

A More Complex Example

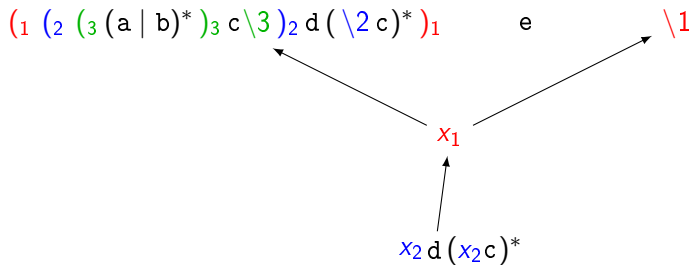
(₁ (₂ (₃ (a | b)*)₃ c \3)₂ d (\2 c)*)₁ e \1

A More Complex Example

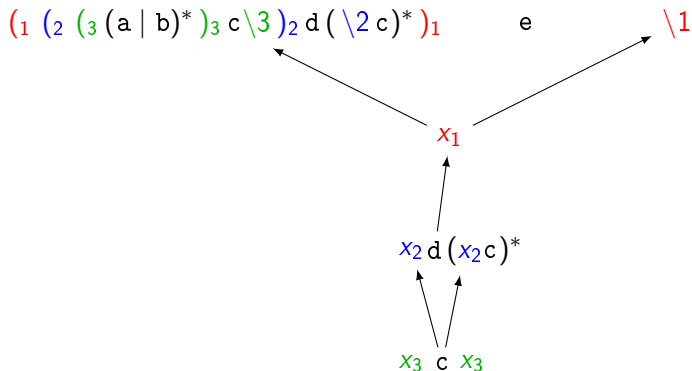
(₁ (₂ (₃ (a | b)*)₃ c \3)₂ d (\2 c)*)₁ e \1

The diagram illustrates a complex regular expression with nested backreferences. The expression is: `(1 (2 (3 (a | b)*)3 c \3)2 d (\2 c)*)1 e \1`. The subexpressions are numbered: ₁ is the outermost group, ₂ is the middle group, and ₃ is the innermost group. The backreference `\3` refers to the innermost group, `\2` refers to the middle group, and `\1` refers to the outermost group. A red `x1` is positioned below the expression, with two arrows pointing to the `\3` and `\1` backreferences, indicating that the same variable `x1` is used to refer to different parts of the expression.

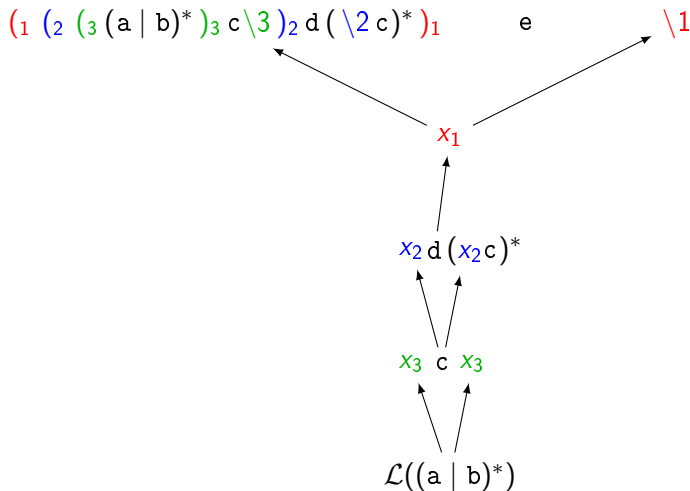
A More Complex Example



A More Complex Example



A More Complex Example



A Fairly Involved Example

$(a^*)_1 c (b | c)_2 c (b | d)^*_2 c$

A Fairly Involved Example

$(\color{red}{1} a^* \color{red}{1})_1$ c $(\color{blue}{2} \color{red}{\backslash 1} (b | c) \color{red}{\backslash 1})_2$ c $(\color{green}{3} (\color{red}{\backslash 1} | d)^* \color{blue}{\backslash 2})_3$ $\color{green}{\backslash 3}$

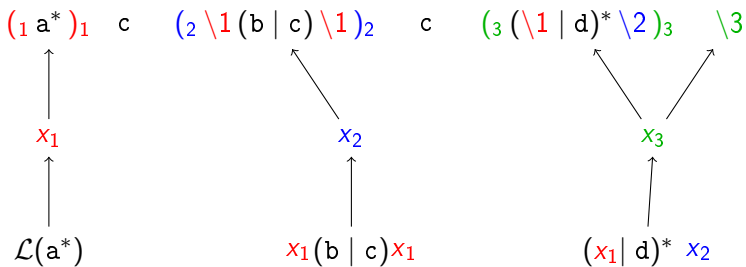
A Fairly Involved Example

$(x_1 a^*)_1 c (x_2 \backslash 1 (b | c) \backslash 1)_2 c (x_3 (\backslash 1 | d)^* \backslash 2)_3 \backslash 3$

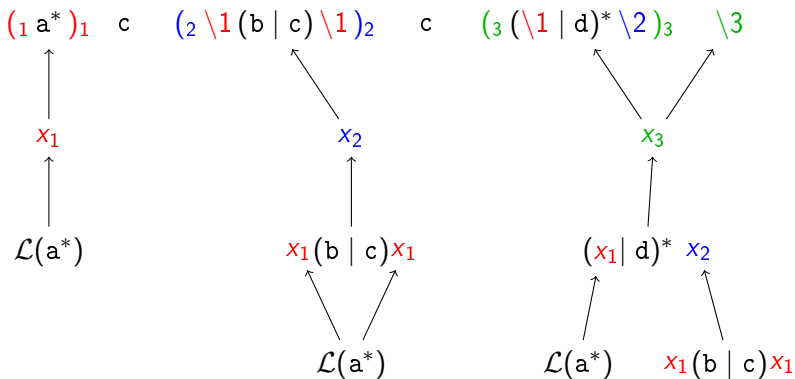
The diagram illustrates three regular expressions with backreferences, each associated with a variable x_i :

- Expression 1: $(x_1 a^*)_1 c$. An arrow points from x_1 to the first 1 .
- Expression 2: $(x_2 \backslash 1 (b | c) \backslash 1)_2 c$. An arrow points from x_2 to the first $\backslash 1$.
- Expression 3: $(x_3 (\backslash 1 | d)^* \backslash 2)_3 \backslash 3$. An arrow points from x_3 to the first $\backslash 1$, and another arrow points from x_3 to $\backslash 2$.

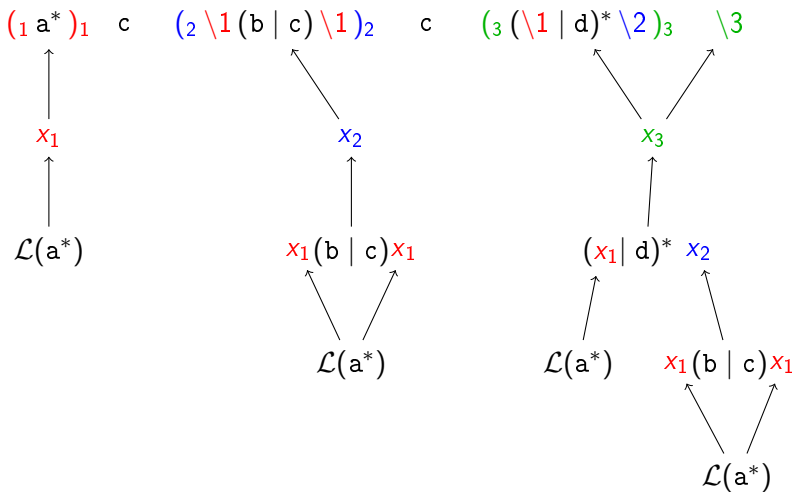
A Fairly Involved Example



A Fairly Involved Example



A Fairly Involved Example



Aim of This Paper

- Regular expressions on the one hand and homomorphic replacement on the other a well understood concepts in language theory.

Aim of This Paper

- Regular expressions on the one hand and homomorphic replacement on the other a well understood concepts in language theory.
- In REGEX, these two concepts seem inherently entangled and it seems difficult to treat them separately.

Aim of This Paper

- Regular expressions on the one hand and homomorphic replacement on the other a well understood concepts in language theory.
- In REGEX, these two concepts seem inherently entangled and it seems difficult to treat them separately.
- Our approach: Study REGEX by investigating alternative ways to combine regular expressions and homomorphic replacement...

Aim of This Paper

- Regular expressions on the one hand and homomorphic replacement on the other a well understood concepts in language theory.
- In REGEX, these two concepts seem inherently entangled and it seems difficult to treat them separately.
- Our approach: Study REGEX by investigating alternative ways to combine regular expressions and homomorphic replacement...
- ...without exceeding the expressive power of REGEX languages.

Aim of This Paper

- Regular expressions on the one hand and homomorphic replacement on the other a well understood concepts in language theory.
- In REGEX, these two concepts seem inherently entangled and it seems difficult to treat them separately.
- Our approach: Study REGEX by investigating alternative ways to combine regular expressions and homomorphic replacement...
- ...without exceeding the expressive power of REGEX languages.
- Informally: Take regular expressions, take some mechanism of homomorphic replacement, combine them and see how much of the class of REGEX languages we actually get.

(Typed) Pattern languages

Pattern: A word containing terminals (e. g. $\Sigma = \{a, b, c\}$) and variables ($X := \{x_1, x_2, x_3, \dots\}$).

(Typed) Pattern languages

Pattern: A word containing terminals (e. g. $\Sigma = \{a, b, c\}$) and variables ($X := \{x_1, x_2, x_3, \dots\}$).

$\alpha = x_1 x_2 b x_1 x_2 x_1$

(Typed) Pattern languages

Pattern: A word containing terminals (e. g. $\Sigma = \{a, b, c\}$) and variables ($X := \{x_1, x_2, x_3, \dots\}$).

$$\alpha = x_1 x_2 b x_1 x_2 x_1$$

$$\mathcal{L}_\Sigma(\alpha) = \{w \mid w = uvbuvu, u, v \in \Sigma^*\}.$$

(Typed) Pattern languages

Pattern: A word containing terminals (e. g. $\Sigma = \{a, b, c\}$) and variables ($X := \{x_1, x_2, x_3, \dots\}$).

$$\alpha = x_1 x_2 b x_1 x_2 x_1$$

$$\mathcal{L}_\Sigma(\alpha) = \{w \mid w = uvbuvu, u, v \in \Sigma^*\}.$$

A *type* for α : $\mathcal{T} := (T_{x_1}, T_{x_2})$

(Typed) Pattern languages

Pattern: A word containing terminals (e. g. $\Sigma = \{a, b, c\}$) and variables ($X := \{x_1, x_2, x_3, \dots\}$).

$$\alpha = x_1 x_2 b x_1 x_2 x_1$$

$$\mathcal{L}_\Sigma(\alpha) = \{w \mid w = uvbuvu, u, v \in \Sigma^*\}.$$

A *type* for α : $\mathcal{T} := (T_{x_1}, T_{x_2})$

$$\mathcal{L}_\mathcal{T}(\alpha) = \{w \mid w = uvbuvu, u \in T_{x_1}, v \in T_{x_2}\}.$$

(Typed) Pattern languages

$$\text{PAT} := (\Sigma \cup X)^+.$$

(Typed) Pattern languages

$\text{PAT} := (\Sigma \cup X)^+$.

$\text{var}(\alpha)$: Set of variables occurring in α .

E. g. $\text{var}(x_1 a b x_2 b a x_1 x_2 c x_3) = \{x_1, x_2, x_3\}$.

(Typed) Pattern languages

$$\text{PAT} := (\Sigma \cup X)^+.$$

$\text{var}(\alpha)$: Set of variables occurring in α .

E. g. $\text{var}(x_1 abx_2 bax_1 x_2 cx_3) = \{x_1, x_2, x_3\}$.

For any language class \mathfrak{L} ,

$$\mathcal{L}_{\mathfrak{L}}(\text{PAT}) := \{\mathcal{L}_{\mathcal{T}}(\alpha) \mid \alpha \in \text{PAT}, \mathcal{T} \in \mathfrak{L}^{|\text{var}(\alpha)|}\}.$$

(Typed) Pattern languages

$\text{PAT} := (\Sigma \cup X)^+$.

$\text{var}(\alpha)$: Set of variables occurring in α .

E. g. $\text{var}(x_1 a b x_2 b a x_1 x_2 c x_3) = \{x_1, x_2, x_3\}$.

For any language class \mathfrak{L} ,

$\mathcal{L}_{\mathfrak{L}}(\text{PAT}) := \{\mathcal{L}_{\mathcal{T}}(\alpha) \mid \alpha \in \text{PAT}, \mathcal{T} \in \mathfrak{L}^{|\text{var}(\alpha)|}\}$.

Proposition

$\mathcal{L}_{\text{REG}}(\text{PAT}) \subseteq \mathcal{L}(\text{REGEX})$.

(Typed) Pattern languages

Idea:

$$\mathcal{L}_1 := \mathcal{L}_{\text{REG}}(\text{PAT}),$$

(Typed) Pattern languages

Idea:

$$\mathcal{L}_1 := \mathcal{L}_{\text{REG}}(\text{PAT}),$$

$$\mathcal{L}_2 := \mathcal{L}_{\mathcal{L}_1}(\text{PAT}),$$

(Typed) Pattern languages

Idea:

$$\mathcal{L}_1 := \mathcal{L}_{\text{REG}}(\text{PAT}),$$

$$\mathcal{L}_2 := \mathcal{L}_{\mathcal{L}_1}(\text{PAT}),$$

$$\mathcal{L}_3 := \mathcal{L}_{\mathcal{L}_2}(\text{PAT}),$$

\vdots

(Typed) Pattern languages

Idea:

$$\mathcal{L}_1 := \mathcal{L}_{\text{REG}}(\text{PAT}),$$

$$\mathcal{L}_2 := \mathcal{L}_{\mathcal{L}_1}(\text{PAT}),$$

$$\mathcal{L}_3 := \mathcal{L}_{\mathcal{L}_2}(\text{PAT}),$$

\vdots

Proposition

For any class of languages \mathcal{L} , $\mathcal{L}_{\mathcal{L}}(\text{PAT}) = \mathcal{L}_{\mathcal{L}_{\mathcal{L}}(\text{PAT})}(\text{PAT})$.

(Typed) Pattern languages

Idea:

$$\mathcal{L}_1 := \mathcal{L}_{\text{REG}}(\text{PAT}),$$

$$\mathcal{L}_2 := \mathcal{L}_{\mathcal{L}_1}(\text{PAT}),$$

$$\mathcal{L}_3 := \mathcal{L}_{\mathcal{L}_2}(\text{PAT}),$$

\vdots

Proposition

For any class of languages \mathcal{L} , $\mathcal{L}_{\mathcal{L}}(\text{PAT}) = \mathcal{L}_{\mathcal{L}_{\mathcal{L}}(\text{PAT})}(\text{PAT})$.

Hence, the aspect of regular expressions cannot be limited to the type languages.

Patterns with Regular Operators

$\text{PAT}_{ro} := \{\alpha \mid \alpha \text{ is a regular expression over } (\Sigma \cup X)\}$. Every $\alpha \in \text{PAT}_{ro}$ is a *pattern with regular operators*.

Patterns with Regular Operators

$\text{PAT}_{ro} := \{\alpha \mid \alpha \text{ is a regular expression over } (\Sigma \cup X)\}$. Every $\alpha \in \text{PAT}_{ro}$ is a *pattern with regular operators*.

$\mathcal{L}_{\mathcal{T}}(\alpha) := \mathcal{L}_{\mathcal{T}}(\beta_1) \cup \mathcal{L}_{\mathcal{T}}(\beta_2) \cup \mathcal{L}_{\mathcal{T}}(\beta_3) \cup \dots$, where

Patterns with Regular Operators

$\text{PAT}_{ro} := \{\alpha \mid \alpha \text{ is a regular expression over } (\Sigma \cup X)\}$. Every $\alpha \in \text{PAT}_{ro}$ is a *pattern with regular operators*.

$\mathcal{L}_{\mathcal{T}}(\alpha) := \mathcal{L}_{\mathcal{T}}(\beta_1) \cup \mathcal{L}_{\mathcal{T}}(\beta_2) \cup \mathcal{L}_{\mathcal{T}}(\beta_3) \cup \dots$, where

- \mathcal{T} is a type for α and

Patterns with Regular Operators

$\text{PAT}_{ro} := \{\alpha \mid \alpha \text{ is a regular expression over } (\Sigma \cup X)\}$. Every $\alpha \in \text{PAT}_{ro}$ is a *pattern with regular operators*.

$\mathcal{L}_{\mathcal{T}}(\alpha) := \mathcal{L}_{\mathcal{T}}(\beta_1) \cup \mathcal{L}_{\mathcal{T}}(\beta_2) \cup \mathcal{L}_{\mathcal{T}}(\beta_3) \cup \dots$, where

- \mathcal{T} is a type for α and
- $\mathcal{L}(\alpha) = \{\beta_1, \beta_2, \beta_3, \dots\}$.

Patterns with Regular Operators

$\text{PAT}_{ro} := \{\alpha \mid \alpha \text{ is a regular expression over } (\Sigma \cup X)\}$. Every $\alpha \in \text{PAT}_{ro}$ is a *pattern with regular operators*.

$\mathcal{L}_{\mathcal{T}}(\alpha) := \mathcal{L}_{\mathcal{T}}(\beta_1) \cup \mathcal{L}_{\mathcal{T}}(\beta_2) \cup \mathcal{L}_{\mathcal{T}}(\beta_3) \cup \dots$, where

- \mathcal{T} is a type for α and
- $\mathcal{L}(\alpha) = \{\beta_1, \beta_2, \beta_3, \dots\}$.

Example: $\mathcal{L}_{(\mathcal{L}(b^*))}((x_1c)^+) =$

Patterns with Regular Operators

$\text{PAT}_{ro} := \{\alpha \mid \alpha \text{ is a regular expression over } (\Sigma \cup X)\}$. Every $\alpha \in \text{PAT}_{ro}$ is a *pattern with regular operators*.

$\mathcal{L}_{\mathcal{T}}(\alpha) := \mathcal{L}_{\mathcal{T}}(\beta_1) \cup \mathcal{L}_{\mathcal{T}}(\beta_2) \cup \mathcal{L}_{\mathcal{T}}(\beta_3) \cup \dots$, where

- \mathcal{T} is a type for α and
- $\mathcal{L}(\alpha) = \{\beta_1, \beta_2, \beta_3, \dots\}$.

Example: $\mathcal{L}_{(\mathcal{L}(b^*))}((x_1c)^+) =$
 $\mathcal{L}_{(\mathcal{L}(b^*))}(x_1c) \cup \mathcal{L}_{(\mathcal{L}(b^*))}(x_1cx_1c) \cup \mathcal{L}_{(\mathcal{L}(b^*))}(x_1cx_1cx_1c) \cup \dots =$

Patterns with Regular Operators

$\text{PAT}_{ro} := \{\alpha \mid \alpha \text{ is a regular expression over } (\Sigma \cup X)\}$. Every $\alpha \in \text{PAT}_{ro}$ is a *pattern with regular operators*.

$\mathcal{L}_{\mathcal{T}}(\alpha) := \mathcal{L}_{\mathcal{T}}(\beta_1) \cup \mathcal{L}_{\mathcal{T}}(\beta_2) \cup \mathcal{L}_{\mathcal{T}}(\beta_3) \cup \dots$, where

- \mathcal{T} is a type for α and
- $\mathcal{L}(\alpha) = \{\beta_1, \beta_2, \beta_3, \dots\}$.

Example: $\mathcal{L}_{(\mathcal{L}(b^*))}((x_1c)^+) =$
 $\mathcal{L}_{(\mathcal{L}(b^*))}(x_1c) \cup \mathcal{L}_{(\mathcal{L}(b^*))}(x_1cx_1c) \cup \mathcal{L}_{(\mathcal{L}(b^*))}(x_1cx_1cx_1c) \cup \dots =$
 $\{(b^n c)^m \mid n \geq 0, m \geq 1\}$.

Expressive Power

Theorem

$$\mathcal{L}_{\{\Sigma^*\}}(\text{PAT}) \subset \mathcal{L}_{\text{REG}}(\text{PAT}) \subset \mathcal{L}_{\text{REG}}(\text{PAT}_{\text{ro}}).$$

Iteratively Typing Patterns with Regular Operators

$\mathcal{L}_{ro,0} := \text{REG},$

Iteratively Typing Patterns with Regular Operators

$$\mathcal{L}_{ro,0} := \text{REG},$$

$$\mathcal{L}_{ro,1} := \mathcal{L}_{\mathcal{L}_{ro,0}}(\text{PAT}_{ro}) = \mathcal{L}_{\text{REG}}(\text{PAT}_{ro}),$$

Iteratively Typing Patterns with Regular Operators

$$\mathcal{L}_{ro,0} := \text{REG},$$

$$\mathcal{L}_{ro,1} := \mathcal{L}_{\mathcal{L}_{ro,0}}(\text{PAT}_{ro}) = \mathcal{L}_{\text{REG}}(\text{PAT}_{ro}),$$

$$\mathcal{L}_{ro,2} := \mathcal{L}_{\mathcal{L}_{ro,1}}(\text{PAT}_{ro}),$$

Iteratively Typing Patterns with Regular Operators

$$\mathcal{L}_{ro,0} := \text{REG},$$

$$\mathcal{L}_{ro,1} := \mathcal{L}_{\mathcal{L}_{ro,0}}(\text{PAT}_{ro}) = \mathcal{L}_{\text{REG}}(\text{PAT}_{ro}),$$

$$\mathcal{L}_{ro,2} := \mathcal{L}_{\mathcal{L}_{ro,1}}(\text{PAT}_{ro}),$$

$$\mathcal{L}_{ro,3} := \mathcal{L}_{\mathcal{L}_{ro,2}}(\text{PAT}_{ro}),$$

⋮

Iteratively Typing Patterns with Regular Operators

$$\mathcal{L}_{ro,0} := \text{REG},$$

$$\mathcal{L}_{ro,1} := \mathcal{L}_{\mathcal{L}_{ro,0}}(\text{PAT}_{ro}) = \mathcal{L}_{\text{REG}}(\text{PAT}_{ro}),$$

$$\mathcal{L}_{ro,2} := \mathcal{L}_{\mathcal{L}_{ro,1}}(\text{PAT}_{ro}),$$

$$\mathcal{L}_{ro,3} := \mathcal{L}_{\mathcal{L}_{ro,2}}(\text{PAT}_{ro}),$$

⋮

$$\mathcal{L}_{ro,\infty} := \bigcup_{i=0}^{\infty} \mathcal{L}_{ro,i}.$$

Iteratively Typing Patterns with Regular Operators

$$\mathcal{L}_{ro,0} := \text{REG},$$

$$\mathcal{L}_{ro,1} := \mathcal{L}_{\mathcal{L}_{ro,0}}(\text{PAT}_{ro}) = \mathcal{L}_{\text{REG}}(\text{PAT}_{ro}),$$

$$\mathcal{L}_{ro,2} := \mathcal{L}_{\mathcal{L}_{ro,1}}(\text{PAT}_{ro}),$$

$$\mathcal{L}_{ro,3} := \mathcal{L}_{\mathcal{L}_{ro,2}}(\text{PAT}_{ro}),$$

⋮

$$\mathcal{L}_{ro,\infty} := \bigcup_{i=0}^{\infty} \mathcal{L}_{ro,i}.$$

Theorem

$$\mathcal{L}_{ro,0} \subset \mathcal{L}_{ro,1} \subset \mathcal{L}_{ro,2} \subseteq \mathcal{L}_{ro,3} \subseteq \mathcal{L}_{ro,4} \subseteq \dots$$

Pattern Expressions

Introduced by Câmpeanu and Yu, 2004.

A *pattern expression* is a tuple

$$(x_1 \rightarrow r_1, x_2 \rightarrow r_2, \dots, x_n \rightarrow r_n),$$

Pattern Expressions

Introduced by Câmpeanu and Yu, 2004.

A *pattern expression* is a tuple

$$(x_1 \rightarrow r_1, x_2 \rightarrow r_2, \dots, x_n \rightarrow r_n),$$

where

- $\text{var}(r_1) = \emptyset,$
- $\text{var}(r_2) \subseteq \{x_1\},$
- $\text{var}(r_3) \subseteq \{x_1, x_2\},$
- $\text{var}(r_4) \subseteq \{x_1, x_2, x_3\},$
- \vdots

Pattern Expressions

Introduced by Câmpeanu and Yu, 2004.

A *pattern expression* is a tuple

$$(x_1 \rightarrow r_1, x_2 \rightarrow r_2, \dots, x_n \rightarrow r_n),$$

where

- $\text{var}(r_1) = \emptyset,$
- $\text{var}(r_2) \subseteq \{x_1\},$
- $\text{var}(r_3) \subseteq \{x_1, x_2\},$
- $\text{var}(r_4) \subseteq \{x_1, x_2, x_3\},$
- \vdots

The set of all pattern expressions is denoted by PE.

Pattern Expressions

Introduced by Câmpeanu and Yu, 2004.

A *pattern expression* is a tuple

$$(x_1 \rightarrow r_1, x_2 \rightarrow r_2, \dots, x_n \rightarrow r_n),$$

where

- $\text{var}(r_1) = \emptyset,$
- $\text{var}(r_2) \subseteq \{x_1\},$
- $\text{var}(r_3) \subseteq \{x_1, x_2\},$
- $\text{var}(r_4) \subseteq \{x_1, x_2, x_3\},$
- \vdots

The set of all pattern expressions is denoted by PE.

Example: $q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c \mid d)x_1, x_3 \rightarrow x_1cx_2).$

Pattern Expression Languages

$$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c \mid d)x_1, x_3 \rightarrow x_1cx_2),$$

Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c \mid d)x_1, x_3 \rightarrow x_1cx_2),$

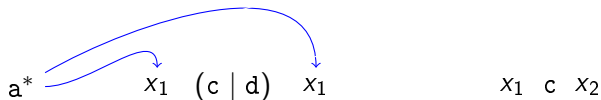
a^*

$x_1 (c \mid d) x_1$

$x_1 c x_2$

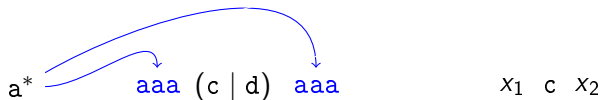
Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c \mid d)x_1, x_3 \rightarrow x_1cx_2),$



Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c | d)x_1, x_3 \rightarrow x_1cx_2),$




Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c | d)x_1, x_3 \rightarrow x_1cx_2),$

a^*

$aaa (c | d) aaa$

$x_1 c x_2$



Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c | d)x_1, x_3 \rightarrow x_1cx_2),$

a^*

aaa (c | d) aaa

x_1 c aaacaaa



Pattern Expression Languages

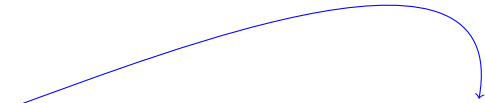
$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c | d)x_1, x_3 \rightarrow x_1cx_2),$

a^* $aaa (c | d) aaa$ $x_1 c aaacaaa$

Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c | d)x_1, x_3 \rightarrow x_1cx_2),$

a^* $aaa (c | d) aaa$ $a c aaacaaa$



Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c | d)x_1, x_3 \rightarrow x_1cx_2),$

a^*

aaa (c | d) aaa

a c aaacaaa

Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c \mid d)x_1, x_3 \rightarrow x_1cx_2),$

a^* $aaa (c \mid d) aaa$ $a c aaacaaa$

$\mathcal{L}_{it}(q) = \{a^k ca^m ua^m \mid k, m \in \mathbb{N}_0, u \in \{c, d\}\}$ is the *language generated by q with respect to iterated substitution.*

Pattern Expression Languages

$$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c \mid d)x_1, x_3 \rightarrow x_1cx_2),$$

Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c \mid d)x_1, x_3 \rightarrow x_1cx_2),$

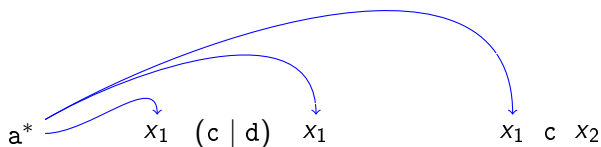
a^*

$x_1 (c \mid d) x_1$

$x_1 c x_2$

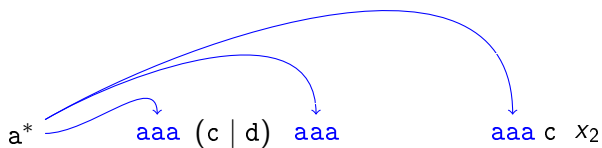
Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c \mid d)x_1, x_3 \rightarrow x_1cx_2),$



Pattern Expression Languages


$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c | d)x_1, x_3 \rightarrow x_1cx_2),$



Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c | d)x_1, x_3 \rightarrow x_1cx_2),$

a^* $aaa (c | d) aaa$ $aaa c x_2$



Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c | d)x_1, x_3 \rightarrow x_1cx_2),$

a^* $aaa (c | d) aaa$ $aaa c aaadaaa$



Pattern Expression Languages

$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c | d)x_1, x_3 \rightarrow x_1cx_2),$

a^*

aaa (c | d) aaa

aaa c aaadaaa

Pattern Expression Languages

$$q := (x_1 \rightarrow a^*, x_2 \rightarrow x_1(c \mid d)x_1, x_3 \rightarrow x_1cx_2),$$

a^* $aaa (c \mid d) aaa$ $aaa c aaadaaa$

$\mathcal{L}_{\text{uni}}(q) = \{a^mca^mua^m \mid m \in \mathbb{N}_0, u \in \{c, d\}\}$ is the *language generated by q with respect to **uniform** substitution.*

Pattern Expression Languages

Proposition

[Campeanu and Yu] For every $p \in \text{PE}$, $\mathcal{L}_{\text{it}}(p)$ is a REGEX language.

Theorem

$$\mathcal{L}_{\text{ro},\infty} = \mathcal{L}_{\text{it}}(\text{PE}).$$

Iterated vs. Uniform Substitution

Proposition

Let $p := (x_1 \rightarrow r_1, \dots, x_m \rightarrow r_m) \in \text{PE}$.

- $\mathcal{L}_{\text{uni}}(p) \subseteq \mathcal{L}_{\text{it}}(p)$,
- if, for every i, j , $1 \leq i < j \leq m$, $\text{var}(r_i) \cap \text{var}(r_j) = \emptyset$, then $\mathcal{L}_{\text{it}}(p) \subseteq \mathcal{L}_{\text{uni}}(p)$.

Iterated vs. Uniform Substitution

Proposition

Let $p := (x_1 \rightarrow r_1, \dots, x_m \rightarrow r_m) \in \text{PE}$.

- $\mathcal{L}_{\text{uni}}(p) \subseteq \mathcal{L}_{\text{it}}(p)$,
- if, for every i, j , $1 \leq i < j \leq m$, $\text{var}(r_i) \cap \text{var}(r_j) = \emptyset$, then $\mathcal{L}_{\text{it}}(p) \subseteq \mathcal{L}_{\text{uni}}(p)$.

Theorem

$\mathcal{L}_{\text{it}}(\text{PE}) \subset \mathcal{L}_{\text{uni}}(\text{PE})$.

Notation

A REGEX r is *star-free initialised* iff every referenced subexpression does not occur under a star.

Notation

A REGEX r is *star-free initialised* iff every referenced subexpression does not occur under a star.

- $((_1(a | b)^*_1 b \backslash 1)^* b \backslash 1$
- $(_1(a | b)^*_1 \backslash 1 ({}_2 c^*_2 (d \backslash 1 \backslash 2)^*$

PE w. r. t. uniform subst. vs. star-free initialised REGEX

Lemma

For every pattern expression p , there exists a star-free initialised REGEX r with $\mathcal{L}_{\text{uni}}(p) = \mathcal{L}(r)$.

PE w. r. t. uniform subst. vs. star-free initialised REGEX

Lemma

For every pattern expression p , there exists a star-free initialised REGEX r with $\mathcal{L}_{\text{uni}}(p) = \mathcal{L}(r)$.

Lemma

For every star-free initialised REGEX r , there exists a pattern expression p with $\mathcal{L}(r) = \mathcal{L}_{\text{uni}}(p)$.

PE w. r. t. uniform subst. vs. star-free initialised REGEX

Lemma

For every pattern expression p , there exists a star-free initialised REGEX r with $\mathcal{L}_{\text{uni}}(p) = \mathcal{L}(r)$.

Lemma

For every star-free initialised REGEX r , there exists a pattern expression p with $\mathcal{L}(r) = \mathcal{L}_{\text{uni}}(p)$.

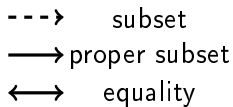
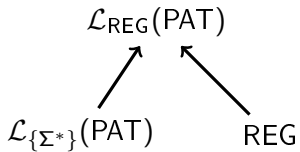
Theorem

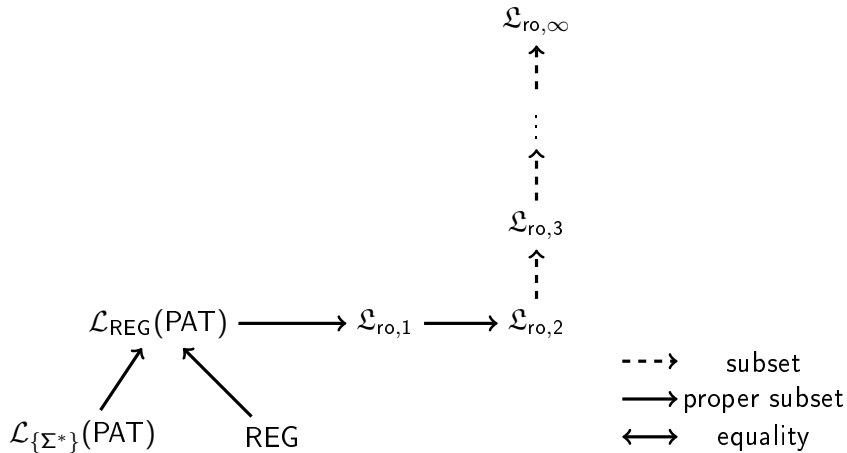
$\mathcal{L}(\text{REGEX}_{\text{sfi}}) = \mathcal{L}_{\text{uni}}(\text{PE})$.

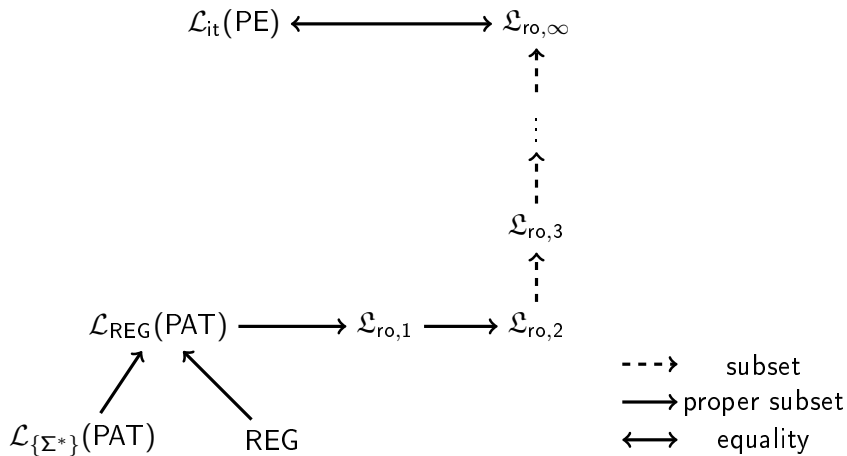
$\mathcal{L}_{\{\Sigma^*\}}(\text{PAT})$

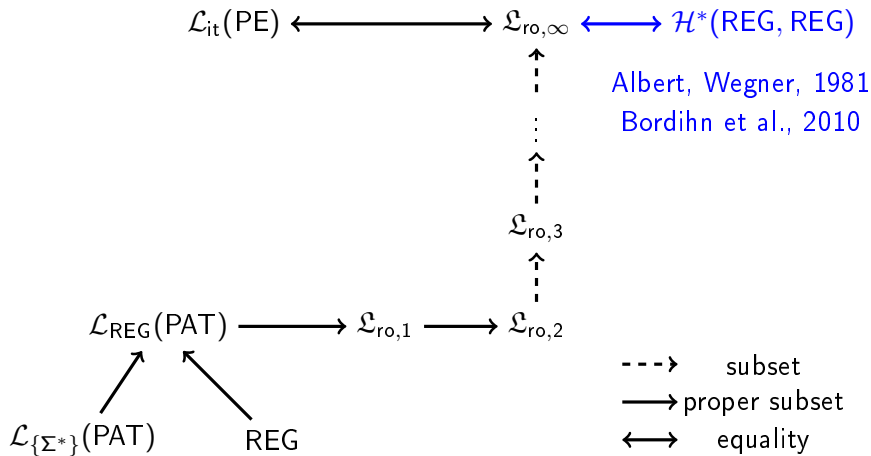
REG

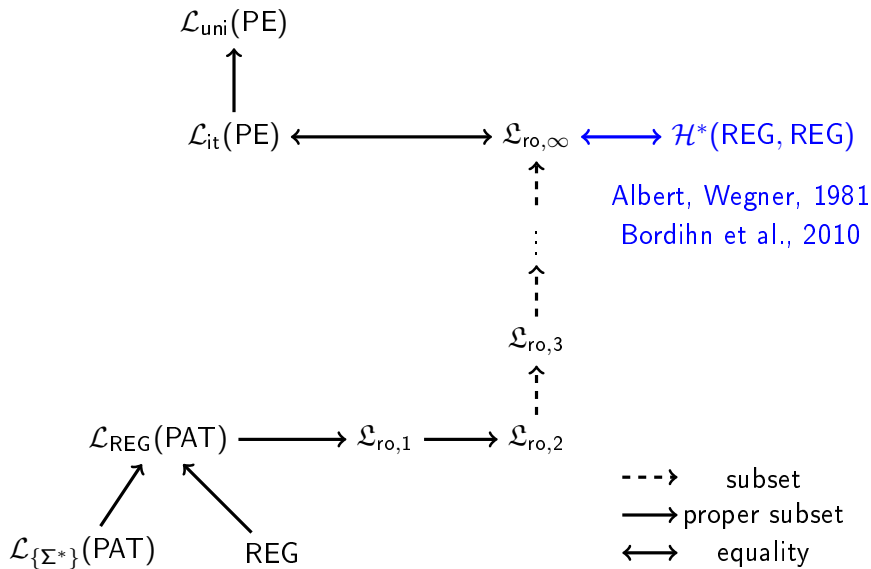
---> subset
-> proper subset
<=> equality

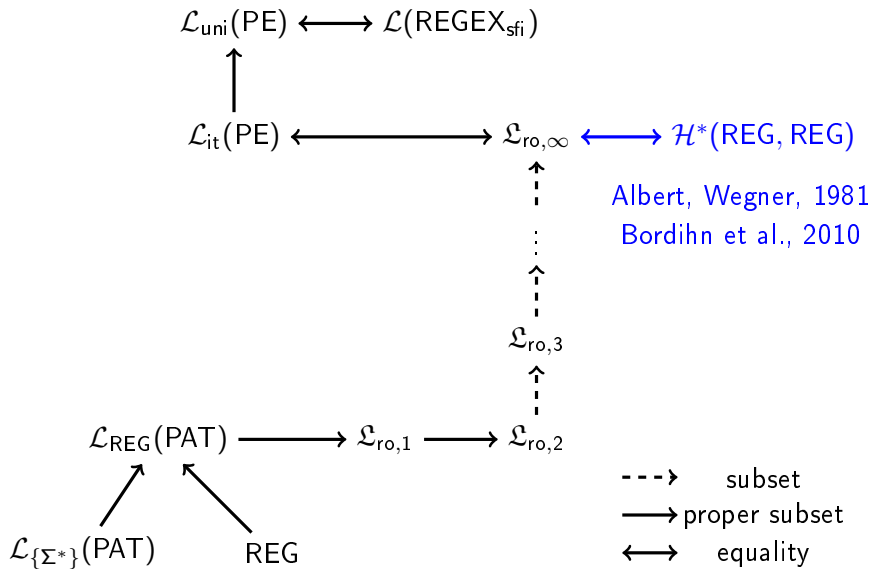












Albert, Wegner, 1981
 Bordihn et al., 2010

