

On expressive power of regular expressions with subroutine calls and lookahead assertions

Ondřej Guth
ondrej.guth@fit.cvut.cz

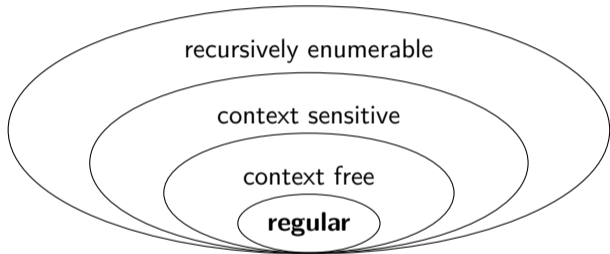
Department of Theoretical Computer Science
Faculty of Information Technology
Czech Technical University in Prague

The Prague Stringology Conference 2023

Regular expressions¹ match regular languages²

Regular expression is the minimal set of operations to express regular language

$(my + \epsilon)(great)^* grandfather$



¹Stephen Cole Kleene. "Representation of Events in Nerve Nets and Finite Automata". In: *Automata Studies*. (AM-34). Ed. by C. E. Shannon and J. McCarthy. Princeton University Press, Dec. 1956, pp. 3–42. ISBN: 978-1-4008-8261-8. DOI: 10.1515/9781400882618-002. (Visited on 08/20/2021).

²Noam Chomsky. "On Certain Formal Properties of Grammars". In: *Information and Control* (1959). ISSN: 00199958. DOI: 10.1016/S0019-9958(59)90362-6.

Practical “regular” expressions are different and have various flavours



PCRE2: $a^n b^n$
(a(?1)?b)

What language class can be expressed by a particular regex flavour?

?

Expressive power of certain combinations of features was already known

Character class, interval quantifier, concatenation, alternative, iteration: RLs

`[Mm]y ((great){1,3}-grand)(father|mother)`

Expressive power of certain combinations of features was already known

Character class, interval quantifier, concatenation, alternative, iteration: RLs

```
[Mm]y ((great){1,3}-grand)(father|mother)
```

Zero-width lookahead assertions: RLs

```
girls? need((?<=s need)|(?<=1 need)s)
```

Expressive power of certain combinations of features was already known

Character class, interval quantifier, concatenation, alternative, iteration: RLS

`[Mm]y ((great){1,3}-grand)(father|mother)`

Zero-width lookahead assertions: RLS

`girls? need((?<=s need)|(?<=1 need)s)`

Backreferences: some CSLs and not some CFLs

`([a-z]{10})\1`

Formalisation of regex: matching relation³

Other formalisms exist

Definition

A matching relation \rightsquigarrow is of the form $(r, x, i) \rightsquigarrow \mathcal{R}$ where $\mathcal{R} = \{i : i \in \mathbb{N} \wedge i \leq |x| + 1\}$ (matching result).

Definition

The language of a regex $r \in \mathbb{E}_{\text{LS}, \mathcal{A}, \mathcal{X}}$ is $L(r) = \{x : (r, x, 1) \rightsquigarrow \mathcal{R} \wedge |x| + 1 \in \mathcal{R}\}$.

³Nariyoshi Chida and Tachio Terauchi. “On Lookaheads in Regular Expressions with Backreferences”. In: *7th International Conference on Formal Structures for Computation and Deduction (FSCD 2022)*. Ed. by Amy P. Felty. Vol. 228. Leibniz International Proceedings in Informatics (LIPIcs). Dagstuhl, Germany: Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2022, 15:1–15:18. ISBN: 978-3-95977-233-4. DOI: 10.4230/LIPIcs.FSCD.2022.15.

Formalisation of regex: matching relation³

Other formalisms exist

Definition

A matching relation \rightsquigarrow is of the form $(r, \mathbf{x}, i) \rightsquigarrow \mathcal{R}$ where $\mathcal{R} = \{i : i \in \mathbb{N} \wedge i \leq |\mathbf{x}| + 1\}$ (matching result).

Definition

The language of a regex $r \in \mathbb{E}_{\text{LS}, \mathcal{A}, \mathcal{X}}$ is $L(r) = \{\mathbf{x} : (r, \mathbf{x}, 1) \rightsquigarrow \mathcal{R} \wedge |\mathbf{x}| + 1 \in \mathcal{R}\}$.

Example

$$\frac{a \in \mathcal{A} \wedge i \leq |\mathbf{x}| \wedge \mathbf{x}[i] = a}{(a, \mathbf{x}, i) \rightsquigarrow \{i + 1\}}$$

Example (Lookahead)

$$\frac{(r, \mathbf{x}, i) \rightsquigarrow \mathcal{R}}{((?=r), \mathbf{x}, i) \rightsquigarrow \{i \wedge \mathcal{R} \neq \emptyset\}}$$

³Nariyoshi Chida and Tachio Terauchi. “On Lookaheads in Regular Expressions with Backreferences”. In: *7th International Conference on Formal Structures for Computation and Deduction (FSCD 2022)*. Ed. by Amy P. Felty. Vol. 228. Leibniz International Proceedings in Informatics (LIPIcs). Dagstuhl, Germany: Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2022, 15:1–15:18. ISBN: 978-3-95977-233-4. DOI: 10.4230/LIPIcs.FSCD.2022.15.

Formalisation of subroutine calls

Definition (Numbered subroutine call)

$$\frac{(\sigma(l), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}{((?l), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}$$

Definition (Named subroutine call)

$$\frac{(\sigma(\nu(N)), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}{((?P>N), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}$$

Formalisation of subroutine calls

Definition (Numbered subroutine call)

$$\frac{(\sigma(l), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}{((?l), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}$$

Definition (Named subroutine call)

$$\frac{(\sigma(\nu(N)), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}{((?P>N), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}$$

Example

$$\frac{\dots}{\frac{\frac{(\mathbf{a}((?1) \mid \varepsilon)\mathbf{b}, \mathbf{ab}, 2) \rightsquigarrow \emptyset}{((?1), \mathbf{ab}, 2) \rightsquigarrow \emptyset} \quad (\varepsilon, \mathbf{ab}, 2) \rightsquigarrow \{2\}}{(\mathbf{a}, \mathbf{ab}, 1) \rightsquigarrow \{2\}} \quad \frac{((?1) \mid \varepsilon), \mathbf{ab}, 2 \rightsquigarrow \{2\}}{((?(\text{DEFINE}))(?<S>\mathbf{a}((?1) \mid \varepsilon)\mathbf{b}), \mathbf{ab}, 1) \rightsquigarrow \{1\}} \quad \frac{(\mathbf{b}, \mathbf{ab}, 2) \rightsquigarrow \{3\}}{(\mathbf{a}((?1) \mid \varepsilon)\mathbf{b}, \mathbf{ab}, 1) \rightsquigarrow \{3\}}}{((?P>S), \mathbf{ab}, 1) \rightsquigarrow \{3\}} \quad \frac{((?(\text{DEFINE}))(?<S>\mathbf{a}((?1) \mid \varepsilon)\mathbf{b}), \mathbf{ab}, 1) \rightsquigarrow \{1\}}{((?(\text{DEFINE}))(?<S>\mathbf{a}((?1) \mid \varepsilon)\mathbf{b})(?P>S), \mathbf{ab}, 1) \rightsquigarrow \{3\}}$$

Example (Context-free languages)

$\{a^n b^n\}$, $\{ww^R\}$

Example (Context-free languages)

$\{a^n b^n\}$, $\{ww^R\}$

Definition (Context-free grammar)

A quadruple $(\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$ where every member of \mathcal{R} is in the form of

$N \rightarrow \mathbf{v}$, $\mathbf{v} \in (\mathcal{A} \cup \mathcal{V})^*$.

Definition (Derivation step in CFG)

If $\mathbf{v}_1 \rightarrow \mathbf{v}_2 \in \mathcal{R}$ then

$\mathbf{pv}_1\mathbf{s} \Rightarrow \mathbf{pv}_2\mathbf{s}$ is possible.

Definition (Language generated by grammar)

$\{\mathbf{x} \in \mathcal{A}^* : S \Rightarrow^* \mathbf{x}\}$

Example (Context-free languages)

$\{a^n b^n\}$, $\{ww^R\}$

Definition (Context-free grammar)

A quadruple $(\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$ where every member of \mathcal{R} is in the form of $N \rightarrow \mathbf{v}$, $\mathbf{v} \in (\mathcal{A} \cup \mathcal{V})^*$.

Example

$(\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \varepsilon\}, S)$

S

Definition (Derivation step in CFG)

If $\mathbf{v}_1 \rightarrow \mathbf{v}_2 \in \mathcal{R}$ then $\mathbf{p}\mathbf{v}_1\mathbf{s} \Rightarrow \mathbf{p}\mathbf{v}_2\mathbf{s}$ is possible.

Definition (Language generated by grammar)

$\{x \in \mathcal{A}^* : S \Rightarrow^* x\}$

Example (Context-free languages)

$\{a^n b^n\}, \{ww^R\}$

Definition (Context-free grammar)

A quadruple $(\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$ where every member of \mathcal{R} is in the form of

$N \rightarrow \mathbf{v}, \mathbf{v} \in (\mathcal{A} \cup \mathcal{V})^*$.

Definition (Derivation step in CFG)

If $\mathbf{v}_1 \rightarrow \mathbf{v}_2 \in \mathcal{R}$ then

$\mathbf{pv}_1\mathbf{s} \Rightarrow \mathbf{pv}_2\mathbf{s}$ is possible.

Definition (Language generated by grammar)

$\{\mathbf{x} \in \mathcal{A}^* : S \Rightarrow^* \mathbf{x}\}$

Example

$(\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \varepsilon\}, S)$

$S \Rightarrow aSb$

Example (Context-free languages)

$\{a^n b^n\}$, $\{ww^R\}$

Definition (Context-free grammar)

A quadruple $(\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$ where every member of \mathcal{R} is in the form of

$N \rightarrow \mathbf{v}$, $\mathbf{v} \in (\mathcal{A} \cup \mathcal{V})^*$.

Definition (Derivation step in CFG)

If $\mathbf{v}_1 \rightarrow \mathbf{v}_2 \in \mathcal{R}$ then

$\mathbf{pv}_1\mathbf{s} \Rightarrow \mathbf{pv}_2\mathbf{s}$ is possible.

Definition (Language generated by grammar)

$\{\mathbf{x} \in \mathcal{A}^* : S \Rightarrow^* \mathbf{x}\}$

Example

$(\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \varepsilon\}, S)$

$S \Rightarrow aSb \Rightarrow aaSbb$

Example (Context-free languages)

$\{a^n b^n\}$, $\{ww^R\}$

Definition (Context-free grammar)

A quadruple $(\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$ where every member of \mathcal{R} is in the form of

$N \rightarrow \mathbf{v}$, $\mathbf{v} \in (\mathcal{A} \cup \mathcal{V})^*$.

Definition (Derivation step in CFG)

If $\mathbf{v}_1 \rightarrow \mathbf{v}_2 \in \mathcal{R}$ then

$\mathbf{pv}_1\mathbf{s} \Rightarrow \mathbf{pv}_2\mathbf{s}$ is possible.

Definition (Language generated by grammar)

$\{\mathbf{x} \in \mathcal{A}^* : S \Rightarrow^* \mathbf{x}\}$

Example

$(\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \varepsilon\}, S)$

$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aa\varepsilon bb$

Every context-free grammar can be expressed by a regex with subroutine calls

Conversion of a CFG to an equivalent regex

$(\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \varepsilon\}, S)$

1. $r_1 = \text{rx}(aSb) = a(?P>S)b$

Every context-free grammar can be expressed by a regex with subroutine calls

Conversion of a CFG to an equivalent regex

$(\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \varepsilon\}, S)$

1. $r_1 = \text{rx}(aSb) = a(?P>S)b$
2. $r_2 = \varepsilon$

Every context-free grammar can be expressed by a regex with subroutine calls

Conversion of a CFG to an equivalent regex

$(\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \varepsilon\}, S)$

1. $r_1 = \text{rx}(aSb) = a(?P>S)b$
2. $r_2 = \varepsilon$
3. $r = (?(\text{DEFINE})(?<S>r_1 \mid r_2))$

Every context-free grammar can be expressed by a regex with subroutine calls

Conversion of a CFG to an equivalent regex

$(\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \varepsilon\}, S)$

1. $r_1 = \text{rx}(aSb) = a(?P>S)b$
2. $r_2 = \varepsilon$
3. $r = (?(\text{DEFINE})(?<S>r_1 \mid r_2))$
4. $r = (?(\text{DEFINE})(?<S>r_1 \mid r_2))(?P>S)$

Every regex with concatenation, alternative, and subroutine call can be expressed by a context-free grammar

Conversion of a regex to an equivalent CFG

$(?P>N)a(?(\text{DEFINE})(?<N>\epsilon \mid b(?P>N)b \mid c(?P>N)c))$

► $\mathfrak{G}_\epsilon = (\{S_\epsilon\}, \{a, b, c\}, \{S_\epsilon \rightarrow \epsilon\}, S_\epsilon)$, $\mathfrak{G}_a = (\{S_a\}, \{a, b, c\}, \{S_a \rightarrow a\}, S_a) \dots$

Every regex with concatenation, alternative, and subroutine call can be expressed by a context-free grammar

Conversion of a regex to an equivalent CFG

$(?P>N)a(?(\text{DEFINE})(?<N>\epsilon \mid b(?P>N)b \mid c(?P>N)c))$

- ▶ $\mathfrak{G}_\epsilon = (\{S_\epsilon\}, \{a, b, c\}, \{S_\epsilon \rightarrow \epsilon\}, S_\epsilon)$, $\mathfrak{G}_a = (\{S_a\}, \{a, b, c\}, \{S_a \rightarrow a\}, S_a) \dots$
- ▶ $\mathfrak{G}_{(?P>N)} = (\{S_{(?P>N)}, N\}, \{a, b, c\}, \{S_{(?P>N)} \rightarrow N\}, S_{(?P>N)})$

Every regex with concatenation, alternative, and subroutine call can be expressed by a context-free grammar

Conversion of a regex to an equivalent CFG

$(?P>N)a(?(\text{DEFINE})(?<N>\epsilon \mid b(?P>N)b \mid c(?P>N)c))$

- ▶ $\mathfrak{G}_\epsilon = (\{S_\epsilon\}, \{a, b, c\}, \{S_\epsilon \rightarrow \epsilon\}, S_\epsilon)$, $\mathfrak{G}_a = (\{S_a\}, \{a, b, c\}, \{S_a \rightarrow a\}, S_a) \dots$
- ▶ $\mathfrak{G}_{(?P>N)} = (\{S_{(?P>N)}, N\}, \{a, b, c\}, \{S_{(?P>N)} \rightarrow N\}, S_{(?P>N)})$
- ▶ after constructing grammars for elementary expressions. . .

Every regex with concatenation, alternative, and subroutine call can be expressed by a context-free grammar

Conversion of a regex to an equivalent CFG

$(?P>N)a(?(\text{DEFINE})(?<N>\epsilon \mid b(?P>N)b \mid c(?P>N)c))$

- ▶ $\mathfrak{G}_\epsilon = (\{S_\epsilon\}, \{a, b, c\}, \{S_\epsilon \rightarrow \epsilon\}, S_\epsilon)$, $\mathfrak{G}_a = (\{S_a\}, \{a, b, c\}, \{S_a \rightarrow a\}, S_a) \dots$
- ▶ $\mathfrak{G}_{(?P>N)} = (\{S_{(?P>N)}, N\}, \{a, b, c\}, \{S_{(?P>N)} \rightarrow N\}, S_{(?P>N)})$
- ▶ after constructing grammars for elementary expressions. . .
- ▶ $\mathfrak{G}_{b(?P>N)} = (\mathcal{V}_b \cup \mathcal{V}_{(?P>N)} \cup \{S_{b(?P>N)}\}, \{a, b, c\}, \mathcal{R}_b \cup \mathcal{R}_{(?P>N)} \cup \{S_{b(?P>N)} \rightarrow S_b S_{(?P>N)}\}, S_{b(?P>N)}) : S_{b(?P>N)} \notin \mathcal{V}_b \cup \mathcal{V}_{(?P>N)} \dots$

Every regex with concatenation, alternative, and subroutine call can be expressed by a context-free grammar

Conversion of a regex to an equivalent CFG

$(?P>N)a(?(\text{DEFINE})(?<N>\epsilon \mid b(?P>N)b \mid c(?P>N)c))$

- ▶ $\mathfrak{G}_\epsilon = (\{S_\epsilon\}, \{a, b, c\}, \{S_\epsilon \rightarrow \epsilon\}, S_\epsilon)$, $\mathfrak{G}_a = (\{S_a\}, \{a, b, c\}, \{S_a \rightarrow a\}, S_a) \dots$
- ▶ $\mathfrak{G}_{(?P>N)} = (\{S_{(?P>N)}, N\}, \{a, b, c\}, \{S_{(?P>N)} \rightarrow N\}, S_{(?P>N)})$
- ▶ after constructing grammars for elementary expressions. . .
- ▶ $\mathfrak{G}_{b(?P>N)} = (\mathcal{V}_b \cup \mathcal{V}_{(?P>N)} \cup \{S_{b(?P>N)}\}, \{a, b, c\}, \mathcal{R}_b \cup \mathcal{R}_{(?P>N)} \cup \{S_{b(?P>N)} \rightarrow S_b S_{(?P>N)}\}, S_{b(?P>N)}) : S_{b(?P>N)} \notin \mathcal{V}_b \cup \mathcal{V}_{(?P>N)} \dots$
- ▶ $\mathfrak{G}_{\epsilon|b(?P>N)b} = (\mathcal{V}_\epsilon \cup \mathcal{V}_{b(?P>N)b} \cup \{S_{\epsilon|b(?P>N)b}\}, \{a, b, c\}, \mathcal{R}_\epsilon \cup \mathcal{R}_{b(?P>N)b} \cup \{S_{\epsilon|b(?P>N)b} \rightarrow S_\epsilon \mid S_{b(?P>N)b}\}, S_{\epsilon|b(?P>N)b}) : S_{\epsilon|b(?P>N)b} \notin \mathcal{V}_\epsilon \cup \mathcal{V}_{b(?P>N)b}$

Every regex with concatenation, alternative, and subroutine call can be expressed by a context-free grammar

Conversion of a regex to an equivalent CFG

$(?P>N)a(?(\text{DEFINE})(?<N>\epsilon \mid b(?P>N)b \mid c(?P>N)c))$

- ▶ $\mathfrak{G}_\epsilon = (\{S_\epsilon\}, \{a, b, c\}, \{S_\epsilon \rightarrow \epsilon\}, S_\epsilon)$, $\mathfrak{G}_a = (\{S_a\}, \{a, b, c\}, \{S_a \rightarrow a\}, S_a)$...
- ▶ $\mathfrak{G}_{(?P>N)} = (\{S_{(?P>N)}, N\}, \{a, b, c\}, \{S_{(?P>N)} \rightarrow N\}, S_{(?P>N)})$
- ▶ after constructing grammars for elementary expressions...
- ▶ $\mathfrak{G}_{b(?P>N)} = (\mathcal{V}_b \cup \mathcal{V}_{(?P>N)} \cup \{S_{b(?P>N)}\}, \{a, b, c\}, \mathcal{R}_b \cup \mathcal{R}_{(?P>N)} \cup \{S_{b(?P>N)} \rightarrow S_b S_{(?P>N)}\}, S_{b(?P>N)}) : S_{b(?P>N)} \notin \mathcal{V}_b \cup \mathcal{V}_{(?P>N)}$...
- ▶ $\mathfrak{G}_{\epsilon|b(?P>N)b} = (\mathcal{V}_\epsilon \cup \mathcal{V}_{b(?P>N)b} \cup \{S_{\epsilon|b(?P>N)b}\}, \{a, b, c\}, \mathcal{R}_\epsilon \cup \mathcal{R}_{b(?P>N)b} \cup \{S_{\epsilon|b(?P>N)b} \rightarrow S_\epsilon \mid S_{b(?P>N)b}\}, S_{\epsilon|b(?P>N)b}) : S_{\epsilon|b(?P>N)b} \notin \mathcal{V}_\epsilon \cup \mathcal{V}_{b(?P>N)b}$
- ▶ ...

Thus it is shown that subroutine calls match exactly context-free languages

Theorem

$$\mathbb{L}_{ES} = \mathbb{L}_{CF}$$

Adding lookahead to a regex with subroutine calls extends its expressive power

Theorem

$$\mathbb{L}_{ES} \subsetneq \mathbb{L}_{ELS}$$

Proof.

$$\{a^g b^g c^g : g \in \mathbb{N}\} =$$
$$L((?= (?<N_1>a(\epsilon | (?P>N_1))b) c)aa^*(?<N_2>b(\epsilon | (?P>N_2))c))$$

□

Summary

- ▶ expressive power of expressions with concatenation, alternative, and subroutine call is equivalent to the class of context-free languages
 - ▶ algorithms to convert between regex and context-free grammar
- ▶ expressive power of expressions with concatenation, alternative, subroutine call, and lookahead is beyond context-free languages

Future work: expressive power of regex with both subroutine call and lookahead.