

# On Expressive Power of Regular Expressions with Subroutine Calls and Lookaround Assertions

Ondřej Guth

Czech Technical University in Prague  
Thákurova 9  
16000 Praha  
Czech Republic  
`ondrej.guth@fit.cvut.cz`

**Abstract.** Many regular expression engines employ syntactical extensions to provide simple, expressive support for real-world needs. These features are subroutine calls, zero-width lookaround assertions, DEFINE rules, and named parenthesised expressions. A subroutine call executes a specified subpattern where the call is placed, possibly recursively. Lookaround assertions are either lookahead or lookbehind: a lookahead is a conditional within a subpattern: when it fails, the match at the current position of the whole subpattern fails, while a lookahead itself does not consume any input; a lookbehind works as a lookahead except it checks the input prior to the current position. A DEFINE rule introduces a subpattern for use by a subroutine call, while not involved in matching where the rule is placed. A named parenthesised expression can be executed by its name in addition to the parenthesis number. This paper presents a formalisation of subroutine calls, DEFINE rules, and named parenthesised expressions using the matching relation while attempting to mimic the behaviour of real-world regular expression engines. Also, we give an alternative constructive proof of equivalence of expressive power of regular expressions extended with subroutine calls and the class of context-free languages: a conversion between such expressions and context-free grammars. Finally, the question of whether regular expressions with operations lookaround assertion combined with subroutine call have greater expressive power than expressions with only subroutine call is answered positively.

## 1 Introduction

Regular expressions were introduced by Kleene[13] as a theoretical concept with expressive power equivalent to regular languages (these are further referred to as classical regular expressions or RE). This concept plays an important role in pattern matching and its variants with multiple (finite or infinite) patterns (see the taxonomy of pattern matching problems by Melichar and Holub[16]). So-called regular expressions have been implemented in many tools (e.g., UNIX text filters, text editors), programming languages, and libraries (these expressions are often referred to as extended regular expressions, practical regular expressions, or regexes). Unlike classical regular expressions, regexes “seem to have been invented entirely on the level of software implementation, without prior theoretical formalisation” (Schmid[22]). Moreover, both the syntax and semantics of the regex flavours used in implementations differ from each other (differences among the flavours were described by Friedl[8]). However, researchers have been exploring the algorithmic and language properties of practical regular expressions. Research results address properties of combinations of particular features used in regexes rather than complete flavours.

One of the fields in which this paper is concerned focuses on the expressive power of regexes and its relation to known language classes. Some syntactic constructs (not

used in classical regular expressions) are known to be mere syntax sugar: they can be rewritten to equivalent REs. Among these constructs are positive iteration (e.g., `a+`), character class (e.g., `[abc]`), or counting constraint (interval quantifier, e.g., `a{3,8}`) as pointed out, among others, by Câmpeanu et al.[4]. However, some features of regexes impact their expressive power: nonregular languages can be matched. A backreference indicates that a substring matched by a corresponding parenthesised subpattern should be matched again at the positions where the backreference is placed. Regexes with backreferences can match a proper subset of the class of context-sensitive languages (the expressive power of backreferences was studied, among others, by Câmpeanu et al.[4], Berglund et al.[2], or Schmid[22]). The expressive power of a practical regular expression with features of RE extended by lookahead stays within regular languages (Berglund et al.[3]). Regexes with subroutine calls describe context-free languages (addressed in master's thesis by Hruša[12]). When a practical regular expression with backreferences is extended by lookaheads, its expressive power supersedes the expressive power of regex with only backreferences (Chida and Terauchi[5]). Therefore, it is natural to wonder whether adding lookahead to a regex with subroutine calls also impacts expressive power. This question is addressed in this paper.

The formalisation of syntax and semantics of features of practical regular expressions is an essential part of proving their expressive power. Aho[1] gave a relatively informal definition of a regex with backreferences: the definition uses named variables, while any variable can be reassigned. Câmpeanu, Salomaa, and Yu[4] precisely formalised the numbered backreferences. Another formalisation of backreferences, factor-referencing, was introduced by Schmid[22]: it uses named variables which can be reassigned, and unlike the first formalisation[1], it deals with details of both syntax and semantics. Regexes with lookahead were originally formalised by Morihata[17] according to Berglund et al.[3]: the definition uses lookahead language. Chida and Terauchi[6] formalised the regexes with lookaheads and numbered backreferences using the matching relation. The syntax and semantics of the regexes with numbered subroutine calls were defined by Hruša[12]. To the author's knowledge, there is no formalisation of DEFINE rules, named parenthesised expressions, or named subroutine calls. This paper fills this gap by extending the notion matching relation.

Finding the expressive power of regex flavours is motivated by more than scientific curiosity. Users of pattern matching tools need to know what can and can not be matched by particular flavours of practical regular expressions<sup>1</sup>.

To the author's knowledge, there has been almost no research on the expressive power of regexes with subroutine calls. The first known text dealing with this gap was published as a blog post by Popov[19] providing a sketch of a reduction of context-free grammars to regexes with DEFINE rules and named subroutine calls. Popov's claim of the equivalence of expressive power of regexes with subroutine calls and context-free languages was later formally proved by Hruša[12] while using numbered-only subroutine calls. This paper gives an alternative proof to the Hruša's while using the matching relation and regexes with DEFINE rules and named subroutine calls. We hope that our approach is more straightforward and extensible in future research.

---

<sup>1</sup> As shown by several discussions at Stack Overflow, for example, <https://stackoverflow.com/q/35449863>, <https://stackoverflow.com/q/2974210>, or <https://stackoverflow.com/q/4840988>.

To the author's knowledge, there is no peer-reviewed publication on the expressive power of practical regular expressions with both subroutine calls and lookahead assertions. The problem seems to have a solution due to Popov's[19] sketch (or, more precisely, an idea) of a reduction from context-sensitive grammars to regexes<sup>2</sup>. However, we show a counterexample. In addition, we present proof that the expressive power of practical regular expressions with lookahead assertions and subroutine calls is greater than the expressive power of expressions with only subroutine calls.

This paper focuses on the expressive power of regex with subroutine calls (sub-pattern recursion) and lookahead assertions. Our contributions are the following:

- We give a formalisation of practical regular expressions by extending the notion of the matching relation. In particular, this paper gives the formalisation of constructs named parenthesised expression, DEFINE rule, and both numbered and named subroutine calls. Our formalisation mimics the syntax and semantics of Perl-Compatible Regular Expressions (PCRE2, as documented on the manual page[10] and as we experimentally verified). At the same time, it also works for Perl regular expressions[18] and Ruby Regexp class[20].
- We prove that the expressive power of regex with concatenation, alternative, DEFINE rule, and subroutine call is equal to the class of context-free languages. This proof is based on the matching relation and works for regex with numbered and named subroutines.
- We show that adding lookahead assertions to regexes with subroutine calls extends their expressive power beyond context-free languages. In addition, we also show that the equivalence of the expressive power of these regexes to the class of context-sensitive languages remains an open problem.

This paper is structured as follows. In section 2, we give notational preliminaries. Section 3 contains the formalisation of practical regular expressions with subroutine calls. In section 4, we present proof that the expressive power of regexes with subroutine calls is equal to context-free languages: a conversion between such a regex and context-free grammar. Section 5 contains proof that adding lookahead assertion extends the expressive power of practical regular expressions with subroutine call. In section 6, we conclude the paper and discuss future work.

## 2 Preliminaries

The set of natural numbers is denoted by  $\mathbb{N}$  and is without zero. The mathematical symbols  $\emptyset$ ,  $\cup$ ,  $\cap$ ,  $\setminus$ ,  $\wedge$ ,  $\vee$ , and  $\nexists$  denote the empty set, set union, set intersection, set difference, logical conjunction, universal quantifier, and negated existential quantifier, respectively. If  $\mathcal{V}$  and  $\mathcal{X}$  are sets,  $\mathcal{V}$  being a subset (or a strict subset) of  $\mathcal{X}$  is denoted by  $\mathcal{V} \subseteq \mathcal{X}$  (or  $\mathcal{V} \subsetneq \mathcal{X}$ , respectively). An *alphabet*, denoted by  $\mathcal{A}$ , is a finite nonempty set whose elements are called *symbols*. A *string* over  $\mathcal{A}$  is a finite sequence of elements of  $\mathcal{A}$ . The empty sequence is called the *empty string* and is denoted by  $\varepsilon$ . The *set of all strings* over  $\mathcal{A}$  is denoted by  $\mathcal{A}^*$  and the *set of all nonempty strings* over  $\mathcal{A}$  is denoted by  $\mathcal{A}^+$ . The *length* of a string  $\mathbf{y}$  is the length of the sequence associated with

---

<sup>2</sup> The reduction of context-sensitive grammars to regexes seem to be trusted: for example, a Stack Exchange contributor claims that regexes with subroutine calls and lookahead assertions can express any context-sensitive language using Popov's argument: <https://cs.stackexchange.com/q/143221>.

$\mathbf{y}$  and is denoted by  $|\mathbf{y}|$ . By  $\mathbf{y}[i]$ , where  $i \in \mathbb{N} \wedge i \in \{1, \dots, |\mathbf{y}|\}$ , we denote the symbol at the index  $i$  of  $\mathbf{y}$ . The concatenation of strings  $\mathbf{y}_1$  and  $\mathbf{y}_2$  is denoted by  $\mathbf{y}_1\mathbf{y}_2$ . Thus,  $\mathbf{y} = \mathbf{y}[1]\mathbf{y}[2] \dots \mathbf{y}[|\mathbf{y}|]$ . The substring of  $\mathbf{y}$  that starts at the index  $i$  and ends at the index  $g$  is denoted by  $\mathbf{y}[i..g]$ ; that is,  $\mathbf{y}[i..g] = \mathbf{y}[i]\mathbf{y}[i+1] \dots \mathbf{y}[g]$ . A *language* over an alphabet  $\mathcal{A}$  is a set of strings over  $\mathcal{A}$ , denoted by  $\mathcal{L} \subseteq \mathcal{A}^*$ . The *concatenation of languages*  $\mathcal{L}_1, \mathcal{L}_2$  is denoted by  $\mathcal{L}_1 \cdot \mathcal{L}_2$  and is defined as  $\mathcal{L}_1 \cdot \mathcal{L}_2 = \{\mathbf{y} = \mathbf{y}_1\mathbf{y}_2 : \mathbf{y}_1 \in \mathcal{L}_1 \wedge \mathbf{y}_2 \in \mathcal{L}_2\}$ . The *closure of a language*  $\mathcal{L}$  is denoted by  $\mathcal{L}^*$  and is defined as  $\bigcup_{g \geq 0} \mathcal{L}^g$  where  $\mathcal{L}^0 = \{\varepsilon\}$ ,  $\mathcal{L}^1 = \mathcal{L}$ , and for  $g > 1$ :  $\mathcal{L}^g = \mathcal{L} \cdot \mathcal{L}^{g-1}$ .

The following grammar-related notions follow the conventions of Hopcroft et al.[11] and Mateescu et al.[15]. A *grammar* is a quadruple  $\mathfrak{G} = (\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$  where  $\mathcal{A} \cap \mathcal{V} = \emptyset$ ,  $S \in \mathcal{V}$ , and  $\mathcal{R}$  is a set of pairs  $(\mathbf{v}_1, \mathbf{v}_2)$  where  $\mathbf{v}_1 \in (\mathcal{A} \cup \mathcal{V})^* \cap \mathcal{V}^+$  and  $\mathbf{v}_2 \in (\mathcal{A} \cup \mathcal{V})^*$ . The sets  $\mathcal{V}$  (nonterminals),  $\mathcal{A}$ , and  $\mathcal{R}$  are finite. We use the following naming and typographic conventions:  $a \in \mathcal{A}$ ,  $N \in \mathcal{V}$ ,  $\mathbf{x}, \mathbf{y} \in \mathcal{A}^*$ , and  $\mathbf{v} \in (\mathcal{A} \cup \mathcal{V})^*$  (bold italic sans serif for a string that can contain a nonterminal). The members of the set  $\mathcal{R}$  are called *productions* and are written with  $\rightarrow$  as a delimiter of the left- and right-hand side. Multiple productions with the same left-hand side can be contracted: for instance, if  $\mathbf{v} \rightarrow \mathbf{v}_1, \mathbf{v} \rightarrow \mathbf{v}_2 \in \mathcal{R}$  then we can write  $\mathbf{v} \rightarrow \mathbf{v}_1 \mid \mathbf{v}_2 \in \mathcal{R}$ . A *derivation step* in grammar  $\mathfrak{G} = (\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$  is denoted by  $\xrightarrow{\mathfrak{G}}$  and defined as follows. If  $\mathbf{v}_1 \in (\mathcal{A} \cup \mathcal{V})^+, \mathbf{p}, \mathbf{s}, \mathbf{v}_2 \in (\mathcal{A} \cup \mathcal{V})^*$ , and  $\mathbf{v}_1 \rightarrow \mathbf{v}_2 \in \mathcal{R}$  then  $\mathbf{p}\mathbf{v}_1\mathbf{s} \xrightarrow{\mathfrak{G}} \mathbf{p}\mathbf{v}_2\mathbf{s}$ . The transitive closure of  $\xrightarrow{\mathfrak{G}}$  is denoted by  $\xrightarrow{\mathfrak{G}}^+$ , and the reflective and transitive closure is denoted by  $\xrightarrow{\mathfrak{G}}^*$ . If  $S \xrightarrow{\mathfrak{G}}^* \mathbf{v}$  then the string  $\mathbf{v}$  is called a *sentential form* and  $S \xrightarrow{\mathfrak{G}}^* \mathbf{v}$  is called a *derivation* of  $\mathbf{v}$  in  $\mathfrak{G}$ . The *language generated by grammar*  $\mathfrak{G} = (\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$  is denoted by  $L(\mathfrak{G})$  and is defined as  $L(\mathfrak{G}) = \{\mathbf{x} \in \mathcal{A}^* : S \xrightarrow{\mathfrak{G}}^* \mathbf{x}\}$ .

We relate language classes of practical regular expressions to the Chomsky hierarchy ([7]). A *context-sensitive grammar* is a grammar  $(\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$  where every member of  $\mathcal{R}$  is of the form  $\mathbf{p}Ns \rightarrow \mathbf{pvs}$  where  $\mathbf{v} \neq \varepsilon$ , and  $\mathbf{p}, \mathbf{s} \in (\mathcal{A} \cup \mathcal{V})^*$ . A *context-free grammar* is a grammar  $(\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$  where all members of  $\mathcal{R}$  are of the form  $N \rightarrow \mathbf{v}$ .

A derivation  $S \xrightarrow{\mathfrak{G}}^* \mathbf{v}$  in a context-free grammar  $\mathfrak{G}$  is called *the leftmost* if at each derivation step we replace the leftmost nonterminal (in a sentential form) by the right-hand side of one of its productions. A context-free grammar  $(\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$  is in *Greibach normal form* if every production is of the form  $N \rightarrow a\mathbf{v}$ . [9],[14, lecture 21]

A *context-sensitive language* is a language that is generated by some context-sensitive grammar. Context-free languages are defined likewise. The class of context-free languages is denoted by  $\mathbb{L}_{CF}$ .

## 2.1 Regular expressions

The set of all classical regular expressions over an alphabet  $\mathcal{A}$  is denoted by  $\mathbb{E}_{C,\mathcal{A}}$ . The syntax of classical regular expressions is given as follows (as defined in the literature[11,13] and adapted to conform conventions used in tools and libraries; operators are ordered by their precedence from the highest):

1.  $\emptyset$  and  $\varepsilon$  are regular expressions,
2. for  $a \in \mathcal{A}$ ,  $a$  is a regular expression,
3. for  $\mathbf{r} \in \mathbb{E}_{C,\mathcal{A}}$ ,  $(\mathbf{r})$  (parenthesised expression) is a regular expression,
4. for  $\mathbf{r} \in \mathbb{E}_{C,\mathcal{A}}$ ,  $\mathbf{r}^*$  (iteration, Kleene star) is a regular expression.
5. for  $\mathbf{r}_1, \mathbf{r}_2 \in \mathbb{E}_{C,\mathcal{A}}$ ,  $\mathbf{r}_1 \cdot \mathbf{r}_2$  or  $\mathbf{r}_1\mathbf{r}_2$  (concatenation) is a regular expression,
6. for  $\mathbf{r}_1, \mathbf{r}_2 \in \mathbb{E}_{C,\mathcal{A}}$ ,  $\mathbf{r}_1 \mid \mathbf{r}_2$  (alternative) is a regular expression.

The semantics of classical regular expressions (where the language matched by RE  $\mathbf{r}$  is denoted by  $L(\mathbf{r})$ ) is given as follows[11]:

- $L(\emptyset) = \emptyset$ ,
- $L(\varepsilon) = \{\varepsilon\}$ ,
- for  $a \in \mathcal{A}$ ,  $L(a) = \{a\}$ ,
- for  $\mathbf{r} \in \mathbb{E}_{C,\mathcal{A}}$ ,  $L((\mathbf{r})) = L(\mathbf{r})$ ,
- for  $\mathbf{r} \in \mathbb{E}_{C,\mathcal{A}}$ ,  $L(\mathbf{r}^*) = (L(\mathbf{r}))^*$ ,
- for  $\mathbf{r}_1, \mathbf{r}_2 \in \mathbb{E}_{C,\mathcal{A}}$ ,  $L(\mathbf{r}_1 \mathbf{r}_2) = L(\mathbf{r}_1) \cdot L(\mathbf{r}_2)$ ,
- for  $\mathbf{r}_1, \mathbf{r}_2 \in \mathbb{E}_{C,\mathcal{A}}$ ,  $L(\mathbf{r}_1 | \mathbf{r}_2) = L(\mathbf{r}_1) \cup L(\mathbf{r}_2)$ .

The set of all practical regular expressions with operations iteration, concatenation, alternative, DEFINE rule, lookaround assertion, subroutine calls, and numbered and named parenthesised subexpressions over alphabet  $\mathcal{A}$  is denoted by  $\mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$  where the set of labels of named parenthesised expressions is denoted by  $\mathcal{X}$  ( $\mathcal{X} \cap \mathcal{A} = \emptyset$ ). The set of all regexes without lookaround assertions is denoted by  $\mathbb{E}_{S,\mathcal{A},\mathcal{X}}$ . Each practical regular expression consists of characters that may occur in the input string (i.e.,  $a \in \mathcal{A}$ ) and metacharacters that cannot occur in the input<sup>3</sup>:  $(,), ?, =, <, > \notin \mathcal{A}$ . For brevity, we write parentheses that denote a parenthesised expression with their assigned number, e.g.,  $(_l)_l$ . Our syntax closely follows the flavour PCRE2:

- the empty string, a character, iteration, concatenation, and alternative are defined the same way as for classical regular expressions,
- $\mathbf{r} \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}} \wedge l \in \mathbb{N} : (_l \mathbf{r})_l \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$ ,
- $\mathbf{r} \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}} \wedge l \in \mathbb{N} \wedge N \in \mathcal{X} : (_l ? < N > \mathbf{r})_l \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$  (parenthesised expression named  $N$  and numbered  $l$ ),
- $\mathbf{r} \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}} : (?(\text{DEFINE})\mathbf{r}) \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$  (DEFINE rule),
- $\mathbf{r} \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}} : (?= \mathbf{r}) \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$  (lookahead),
- $\mathbf{r} \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}} : (?<= \mathbf{r}) \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$  (lookbehind),
- $l \in \mathbb{N} : (?l) \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$  (numbered subroutine call),
- $N \in \mathcal{X} : (?P>N) \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$  (named subroutine call).

The numbering of parenthesised expressions, both named and unnamed, is unique. (Note that neither parentheses around lookahead, lookbehind, nor subroutine call delimit a parenthesised expression.)

The semantics of regexes with numbered backreferences and lookaround assertions was defined using the matching relation by Chida and Terauchi[5,6]. We closely follow their definition<sup>4</sup>. A matching relation  $\rightsquigarrow$  is of the form  $(\mathbf{r}, \mathbf{x}, i) \rightsquigarrow \mathcal{R}$  where  $\mathbf{r} \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$ ,  $\mathbf{x} \in \mathcal{A}^*$ ,  $i \in \mathbb{N} \wedge i \leq |\mathbf{x}|$ , and  $\mathcal{R} = \{i : i \in \mathbb{N} \wedge i \leq |\mathbf{x}| + 1\}$  (matching result). The rules for deriving the matching relation for practical regular expressions with lookaround assertions are as follows:

$$\overline{(\emptyset, \mathbf{x}, i) \rightsquigarrow \emptyset}$$

$$\overline{(\varepsilon, \mathbf{x}, i) \rightsquigarrow \{i\}}$$

<sup>3</sup> For brevity, we deviate from the way real-world engines treat metacharacters: they can occur in the input and can be matched in a regex when following a special escaping metacharacter (some flavours can match some metacharacters even without escaping), mostly the backslash. We refer the reader to Friedl[8].

<sup>4</sup> We omit capturing environment as this paper does not deal with backreferences.

$$\begin{aligned}
& \frac{a \in \mathcal{A} \wedge i \leq |\mathbf{x}| \wedge \mathbf{x}[i] = a}{(\mathbf{a}, \mathbf{x}, i) \rightsquigarrow \{i+1\}}, \frac{a \in \mathcal{A} \wedge (i > |\mathbf{x}| \vee \mathbf{x}[i] \neq a)}{(\mathbf{a}, \mathbf{x}, i) \rightsquigarrow \emptyset} \\
& \frac{(\mathbf{r}, \mathbf{x}, i) \rightsquigarrow \mathcal{R} \wedge \forall i_h \in \mathcal{R} \setminus \{i\} : (\mathbf{r}^*, \mathbf{x}, i_h) \rightsquigarrow \mathcal{R}_h}{(\mathbf{r}^*, \mathbf{x}, i) \rightsquigarrow \{i\} \cup \bigcup_{1 \leq h \leq |\mathcal{R} \setminus \{i\}|} \mathcal{R}_h} \tag{1} \\
& \frac{(\mathbf{r}_1, \mathbf{x}, i) \rightsquigarrow \mathcal{R} \wedge \forall (i_h) \in \mathcal{R} : (\mathbf{r}_2, \mathbf{x}, i_h) \rightsquigarrow \mathcal{R}_h}{(\mathbf{r}_1 \mathbf{r}_2, \mathbf{x}, i) \rightsquigarrow \bigcup_{1 \leq h \leq |\mathcal{R}|} \mathcal{R}_h} \tag{2} \\
& \frac{(\mathbf{r}_1, \mathbf{x}, i) \rightsquigarrow \mathcal{R}_1 \wedge (\mathbf{r}_2, \mathbf{x}, i) \rightsquigarrow \mathcal{R}_2}{(\mathbf{r}_1 | \mathbf{r}_2, \mathbf{x}, i) \rightsquigarrow \mathcal{R}_1 \cup \mathcal{R}_2} \tag{3} \\
& \frac{(\mathbf{r}, \mathbf{x}, i) \rightsquigarrow \mathcal{R}}{((\mathbf{r}), \mathbf{x}, i) \rightsquigarrow \mathcal{R}} \tag{4} \\
& \frac{(\mathbf{r}, \mathbf{x}, i) \rightsquigarrow \mathcal{R}}{((?=r), \mathbf{x}, i) \rightsquigarrow \{i \wedge \mathcal{R} \neq \emptyset\}} \tag{5} \\
& \frac{\mathbf{y} \in \mathcal{A}^* \wedge (\mathbf{y}, \mathbf{x}[i - |\mathbf{y}|..i-1], 1) \rightsquigarrow \mathcal{R}}{((?<=\mathbf{y}), \mathbf{x}, i) \rightsquigarrow \{i \wedge \mathcal{R} \neq \emptyset\}} \tag{6}
\end{aligned}$$

The language of a regex  $\mathbf{r} \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$  is  $L(\mathbf{r}) = \{\mathbf{x} : (\mathbf{r}, \mathbf{x}, 1) \rightsquigarrow \mathcal{R} \wedge |\mathbf{x}| + 1 \in \mathcal{R}\}$ . We also say that a string  $\mathbf{x} \in L(\mathbf{r})$  *matches* a regex  $\mathbf{r}$  (similarly,  $L(\mathbf{r})$  is the language matched by  $\mathbf{r}$ ). The class of all languages that can be matched by the regexes of  $\mathbb{E}_{LS,\mathcal{A},\mathcal{X}}$  is denoted by  $\mathbb{L}_{E_{LS}}$ . The class of all languages that can be matched by the regexes of  $\mathbb{E}_{S,\mathcal{A},\mathcal{X}}$  is denoted by  $\mathbb{L}_{E_S}$ .

### 3 Formalizing expressions with subroutine calls and lookahead assertions

We now formally define the semantics of practical regular expressions with named parenthesised expressions, DEFINE rules, and numbered and named subroutine calls. Our definition is an extension of the matching relation in the previous section. In this section, the following notation is used:  $i, l \in \mathbb{N}; N \in \mathcal{X}; \mathbf{r}, \mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3 \in \mathbb{E}_{LS,\mathcal{A},\mathcal{X}}; \mathbf{x} \in \mathcal{A}^*$ .

The subroutine call attempts to match a given parenthesised expression at a current position. To be able to use the subexpression given the parenthesis number, the partial function  $\sigma_{\mathbf{r}}$  is computed before the matching of regex  $\mathbf{r}$  starts; it is defined as follows:

$$\sigma_{\mathbf{r}} : \mathbb{N} \rightarrow \mathbb{E}_{LS,\mathcal{A},\mathcal{X}} \wedge \mathbf{r} = \mathbf{r}_1({}_l \mathbf{r}_2)_l \mathbf{r}_3 \text{ implies } \sigma_{\mathbf{r}}(l) = \mathbf{r}_2$$

If no confusion can arise, we use  $\sigma$  for simplicity. In addition to being unambiguous, the numbering of parenthesised expressions (both named and unnamed) is not important for our results. Our definition conforms to some flavours of practical regular expressions.<sup>5</sup>

Matching a numbered subroutine call means matching the subpattern given in  $l$ -th parentheses from the current position. The subroutine call can be located anywhere related to the referred subpattern (i.e., both forward and recursive calls are valid).

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

<sup>5</sup> Namely PCRE2 and Perl. Both PCRE2 and Perl even support enabling duplicate parenthesis numbers (it is not the default). Our definition needs to be modified for the Ruby Regexp class: numbers cannot be used when at least one named expression is present.

Matching a named subroutine call uses the parenthesised expression assigned to the given name. Named parenthesised expressions can be identified by either their name or number. Thus, a name is just an alias for the parenthesis number. The partial function  $\nu_r$  is computed before the matching of  $r$  starts. If the name is used for multiple parenthesised expressions,  $\nu_r$  assigns the name to the leftmost parentheses.

$$\nu_r : \mathcal{X} \rightarrow \mathbb{N} \wedge r = r_1(l?< N >r_2)_l r_3 \text{ implies } \nu_r(N) = l \wedge \nexists l' < l : r = r'_1(\nu?< N >r'_2)_{l'} r'_3$$

If no confusion can arise, we use  $\nu$  for simplicity.

$$\frac{(\sigma(\nu(N)), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}{((?P>N), \mathbf{x}, i) \rightsquigarrow \mathcal{R}} \quad (7)$$

This matching relation for named subroutine call closely mimics the semantics of PCRE2, Perl, and Ruby.

Any regex can be wrapped in a DEFINE rule. In addition to the possibility of extending  $\nu$  or  $\sigma$ , the DEFINE rule does not affect the matching.

$$\overline{((?(\text{DEFINE})r), \mathbf{x}, i) \rightsquigarrow \{i\}} \quad (8)$$

DEFINE rules with the above-defined semantics are supported by PCRE2 and Perl.

## 4 Expressive power of subroutine call

We give a rigorous proof of the equivalence of expressive power of context-free languages and practical regular expressions with subroutine calls. Our proof is based on the matching relation and extends Hruša's work[12], which is based on Popov's claim[19].

**Theorem 1.**  $\mathbb{L}_{E_S} = \mathbb{L}_{CF}$

To prove the class equivalence, we first show that every context-free grammar can be converted to a regex with subroutine calls. Later, we show that every such regex can be converted into context-free grammar.

**Lemma 2.**  $\mathbb{L}_{CF} \subseteq \mathbb{L}_{E_S}$ .

We show that every context-free grammar can be expressed by a practical regular expression with the following sufficient operations: concatenation, alternative, DEFINE rule with named parenthesised expression, and named subroutine call. Intuitively, such a regex contains all the building blocks of context-free grammar: concatenation, alternative, and the ability to reuse a subexpression, even recursively. The conversion is formally defined by algorithm 1 and definition 3. The restriction of conversion to Greibach normal form grammar does not change the expressive power: Any context-free grammar (and therefore any context-free language) can be expressed by a Greibach normal form grammar using a known transformation.[9][14, lecture 21]

**Definition 3.** Let us have a context-free grammar  $\mathfrak{G} = (\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$  and a regex  $r \in \mathbb{E}_{S, \mathcal{A}, \mathcal{X}}$  where  $\mathcal{V} = \mathcal{X}$ . The function  $rx : (\mathcal{V} \cup \mathcal{A})^* \rightarrow \mathbb{E}_{S, \mathcal{A}, \mathcal{X}}$  is defined as follows: let  $\mathbf{v}_1, \mathbf{v}_2 \in (\mathcal{A} \cup \mathcal{V})^*$ ,  $a \in \mathcal{A}$ , and  $N \in \mathcal{V}$  then  $rx(\varepsilon) = \varepsilon$ ,  $rx(a) = a$ ,  $rx(N) = (?P>N)$ , and  $rx(\mathbf{v}_1 \mathbf{v}_2) = rx(\mathbf{v}_1) rx(\mathbf{v}_2)$ .

**Algorithm 1** Conversion of a context-free grammar to a regex

---

**Input:** a context-free grammar  $\mathfrak{G} = (\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$  in Greibach normal form**Output:** a regex  $r \in \mathbb{E}_{\mathcal{S}, \mathcal{A}, \mathcal{X}}$  such that  $L(\mathfrak{G}) = L(r)$ 

1. initialize  $r = \varepsilon$  and consider  $\mathcal{X} = \mathcal{V}$
  2. for all productions with a non-terminal  $N$  on the left-hand side ( $N \rightarrow \mathbf{v}_{N1} \mid \dots \mid \mathbf{v}_{Nm_N} \in \mathcal{R}$ ):  
 (a) let  $r_{Ng}$  ( $1 \leq g \leq m_N$ ) be constructed from the  $g$ -th right-hand side of the production for  $N$  ( $\mathbf{v}_{Ng}$ ) by replacing the non-terminals with subroutine call:  $r_{Ng} = \text{rx}(\mathbf{v}_{Ng})$   
 (b) add a DEFINE rule with parenthesised expression named  $N$  containing the strings  $r_{Ng}$  constructed from right-hand side of these productions, i.e., let  $r = r(?(\text{DEFINE})(?<N>r_{N1} \mid r_{N2} \mid \dots \mid r_{Nm_N}))$
  3. add the matching of the initial symbol, i.e., let  $r = r(?P>S)$
- 

**Lemma 4.** Let us have a Greibach normal form grammar  $\mathfrak{G}$  and a practical regular expression  $r$  such that  $r$  is constructed by the algorithm 1 from  $\mathfrak{G}$ . For any possible derivation step in  $\mathfrak{G}$ , a possible step exists in the matching relation for  $r$ .

*Proof.* Let the grammar be  $\mathfrak{G} = (\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$ . Without loss of generality, we can assume that only the leftmost derivations are used to derive any sentential form. Let  $\mathbf{p}Ns \xrightarrow[\mathfrak{G}]{} \mathbf{pvs}$  be a derivation step (recall that this derivation step is possible only if  $N \rightarrow \mathbf{v} \in \mathcal{R}$ ) where  $\mathbf{p} \in \mathcal{A}^*$ ,  $\mathbf{s} \in \mathcal{V}^*$  and  $\mathbf{v} \in (\mathcal{V} \cup \mathcal{A})^+$ . Recall that due to algorithm 1,  $r$  is in the following form:

$$(?(\text{DEFINE}) \dots) \cdots (?(\text{DEFINE})(?<N>\dots \mid \text{rx}(\mathbf{v}) \mid \dots)) \cdots (?P>S) \quad (9)$$

The following steps of the matching relation show that an expression of the form  $\mathbf{p}(?P>N) \text{rx}(\mathbf{s})$  (from step 2a follows that  $\text{rx}(\mathbf{p}) = \mathbf{p}$ ) can always be resolved by the concatenation of  $\mathbf{p}$ ,  $\text{rx}(\mathbf{v})$ , and  $\text{rx}(\mathbf{s})$ .

$$\begin{array}{c} \dots \\ \hline \frac{\dots}{(\text{rx}(\mathbf{v}_{N1}), \mathbf{x}, i_h) \rightsquigarrow \mathcal{R}_{hN1}} \cdots \frac{(\text{rx}(\mathbf{v}), \mathbf{x}, i_h) \rightsquigarrow \mathcal{R}_{hNg}}{(\text{rx}(\mathbf{v}_{N1}) \mid \dots \mid \text{rx}(\mathbf{v}) \mid \dots \mid \text{rx}(\mathbf{v}_{Nm_N}), \mathbf{x}, i_h) \rightsquigarrow \mathcal{R}_{hN} = \bigcup_{1 \leq h' \leq m_N} \mathcal{R}_{hNh'}} \cdots \\ \frac{(\sigma(l), \mathbf{x}, i_h) \rightsquigarrow \mathcal{R}_{hN}}{\frac{\dots}{\frac{(\sigma(\nu(N)), \mathbf{x}, i_h) \rightsquigarrow \mathcal{R}_{hN}}{(\mathbf{p}, \mathbf{x}, i) \rightsquigarrow \mathcal{R}'}} \cdots} \frac{(\mathbf{p}(?P>N), \mathbf{x}, i_h) \rightsquigarrow \mathcal{R}_{hN} \quad \forall i_{hN} \in \mathcal{R}_{hN} : (\text{rx}(\mathbf{s}), \mathbf{x}, i_{hN}) \rightsquigarrow \mathcal{R}_{hsN}}{\frac{\dots}{\frac{(\mathbf{p}(?P>N) \text{rx}(\mathbf{s}), \mathbf{x}, i) \rightsquigarrow \mathcal{R}_h = \bigcup_{1 \leq h \leq |\mathcal{R}'|} \mathcal{R}_h}{(\mathbf{p}(?P>N) \text{rx}(\mathbf{s}), \mathbf{x}, i) \rightsquigarrow \bigcup_{1 \leq h \leq |\mathcal{R}'|} \mathcal{R}_h}}} \end{array}$$

Note that  $\mathcal{R}' = \{i + |\mathbf{p}|\}$  if  $\mathbf{x}[i..i + |\mathbf{p}| - 1] = \mathbf{p}$  and  $\mathcal{R}' = \emptyset$  otherwise. Furthermore, the matching result for  $\text{rx}(\mathbf{v})$ ,  $\mathcal{R}_{hNg}$ , is involved in the expression matching.  $\square$

The following lemma holds due to lemma 4.

**Lemma 5.** Let us have a Greibach normal form grammar  $\mathfrak{G}$  and a practical regular expression  $r$  such that  $r$  is constructed by the algorithm 1 from  $\mathfrak{G}$ . Then  $L(\mathfrak{G}) \subseteq L(r)$ .

**Lemma 6.** Let us have a Greibach normal form grammar  $\mathfrak{G}$  and a practical regular expression  $r$  such that  $r$  is constructed by the algorithm 1 from  $\mathfrak{G}$ . Then  $L(r) \subseteq L(\mathfrak{G})$ .

*Proof.* Let the grammar be  $\mathfrak{G} = (\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$  and  $\mathbf{x} \in L(\mathfrak{G})$ . The general form of a regex constructed by the algorithm is (9) while all subpatterns  $\text{rx}(\mathbf{v})$  of  $r$  are of

the form  $a_1 \text{rx}(\mathbf{s})$  where  $a_1 \in \mathcal{A}$  and  $\mathbf{s} \in \mathcal{V}^*$  (due to the normal Greibach form). Such a subpattern can exist only if  $N \rightarrow a_1 \mathbf{s} \in \mathcal{R}$  (due to step 2a). Following the matching relation, the initial step in matching, that is,  $(?P>S)$  which applies some subpattern of named parentheses  $S$ ,  $a_2 \text{rx}(\mathbf{s}')$ , is possible only when  $\mathbf{x}[1] = a_2$ ; such a subpattern is constructed only if  $S \rightarrow a_2 \mathbf{s}' \in \mathcal{R}$ . Suppose, for contradiction, that  $a_2 \text{rx}(\mathbf{s}')$  matches  $\mathbf{x}$  and  $\mathbf{x} \notin L(\mathfrak{G})$ . As  $\mathbf{s}' = N_1 N_2 \dots N_{|\mathbf{s}'|}$ , the only possibility is that any of  $(?P>N_g)$  matches a substring  $\mathbf{y}$  of  $\mathbf{x}$  and  $N_g \xrightarrow[\mathfrak{G}]{}^* \mathbf{y}$  is not possible. However, due to the algorithm 1, the named parenthesised expressions in  $\mathbf{r}$  contain only subpatterns corresponding to the right-hand sides of the productions in  $\mathfrak{G}$ . Due to the Greibach normal form,  $\mathbf{r}$  can only match  $\mathbf{x}$  using subpatterns in the same order as productions of  $\mathfrak{G}$  are applied when generating  $\mathbf{x}$ . In any situation when a subroutine call  $(?P>N_g)$  occurs,  $\mathfrak{G}$  can use any production for  $N_g$ , because the application of productions in context-free grammar is not restricted by their order or context.  $\square$

The validity of lemma 2 was shown by lemmas 5 and 6: any context-free language can be expressed by a regex with subroutine calls.

**Lemma 7.**  $\mathbb{L}_{E_S} \subseteq \mathbb{L}_{CF}$ .

We show that an equivalent context-free grammar can express any practical regular expression with operations concatenation, alternative, DEFINE rule with named parenthesised expression, and named subroutine call. We begin by showing that removing the other valid operations in  $\mathbb{E}_{S,\mathcal{A},\mathcal{X}}$  does not change the expressive power.

**Lemma 8 (Redundancy of Kleene star).** *Every regex of the form  $\mathbf{r}^* \in \mathbb{E}_{S,\mathcal{A},\mathcal{X}}$  can be expressed as*

$$((?(\text{DEFINE})(?<N>\varepsilon \mid \mathbf{r} (?P>N))) (?P>N)) \quad (10)$$

where  $N \notin \mathcal{X}$ .

*Proof.* The matching relation can be derived as follows:

$$\begin{array}{c} (\mathbf{r}, \mathbf{x}, i) \rightsquigarrow \mathcal{R}' \wedge \forall i_h \in \mathcal{R}' : ((?P>N), \mathbf{x}, i_h) \rightsquigarrow \mathcal{R}_h \\ \hline (\varepsilon, \mathbf{x}, i) \rightsquigarrow \{i\} \quad \frac{(\mathbf{r} (?P>N), \mathbf{x}, i) \rightsquigarrow \mathcal{R} = \bigcup_{1 \leq h \leq |\mathcal{R}'|} \mathcal{R}_h}{(\varepsilon \mid \mathbf{r} (?P>N), \mathbf{x}, i) \rightsquigarrow \mathcal{R} \cup \{i\}} \\ \hline \frac{\frac{((?(\text{DEFINE})(?<N>\varepsilon \mid \mathbf{r} (?P>N))), \mathbf{x}, i) \rightsquigarrow \{i\}}{((?P>N), \mathbf{x}, i) \rightsquigarrow \mathcal{R} \cup \{i\}}} {((10), \mathbf{x}, i) \rightsquigarrow \mathcal{R} \cup \{i\}} \end{array}$$

It is clear that the matching result is the same as that of the Kleene star (1). Although the matching relation of the Kleene star (1) excludes  $i$  from rematching  $\mathbf{r}^*$ , this exclusion does not affect the positions in its matching result.  $\square$

The following two redundancies of a standalone parenthesised expression are straightforward and thus are left without proof. Each occurrence of parenthesised expression put inside a DEFINE rule retains its parenthesis number and name, and thus it does not affect any subroutine call.

**Lemma 9 (Redundancy of named parenthesised expression).** *A regex in the form  $(?<N>\mathbf{r}) \in \mathbb{E}_{S,\mathcal{A},\mathcal{X}}$  (named parenthesised expression outside the DEFINE rule) can be expressed as*

$$((?(\text{DEFINE})(?<N>\mathbf{r})) (?P>N))$$

**Lemma 10 (Redundancy of numbered parenthesised expression).** *A regex in the form  $(_l \mathbf{r})_l \in \mathbb{E}_{S,A,\mathcal{X}}$  (numbered parenthesised expression outside the DEFINE rule) can be expressed as*

$$(\text{?(DEFINE})(_l ? < N > \mathbf{r})_l)(\text{?P} > N)$$

where  $N \notin \mathcal{X}$ .

The redundancy of the numbered subroutine call (and its replacement with the named subroutine call) is straightforward and thus left without proof.

**Lemma 11 (Redundancy of numbered subroutine call).** *Let  $\mathbf{r} \in \mathbb{E}_{S,A,\mathcal{X}}$  be a regex with operations concatenation, alternative, DEFINE rule with named parenthesised expression, named subroutine call, and numbered subroutine call. To construct a regex  $\mathbf{r}' \in \mathbb{E}_{S,A,\mathcal{X}}$  such that  $L(\mathbf{r}) = L(\mathbf{r}')$  and  $\mathbf{r}'$  does not contain numbered subroutine call,  $\mathbf{r}'$  is the same as  $\mathbf{r}$  with the following modifications: any occurrence of  $(?l)$  from  $\mathbf{r}$  is replaced with  $(\text{?P} > N)$  in  $\mathbf{r}'$  where  $l$  refers to  $(\text{?(DEFINE})(_l ? < N > \mathbf{r}_l)_l)$ .*

The conversion of a regex to a context-free grammar is formally defined in algorithm 2, which is inspired by Thompson's[23] pattern matching algorithm as presented by Hopcroft et al.[11, theorem 3.7] and closure properties of context-free languages studied by Scheinberg[21] as presented by Kozen[14].

---

### Algorithm 2 Conversion of a regex to a context-free grammar

**Input:** a practical regular expression  $\mathbf{r} \in \mathbb{E}_{S,A,\mathcal{X}}$  with operations concatenation, alternative, DEFINE rule with named parenthesised expression, and named subroutine call

**Output:** a context-free grammar  $\mathfrak{G}_r = (\mathcal{V}_r, \mathcal{A}, \mathcal{R}_r, S_r)$  such that  $L(\mathbf{r}) = L(\mathfrak{G}_r)$

1. construct grammars for elementary expressions
    - $\mathfrak{G}_\emptyset = (\{S_\emptyset\}, \mathcal{A}, \emptyset, S_\emptyset)$
    - $\mathfrak{G}_\varepsilon = (\{S_\varepsilon\}, \mathcal{A}, \{S_\varepsilon \rightarrow \varepsilon\}, S_\varepsilon)$
    - $\mathfrak{G}_a = (\{S_a\}, \mathcal{A}, \{S_a \rightarrow a\}, S_a) : a \in \mathcal{A}$
    - $\mathfrak{G}_{(\text{?P} > N)} = (\{S_{(\text{?P} > N)}\}, \mathcal{A}, \{S_{(\text{?P} > N)} \rightarrow N\}, S_{(\text{?P} > N)}) : N \in \mathcal{X}$
  2. iteratively construct grammars for operations in a regex (suppose  $\mathfrak{G}_{r_1} = (\mathcal{V}_{r_1}, \mathcal{A}, \mathcal{R}_{r_1}, S_{r_1})$  and  $\mathfrak{G}_{r_2} = (\mathcal{V}_{r_2}, \mathcal{A}, \mathcal{R}_{r_2}, S_{r_2})$  are already constructed for  $r_1$  and  $r_2$ , respectively)
    - $\mathfrak{G}_{r_1 r_2} = (\mathcal{V}_{r_1} \cup \mathcal{V}_{r_2} \cup \{S_{r_1 r_2}\}, \mathcal{A}, \mathcal{R}_{r_1} \cup \mathcal{R}_{r_2} \cup \{S_{r_1 r_2} \rightarrow S_{r_1} S_{r_2}\}, S_{r_1 r_2}) : S_{r_1 r_2} \notin \mathcal{V}_{r_1} \cup \mathcal{V}_{r_2}$
    - $\mathfrak{G}_{r_1 | r_2} = (\mathcal{V}_{r_1} \cup \mathcal{V}_{r_2} \cup \{S_{r_1 | r_2}\}, \mathcal{A}, \mathcal{R}_{r_1} \cup \mathcal{R}_{r_2} \cup \{S_{r_1 | r_2} \rightarrow S_{r_1} \mid S_{r_2}\}, S_{r_1 | r_2}) : S_{r_1 | r_2} \notin \mathcal{V}_{r_1} \cup \mathcal{V}_{r_2}$
    - $\mathfrak{G}_{(\text{?(DEFINE})(_l ? < N > r_1)_l)} = (\mathcal{V}_{r_1} \cup \{S_{(\text{?(DEFINE})(_l ? < N > r_1)_l)}\}, \mathcal{A}, \mathcal{R}_{r_1} \cup \{N_l \rightarrow S_{r_1}, S_{(\text{?(DEFINE})(_l ? < N > r_1)_l)} \rightarrow \varepsilon\}, S_{(\text{?(DEFINE})(_l ? < N > r_1)_l)}) : S_{(\text{?(DEFINE})(_l ? < N > r_1)_l)}, N_l \notin \mathcal{V}_{r_1}$
  3. having grammar  $\mathfrak{G}_r = (\mathcal{V}_r, \mathcal{A}, \mathcal{R}_r, S_r)$ , for every nonterminal  $N \in \mathcal{X} \cap \mathcal{V}_r$ : if  $N_{\nu(N)} \in \mathcal{V}_r$  then replace all occurrences of  $N$  in the right-hand sides of productions  $\mathcal{R}_r$  with  $N_{\nu(N)}$
  4. return  $\mathfrak{G}_r$
- 

**Lemma 12.** *Let  $r_1, r_2 \in \mathbb{E}_{S,A,\mathcal{X}}$ . Let  $\mathfrak{G}_{r_1}, \mathfrak{G}_{r_2}$  be the grammars constructed by the algorithm 2 from  $r_1, r_2$ , respectively. If  $L(r_1) = L(\mathfrak{G}_{r_1})$  and  $L(r_2) = L(\mathfrak{G}_{r_2})$ , then  $L(r_1 \mid r_2) = L(\mathfrak{G}_{r_1 | r_2})$  and  $L(r_1 r_2) = L(\mathfrak{G}_{r_1 r_2})$ .*

*Proof.* For the alternative of regexes  $r_1, r_2$ , following the matching relation (3),  $r_1 \mid r_2$  matches some  $x$  if at least one of  $r_1, r_2$  matches  $x$ . Clearly,  $x \in L(\mathfrak{G}_{r_1 | r_2})$  if  $x \in L(\mathfrak{G}_{r_1}) \vee x \in L(\mathfrak{G}_{r_2})$ . Following the well-known construction of a context-free grammar for the union of languages[21],  $L(r_1 \mid r_2) \subseteq L(\mathfrak{G}_{r_1 | r_2})$  because  $\mathcal{V}_{r_1} \cap \mathcal{V}_{r_2}$  can be nonempty. Suppose, for contradiction, that  $x \notin L(r_1 \mid r_2) \wedge x \in L(\mathfrak{G}_{r_1 | r_2})$  while

$L(\mathbf{r}_1) = L(\mathfrak{G}_{\mathbf{r}_1})$  and  $L(\mathbf{r}_2) = L(\mathfrak{G}_{\mathbf{r}_2})$ . The only way it can happen is  $S_{\mathbf{r}_1|\mathbf{r}_2} \xrightarrow[\mathfrak{G}_{\mathbf{r}_1|\mathbf{r}_2}]{}^+$   $\mathbf{p}_1 N_1 \mathbf{s}_1 \xrightarrow[\mathfrak{G}_{\mathbf{r}_1|\mathbf{r}_2}]{}^* \mathbf{p}_2 N_2 \mathbf{s}_2 \xrightarrow[\mathfrak{G}_{\mathbf{r}_1|\mathbf{r}_2}]{}^* \mathbf{p}_2 \mathbf{v}_2 \mathbf{s}_2 \xrightarrow[\mathfrak{G}_{\mathbf{r}_1|\mathbf{r}_2}]{}^* \mathbf{x}$  where  $N_1 \rightarrow \mathbf{v}_1$  and  $N_2 \rightarrow \mathbf{v}_2$  are not from the same grammar  $\mathfrak{G}_{\mathbf{r}_1}, \mathfrak{G}_{\mathbf{r}_2}$ ; in other words,  $N_2 \in \mathcal{V}_{\mathbf{r}_1} \cap \mathcal{V}_{\mathbf{r}_2} \wedge N_2 \rightarrow \mathbf{v}_2 \notin \mathcal{R}_{\mathbf{r}_1} \cap \mathcal{R}_{\mathbf{r}_2}$ . All grammars from step 2 introduce unique nonterminals that cannot appear in both  $\mathcal{V}_{\mathbf{r}_1}$  and  $\mathcal{V}_{\mathbf{r}_2}$ . The nonterminals of grammars  $\mathfrak{G}_\emptyset, \mathfrak{G}_\epsilon$ , and  $\mathfrak{G}_a$  clearly cannot have different right-hand sides of productions in different grammars  $\mathfrak{G}_{\mathbf{r}_1}, \mathfrak{G}_{\mathbf{r}_2}$ . Both nonterminals introduced by  $\mathfrak{G}_{(\text{?P}>N)}$  can appear in both  $\mathcal{V}_{\mathbf{r}_1}$  and  $\mathcal{V}_{\mathbf{r}_2}$ , however, due to step 3, every  $N$  from  $\mathfrak{G}_{(\text{?P}>N)}$  is replaced by a single  $N_l$  that rewrites to a unique nonterminal determined by a single DEFINE rule. Thus, it is not possible to achieve  $N_2 \in \mathcal{V}_{\mathbf{r}_1} \cap \mathcal{V}_{\mathbf{r}_2} \wedge N_2 \rightarrow \mathbf{v}_2 \notin \mathcal{R}_{\mathbf{r}_1} \cap \mathcal{R}_{\mathbf{r}_2}$ .

Similar arguments can be used to prove that  $L(\mathbf{r}_1 \mathbf{r}_2) = L(\mathfrak{G}_{\mathbf{r}_1 \mathbf{r}_2})$ .  $\square$

**Lemma 13.** *If  $\mathfrak{G}_r = (\mathcal{V}_r, \mathcal{A}, \mathcal{R}_r, S_r)$  is constructed from  $r$  by algorithm 2 then  $L(r) = L(\mathfrak{G}_r)$  for any  $r \in \mathbb{E}_{S, A, X}$  (with operations concatenation, alternative, DEFINE rule with named parenthesised expression, and named subroutine call) and  $\mathfrak{G}_r$  is context-free.*

*Proof.* The grammars are clearly correct for the cases of elementary expressions  $\emptyset, \epsilon, a \in \mathcal{A}$ . Assume that for  $\mathbf{r}_1$  and  $\mathbf{r}_2$ ,  $L(\mathbf{r}_1) = L(\mathfrak{G}_{\mathbf{r}_1})$  and  $L(\mathbf{r}_2) = L(\mathfrak{G}_{\mathbf{r}_2})$ , respectively. For a subroutine call  $(\text{?P}>N)$ , the matching relation (7) is

$$\frac{\frac{(r_1, \mathbf{x}, i) \rightsquigarrow \mathcal{R}}{(\sigma(l), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}}{\frac{(\sigma(\nu(N)), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}{((\text{?P}>N}), \mathbf{x}, i) \rightsquigarrow \mathcal{R}}}$$

where  $\mathcal{R}$  contains all  $i'$  such that  $\mathbf{x}[i..i'-1]$  matches  $\mathbf{r}_1$  and  $l$  identifies the leftmost parenthesised expression named  $N$ . Production  $S_{(\text{?P}>N)} \rightarrow N$  of  $\mathfrak{G}_{(\text{?P}>N)}$  is effectively  $S_{(\text{?P}>N)} \rightarrow N_{\nu(N)}$  due to step 3. Nonterminal  $N_{\nu(N)}$  rewrites to  $S_{\mathbf{r}_1}$ . Therefore, the grammar  $\mathfrak{G}_{(\text{?P}>N)}$  generates the same language as is matched by  $(\text{?P}>N)$ .

The grammar  $\mathfrak{G}_{(\text{?DEFINE})(l(\text{?<}N\text{>}\mathbf{r}_1)_l)}$  follows the matching relation for the DEFINE rule (8). The correctness of both  $\mathfrak{G}_{\mathbf{r}_1|\mathbf{r}_2}$  and  $\mathfrak{G}_{\mathbf{r}_1 \mathbf{r}_2}$  follows from lemma 12. Therefore, step 2 constructs the correct grammars.

All grammars add only productions with a single nonterminal on the left-hand side; therefore, all grammars constructed by the algorithm 2 are context-free.  $\square$

Lemma 2 is proved by lemmas 5 and 6; also, lemma 7 is proved by lemmas 8, 9, 10, 11, and 13. Therefore, theorem 1 is proved.

## 5 Expressive power of subroutine call combined with lookahead assertions

We show that lookahead assertion combined with subroutine call has greater expressive power than subroutine call alone. We use an example of such regex inspired by Popov's blog post[19] and arguments by Scheinberg[21].

**Theorem 14.**  $\mathbb{L}_{E_S} \subsetneq \mathbb{L}_{E_{LS}}$

*Proof.* We show a regex  $\mathbf{r} \in \mathbb{E}_{\text{LS}, \mathcal{A}, \mathcal{X}}$  that matches a language that is not context-free (the equality of  $\mathbb{L}_{\text{CF}}$  and  $\mathbb{L}_{\text{ES}}$  is shown by theorem 1). Language  $\mathcal{L} = \{\mathbf{a}^g \mathbf{b}^g \mathbf{c}^g : g \in \mathbb{N}\}$  is a well-known language that is not context-free[11, example 7.19]. We show that for  $\mathbf{r} = (?=(?<N_1>\mathbf{a}(\varepsilon \mid (?P>N_1))\mathbf{b})\mathbf{c})\mathbf{a}\mathbf{a}^*(?<N_2>\mathbf{b}(\varepsilon \mid (?P>N_2))\mathbf{c})$ ,  $\mathcal{L} = L(\mathbf{r})$ . Let  $\mathbf{r}_1 = (?<N_1>\mathbf{a}(\varepsilon \mid (?P>N_1))\mathbf{b})\mathbf{c}$ ,  $\mathbf{r}_2 = \mathbf{a}\mathbf{a}^*$ ,  $\mathbf{r}_{N_2} = \mathbf{b}(\varepsilon \mid (?P>N_2))\mathbf{c}$ , and thus  $\mathbf{r} = (?=\mathbf{r}_1)\mathbf{r}_2(?<N_2>\mathbf{r}_{N_2})$ . Let  $\mathbf{r}_{N_1} = \mathbf{a}(\varepsilon \mid (?P>N_1))\mathbf{b}$ . The matching relation for  $\mathbf{r}_1$  can be derived as follows:

$$\frac{\dots}{\overline{(\mathbf{r}_{N_1}, \mathbf{x}, 1) \rightsquigarrow \mathcal{R}}} \quad \mathbf{x}[i] = \mathbf{c}$$

$$\frac{\overline{((?<N_1>\mathbf{r}_{N_1}), \mathbf{x}, 1) \rightsquigarrow \mathcal{R}} \quad \forall i \in \mathcal{R} : (\mathbf{c}, \mathbf{x}, i) \rightsquigarrow \{i+1 : \mathbf{x}[i] = \mathbf{c}\}}{\overline{((?<N_1>\mathbf{r}_{N_1})\mathbf{c}, \mathbf{x}, 1) \rightsquigarrow \bigcup_{i \in \mathcal{R}} \{i+1 : \mathbf{x}[i] = \mathbf{c}\}}}$$

Therefore,  $L(\mathbf{r}_1) = L(\mathbf{r}_{N_1}) \cdot \{\mathbf{c}\}$ . Let us derive the matching relation for  $\mathbf{r}_{N_1}$ :

$$\frac{\dots}{\overline{(\mathbf{r}_{N_1}, \mathbf{x}, 2) \rightsquigarrow \mathcal{R}_{N_1}}}$$

$$\frac{\overline{(\varepsilon, \mathbf{x}, 2) \rightsquigarrow \{2\}} \quad \overline{((?P>N_1), \mathbf{x}, 2) \rightsquigarrow \mathcal{R}_{N_1}}}{\overline{(\varepsilon \mid (?P>N_1), \mathbf{x}, 2) \rightsquigarrow \mathcal{R}_{N_1} = \{2\} \cup \mathcal{R}_{N_1}}}$$

$$\frac{\mathbf{x}[1] = \mathbf{a}}{\overline{(\mathbf{a}, \mathbf{x}, 1) \rightsquigarrow \mathcal{R}_{\mathbf{a}}}} \quad \frac{\overline{((\varepsilon \mid (?P>N_1)), \mathbf{x}, 2) \rightsquigarrow \mathcal{R}_{N_1}} \quad \forall i \in \mathcal{R}_{N_1} : (\mathbf{b}, \mathbf{x}, i) \rightsquigarrow \mathcal{R}_{\mathbf{b}_i}}{\overline{\text{if } \mathbf{x}[1] = \mathbf{a} : ((\varepsilon \mid (?P>N_1))\mathbf{b}, \mathbf{x}, 2) \rightsquigarrow \mathcal{R} = \bigcup_{i \in \mathcal{R}_{N_1}} \mathcal{R}_{\mathbf{b}_i}}}$$

$$\overline{(\mathbf{a}(\varepsilon \mid (?P>N_1))\mathbf{b}, \mathbf{x}, 1) \rightsquigarrow \mathcal{R} \text{ if } \mathbf{x}[1] = \mathbf{a}}$$

The regex  $\mathbf{r}_{N_1}$  matches only if a prefix of  $\mathbf{x}$  has the form  $\mathbf{a} \dots \mathbf{ab} \dots \mathbf{b}$ . Furthermore, because  $\mathbf{as}$  and  $\mathbf{bs}$  are matched only within the same parenthesised expression and the same number of times,  $L(\mathbf{r}_{N_1}) = \{\mathbf{a}^g \mathbf{b}^g : g \in \mathbb{N}\}$ . Clearly,  $L(\mathbf{r}_2) = \{\mathbf{a}\}^+$ . Following similar arguments as for  $L(\mathbf{r}_{N_1})$ ,  $L(\mathbf{r}_{N_2}) = \{\mathbf{b}^g \mathbf{c}^g : g \in \mathbb{N}\}$ . The matching relation for  $\mathbf{r}$  can be derived as follows:

$$\frac{\dots}{\overline{(\mathbf{r}_1, \mathbf{x}, 1) \rightsquigarrow \mathcal{R}_1}} \quad \dots$$

$$\frac{\overline{((?=\mathbf{r}_1), \mathbf{x}, 1) \rightsquigarrow \{1 : \mathcal{R}_1 \neq \emptyset\}} \text{ if } \mathcal{R}_1 \neq \emptyset : \overline{(\mathbf{r}_2(?<N_2>\mathbf{r}_{N_2}), \mathbf{x}, 1) \rightsquigarrow \mathcal{R}}}{\overline{((?=(?<N_1>\mathbf{a}(\varepsilon \mid (?P>N_1))\mathbf{b})\mathbf{c})\mathbf{a}\mathbf{a}^*(?<N_2>\mathbf{b}(\varepsilon \mid (?P>N_2))\mathbf{c}), \mathbf{x}, 1) \rightsquigarrow \mathcal{R}}}$$

The lookahead matches, following the matching relation (5), only if the regex  $\mathbf{r}_1$  matches, while the current position in  $\mathbf{x}$  is unchanged. Thus,  $\mathbf{r}$  matches  $\mathbf{x}$  if  $\mathbf{x}$  starts with  $\mathbf{a}^g \mathbf{b}^g \mathbf{c}$  and also has the form  $\mathbf{a}^g \mathbf{b}^g \mathbf{c}^g$ . In other words,  $L(\mathbf{r}) = (\{\mathbf{a}^g \mathbf{b}^g \mathbf{c} : g \in \mathbb{N}\} \cdot \{\mathbf{a}, \mathbf{b}, \mathbf{c}\}^*) \cap (\{\mathbf{a}^g : g \in \mathbb{N}\} \cdot \{\mathbf{b}^g \mathbf{c}^g : g \in \mathbb{N}\})$ .  $\square$

## 5.1 Relation with context-sensitive languages

To the author's knowledge, there is no peer-reviewed or academic publication concerning the expressive power of practical regular expressions with both lookaround assertions and subroutine calls. The only known text on this topic is due to Popov[19]: an idea of what a reduction of context-sensitive grammars to regexes might look like.

Popov claims that having a context-sensitive grammar  $\mathfrak{G} = (\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$ , any production in the form  $\mathbf{pNs} \rightarrow \mathbf{pvs}$  can be converted into a DEFINE rule of the form  $(?(DEFINE)(?<N>(?<= \text{rx}(\mathbf{p})) \text{rx}(\mathbf{v})(?= \text{rx}(\mathbf{s}))))$ . However, this alone does not work for all context-sensitive grammars. Let us attempt to formalize the conversion in algorithm 3 as a modification of algorithm 1:

Although a regex  $\mathbf{r}$  matches any string generated by grammar  $\mathfrak{G}$ , it is still not correct because, in general, it can match more. To apply a production of the form

**Algorithm 3** Conversion of a context-sensitive grammar to a regex ([19], incorrect)

**Input:** a context-sensitive grammar  $\mathfrak{G} = (\mathcal{V}, \mathcal{A}, \mathcal{R}, S)$

**Output:** a regex  $r \in \mathbb{E}_{\mathcal{S}, \mathcal{A}, \mathcal{X}}$  such that  $L(\mathfrak{G}) = L(r)$

1. initialize  $r = \epsilon$  and consider  $\mathcal{X} = \mathcal{V}$
2. for all  $N \in \mathcal{V}$ :
  - (a) let  $r_N = \epsilon$
  - (b) for all productions with  $N$  on the left-hand side and particular left and right context ( $\mathbf{p}Ns \rightarrow \mathbf{p}\mathbf{v}_{N1}s \mid \dots \mid \mathbf{p}\mathbf{v}_{Nm_N}s \in \mathcal{R}$ ):
    - i. let  $r_{Ng}$  ( $1 \leq g \leq m_N$ ) be constructed from  $\mathbf{v}_{Ng}$  for  $\mathbf{p}Ns$  ( $\mathbf{v}_{Ng}$ ):  $r_{Ng} = rx(\mathbf{v}_{Ng})$
    - ii. if  $r_N = \epsilon$  then  $r_N = (?<=\mathbf{p})r_{Ng}(?=s)$  else  $r_N = r_N \mid (?<=\mathbf{p})r_{Ng}(?=s)$
  - (c) let  $r = r(?(\text{DEFINE})(?<N>r_N))$
3. add the matching of the initial symbol, i.e., let  $r = r((?P>S))$

$\mathbf{p}Ns \rightarrow \mathbf{p}\mathbf{v}s$  in the generation of string  $x$ , the left-hand side of the production must appear in a sentential form, that is,  $S \xrightarrow[\mathfrak{G}]{}^* \mathbf{p}'\mathbf{p}Ns s' \xrightarrow[\mathfrak{G}]{}^* \mathbf{p}'\mathbf{p}\mathbf{v}ss' \xrightarrow[\mathfrak{G}]{}^* x$ . In other words, the context  $(\mathbf{p}, s)$  of the production must already be present in the sentential form. Similarly to  $\mathbf{p}Ns \rightarrow \mathbf{p}\mathbf{v}s$  being part of generating a substring of  $x$ , the subpattern  $(?<=\mathbf{p})\mathbf{v}(?=s)$  matches a substring of  $x$ . However,  $(?<=\mathbf{p})\mathbf{v}(?=s)$  can match a substring of another  $x' \notin L(\mathfrak{G})$  because during the matching of a regex, the order of the use of subpatterns is independent of the derivations of sentential forms by  $\mathfrak{G}$ . The following example illustrates this.

*Example 15.* Let us have  $\mathfrak{G} = (\{N_1, N_2, S\}, \{a, c\}, \{S \rightarrow N_2N_1N_1, N_1 \rightarrow a, aN_1 \rightarrow aaa, N_2a \rightarrow caa\}, S)$ . Clearly,  $L(\mathfrak{G}) = \{caaa, caaaa\}$ . After applying algorithm 3,

$$\begin{aligned} r = & (?(\text{DEFINE})(?<N_1>a \mid (?<=a)aa)) (?(\text{DEFINE})(?<N_2>ca(?=a))) \cdot \\ & \cdot (?(\text{DEFINE})(?<S>(?P>N_2)(?P>N_1)(?P>N_1)))(?P>S) \end{aligned}$$

and  $L(\mathfrak{G}) \subsetneq L(r)$ , as  $caaaaa \in L(r)$ :  $caaaaa[1..2]$  is matched by the subpattern  $ca(?=a)$ , and both  $caaaaa[3..4]$  and  $caaaaa[5..6]$  are matched by the subpattern  $(?<=a)aa$ . However, in  $\mathfrak{G}$ , there is no way to apply production  $aN_1 \rightarrow aaa$  twice, as there must first be symbol  $a$  in a sentential form (which consumes one  $N_1$ ).

As a result, the relation between the class of context-sensitive languages and the class of languages expressed by regexes with both subroutine calls and lookahead assertions remains an open problem.

## 6 Conclusions

We presented a formalisation of syntax and semantics of certain features of practical regular expressions using the matching relation: subroutine call, named parenthesised expression, and DEFINE rule. We attempted to mimic documented (and real) behaviour of certain flavours of practical regular expressions: Perl-compatible regular expressions, Perl, and Ruby Regexp class.

This paper showed the equivalence of context-free languages and languages expressed by practical regular expressions with concatenation, alternative, and subroutine call. This result applies to flavours that support subroutine calls. We presented an alternative constructive proof employing named subroutine calls, DEFINE rules, and the matching relation: a conversion between such practical regular expressions and context-free grammar.

We showed that adding zero-width lookahead assertions to practical regular expressions with operations concatenation, alternative, and subroutine call extends their expressive power beyond context-free languages. However, the relation of the language class expressed by such expressions to some non-context-free languages, particularly the class of context-sensitive languages, remains an open problem.

We hope that our results stimulate more work on the expressive power of specific combinations of operations used in practical regular expressions, such as backreferences, subroutine calls, lookahead assertions, or atomic groups.

## References

1. A. V. AHO: *Algorithms for Finding Patterns in Strings*, in Algorithms and Complexity, J. van Leeuwen, ed., Handbook of Theoretical Computer Science, Elsevier, 1990, pp. 255–300.
2. M. BERGLUND AND B. VAN DER MERWE: *Regular Expressions with Backreferences Re-examined*, in Proceedings of the Prague Stringology Conference 2017, Czech Technical University in Prague, 2017, pp. 30–41.
3. M. BERGLUND, B. VAN DER MERWE, AND S. VAN LITSENBORGH: *Regular Expressions with Lookahead*. Journal of Universal Computer Science, 27(4) 2021, pp. 324–340.
4. C. CÂMPEANU, K. SALOMAA, AND S. YU: *A formal study of practical regular expressions*. International Journal of Foundations of Computer Science, 14(06) Dec. 2003, pp. 1007–1018.
5. N. CHIDA AND T. TERAUCHI: *On lookaheads in regular expressions with backreferences*, in 7th International Conference on Formal Structures for Computation and Deduction (FSCD 2022), A. P. Felty, ed., vol. 228 of Leibniz International Proceedings in Informatics (LIPIcs), Dagstuhl, Germany, 2022, Schloss Dagstuhl – Leibniz-Zentrum für Informatik, pp. 15:1–15:18.
6. N. CHIDA AND T. TERAUCHI: *Repairing DoS vulnerability of real-world regexes*, in 2022 IEEE Symposium on Security and Privacy (SP), Los Alamitos, CA, USA, May 2022, IEEE Computer Society, pp. 1049–1066.
7. N. CHOMSKY: *Three models for the description of language*. IEEE Transactions on Information Theory, 2(3) Sept. 1956, pp. 113–124.
8. J. E. F. FRIEDL: *Mastering Regular Expressions*, O'Reilly, Sebastopol, CA, 3rd ed., 2006.
9. S. A. GREIBACH: *A New Normal-Form Theorem for Context-Free Phrase Structure Grammars*. Journal of the ACM, 12(1) Jan. 1965, pp. 42–52.
10. P. HAZEL: *PCRE2 — Perl-compatible regular expressions (revised API)*, University of Cambridge, Cambridge CB2 3QH, England, 2022, UNIX manual page PCRE2PATTERN(3) referring to PCRE2 10.40.
11. J. E. HOPCROFT, R. MOTWANI, AND J. D. ULLMAN: *Introduction to Automata Theory, Languages, and Computation*, Always Learning / Pearson, Pearson Education, Harlow, 3rd ed., 2014.
12. V. HRUŠA: *Regular Expressions with Subpattern Recursion*, Master's thesis, Czech Technical University in Prague, 2021.
13. S. C. KLEENE: *Representation of Events in Nerve Nets and Finite Automata*, in Automata Studies. (AM-34), C. E. Shannon and J. McCarthy, eds., Princeton University Press, Dec. 1956, pp. 3–42.
14. D. C. KOZEN: *Automata and Computability*, Undergraduate Texts in Computer Science, Springer, New York, Oct. 2012.
15. A. MATEESCU AND A. SALOMAA: *Aspects of Classical Language Theory*, in Handbook of Formal Languages, G. Rozenberg and A. Salomaa, eds., Springer Berlin Heidelberg, 1997, pp. 175–251.
16. B. MELICHAR AND J. HOLUB: *6D Classification of Pattern Matching Problems*, in Proceedings of the Prague Stringology Club Workshop '97, Czech Technical University in Prague, 1997.
17. A. MORIHATA: *Translation of regular expression with lookahead into finite state automaton*. Computer Software, 12(1) 2012, pp. 148–158.
18. *perlre — Perl regular expressions*, Perl Programmers Reference Guide, 2022, UNIX manual page PERLRE(1) referring to perl v5.36.0.
19. N. POPOV: *The true power of regular expressions*. <https://www.npopov.com/2012/06/15/The-true-power-of-regular-expressions.html>, June 2012.

20. *class Regexp*, <https://docs.ruby-lang.org/en/3.1/Regexp.html>, The API documentation for Ruby 3.1.
21. S. SCHEINBERG: *Note on the boolean properties of context free languages*. Information and Control, 3(4) Dec. 1960, pp. 372–375.
22. M. L. SCHMID: *Characterising REGEX languages by regular languages equipped with factor-referencing*. Information and Computation, 249 Oct. 2016.
23. K. THOMPSON: *Programming Techniques: Regular expression search algorithm*. Communications of the ACM, 11(6) June 1968, pp. 419–422.