

On Lookaheads in Regular Expressions with Backreferences

Nariyoshi Chida  

NTT Corporation, Japan

Waseda University, Japan

Tachio Terauchi 

Waseda University, Japan

Abstract

Many modern regular expression engines employ various extensions to give more expressive support for real-world usages. Among the major extensions employed by many of the modern regular expression engines are *backreferences* and *lookaheads*. A question of interest about these extended regular expressions is their expressive power. Previous works have shown that (i) the extension by lookaheads does not enhance the expressive power, i.e., the expressive power of regular expressions with lookaheads is still regular, and that (ii) the extension by backreferences enhances the expressive power, i.e., the expressive power of regular expressions with backreferences (abbreviated as *rewb*) is no longer regular. This raises the following natural question: Does the extension of regular expressions with backreferences by lookaheads enhance the expressive power of regular expressions with backreferences? This paper answers the question positively by proving that adding either positive lookaheads or negative lookaheads increases the expressive power of *rewb* (the former abbreviated as *rewb_p* and the latter as *rewb_n*). A consequence of our result is that neither the class of finite state automata nor that of memory automata (MFA) of Schmid [14] (which corresponds to regular expressions with backreferences but without lookaheads) corresponds to *rewb_p* or *rewb_n*. To fill the void, as a first step toward building such automata, we propose a new class of automata called *memory automata with positive lookaheads* (PLMFA) that corresponds to *rewb_p*. The key idea of PLMFA is to extend MFA with a new kind of memories, called *positive-lookahead memory*, that is used to simulate the backtracking behavior of positive lookaheads. Interestingly, our positive-lookahead memories are almost perfectly *symmetric* to the capturing-group memories of MFA. Therefore, our PLMFA can be seen as a natural extension of MFA that can be obtained independently of its original intended purpose of simulating *rewb_p*.

2012 ACM Subject Classification Theory of computation → Formal languages and automata theory

Keywords and phrases Regular expressions, Lookaheads, Backreferences, Memory automata.

Digital Object Identifier 10.4230/LIPIcs.FSCD.2022.13

Funding This work was supported by JSPS KAKENHI Grant Numbers 17H01720, 18K19787, 20H04162, 20K20625, and 22H03570.

1 Introduction

Regular expressions, introduced by Kleene [9], and the extensions employed by many of the modern regular expression engines are widely studied in formal language theory. Among the major extensions are *backreferences* and *lookaheads*. Previous works on formal language theory have studied the two features mostly in isolation. Morihata [12] and Berglund et al. [3] showed that extending regular expressions by lookaheads does not enhance their expressive power. Their proofs are by a translation to boolean finite automata [4] whose expressive power is regular. The formal study of regular expressions with backreferences (*rewb*) dates back to the seminal work by Aho [1]. More recently, a formal semantics and a pumping lemma were given by Câmpeanu et al. [5], and Berglund and van der Merwe [2] showed that different variants of backreference semantics give rise to differences in expressive powers.



© Nariyoshi Chida and Tachio Terauchi;
licensed under Creative Commons License CC-BY 4.0

7th International Conference on Formal Structures for Computation and Deduction (FSCD 2022).

Editor: Amy P. Felty; Article No. 13; pp. 13:1–13:18

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

46 Schmid [14] proposed *memory automata* (MFA) and showed that the expressive power of the
47 automata is equivalent to that of *rewb*.

48 In this paper, we initiate a formal study of *regular expression with backreferences and*
49 *lookaheads* (*rewbl* for short). We call the fragment containing only positive (resp. negative)
50 lookaheads *rewbl_p* (resp. *rewbl_n*). We show that both *rewbl_p* and *rewbl_n* are more expressive
51 than *rewb*, and also prove some language-theoretic properties of *rewbl*. One consequence of
52 the results is the undecidability of a problem tackled in a recent work [11].

53 Another consequence of our results is that neither the class of finite state automata nor that
54 of memory automata (MFA) of Schmid [14] (which corresponds to regular expressions with
55 backreferences but without lookaheads) corresponds to *rewbl_p* or *rewbl_n*. As remarked above,
56 prior works [3, 12] have applied translation to boolean finite automata [4] (or alternating
57 finite automata [7]) to build automata equivalent to regular expressions with lookaheads.
58 They simulate lookaheads by executing multiple runs simultaneously without backtracking.
59 Unfortunately, the interaction of lookaheads with backreferences prevents us from applying
60 the approaches to *rewbl*. Namely, *rewbl* permits *cross-lookahead backreferences* whereby a
61 string captured outside of a lookahead is referred from inside of the lookahead, or vice versa
62 (only the former is allowed for negative lookaheads whereas both are allowed for positive
63 lookaheads). Such cross-lookahead backreferences intrinsically require backtracking. In our
64 work, as a first step toward building automata equivalent to *rewbl*, we introduce a new class
65 of automata called *memory automata with positive lookaheads* (PLMFAs). We prove that
66 PLMFAs are equivalent to *rewbl_p* in expressive power. A key component of PLMFAs is
67 a new kind of memories, called a *positive-lookahead memory*, that is used to simulate the
68 backtracking behavior of positive lookaheads. Interestingly, our positive-lookahead memories
69 are almost perfectly *symmetric* to the capturing-group memories of MFA. Therefore, our
70 PLMFA can be seen as a natural extension of MFA that can be obtained independently of
71 its original intended purpose of simulating *rewbl_p*.

72 In summary, this paper makes the following contributions:

- 73 ■ We show that the extension of *rewb* by either positive or negative lookaheads enhances
74 the expressive power. Additionally, we prove some language-theoretic properties of *rewbl*.
75 (Sec. 3)
- 76 ■ We introduce memory automata with positive lookaheads (PLMFAs), a new class of
77 automata that we prove to be equivalent in expressive power to *rewbl_p*. A key component
78 of PLMFAs is a new kind of memories called positive-lookahead memory, which is almost
79 perfectly symmetric to capturing-group memory of MFA. (Sec. 4)

80 We believe that our work leads to interesting future developments in both theoretical and
81 practical fronts: interesting practically because backreferences and lookaheads are practically
82 motivated by real-world needs, and interesting theoretically because, as we shall show, *rewbl*
83 does not appear to correspond to any known formal language classes.

84 **2 Preliminaries**

85 In this section, we introduce the preliminary notations (Sec. 2.1) and present the syntax and
86 the semantics of *rewbl* (Sec. 2.2).

87 **2.1 Notation**

88 We write \mathbb{N} for the set of natural numbers and $[i]$ for the set $\{1, 2, \dots, i\}$ where $i \in \mathbb{N}$. For
89 a sequence l , we write $|l|$ for its length, $l[i]$ (for $1 \leq i \leq |l|$) for its i th element, $l[i..j]$
90 for the sub-sequence from the i th element to the j th element (for $1 \leq i \leq j \leq |l|$). We write

$r ::= a$	character		r^*	repetition
\emptyset	empty set		$(_i r)_i$	capturing group
ϵ	empty string		$\backslash i$	backreference
rr	concatenation		$(?=r)$	pos-lookahead
$r r$	union		$(?!r)$	neg-lookahead

■ **Figure 1** The syntax of rewbl expressions

91 $l_1 :: l_2$ for the concatenation of l_1 and l_2 . We abbreviate it as $l_1 l_2$ if clear from the context.
 92 We write $v \in l$ to denote that l contains v . We write Σ for a finite alphabet; $a, b \in \Sigma$ for a
 93 character; $x, y \in \Sigma^*$ for a sequence of characters (i.e., *string*); ϵ for the empty string; Σ_ϵ for
 94 $\Sigma \cup \{\epsilon\}$; In what follows, we fix a finite alphabet Σ . For $1 \leq i < j \leq |x|$, we define $x[i..j]$ to
 95 be $x[i..j-1]$. For $x, y \in \Sigma^*$, we define $x \backslash y$ to be the left quotient of x divided by y , i.e.,
 96 v where $yv = x$. Dually, the right quotient of x divided y , x/y , is v where $vy = x$. $S \subset U$
 97 denotes that S is a proper subset of U , i.e., $S \subseteq U \wedge S \neq U$. For a partial map f , we write
 98 $dom(f)$ for the domain of f . $\mathcal{P}(S)$ denotes that the power set of a set S . For f a (partial)
 99 function, $f[\alpha \mapsto \beta]$ denotes the (partial) function that maps α to β and behaves as f for all
 100 other arguments. We write $f(\alpha) = \perp$ if f is undefined at α .

101 2.2 Regular Expressions with Backreferences and Lookaheads

102 The syntax of *regular expressions with backreferences and lookaheads (rewbl)* is given by
 103 Fig. 1. The semantics of the pure regular expression constructs (i.e., the first six constructs
 104 of Fig. 1) is standard. We write \cdot for the rewbl that matches any character (i.e., $a_1 | \dots | a_n$
 105 where $\Sigma = \{a_1, \dots, a_n\}$). The precedence order of the operators is as follows: Kleene-*,
 106 concatenation, and union. The left has a higher precedence. For example, the expression
 107 $a^* | bc^*$ means $((a^*) | (bc^*))$ due to the priority.

108 The remaining constructs, i.e., capturing groups, backreferences, and lookaheads, are
 109 the extensions considered in this paper. In conformance with the nomenclature from the
 110 literature [10], we call the fragment of rewbl without lookaheads *rewb*. We call the fragment
 111 of rewbl without negative (resp. positive) lookaheads *rewbl_p* (resp. *rewbl_n*). In what follows,
 112 we explain the semantics of the extended features informally in terms of the standard
 113 backtracking-based matching algorithm which attempts to match the given regular expression
 114 with the given string and backtracks when the attempt fails. A *capturing group* $(_i r)_i$ (or $(r)_i$
 115 if no ambiguity arises) attempts to match r , and if successful, stores the matched substring in
 116 the storage identified by the index i . Otherwise, the match fails and the algorithm backtracks.
 117 A *backreference* $\backslash i$ refers to the substring matched to the corresponding capturing group
 118 $(_i r)_i$, and attempts to match the same substring if the capture had succeeded. If the capture
 119 had not succeeded, i.e., is an *unassigned backreference*, or the matching against the captured
 120 substring fails, then the algorithm backtracks. Capturing groups in practice often do not have
 121 explicit indexes, but we write them here for readability. A *positive* (resp. *negative*) *lookahead*
 122 $(?=r)$ (resp. $(?!r)$) attempts to match r without any character consumption, proceeds if the
 123 match succeeds (resp. fails), and backtracks otherwise.

124 More formally, the semantics is defined by the *matching relation* \rightsquigarrow that models the
 125 behavior of backtracking matching algorithms. The full rules for deriving the matching
 126 relation \rightsquigarrow is shown in Fig. 2. The semantics is same as the one defined in our recent work [8]
 127 except for specializing the set-of-characters rules to the rules for a single character and the
 128 empty set. We define $ite(true, A, B) = A$ and $ite(false, A, B) = B$.

13:4 On Lookaheads in Regular Expressions with Backreferences

$$\begin{array}{c}
\frac{p \leq |w| \quad w[p] = a}{(a, w, p, \Lambda) \rightsquigarrow \{(p+1, \Lambda)\}} \text{ (CHARACTER)} \\
\frac{p > |w| \vee w[p] \neq a}{(a, w, p, \Lambda) \rightsquigarrow \emptyset} \text{ (CHARACTER FAILURE)} \\
\frac{}{(\emptyset, w, p, \Lambda) \rightsquigarrow \emptyset} \text{ (EMPTY SET)} \\
\frac{}{(\epsilon, w, p, \Lambda) \rightsquigarrow \{(p, \Lambda)\}} \text{ (EMPTY STRING)} \\
\frac{(r_1, w, p, \Lambda) \rightsquigarrow \mathcal{N} \quad \forall (p_i, \Lambda_i) \in \mathcal{N}, (r_2, w, p_i, \Lambda_i) \rightsquigarrow \mathcal{N}_i}{(r_1 r_2, w, p, \Lambda) \rightsquigarrow \bigcup_{0 \leq i < |\mathcal{N}|} \mathcal{N}_i} \text{ (CONCATENATION)} \\
\frac{(r_1, w, p, \Lambda) \rightsquigarrow \mathcal{N} \quad (r_2, w, p, \Lambda) \rightsquigarrow \mathcal{N}'}{(r_1 | r_2, w, p, \Lambda) \rightsquigarrow \mathcal{N} \cup \mathcal{N}'} \text{ (UNION)} \\
\frac{\forall (p_i, \Lambda_i) \in (\mathcal{N} \setminus \{(p, \Lambda)\}), (r^*, w, p_i, \Lambda_i) \rightsquigarrow \mathcal{N}_i}{(r^*, w, p, \Lambda) \rightsquigarrow \{(p, \Lambda)\} \cup \bigcup_{0 \leq i < |\mathcal{N} \setminus \{(p, \Lambda)\}|} \mathcal{N}_i} \text{ (REPETITION)} \\
\frac{(r, w, p, \Lambda) \rightsquigarrow \mathcal{N}}{((r)_j, w, p, \Lambda) \rightsquigarrow \{(p_i, \Lambda_i [j \mapsto w[p..p_i]]) \mid (p_i, \Lambda_i) \in \mathcal{N}\}} \text{ (CAPTURING GROUP)} \\
\frac{\Lambda(i) \neq \perp \quad (\Lambda(i), w, p, \Lambda) \rightsquigarrow \mathcal{N}}{(\backslash i, w, p, \Lambda) \rightsquigarrow \mathcal{N}} \text{ (BACKREFERENCE)} \\
\frac{\Lambda(i) = \perp}{(\backslash i, w, p, \Lambda) \rightsquigarrow \emptyset} \text{ (BACKREFERENCE FAILURE)} \\
\frac{(r, w, p, \Lambda) \rightsquigarrow \mathcal{N}}{((?=r), w, p, \Lambda) \rightsquigarrow \{(p, \Lambda') \mid (_, \Lambda') \in \mathcal{N}\}} \text{ (POSITIVE LOOKAHEAD)} \\
\frac{(r, w, p, \Lambda) \rightsquigarrow \mathcal{N} \quad \mathcal{N}' = \text{ite}(\mathcal{N} \neq \emptyset, \emptyset, \{(p, \Lambda)\})}{((?!r), w, p, \Lambda) \rightsquigarrow \mathcal{N}'} \text{ (NEGATIVE LOOKAHEAD)}
\end{array}$$

■ **Figure 2** Rules of the matching relation \rightsquigarrow

129 A matching relation is of the form $(r, w, p, \Lambda) \rightsquigarrow \mathcal{N}$ where p is a position on the string w
130 such that $1 \leq p \leq |w| + 1$, Λ , called an *environment*, is a function that maps each capturing
131 group index to a string captured by the corresponding capturing group, and \mathcal{N} is a set of
132 matching results. A *matching result* is a pair of a position and an environment. Roughly,
133 $(r, w, p, \Lambda) \rightsquigarrow \mathcal{N}$ is read: a rewbl expression r tries to match the string w from the position p ,
134 with the environment Λ and, if $(p', \Lambda') \in \mathcal{N}$, r consumed $p' - p$ characters and updated the
135 environment to Λ' . Additionally, if $\mathcal{N} = \emptyset$, it means that the matching failed.

136 In the two rules for a character, the rewbl a tries to match the string w at the position p
137 with the function capturing Λ . If the p th character $w[p]$ is a , then the matching succeeds
138 returning the matching result $(p+1, \Lambda)$ (CHARACTER). Otherwise, the character $w[p]$ does
139 not match or the position is at the end of the string, and \emptyset is returned as the matching result
140 indicating the match failure (CHARACTER FAILURE).

141 The rules (EMPTY SET), (EMPTY STRING), (CONCATENATION), (UNION) and (REPETITION)
142 are self explanatory. Note that we avoid self looping in (REPETITION) by not repeating the
143 match from the same position.

144 In the rule (CAPTURING GROUP), we first get the matching result \mathcal{N} from matching w
145 against r at the current position p . And for each matching result $(p_i, \Lambda_i) \in \mathcal{N}$ (if any), we
146 record the matched substring $w[p..p_i]$ in the corresponding environment Λ_i at the index i .
147 The rule (BACKREFERENCE) looks up the captured substring and tries to match it with the
148 input at the current position. The match fails if the corresponding capture has failed as
149 stipulated by the rule (BACKREFERENCE FAILURE).

150 In the rule (POSITIVE LOOKAHEAD), the expression r is matched against the given
151 string w at the current position p to obtain the matching results \mathcal{N} . Then, for every match
152 result $(p', \Lambda') \in \mathcal{N}$ (if any), we reset the position from p' to p . This models the behavior of
153 lookaheads which does not consume the string. The rule (NEGATIVE LOOKAHEAD) is similar,
154 except that we reset and proceed when there is no match. Note that captures made inside
155 of a negative lookahead cannot be referred outside of the lookahead, which agrees with the
156 behavior of regular expression engines in practice.

157 ► **Definition 1 (Language).** The *language* of a rewbl r is defined as $L(r) = \{w \mid (r, w, 1, \emptyset) \rightsquigarrow$
158 $\mathcal{N} \wedge \exists \Lambda. (|w| + 1, \Lambda) \in \mathcal{N}\}$.

159 Recall that one subtle aspect of rewbl is that backreferences can cross lookahead boundaries
160 (cf. Sec. 1). We next show some examples of cross-lookahead backreferences.

161 ► **Example 2.** Consider the expression $(\cdot^*z)_1(=?\setminus 1)^*$. Its language is $\{xzxzy \mid x, y \in \Sigma^*\}$.
162 For example, when the input string is **azazbc**, the expression captures the prefix **az** and
163 refers it from the inside of the positive lookahead $(=?\setminus 1)$.

164 ► **Example 3.** Consider the expression $(\cdot)_1(?!\setminus 1)^*$. The language is $\{a \mid a \in \Sigma\} \cup \{abx \mid$
165 $a, b \in \Sigma \wedge a \neq b \wedge x \in \Sigma^*\}$.

166 ► **Example 4.** Consider the expression $(?=(\cdot^*)_1z)\setminus 1$. The language is $\{zx \mid x \in \{z\}^*\}$.

167 Example 2 (resp. 3) shows an example where a string captured outside of a positive (resp. neg-
168 ative) lookahead is backreferenced in the lookahead. Example 4 shows an example where a
169 string captured inside of a positive lookahead is backreferenced from outside of the lookahead.

170 2.2.1 Conventions on Syntax and Semantics

171 We review the conventions regarding capturing groups and unassigned references. The
172 conventions are proposed in prior works on rewbl, and as shown by [2], they affect the
173 expressive power of rewbl. Here, we simply present the conventions and refer interested
174 readers to [2] for the expressive power differences.

175 There are two conventions regarding capturing groups: *no label repetitions (NLR)* and
176 *may repeat labels (MRL)*. NLR requires the indexes of capturing groups to be distinct, whereas
177 MRL imposes no such restrictions. For example, $(\cdot^*)_1\setminus 1$ satisfies NLR, but $((\cdot^*)_1\mid(\cdot^*)_1)\setminus 1$
178 does not because the capturing group with index 1 appears twice. NLR is assumed in the
179 prior works by Câmpeanu et al. [5] and Carle and Narendran [6] on the expressive power of
180 rewbl.

181 There are two conventions regarding unassigned references: the ϵ *semantics* and the \emptyset
182 *semantics*. The ϵ (resp. \emptyset) semantics defines that unassigned references are handled as an
183 empty string ϵ (resp. a failure \emptyset). For example, for $r = a\setminus 1$, $L(r) = \{a\}$ with the ϵ semantics
184 but $L(r) = \emptyset$ with the \emptyset semantics. Additionally, prior works have proposed a condition called
185 *no unassigned reference (NUR)*. The NUR condition does not allow unassigned references in
186 expressions, i.e., all expressions with unassigned references are to be excluded (see below for
187 the formal definition). For example, $r = a\setminus 1$ does not satisfy the NUR condition because $\setminus 1$
188 is an unassigned references. Note that the ϵ semantics and the \emptyset semantics coincide under
189 the NUR condition because there would be no unassigned references. The condition is also
190 assumed in [5, 6] ([6] incorrectly remarks that [5] does not assume the condition).

191 In the rest of this section, we give a formal definition of the NUR condition that we shall
192 also use later in our proofs. We note that prior works that proposed the condition did not
193 provide a formal definition of it [5, 6]. First, we define the function *Capture* from rewbls to
194 the set of capturing group indexes that can be referred from their continuations:

$$195 \text{Capture}(r) = \begin{cases} \emptyset & (\text{if } r = a, \emptyset, \epsilon, r_1^*, \setminus i, \text{ or } (?!r_1)) \\ \text{Capture}(r_1) \cup \text{Capture}(r_2) & (\text{if } r = r_1r_2) \\ \text{Capture}(r_1) \cap \text{Capture}(r_2) & (\text{if } r = r_1|r_2) \\ \text{Capture}(r_1) \cup \{i\} & (\text{if } r = (i r_1)_i) \\ \text{Capture}(r_1) & (\text{if } r = (?=r_1)) \end{cases}$$

196 With this, we can define the predicate $\text{NUR}(S, r)$ that says that r satisfies NUR condition if
197 it occurs in a context where the capturing group indexes in S can be referred:

$$\begin{aligned}
198 \quad \text{NUR}(S, r) = \begin{cases} \text{true} & (\text{if } r = a, \emptyset, \text{ or } \epsilon) \\ \text{NUR}(S, r_1) \wedge \text{NUR}(S \cup \text{Capture}(r_1), r_2) & (\text{if } r = r_1 r_2) \\ \text{NUR}(S, r_1) \wedge \text{NUR}(S, r_2) & (\text{if } r = r_1 | r_2) \\ \text{NUR}(S, r_1) & (\text{if } r = r_1^*, (i r_1)_i, (?=r_1), \text{ or } (?!r_1)) \\ i \in S & (\text{if } r = \setminus i) \end{cases}
\end{aligned}$$

199 Then, r can be said to satisfy the NUR condition iff $\text{NUR}(\emptyset, r) = \text{true}$. For example, the
200 expression $r = ({}_1\mathbf{a})_1 \setminus 1$ satisfies the NUR condition because $\text{NUR}(\emptyset, r) = \text{NUR}(\emptyset, ({}_1\mathbf{a})_1) \wedge$
201 $\text{NUR}(\emptyset \cup \text{Capture}(({}_1\mathbf{a})_1), \setminus 1) = \text{NUR}(\emptyset, \mathbf{a}) \wedge \text{NUR}(\{1\}, \setminus 1) = \text{true}$. As another example, $r =$
202 $({}_1\mathbf{a} \setminus 1)_1$ does not satisfy the NUR condition because $\text{NUR}(\emptyset, r) = \text{NUR}(\emptyset, \mathbf{a} \setminus 1) = \text{NUR}(\emptyset,$
203 $\mathbf{a}) \wedge \text{NUR}(\emptyset \cup \text{Capture}(\mathbf{a}), \setminus 1) = \text{false}$. In what follows, unless explicitly stated otherwise, we
204 assume that a rewbl that satisfies NLR and NUR.

205 3 Language Properties of Rewbl

206 In this section, we prove some salient language properties of rewbl. Importantly, we show
207 that both rewbl_p and rewbl_n is strictly more expressive than rewbl, thus showing that the
208 extension by either positive or negative lookaheads changes the expressive power of rewbl.
209 In the following, we denote by $\mathbb{L}_B, \mathbb{L}_{BL}, \mathbb{L}_{BL_n},$ and \mathbb{L}_{BL_p} the class of languages matched by
210 rewbl, rewbl, rewbl_n , and rewbl_p , respectively.

211 Our first result states that \mathbb{L}_{BL_n} is closed under union, intersection, and complement.

212 ► **Theorem 5.** \mathbb{L}_{BL_n} is closed under union, intersection, and complement.

213 **Proof.** Suppose we have rewbl expressions r_1 and r_2 . Then, rewbl expressions that accept
214 the union of $L(r_1)$ and $L(r_2)$, the intersection of $L(r_1)$ and $L(r_2)$, and the complement of
215 $L(r_1)$ can be constructed respectively as follows.

- 216 ■ **Union:** $L(r_1) \cup L(r_2) = L(r_1 | r_2)$;
- 217 ■ **Intersection:** $L(r_1) \cap L(r_2) = L((?!r_1(?!))r_2)$; and
- 218 ■ **Complement:** $L(r_1)^c \triangleq \Sigma^* \setminus L(r_1) = L((?!r_1(?!)).*)$.

219 In **Intersection** and **Complement**, a subtle point is that a negative lookahead $(?!r)$
220 accepts a string even if the expression r rejects only a prefix of the string. For example,
221 $L((?!r_1)r_2)$ is the set of strings in $L(r_2)$ that have a prefix that belongs to $L(r_1)$, rather
222 than the intersection of $L(r_1)$ and $L(r_2)$. To force whole matching, the negative lookahead
223 $(?!.)$ is appended. ◀

224 Of course, we could alternatively show **Intersection** from **Union** and **Complement** by
225 applying De Morgan's laws: $L(r_1) \cap L(r_2) = (L(r_1)^c \cup L(r_2)^c)^c$. The above proof gives a
226 direct construction which shows that the intersection can be obtained by a short rewbl_n
227 expression.

228 We next show that \mathbb{L}_{BL} is also closed under union, intersection, and complement.

229 ► **Theorem 6.** \mathbb{L}_{BL} is closed under union, intersection, and complement.

230 **Proof.** The proof is the same as Theorem 5. Or, for **Intersection**, an even shorter proof is
231 possible: $L(r_1) \cap L(r_2) = L((?=r_1(?!))r_2)$. ◀

232 A consequence of Theorem 5 is that rewbl and rewbl_n are more expressive than rewbl.

233 ► **Corollary 7.** $\mathbb{L}_B \subset \mathbb{L}_{BL_n} \subseteq \mathbb{L}_{BL}$.

234 **Proof.** Immediate from Theorem 5 and Lemma 3 of [6] which showed that rewbl is not closed
235 under intersection. ◀

	Closure under				
	\cup	\cap	Complement	Concatenation	Kleene-*
Regular	Yes	Yes	Yes	Yes	Yes
\mathbb{L}_B	Yes	No	No	Yes	Yes
\mathbb{L}_{BL_p}	Yes	?	?	Yes	Yes
\mathbb{L}_{BL_n}	Yes	Yes	Yes	Yes	Yes
\mathbb{L}_{BL}	Yes	Yes	Yes	Yes	Yes

Table 1 Summary of closure properties. The rows highlighted in gray present our new results. ? indicates that the problem is open.

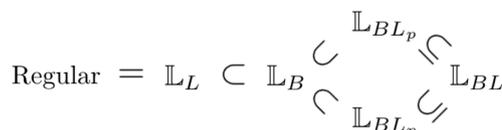
236 We show that adding just positive lookaheads also increases the expressive power of rewbl.
 237

238 ▶ **Theorem 8.** $\mathbb{L}_B \subset \mathbb{L}_{BL_p}$.

239 **Proof.** From [6], the language $S = \{a^i b a^{i+1} b a^k \mid k = i(i+1)k', k' > 0, \text{ and } i > 0\}$ is not in
 240 \mathbb{L}_B . Let $S' = \{a^i b a^{i+1} b a^k c \mid k = i(i+1)k', k' > 0, \text{ and } i > 0\}$. We can prove that S' is also
 241 not in \mathbb{L}_B in a manner similar to the proof that S is not in \mathbb{L}_B [6].

242 Now, S' is the intersection of $L(r_1 c)$ and $L(r_2 c)$ where r_1 and r_2 are $(aa^*)_1 b \setminus 1 a b \setminus 1 \setminus 1^*$
 243 and $(aa^*)_1 b \setminus (1a)_2 b \setminus 2 \setminus 2^*$, respectively. Then, the intersection of the languages of $r_1 c$ and
 244 $r_2 c$ is $L((?=r_1 c)r_2 c) = S' \in \mathbb{L}_{BL_p}$. ◀

245 A summary of the results of closure properties can be found in Table 1. Additionally,
 246 the diagram below summarizes the results of the hierarchy of the language classes. Here,
 247 \mathbb{L}_L denotes the class of languages matched by regular expressions with (both positive and
 248 negative) lookaheads, which is known to be equivalent in expressive power to the set of
 249 regular languages [3, 12].



250

251 From Theorems 5 and 6, we also obtain the following result.

252 ▶ **Theorem 9.** *The emptiness problems of rewbl and rewbl_n are undecidable.*

253 **Proof.** Suppose for contradiction that the emptiness problem of rewbl is decidable. By
 254 Theorem 6, we know that the language of rewbl is closed under intersection. Therefore, the
 255 emptiness problem of intersection of two rewbl expressions is also decidable. However, this
 256 contradicts a result from [6] which states that the latter problem is undecidable. The proof
 257 for rewbl_n is similar by using Theorem 5. ◀

258 A recent work [11] has proposed a method for symbolically executing programs containing
 259 rewbl. Their method generates and tries to solve constraints of the form $x \in L(r)$ where r
 260 is a rewbl expression and x is a variable for which the method tries to find an assignment
 261 that satisfies the constraint. Theorem 9 implies that their constraint solving problem is
 262 undecidable.

263 ▶ **Corollary 10.** *The constraint solving problem of [11] is undecidable.*

264 ► Remark 11. The constructions used to show Cor. 7 and Theorem 8 do not contain
 265 backreferences that cross lookahead boundaries (recall the discussion from Sec. 1 and Sec. 2.2).
 266 Thus, our results show that lookaheads enhance the expressive power of rewbl even without
 267 cross-lookahead backreferences. We leave for future work to investigate whether there are
 268 expressive power changes from allowing or disallowing cross-lookahead backreferences.

269 3.1 Restricted Label Repetitions

270 As pointed out in [2], allowing rewbls to repeat labels of backreferences affects their expressive
 271 powers. In this sub-section, we introduce new conditions called *restricted may repeat labels*
 272 (RMRL) and *no self-capturing reference* (NSR). RMRL allows repeating the labels but only
 273 in a restricted way. NSR enforces that there is no reference that is nested by the capturing
 274 group of the same index. We use RMRL and NSR as an intermediary in the construction of
 275 automata equivalent in expressive power to rewbl_p (with NUR and NLR) in the next section.
 276 But the new conditions may also be of independent interest.

277 Informally, RMRL requires the capturing group that is referred to by any reference is
 278 uniquely determined. For example, $(_1a)_1(_1b)_1 \setminus 1$ satisfies RMRL because the reference $\setminus 1$
 279 refers to the capturing group $(_1b)_1$ while $((_1a)_1|(_1b)_1) \setminus 1$ does not satisfy RMRL because the
 280 reference $\setminus 1$ can refer to the capturing groups $(_1a)_1$ and $(_1b)_1$. To represent the repetitions
 281 of indexes, we use a *multiset*. A multiset, denoted by $\{\!\{ \dots \}\!\}$, is a collection of elements
 282 with repetitions. For example, $\{\!\{ a, a, b \}\!\}$ is a multiset that has two a 's and one b . We use
 283 the notation for sets as that of multisets, e.g., we use \cup for the union of multisets, e.g.,
 284 $\{\!\{ a \}\!\} \cup \{\!\{ a, b \}\!\} = \{\!\{ a, a, b \}\!\}$, and $|\cdot|$ for the number of elements, e.g., $|\{\!\{ a, a, b \}\!\}| = 3$. Formally,
 285 we define RMRL by first defining *NumCaps* that takes a multiset \mathbb{S} and a rewbl expression r
 286 and returns a multiset that represents the number of ways to capture for each index.

$$287 \text{NumCaps}(\mathbb{S}, r) = \begin{cases} \mathbb{S} & (\text{if } r = a, \emptyset, \epsilon, \setminus i, \text{ or } (?!r_1)) \\ \text{NumCaps}(\text{NumCaps}(\mathbb{S}, r_1), r_2) & (\text{if } r = r_1r_2) \\ \text{NumCaps}(\mathbb{S}, r_1) \cup \text{NumCaps}(\mathbb{S}, r_2) & (\text{if } r = r_1|r_2) \\ \{\!\{ j \mid j \in \text{NumCaps}(\mathbb{S}, r_1) \wedge j \neq i \}\!\} \cup \{\!\{ i \}\!\} & (\text{if } r = (_i r_1)_i) \\ \text{NumCaps}(\mathbb{S}, r_1) & (\text{if } r = (?=r_1) \text{ or } r_1^*) \end{cases}$$

288 With this, we define $\text{RMRL}(\mathbb{S}, r)$ that says that r satisfies the RMRL condition if it occurs in
 289 a context where the capturing group indexes in \mathbb{S} can be referred:

$$290 \text{RMRL}(\mathbb{S}, r) = \begin{cases} \text{true} & (\text{if } r = a, \emptyset, \text{ or } \epsilon) \\ \text{RMRL}(\mathbb{S}, r_1) \wedge \text{RMRL}(\text{NumCaps}(\mathbb{S}, r_1), r_2) & (\text{if } r = r_1r_2) \\ \text{RMRL}(\mathbb{S}, r_1) \wedge \text{RMRL}(\mathbb{S}, r_2) & (\text{if } r = r_1|r_2) \\ \text{RMRL}(\mathbb{S}, r_1) & (\text{if } r = r_1^*, (_i r_1)_i, (?=r_1), \text{ or } (?!r_1)) \\ |\{\!\{ i \mid i \in \mathbb{S} \}\!\}| \leq 1 & (\text{if } r = \setminus i) \end{cases}$$

291 We say that a rewbl expression r satisfies RMRL iff $\text{RMRL}(\emptyset, r) = \text{true}$.

292 Next, we explain NSR. NSR requires that for every reference $\setminus i$, the reference is not
 293 nested by the capturing group whose index is i . For example, $(_1a \setminus 1)_1(_2 \setminus 1)_2$ does not satisfy
 294 NSR because $\setminus 1$ appears in its capturing group. Formally, we define NSR as follows.

$$295 \text{NSR}(S, r) = \begin{cases} \text{true} & (\text{if } r = a, \emptyset, \text{ or } \epsilon) \\ \text{NSR}(S, r_1) \wedge \text{NSR}(S, r_2) & (\text{if } r = r_1r_2 \text{ or } r_1|r_2) \\ \text{NSR}(S, r_1) & (\text{if } r = r_1^*, (?=r_1), \text{ or } (?!r_1)) \\ \text{NSR}(S \cup \{i\}, r_1) & (\text{if } r = (_i r_1)_i) \\ i \notin S & (\text{if } r = \setminus i) \end{cases}$$

296 We say that a rewbl expression r satisfies NSR iff $\text{NSR}(\emptyset, r) = \text{true}$. We show that
 297 $\text{RMRL} \wedge \text{NSR}$ is equivalent to NLR in expressive powers.

298 ▶ **Lemma 12.** (1) For any NLR r there exists a NSR and RMRL r' such that $L(r) = L(r')$,
 299 and (2) for any NSR and RMRL r there exists a NLR r' such that $L(r) = L(r')$.

300 **Proof.** (1) is immediate since NLR implies both RMRL and NSR (under the NUR assumption).
 301 To see (2), if r satisfies RMRL and NSR, then capturing groups that are referred
 302 are uniquely determined and are closed when they are referred. Thus, we can construct r'
 303 by replacing indexes of reference i in r and the capturing group referred to i with unique
 304 indexes and removing all unreferred capturing groups. ◀

305 4 Memory Automata with Positive Lookaheads

306 This section present PLMFA, a new class of automata that we prove to be equivalent to
 307 rewbl_p . PLMFA is obtained by extending MFA of Schmid [14] that is equivalent to rewb . The
 308 key extension is the addition of a new kind of memories called *positive-lookahead memories*.
 309 Roughly, a PLMFA is a non-deterministic finite state automata augmented with a list
 310 of capturing-group memories and a list of positive-lookahead memories. The former also
 311 exists in MFA and stores strings captured by capturing groups to simulate the behavior
 312 of backreferences. The latter stores strings matched by positive lookaheads and is used to
 313 simulate the behavior of positive lookaheads.

314 4.1 Formal Definition

315 A *memory* is a tuple (x, \mathbf{s}) of a string $x \in \Sigma^*$ and a *status* \mathbf{s} . A status is either *open* (\mathbf{O})
 316 or *close* (\mathbf{C}). The statuses are changed by *memory instructions* (*instructions* for short)
 317 $\Theta = \{\circ, \mathbf{c}, \diamond\}$ as follows: $\mathbf{s} \oplus \circ = \mathbf{O}$, $\mathbf{s} \oplus \mathbf{c} = \mathbf{C}$, and $\mathbf{s} \oplus \diamond = \mathbf{s}$. Roughly, \mathbf{O} means that
 318 the string in the memory is modified by appending consumed strings, while \mathbf{C} means that
 319 the string in the memory is unmodified. Changing the status from \mathbf{C} to \mathbf{O} (resp. from
 320 \mathbf{O} to \mathbf{C}) representing to an entering (resp. exiting) a capturing group if the memory is
 321 a capturing-group memory and otherwise (i.e., if positive-lookahead memory) a positive
 322 lookahead. At computation steps corresponding to backreferences, the strings in capturing
 323 group memories are used to left-divide the input string and appended to strings stored in
 324 any open memories. Symmetrically, when the strings in positively lookahead memories are
 325 used, they are prepended to the input string and used to right-divide strings stored in any
 326 open memories. A positive lookahead memory is used when it gets closed. For a memory
 327 $t = (x, \mathbf{s})$, we write $t.\text{word}$ for x and $t.\text{status}$ for \mathbf{s} . We define PLMFAs as follows.

328 ▶ **Definition 13** (PLMFA). For $(k_c, k_p) \in \mathbb{N}^2$, a (k_c, k_p) -*memory automaton with positive*
 329 *lookaheads*, $\text{PLMFA}(k_c, k_p)$, is a tuple (Q, δ, q_0, F) such that

- 330 1. Q is a finite set of *states*,
- 331 2. $\delta : Q \times (\Sigma_\epsilon \cup [k_c]) \rightarrow \mathcal{P}(Q \times \Theta^{k_c} \times \Theta^{k_p})$ is the *transition function*,
- 332 3. $q_0 \in Q$ is the *initial state*, and
- 333 4. $F \subseteq Q$ is the set of *accepting states*.

334 Here, k_c and k_p represent the number of capturing-group memories and positive-lookahead
 335 memories, respectively. Next, we define *configurations* of PLMFAs.

336 ▶ **Definition 14** (Configuration). A *configuration* of a PLMFA M is a tuple $(q, w, \sigma_c, \sigma_p)$
 337 where q is a state of M , w is an input string, and σ_c (resp. σ_p) is a list of memories that
 338 represents a list of capturing-group (resp. positive-lookahead) memories.

339 ▶ **Definition 15** (Computation step). For a $\text{PLMFA}(k_c, k_p)$ M and $\ell \in \Sigma_\epsilon \cup [k_c]$, a *step*
 340 *of computation* of M is a binary relation on configurations $\xrightarrow[M]{\ell}$ (or $\xrightarrow{\ell}$ or \rightarrow if irrelevant),

13:10 On Lookaheads in Regular Expressions with Backreferences

341 defined as follows: $(q, w, \sigma_c, \sigma_p) \xrightarrow[M]{\ell} (q', w', \sigma'_c, \sigma'_p)$ iff there is a transition $\delta(q, \ell) \ni (q', ir_c, ir_p)$
 342 satisfying the following conditions.

343 (1) If $\ell \in \Sigma_\epsilon$, $v = \ell$ and otherwise (i.e., if $\ell \in [k_c]$) if $\sigma_c[\ell].status = \mathbf{C}$ then $v = \sigma_c[\ell].word$.

344 (2) $\forall \tau \in \{c, p\}. \forall i \in [k_\tau]. \sigma'_\tau[i].status = \sigma_\tau[i].status \oplus ir_\tau[i]$.

345 (3) $\forall \tau \in \{c, p\}. \forall i \in [k_\tau]$.

$$346 \quad \sigma'_\tau[i].word = \begin{cases} u :: v & (\text{if } \mathbf{O} \Rightarrow_{\tau, i} \mathbf{O}) \\ \sigma_\tau[i].word & (\text{if } _ \Rightarrow_{\tau, i} \mathbf{C}) \\ v & (\text{if } \mathbf{C} \Rightarrow_{\tau, i} \mathbf{O}) \end{cases}$$

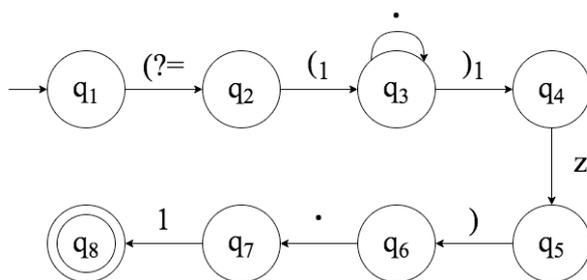
347 where $\sigma_\tau[i].status \Rightarrow_{\tau, i} \sigma'_\tau[i].status$, $bts = \sigma_p[j].word$ where $j = \operatorname{argmax}_{j \in J} |\sigma_p[j].word|$ if
 348 $J = \{j \in [k_p] \mid \mathbf{O} \Rightarrow_{p, j} \mathbf{C}\} \neq \emptyset$ and otherwise $bts = \epsilon$, and $u = \sigma_\tau[i].word / bts$ if bts is a
 349 suffix of $\sigma_\tau[i].word$ and otherwise $u = \epsilon$.

350 (4) $w' = (bts :: w) \setminus v$.

351 We write $(q, w, \sigma_c, \sigma_p) \rightarrow^* (q', w', \sigma'_c, \sigma'_p)$ for $(q, w, \sigma_c, \sigma_p) \rightarrow \dots \rightarrow (q', w', \sigma'_c, \sigma'_p)$. We give
 352 an intuitive reading of the definition. Roughly, v is the string to be consumed by the step,
 353 and (1) says that if $\ell \in \Sigma_\epsilon$ then ℓ is consumed, and otherwise ℓ is a capturing-group index
 354 (i.e., $\ell \in [k_c]$) and the string stored at the corresponding memory, i.e., $\sigma_c[\ell].word$, is consumed
 355 provided that the memory is closed. (2) and (3) stipulate how the statuses and strings of
 356 the memories are updated, respectively. Importantly, (3) defines the *backtrack string* bts
 357 to be used for backtracking caused by a closure of a positive lookahead memory (if any
 358 happens in the step). Namely, bts is set to be the longest string stored in positive lookahead
 359 memories closed by the step ($bts = \epsilon$ if no such closures happen), and is used in (3) to reset
 360 the content of the memories that are open and remain so (i.e., those satisfying $\mathbf{O} \Rightarrow_{\tau, i} \mathbf{O}$) by
 361 right-division (cf. the definition of u). The string contents remain unchanged for memories
 362 that are or remain closed by the step (i.e., those satisfying $_ \Rightarrow_{\tau, i} \mathbf{C}$), and for the rest of the
 363 memories, the consumed string v is appended to their strings. (4) defines w' , which is the
 364 input string in the post configuration (i.e., the string to be consumed in the continuation of
 365 the step), by prepending bts to the previous configuration's input string w to account for any
 366 backtracking that happens in the step, and consuming v by left-division. We remark that, for
 367 any configuration reachable from an initial configuration (see below), $\sigma_p[i].word$ is a prefix of
 368 $\sigma_p[j].word$ or vice versa for any $i, j \in [k_c]$, thus ensuring that bts is uniquely determined.

369 An *initial configuration* is $(q_0, w, \sigma_{c,0}, \sigma_{p,0})$, where $\sigma_{\tau,0}[i] = (\epsilon, \mathbf{C})$ for all $i \in [k_\tau]$ and
 370 $\tau \in \{c, p\}$. That is, every memory is initially closed and stores the empty string. A *run* of a
 371 PLMFA M is a sequence π such that $\pi[1]$ is an initial configuration and $\pi[i] \xrightarrow[M]{\ell} \pi[i+1]$ for all
 372 $1 \leq i < |\pi|$. A run π is *accepting* if $\pi[|\pi|] = (q, \epsilon, _, _)$ for some $q \in F$. A string w is *accepted*
 373 by M if M has an accepting run. The language of M , denoted by $L(M)$, is the set of strings
 374 accepted by M . That is, $L(M) = \{w \in \Sigma^* \mid (q_0, w, \sigma_{c,0}, \sigma_{p,0}) \rightarrow^* (q, \epsilon, \sigma_c, \sigma_p) \wedge q \in F\}$. We
 375 note that when $k_p = 0$, a PLMFA(k_c, k_p) is a k_c -*memory automaton* (MFA(k_c)) or simply
 376 MFA if k_c is irrelevant) introduced by Schmid [14].

377 ► **Example 16.** As an example, consider a run of the PLMFA M shown in Fig. 3, which
 378 is equivalent to the rewbl_p $(?=({}_1 \cdot *)_1 \mathbf{z}) \cdot \setminus 1$ described in Example 4. In the figure, the
 379 (resp. double) circles represent (resp. accepting) states. The arrows represent transitions
 380 and the words on the labels represent labels on the transitions except for $(?=, ({}_1,)_1,$ and
 381 $)$. The arrow with the word $(?=$ (resp. $)$) represents the transition $\delta(q_1, \epsilon) \ni (q_2, \diamond, \circ)$
 382 (resp. $\delta(q_5, \epsilon) \ni (q_6, \diamond, \mathbf{c})$). Additionally, the arrow with the word $({}_1$ (resp. $)_1$) represents the
 383 transition $\delta(q_2, \epsilon) \ni (q_3, \circ, \diamond)$ (resp. $\delta(q_3, \epsilon) \ni (q_4, \mathbf{c}, \diamond)$). The rewbl_p expression contains just
 384 one backreference and positive lookahead. For simplicity, we abbreviate the lists of memories
 385 $[(x, \mathbf{s})]$ as (x, \mathbf{s}) .



■ **Figure 3** A PLMFA equivalent to the rewbl_p expression $(?=(1.*_1z)·\1.$

386 Given an input string $w=zz$, the run of M is as follows: The run begins with the initial con-
 387 figuration $(q_0, zz, (\epsilon, \mathbf{C}), (\epsilon, \mathbf{C}))$. First, the initial configuration changes to $(q_3, zz, (\epsilon, \mathbf{O}), (\epsilon, \mathbf{O}))$
 388 by applying the transitions $\delta(q_1, \epsilon) \ni (q_2, \diamond, \circ)$ and $\delta(q_2, \epsilon) \ni (q_3, \circ, \diamond)$ in this order. The
 389 transition $\delta(q_1, \epsilon) \ni (q_2, \diamond, \circ)$ opens the positive-lookahead memory and the configuration
 390 changes to $(q_2, zz, (\epsilon, \mathbf{C}), (\epsilon, \mathbf{O}))$. The transition $\delta(q_2, \epsilon) \ni (q_3, \circ, \diamond)$ opens the capturing-
 391 group memory and the configuration changes to $(q_3, zz, (\epsilon, \mathbf{O}), (\epsilon, \mathbf{O}))$. Next, the transitions
 392 $\delta(q_3, \cdot) \ni (q_3, \circ, \diamond)$, $\delta(q_3, \epsilon) \ni (q_4, \mathbf{c}, \diamond)$, and $\delta(q_4, \mathbf{z}) \ni (q_5, \diamond, \diamond)$ are applied in this order.
 393 For the first transition, the configuration changes to the configuration $(q_3, \mathbf{z}, (\mathbf{z}, \mathbf{O}), (\mathbf{z}, \mathbf{O}))$
 394 by consuming a character \mathbf{z} . For the second transition, the configuration changes to the
 395 configuration $(q_4, \mathbf{z}, (\mathbf{z}, \mathbf{C}), (\mathbf{z}, \mathbf{O}))$ by closing the capturing-group memory. For the third
 396 transition, the configuration changes to the configuration $(q_5, \epsilon, (\mathbf{z}, \mathbf{C}), (\mathbf{z}, \mathbf{O}))$ by consuming
 397 a character \mathbf{z} .

398 Then, the configuration changes to $(q_6, zz, (\mathbf{z}, \mathbf{C}), (\mathbf{z}, \mathbf{C}))$ by applying the transition
 399 $(q_5, \epsilon) \ni (q_6, \diamond, \mathbf{c})$. The transition closes the positive-lookahead memory and therefore it
 400 simulates the backtracking behavior of the positive lookahead, i.e., it prepends the word of
 401 the positive-lookahead memory zz to the input string. Finally, the configuration changes to
 402 $(q_8, \epsilon, (\mathbf{z}, \mathbf{C}), (\mathbf{z}, \mathbf{C}))$ by applying the transitions $\delta(q_6, \cdot) \ni (q_7, \diamond, \diamond)$ and $\delta(q_7, \mathbf{1}) \ni (q_8, \diamond, \diamond)$
 403 in this order. The transition $\delta(q_6, \cdot) \ni (q_7, \diamond, \diamond)$ consumes a character \mathbf{z} and the configuration
 404 changes to $(q_7, \mathbf{z}, (\mathbf{z}, \mathbf{C}), (\mathbf{z}, \mathbf{C}))$. Next, the current state q_7 has the transition of the
 405 backreference $\backslash 1$, i.e., $\delta(q_7, \mathbf{1}) \ni (q_8, \diamond, \diamond)$. The transition tries to match the captured string,
 406 i.e., \mathbf{z} , with the current input string. Since the match succeeds, the configuration changes to
 407 $(q_8, \epsilon, (\mathbf{z}, \mathbf{C}), (\mathbf{z}, \mathbf{C}))$. Now, the current state q_8 is an accepting state and the current input
 408 string is ϵ , w is accepted by M .

409 ► **Remark 17.** As seen above, capturing-group memories and positive-lookahead memories
 410 exhibit an interesting *symmetry*: at their use, the content of a capturing-group (resp. positive-
 411 lookahead) memory is left-divided from (resp. prepended to) the input string, and appended
 412 to (resp. right-divided from) the strings stored in memories. The symmetry is imperfect
 413 because positive-lookahead memories do not have "triggers" corresponding to backreferences
 414 of capturing-group memories and a use of a positive-lookahead memory is always synchronous
 415 with its closure. A perfect symmetry can be obtained by extending PLMFA with a new kind
 416 of transitions that trigger positive-lookahead memory uses, disassociating them from closures.
 417 The extension certainly does not decrease the expressive power of PLMFA and we conjecture
 418 that it will strictly increase the expressive power.

419 We define conditions on PLMFA that correspond to RMRL and NUR of rewbl (cf. Sec. 2).
 420 Note that we do not define conditions on PLMFA that correspond to NSR of rewbl because
 421 PLMFAs already satisfy such a condition, i.e., PLMFAs do not allow to refer to the memory

13:12 On Lookaheads in Regular Expressions with Backreferences

422 whose status is open. For convenience, we simply call these conditions RMRL and NUR.
 423 Informally, a PLMFA satisfies RMRL if for all capturing-group memory (index) i and a state
 424 q from which the i th memory can be backreferenced, there exist a unique pair of transitions
 425 a and b that opened and closed the i th memory respectively so that the content of the i th
 426 memory when the computation reaches q is what was recorded between a and b . Intuitively,
 427 the pair of transitions correspond to the capturing group opening (i and closing $)_i$ of rewbl.
 428 Next, we formalize RMRL for PLMFA. For a configuration ϖ and $i \in [k_c]$, let us write
 429 $status(\varpi, i)$ for the status of the i th capturing-group memory of ϖ , i.e., $\sigma_c[i].status$ where
 430 $\varpi = (_, _, \sigma_c, _)$.

431 ► **Definition 18** (Opening and Closing Transitions Pair). For $q \in Q$ and $i \in [k_c]$, a pair of
 432 transitions $(\delta(p', \ell') \ni (q', ir'_c, ir'_p), \delta(p'', \ell'') \ni (q'', ir''_c, ir''_p)) = (a, b)$ is called an *opening-and-*
 433 *closing-transitions pair* of index i at state q if there exist a run π and $1 \leq j_1 < j_2 < |\pi|$ such
 434 that (1) the step from $\pi[j_1]$ to $\pi[j_1 + 1]$ takes the transition a and $status(\pi[j_1], i) = \mathbf{C}$, (2)
 435 the step from $\pi[j_2]$ to $\pi[j_2 + 1]$ takes the transition b and $status(\pi[j_2 + 1], i) = \mathbf{C}$, (3) for all
 436 $j_1 < l \leq j_2$, $status(\pi[l], i) = \mathbf{O}$, and for all $j_2 < l \leq |\pi|$, $status(\pi[l], i) = \mathbf{C}$, and (4) the state
 437 of $\pi[|\pi|]$ is q . We define $RefSet_{M,i}(q)$ (or $RefSet_i(q)$ if there is no danger of ambiguity) as
 438 the set of opening-and-closing-transitions pairs of i at q on M .

439 ► **Definition 19** (RMRL-PLMFA). A PLMFA(k_c, k_p) $M = (Q, \delta, q_0, F)$ is called *restricted*
 440 *may repeat labels* (RMRL) if for all $(q, i) \in Q \times [k_c]$ such that $\delta(q, i) \neq \emptyset$, $|RefSet_i(q)| \leq 1$.

441 Next, we define NUR for PLMFA. Informally, a PLMFA satisfies NUR if no capturing-
 442 group memory can be backreferenced without capturing a word. Formally, for $(q, i) \in$
 443 $Q \times [k_c]$, we say that q is *assigned* with respect to index i on M , written $Assigned_M(q, i)$ (or
 444 $Assigned(q, i)$ if there is no danger of ambiguity), if for all runs π such that the state of $\pi[|\pi|]$
 445 is q , there exists $1 \leq j < |\pi|$ such that $status(\pi[j], i) = \mathbf{O}$ and $status(\pi[j + 1], i) = \mathbf{C}$.

446 ► **Definition 20** (NUR-PLMFA). A PLMFA(k_c, k_p) $M = (Q, \delta, q_0, F)$ is *no unassigned*
 447 *reference* (NUR) if for all $(q, i) \in Q \times [k_c]$ such that $\delta(q, i) \neq \emptyset$, $Assigned(q, i) = true$.

448 In what follows, we assume that PLMFAs satisfy RMRL and NUR.

4.2 Normal Forms and Nested Forms

449 We show that a PLMFA can be converted into certain forms. The *normal form* enforces two
 450 restrictions: (1) only ϵ transitions can change memory statuses and at most one status of
 451 the memory at a time, and (2) no transitions open (resp. close) a memory that is already
 452 opened (resp. closed).
 453

454 ► **Definition 21** (Normal Form). A PLMFA is in *normal form* if the following properties are
 455 satisfied. For every transition $\delta(q, \ell) \ni (q', ir'_c, ir'_p)$, $\tau \in \{c, p\}$, and $j \in [k_\tau]$, (1) if $ir_\tau[j] \neq \diamond$,
 456 then $\ell = \epsilon$ and $ir_{\tau'}[l].status = \diamond$ for all $(\tau', l) \in \{c, p\} \times [k_{\tau'}]$ such that $(\tau', l) \neq (\tau, j)$, and
 457 (2) there is no run π such that $\pi[|\pi|] = (_, _, \sigma_c, \sigma_p)$ where $ir_\tau[j] = \circ$ and $\sigma_\tau[j].status = \mathbf{O}$
 458 or $ir_\tau[j] = \mathbf{c}$ and $\sigma_\tau[j].status = \mathbf{C}$.

459 ► **Lemma 22.** *Any PLMFA M can be converted to a normal form PLMFA M' such that*
 460 $L(M) = L(M')$.

461 The proof is by adopting an analogous conversion of [14] and works by extending the states
 462 of M to record memory statuses and splitting simultaneous memory updates to multiple
 463 transitions. We remark that the conversion preserves RMRL and NUR.

464 Next, we define *nested form* which enforces that there are no *overlaps* in any runs. A
 465 run π is said to have an overlap if there exist $1 \leq j_1 < j_2 < j_3 < j_4 < |\pi|$ such that
 466 some memory is opened and closed respectively at the j_1 th and the j_3 th steps and some
 467 memory (possibly the same) is opened and closed respectively at the j_2 th step and the j_4 th
 468 step. There are four types of overlaps, *cc*, *cp*, *pc*, and *pp*, depending on the types of the
 469 first and the second memories (e.g., *cp*-overlap is when the first memory is capturing-group
 470 and the second is positive-lookahead). Intuitively, an overlap corresponds to an invalid
 471 expression that has an overlap of capturing groups or positive lookaheads. For example,
 472 *cc*-overlap corresponds to invalid expressions $({}_i r_1 ({}_j r_2) {}_i r_3)_j$ and *pc*-overlap corresponds to
 473 invalid expressions $(? = r_1 ({}_i r_2) r_3)_i$.

474 ► **Definition 23** (Nested Form). A PLMFA M is in *nested form* if there are no overlaps in
 475 any runs on M .

476 We show that we can transform a PLMFA to the nested form.

477 ► **Lemma 24.** *Any PLMFA M can be converted to a normal and nested form PLMFA M'
 478 such that $L(M) = L(M')$.*

479 The proof is by adding new transitions that close the fragments of the memories which
 480 are open before opening the transitions that cause overlaps and open the next fragments of
 481 the memories after that. For *cc*-overlaps, the conversion coincides with the analogous one
 482 for MFAs [14]. In what follows, without loss of generality, we assume that PLMFAs are in
 483 normal and nested form.

484 4.3 From rewbl_p to PLMFA

485 We show that given a rewbl_p expression r , we can construct a PLMFA M such that
 486 $L(r) = L(M)$. The construction of the PLMFA extends the standard Thompson construction
 487 from a pure regular expression to an NFA [15] with backreferences and lookaheads in a
 488 mostly straightforward manner. For space, the construction and the proof of correctness is
 489 omitted.

490 ► **Theorem 25.** *Let a rewbl_p r include k_c capturing groups and k_p positive lookaheads. Then,
 491 there exists a PLMFA (k_c, k_p) M such that $L(r) = L(M)$.*

492 4.4 From PLMFA to rewbl_p

493 Now, we show that given a PLMFA M , we can construct a rewbl_p expression r such that
 494 $L(M) = L(r)$. We first give the conversion and then show the correctness. For space, we
 495 only show the correctness of the language equivalence and omit the fact that the conditions
 496 NUR, NSR, and RMRL are satisfied by the resulting rewbl_p . However, they can be proved
 497 similarly to the language equivalence.

498 The conversion, referred to as *PtoR*, is inspired by that of MFAs to rewbs [14]. The idea
 499 of the conversion is to use the nested relation of PLMFAs. Since PLMFAs are nested form,
 500 the transitions of the opening-and-closing-transitions pairs have a nested structure, i.e., they
 501 form a directed acyclic graph (DAG). Therefore, we can iteratively convert (sub)automaton
 502 corresponding to the part of the given PLMFA delimited by each such pairs to the rewbl_p
 503 expression in a topological order starting from the pairs that nest nothing. Each step of
 504 the conversion makes an *extended PLMFA* (ePLMFA) whose labels are rewbl_p expressions.
 505 That is, labels ℓ on transitions of PLMFAs are treated as rewbl_p expressions ℓ of ePLMFAs if

13:14 On Lookaheads in Regular Expressions with Backreferences

506 $\ell \in \Sigma_\epsilon$. Additionally, if $\ell = i \in [k_c]$ on the transitions of PLMFAs, $\ell = \setminus i$ on the transitions
 507 of ePLMFAs. For an ePLMFA M and a rewbl_p expression ℓ , a step of computation of
 508 M reading ℓ is defined as follows: $(q, w, \sigma_c, \sigma_p) \xrightarrow{\ell} (q', w', \sigma'_c, \sigma'_p)$ iff either $\ell = \epsilon$ and
 509 $(q, w, \sigma_c, \sigma_p) \xrightarrow{\epsilon} (q', w', \sigma'_c, \sigma'_p)$ according to Def. 15, or $\ell \neq \epsilon$ and there exist a transition
 510 $\delta(q, \ell) \ni (q', \diamond^{k_c+k_p})$ and steps of computations from $(q'_0, w, \sigma_c, \sigma_p)$ to $(q'_f, w', \sigma'_c, \sigma'_p)$ of
 511 $M' = (_, _, q'_0, \{q'_f\})$ obtained from ℓ by the construction from a rewbl_p expression to a
 512 PLMFA mentioned in Sec. 4.3.

513 We also extend the NUR and RMRL conditions for ePLMFAs as follows. For all transition,
 514 we reconstruct PLMFA from the rewbl_p label by applying the construction mentioned in
 515 Sec. 4.3 and replace the transition with the PLMFA by replacing it with ϵ -labeled \diamond -only
 516 transitions to and from the initial and the final state of the PLMFA. We say that an ePLMFA
 517 satisfies NUR and RMRL if the reconstructed PLMFA satisfies NUR and RMRL, respectively.

518 We define the nesting relation. Firstly, for an ePLMFA and $\text{inst} \in \{\text{o}, \text{c}\}$, we define
 519 $Q_{\tau, i, \text{inst}}$ as the set of states q such that $\delta(q, _) \ni (_, ir_c, ir_p)$ where $ir_\tau[i] = \text{inst}$. Let
 520 $\Phi = \{(\tau, i, q, q') \mid \tau \in \{c, p\} \wedge i \in [k_\tau] \wedge q \in Q_{\tau, i, \text{o}} \wedge q' \in Q_{\tau, i, \text{c}}\}$. The *nesting relation*
 521 $\prec \subseteq \Phi \times \Phi$ is defined as follows: $(\tau_1, i_1, q_1, q'_1) \prec (\tau_2, i_2, q_2, q'_2)$ iff there exist a run π
 522 and $s < t < u < v < \pi[|\pi|]$ such that the step from $\pi[s]$ to $\pi[s+1]$ takes a transition
 523 $\delta(q_2, _) \ni (_, ir_c, ir_p)$ where $ir_{\tau_2}[i_2] = \text{o}$, the step from $\pi[t]$ to $\pi[t+1]$ takes a transition
 524 $\delta(q_1, _) \ni (_, ir_c, ir_p)$ where $ir_{\tau_1}[i_1] = \text{o}$, the step from $\pi[u]$ to $\pi[u+1]$ takes a transition
 525 $\delta(q'_1, _) \ni (_, ir_c, ir_p)$ where $ir_{\tau_1}[i_1] = \text{c}$, and the step from $\pi[v]$ to $\pi[v+1]$ takes a transition
 526 $\delta(q'_2, _) \ni (_, ir_c, ir_p)$ where $ir_{\tau_2}[i_2] = \text{c}$.

527 For an ePLMFA M , the procedure of $PtoR(M)$ is defined as follows. Let us initialize
 528 $\Delta : \Phi \rightarrow \mathcal{P}(\Phi)$ as follows: $\Delta(\tau_1, i_1, q_1, q'_1) = \{(\tau_2, i_2, q_2, q'_2) \mid (\tau_2, i_2, q_2, q'_2) \prec (\tau_1, i_1, q_1, q'_1)\}$.
 529 We iteratively update M , Φ , and Δ by the following steps.

- 530 1. Find $(\tau, i, q, q') \in \Phi$ such that $\Delta(\tau, i, q, q') = \emptyset$. Let $\delta(q, _) \ni (q_s, _)$ (resp. $\delta(q', _) \ni$
 531 $(q_e, _)$) be the opening (resp. closing) transition of (τ, i, q, q') . Then, construct an ePLMFA
 532 M' from M by replacing the initial state and set of accepting states with q_s and $\{q'\}$,
 533 respectively, and deleting all transitions that open or close memories.
- 534 2. Convert M' to a rewbl_p expression r using the standard state elimination method.
- 535 3. Delete (τ, i, q, q') from Φ and (the domain and the range of) Δ , and add the transition
 536 $\delta(q, (ir)_i) \ni (q_e, \diamond^{k_c+k_p})$ if $\tau = c$ and otherwise $\delta(q, (?=r)) \ni (q_e, \diamond^{k_c+k_p})$ to M .
- 537 4. Repeat steps 1 to 3 until convergence.
- 538 5. Delete all transitions that open or close memories from M .
- 539 6. Convert M to a rewbl_p expression r by the state elimination method and return r .

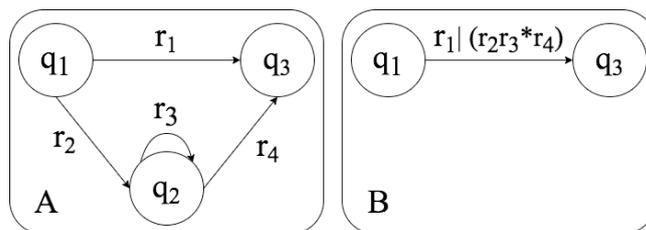
540 At step 1, we can find such a tuple (τ, i, q, q') since it is in nested form. The state elimination
 541 method used in steps 2 and 6 is a straightforward adoption of the standard state elimination
 542 method (see, e.g., [13]) that interprets the rewbl_p expression that appear as labels as ordinary
 543 regular expressions.

544 We proceed to the proof of correctness.

545 ► **Definition 26.** The *language* of an ePLMFA $M = (Q, \delta, q_0, F)$ parameterized by a string
 546 $w \in \Sigma^*$, a capturing-group memory σ_c , and a positive-lookahed memory σ_p , denoted by
 547 $L(M, w, \sigma_c, \sigma_p)$, is defined as

$$548 \quad L(M, w, \sigma_c, \sigma_p) = \{(w', \sigma'_c, \sigma'_p) \mid \exists q' \in F. (q_0, w, \sigma_c, \sigma_p) \rightarrow^* (q', w', \sigma'_c, \sigma'_p)\}.$$

549 For an ePLMFA $M = (Q, \delta, q_0, F)$, we write $M(q', F')$ for the ePLMFA $M' = (Q, \delta, q', F')$
 550 where $q' \in Q$ and $F' \subseteq Q$. We show that eliminating a state from an ePLMFA by applying
 551 one step of the state elimination method does not change the parameterized language of the



■ **Figure 4** (A) Before eliminating the state q_2 . (B) After eliminating the state q_2 .

552 ePLMFA. The state elimination method eliminates a state by deleting and adding transitions
553 as shown in Fig. 4.

554 ► **Lemma 27.** *Let $M = (Q, \delta, q_0, F)$ be an ePLMFA (k_c, k_p) such that for all transitions
555 $\delta(q, \ell) \ni (q', ir_c, ir_p)$, $i \in [k_c]$, and $j \in [k_p]$, $ir_c[i] = \diamond$ and $ir_p[j] = \diamond$. Additionally, let
556 $M' = (Q', \delta', q'_0, F')$ be the ePLMFA obtained by eliminating a state in Q from M by the state
557 elimination method. Then, for all $w \in \Sigma^*$, σ_c , and σ_p , $L(M, w, \sigma_c, \sigma_p) = L(M', w, \sigma_c, \sigma_p)$.*

558 **Proof.** Let us assume that q_1 and q_3 are in Q and Q' , $q_2 \in Q$ is the state eliminated
559 by the state elimination method as shown in Fig. 4. Since the only difference between
560 M and M' comes from the state q_2 , it suffices to show that $L(M(q_1, \{q_3\}), w, \sigma_c, \sigma_p) =$
561 $L(M'(q_1, \{q_3\}), w, \sigma_c, \sigma_p)$. It is immediate from the construction. ◀

562 ► **Theorem 28.** *For an ePLMFA M , let M_j be the ePLMFA obtained after the j th iteration
563 of steps 1 to 3 of PtoR(M). We assume $M_j = M$ if $j = 0$. Then, for all $w \in \Sigma^*$, σ_c , and
564 σ_p , $L(M, w, \sigma_c, \sigma_p) = L(M_j, w, \sigma_c, \sigma_p)$.*

565 **Proof.** The proof is by induction on the number of the iteration j of step 1 to 3. The
566 base step $j = 0$ is immediate since $M_j = M$. For the induction step, let $j > 0$. The
567 inductive hypothesis is that $L(M, w, \sigma_c, \sigma_p) = L(M_j, w, \sigma_c, \sigma_p)$ holds for all j , w , σ_c , and
568 σ_p . We then consider M_{j+1} . We show $L(M_j, w, \sigma_c, \sigma_p) = L(M_{j+1}, w, \sigma_c, \sigma_p)$. As described
569 in the step 3, for the ePLMFA $M_j = (Q, \delta, q_0, F)$, $M_{j+1} = (Q, \delta \cup \{t\}, q_0, F)$ where $t =$
570 $\delta(q, r') \ni (q_e, \diamond^{k_c+k_p})$. Recall that $r' = (ir)_i$ or $(?=r)$ and r is the rewbl $_p$ expression
571 obtained at step 2. For this, $L(M_j, w, \sigma_c, \sigma_p) \subseteq L(M_{j+1}, w, \sigma_c, \sigma_p)$ is immediate. For
572 $L(M_{j+1}, w, \sigma_c, \sigma_p) \subseteq L(M_j, w, \sigma_c, \sigma_p)$, it suffices to show that for all w , σ_c , and σ_p , if
573 there exists $(q, w, \sigma_c, \sigma_p) \xrightarrow{r'} (q_e, w', \sigma'_c, \sigma'_p)$ on M_{j+1} , then there exists $(q, w, \sigma_c, \sigma_p) \rightarrow^*$

574 $(q_e, w', \sigma'_c, \sigma'_p)$ on M_j . We assume that there exists $\pi = (q, w, \sigma_c, \sigma_p) \xrightarrow{r'} (q_e, w', \sigma'_c, \sigma'_p)$
575 on M_{j+1} . The computation of the transition whose label is r' is defined by that of the
576 ePLMFA $M_{r'}$ obtained from r' by the construction mentioned in Sec. 4.3. If $r' = (ir)_i$,
577 $M_{r'} = (Q_{r'}, \delta_{r'}, q, \{q_e\})$ where

- 578 ■ $Q_{r'} = \{q, q_e\} \cup Q_r$; and
- 579 ■ $\delta_{r'} = \delta_r \cup \{((q, \epsilon), \{(q_{0,r}, ir_c, \diamond^{k_p})\})\} \cup \{((q_{F,r}, \epsilon), \{(q_e, ir'_c, \diamond^{k_p})\})\}$ where $ir_c[k] = \circ$ and
580 $ir'_c[k] = \circ$ if $k = i$ and otherwise $ir_c[k] = \diamond$ and $ir'_c[k] = \diamond$ for $k \in [k_c]$

581 with the ePLMFA for r $M_r = (Q_r, \delta_r, q_{0,r}, \{q_{F,r}\})$. By the transitions in $\delta_{r'}$, $\pi = (q, w, \sigma_c, \sigma_p) \rightarrow^*$
582 $(q_{0,r}, w, \sigma''_c, \sigma_p) \xrightarrow{r} (q_{F,r}, w', \sigma'''_c, \sigma'_p) \rightarrow (q_e, w', \sigma'_c, \sigma'_p)$ where $\sigma''_c[k] = (\epsilon, \mathbf{O})$ and $\sigma'''_c[k] =$
583 $(\sigma'_c[k].word, \mathbf{O})$ if $k = i$ and otherwise $\sigma''_c[k] = \sigma_c[k]$ and $\sigma'''_c[k] = \sigma'_c[k]$. We focus on
584 the computation $(q_{0,r}, w, \sigma''_c, \sigma_p) \xrightarrow{r} (q_{F,r}, w', \sigma'''_c, \sigma'_p)$. The label r is constructed from M'
585 by the state elimination method at step 2. Additionally, by Lemma 27, the state elim-
586 ination method preserves the parameterized language equivalence. For this, there exists

13:16 On Lookaheads in Regular Expressions with Backreferences

587 $(q_s, w, \sigma_c'', \sigma_p) \rightarrow^* (q', w', \sigma_c''', \sigma_p')$ on M' . Since the transition function of M_j includes all
588 transitions of M' , there also exists $(q_s, w, \sigma_c'', \sigma_p) \rightarrow^* (q', w', \sigma_c''', \sigma_p')$ on M_j . As described
589 at step 1, M_j has a transition from q to q_s that opens the i th capturing-group memory
590 and a transition from q' to q_e that closes the i th capturing-group memory. Therefore,
591 there exists $(q, w, \sigma_c, \sigma_p) \rightarrow (q_s, w, \sigma_c'', \sigma_p) \rightarrow^* (q', w', \sigma_c''', \sigma_p') \rightarrow (q_e, w', \sigma_c', \sigma_p')$. For this,
592 $L(M_j, w, \sigma_c, \sigma_p) = L(M_{j+1}, w, \sigma_c, \sigma_p)$. By the inductive hypothesis, $L(M_j, w, \sigma_c, \sigma_p) =$
593 $L(M, w, \sigma_c, \sigma_p)$. Thus, $L(M, w, \sigma_c, \sigma_p) = L(M_{j+1}, w, \sigma_c, \sigma_p)$ for the case $r' = ({}_i r)_i$. The case
594 of $r' = (?=r)$ is analogous. ◀

595 Now, we obtain our main result.

596 ▶ **Theorem 29.** *Let M be an ePLMFA. Then, $L(M) = L(r)$ where $r = PtoR(M)$.*

597 **Proof.** Let M' be the ePLMFA at the last step of the state elimination method at step 6 of
598 $PtoR$. Then, $M' = (\{q, q'\}, \{((q, r), \{(q', \diamond^{k_c+k_p})\})\}, q, \{q'\})$. By Theorem 25, $L(M') = L(r)$.
599 By Theorem 28, for all $w \in \Sigma^*$, $L(M', w, \sigma_{c,0}, \sigma_{p,0}) = L(M, w, \sigma_{c,0}, \sigma_{p,0})$, i.e., $L(M') = L(M)$.
600 Thus, $L(M) = L(M') = L(r)$. ◀

601 As a corollary of Theorem 25 and Theorem 29, we obtain the following result.

602 ▶ **Corollary 30.** *The expressive power of PLMFA is equivalent to that of $rewbl_p$.*

5 Related Work

604 Among the major extensions employed by modern real-world regular expression engines are
605 backreferences and lookaheads. However, the previous works on formal language theory have
606 studied the two features mostly in isolation, and to the best of our knowledge, our work is
607 the first formal study of regular expressions extended with both features. Next, we discuss
608 previous works that studied the features in isolation.

609 Prior works by Morihata and Berglund et al. [12, 3] showed that extending regular
610 expressions by lookaheads does not enhance their expressive power. Their proofs are by
611 a translation to boolean finite automata [4] whose expressive power is regular. However,
612 adopting such an approach to defining an equivalent automata is difficult in the presence
613 of backreferences because boolean automata express lookaheads by running several states
614 simultaneously without backtracking, while the combination of lookaheads and backreferences
615 intrinsically requires backtracking. For example, to match against $(?= (\cdot)_1) \setminus 1$, a boolean
616 approach would run $(?= (\cdot)_1)$ and $\setminus 1$ simultaneously, but then the automaton would get
617 stuck while trying to process the backreference $\setminus 1$ as it is unassigned at that point. By
618 contrast, our PLMFA uses the novel positive-lookahead memories to store enough information
619 to simulate the backtracking behavior of positive lookaheads.

620 A formal study of regular expressions with backreferences (rewb) dates back to the
621 seminal work by Aho [1]. More recently, a formal semantics and a pumping lemma were
622 given by Câmpeanu et al. [5]. Berglund and van der Merwe [2] showed that different variants
623 of backreference semantics give rise to differences in expressive powers. Our work adopts
624 and formalizes the no-label-repetitions (NLR) and no-unassigned-reference (NUR) semantics
625 which is also used in [5, 6]. Schmid [14] proposed MFA, and showed that the expressive power
626 of the automata is equivalent to that of rew b. Our PLMFA builds on MFA and extends
627 it with positive-lookahead memories to handle positive lookaheads. As remarked before
628 (cf. Remark 17), our positive-lookahead memory exhibits an interesting symmetry to the
629 capturing-group memory of MFA.

6 Conclusion

We have studied the expressive powers of regular expressions with the popular backreferences and lookaheads extensions. We have shown that extending rewb by positive or negative lookaheads enhance their expressive power. Additionally, we have presented language-theoretic properties of rewb extended by the two forms of lookaheads, and have presented a new class of automata called PLMFA that is equivalent in expressive power to rewbl_p . We have introduced a new kind of memories called a positive-lookahead memory, which is almost perfectly symmetric to capturing-group memory of MFA, as a key component of PLMFA.

Despite the popularity of the backreference and lookaheads extensions in practice, to our knowledge, our work is the first formal study on regular expressions with both extensions. We hope that our results pave the way for more work on the topic.

References

- 1 Alfred V. Aho. Chapter 5 - algorithms for finding patterns in strings. In Jan van Leeuwen, editor, *Algorithms and Complexity*, Handbook of Theoretical Computer Science, pages 255–300. Elsevier, Amsterdam, 1990. URL: <https://www.sciencedirect.com/science/article/pii/B9780444880710500102>, doi:<https://doi.org/10.1016/B978-0-444-88071-0.50010-2>.
- 2 Martin Berglund and Brink van der Merwe. Regular expressions with backreferences re-examined. In Jan Holub and Jan Zdárek, editors, *Proceedings of the Prague Stringology Conference 2017, Prague, Czech Republic, August 28-30, 2017*, pages 30–41. Department of Theoretical Computer Science, Faculty of Information Technology, Czech Technical University in Prague, 2017. URL: <http://www.stringology.org/event/2017/p04.html>.
- 3 Martin Berglund, Brink van der Merwe, and Steyn van Litsenborgh. Regular expressions with lookahead. *J. Univers. Comput. Sci.*, 27(1):324–340, 2021. doi:10.3897/jucs.66330.
- 4 Janusz A. Brzozowski and Ernst L. Leiss. On equations for regular languages, finite automata, and sequential networks. *Theor. Comput. Sci.*, 10:19–35, 1980. URL: <http://dblp.uni-trier.de/db/journals/tcs/tcs10.html#BrzozowskiL80>.
- 5 Cezar Câmpeanu, Kai Salomaa, and Sheng Yu. A formal study of practical regular expressions. *International Journal of Foundations of Computer Science*, 14(06):1007–1018, 2003. arXiv: <https://doi.org/10.1142/S012905410300214X>, doi:10.1142/S012905410300214X.
- 6 Benjamin Carle and Paliath Narendran. On extended regular expressions. In Adrian-Horia Dediu, Armand-Mihai Ionescu, and Carlos Martín-Vide, editors, *Language and Automata Theory and Applications, Third International Conference, LATA 2009, Tarragona, Spain, April 2-8, 2009. Proceedings*, volume 5457 of *Lecture Notes in Computer Science*, pages 279–289. Springer, 2009. doi:10.1007/978-3-642-00982-2_24.
- 7 Ashok K. Chandra and Larry J. Stockmeyer. Alternation. In *17th Annual Symposium on Foundations of Computer Science (sfcs 1976)*, pages 98–108, 1976. doi:10.1109/SFCS.1976.4.
- 8 N. Chida and T. Terauchi. Repairing dos vulnerability of real-world regexes. In *Proceedings of the 43rd IEEE Symposium on Security and Privacy, SP 2022*, pages 1049–1066, Los Alamitos, CA, USA, may 2022. IEEE Computer Society. URL: <https://doi.ieeecomputersociety.org/10.1109/SP46214.2022.00061>, doi:10.1109/SP46214.2022.00061.
- 9 S. C. Kleene. Representation of events in nerve nets and finite automata. In *Automata Studies*. 1956.
- 10 Vladimir Komendantsky. Matching problem for regular expressions with variables. In Hans-Wolfgang Loidl and Ricardo Peña, editors, *Trends in Functional Programming - 13th International Symposium, TFP 2012, St. Andrews, UK, June 12-14, 2012, Revised Selected Papers*, volume 7829 of *Lecture Notes in Computer Science*, pages 149–166. Springer, 2012. doi:10.1007/978-3-642-40447-4_10.
- 11 Blake Loring, Duncan Mitchell, and Johannes Kinder. Sound regular expression semantics for dynamic symbolic execution of javascript. In *Proceedings of the 40th ACM SIGPLAN*

13:18 On Lookaheads in Regular Expressions with Backreferences

- 679 *Conference on Programming Language Design and Implementation, PLDI 2019*, pages 425–438,
680 New York, NY, USA, 2019. Association for Computing Machinery. doi:10.1145/3314221.
681 3314645.
- 682 **12** Akimasa Morihata. Translation of regular expression with lookahead into finite state automaton.
683 *Computer Software*, 29(1):1_147–1_158, 2012. doi:10.11309/jssst.29.1_147.
- 684 **13** Jacques Sakarovitch. *Elements of Automata Theory*. Cambridge University Press, USA, 2009.
- 685 **14** Markus L. Schmid. Characterising REGEX languages by regular languages equipped with
686 factor-referencing. *Inf. Comput.*, 249:1–17, 2016. doi:10.1016/j.ic.2016.02.003.
- 687 **15** Ken Thompson. Programming techniques: Regular expression search algorithm. *Commun.*
688 *ACM*, 11(6):419–422, June 1968. doi:10.1145/363347.363387.