

# Measuring the Expressive Power of Practical Regular Expressions by Classical Stacking Automata Models<sup>★</sup>

Taisei Nogami<sup>a,\*</sup>, Tachio Terauchi<sup>a,\*</sup>

<sup>a</sup>Waseda University, 3-4-1 Okubo, Shinjuku-ku, Tokyo 169-8555, Japan

---

## Abstract

A *rewb* is a regular expression extended with a feature called backreference. It is broadly known that backreference is a practical extension of regular expressions, and is supported by most modern regular expression engines, such as those in the standard libraries of Java, Python, and more. Meanwhile, *indexed languages* are the languages generated by indexed grammars, a formal grammar class proposed by A.V.Aho. We show that these two models' expressive powers are related in the following way: every language described by a rew b is an indexed language. As the smallest formal grammar class previously known to contain rewbs is the class of context sensitive languages, our result strictly improves the known upper-bound. Moreover, we prove the following four claims: (1) there exists a rew b whose language does not belong to the class of stack languages, which is a proper subclass of indexed languages, (2) the language described by a rew b without a captured reference is in the class of nonerasing stack languages, which is a proper subclass of stack languages, (3) there exists a rew b that describes a stack language but not a nonerasing stack language, and (4) a rew b extended with another practical extension called lookaheads can describe a non-indexed language. Finally, we show that the hierarchy investigated in a prior study, which separates the expressive power of rewbs by the notion of nested levels, is within the class of nonerasing stack languages.

**Keywords:** Regular expressions, Backreferences, Lookaheads, Expressive power

---

## 1. Introduction

A *rewb* is a regular expression empowered with a certain extension, called *backreference*, that allows preceding substrings to be used later. It is closer to practical regular expressions than the pure ones, and supported by the standard libraries of most modern programming languages. A typical example of a rew b follows:

**Example 1.** Let  $\Sigma$  be the alphabet  $\{a, b\}$ . The language  $L(\alpha)$  described by the rew b  $\alpha = ({}_1(a+b)^*)_1 \setminus 1$  is  $\{ww \mid w \in \Sigma^*\}$ . Intuitively,  $\alpha$  first *captures* a preceding string  $w \in L((a+b)^*)$  by  $({}_1)_1$ , and second *references* that  $w$  by following  $\setminus 1$ . Therefore,  $\alpha$  matches  $ww$ . Because this  $L(\alpha)$  is a textbook example of a non-context-free language (and therefore non-regular), the expressive power of rewbs exceeds that of the pure ones.

In 1968, A.V.Aho discovered indexed languages with characterizations by two equivalent models: indexed grammars and  $(1N^1)$  nested stack automata (Nested-SA) [1, 2]. The class of indexed languages is a proper superclass of context free languages (CFL), and a proper subclass of context sensitive languages (CSL) [1].

Berglund and van der Merwe [4], and Câmpeanu et al. [6] have shown that the class of rewbs is incomparable with the class of CFLs and is a proper subclass of CSLs. As the first main contribution of this paper, we prove that the language described by a rew b is an indexed language. Since the class of CSLs was the previously known best upper-bound of rewbs, our result gives a novel and strictly tighter upper-bound.

---

<sup>★</sup>This work was supported by JSPS KAKENHI Grant Numbers JP20H04162, JP20K20625, and JP22H03570.

<sup>\*</sup>Corresponding authors.

<sup>1</sup>One-way nondeterministic. “One-way” means that the input cursor will not move back to left. The antonym is “two-way.”

Meanwhile, there is a class of the languages called stack languages [13, 12]. This class corresponds to the model (1N) stack automata (SA), a restriction of Nested-SA. Hence, it trivially follows that the class of stack languages is a subclass of indexed languages. Actually, this containment is known to be proper [2]. Furthermore, a model called nonerasing stack automata (NESA) has been studied in papers such as [13, 16, 20], and its language class is known to be a proper subclass of stack languages [20].

In this paper, we show that every rew without a captured reference (that is, one in which no reference  $\backslash i$  appears as a subexpression of an expression of the form  $(_j \alpha)_j$ ) describes a nonerasing stack language. Given our result, the following question is natural: does every rew describe a (nonerasing) stack language? We show that the answer is no. Namely, we show a rew that describes a non-stack language. Simultaneously, we also show a rew that describes a stack language but not a nonerasing stack language.<sup>2</sup>

Additionally, we investigate the relationship between stacking automata models and rew extended with another popular practical extension called *lookaheads*.<sup>3</sup> Namely, does every rew with lookaheads describe an indexed language? We again answer the question negatively.

Finally, Larsen [17] has proposed a notion called *nested levels* of a rew and showed that they give rise to a concrete increasing hierarchy of expressive powers of rews by exhibiting, for each nested level  $i \in \mathbb{N}$ , a language  $L_i$  that is expressible by a rew at level  $i$  but not at any levels below  $i$ . We show that this hierarchy is within the class of nonerasing stack languages, that is, there exists an NESA  $A_i$  recognizing  $L_i$  for every nested level  $i$ . Below, we summarize the main contributions of the paper.

- (a) Every rew describes an indexed language. (Section 4, Corollary 24)
- (b) Every rew without a captured reference describes a nonerasing stack language. (Section 4, Corollary 25)
- (c) There exists a rew that describes a non-stack language. (Section 5, Theorem 26 and Corollary 31)
- (d) There exists a rew that describes a stack language but not a nonerasing stack language. (Section 5, Corollary 33)
- (e) A rew with lookaheads can describe a non-indexed language. (Section 6, Corollary 39)
- (f) The hierarchy given by Larsen [17] is within the class of nonerasing stack languages. (Section 7, Theorem 40)

Note that by (b) and (c), it follows that there is a rew language that no rew without a captured reference can describe (Section 5, Corollaries 27 and 32). See also Figure 4 for a summary of the results.

The rest of the paper is organized as follows. Section 2 discusses related work. Section 3 defines preliminary notions used in the paper such as the syntax and semantics of rew, SA, NESA, and Nested-SA. Sections 4, 5, 6, and 7 formally state and prove the paper's main contributions listed above. Section 8 concludes the paper with a discussion on future work.

This paper is an extended version of the conference paper [19] and contains the following additional contributions:

1. We present a more natural example of a rew that describes a non-stack language. The non-stack rew  $\alpha_0 = ((_1 \backslash 4 a)_1 (_2 \backslash 3)_2 (_3 \backslash 2 a)_3 (_4 \backslash 1 \backslash 3)_4)^*$  given in the conference version deeply relied on mutual backreferences and was somewhat artificial. We present a new non-stack rew that does not depend on such a complex structure. Furthermore, we show that a slight modification of the rew yields the *first* rew whose language belongs to the stack languages but not to the nonerasing stack languages.
2. We have shown that the language class of regular expressions extended with backreferences is contained in the class of indexed languages, but does the containment still hold if they are further extended with lookaheads? We settle this question negatively, that is, the languages described by rews extended with lookaheads are not always indexed languages.

## 2. Related work

First, we discuss related work on rews. There are several variants of the syntax and semantics of rews since they first appeared in the seminal work by Aho [3]. A recent study by Berglund and van der Merwe [4] summarizes the variants and the relations between them. In sum, there are two variants of the syntax, whether or not a

<sup>2</sup>Note that this language witnesses the separation between the classes of stack languages and nonerasing stack languages, but this separation has been already shown by Ogden [20].

<sup>3</sup>See Example 36 for an example of lookaheads.

same label may appear as the index of more than one capture (“may repeat labels”, “no label repetitions”), and two variants of the semantics, whether an unbound reference is interpreted as the empty string or an undefined factor ( $\varepsilon$ -semantics,  $\emptyset$ -semantics). As shown in [4], there is no difference in the expressive powers between these two semantics under the “may repeat labels” syntax (therefore, there are three classes with different expressive powers, namely “no label repetitions” with  $\emptyset$ -semantics, “no label repetitions” with  $\varepsilon$ -semantics, and “may repeat labels”). In this paper, we focus on the “may repeat labels” formalization, which has the highest expressive power of the three and is often studied in formal language theory. We adopt the  $\varepsilon$ -semantics as the semantics of rewbs. Note that the pioneering formalization of rewbs given by Aho [3] has the equivalent expressive power as this class. The rewbs with “may repeat labels” with  $\varepsilon$ -semantics was recently proposed by Schmid with the notions of a ref-word and a dereferencing function [21]. Simultaneously, he proposed a class of automata called *memory automata* (MFA), and showed that its expressive power is equivalent to that of rewbs. Freydenberger and Schmid extended MFA to *MFA with trap-state* [10]. Berglund and van der Merwe [4] showed that the class of Schmid’s rewbs is a proper subclass of CSLs, and is incomparable with the class of CFLs. Note that there is a pumping lemma for the formalization given by Cămpăneanu et al. [6] but it is known not to work for Schmid’s rewbs. As mentioned above, Larsen introduced the notion of nested levels and showed that increase in the levels increases the expressive powers of rewbs [17]. As for lookaheads, Morihata [18] and Berglund et al. [5] showed that regular expressions with lookaheads are still regular, but Chida and Terauchi [8] showed that rewbs with lookaheads are strictly more powerful than rewbs.

Next, we discuss related work on the three automata models used throughout the paper, namely SA, NESA, and Nested-SA. Ginsburg et al. introduced SA as a mathematical model that is more powerful than pushdown automaton (PDA), and NESA as a restricted version of SA [13]. Hopcroft and Ullman discovered a type of Turing machine corresponding to the class of two-way NESA [16]. Ogden proposed a pumping lemma for stack languages and non-erasing stack languages [20]. Aho proposed Nested-SA with a proof of the fact that (1N)Nested-SA and indexed grammars given by himself in [1] are equivalent in their expressive powers, and recognized PDA and SA as special cases of Nested-SA [2]. Aho also showed that the class of indexed languages is a proper superclass of CFLs, and a proper subclass of CSLs [1]. Hayashi proposed a pumping lemma for indexed languages [14].

### 3. Preliminaries

#### 3.1. Syntax and semantics of rewbs

In this section, we formalize the syntax and the semantics of rewbs following the formalization given in [10]. We begin with the syntax. Let  $\Sigma_\varepsilon = \Sigma \uplus \{\varepsilon\}$  and  $[k] = \{1, 2, \dots, k\}$ , where the symbol  $\uplus$  denotes a disjoint union.

**Definition 2.** For each natural number  $k \geq 1$ , the set of  $k$ -rewbs over  $\Sigma$ , also written  $k$ -rewb by abuse of notation, and the mapping  $\text{var} : k\text{-rewb} \rightarrow \mathcal{P}([k])$  are defined as follows, where  $a \in \Sigma_\varepsilon$  and  $i \in [k]$ :

$$\begin{aligned} (\alpha, \text{var}(\alpha)) ::= & (a, \emptyset) \mid (\backslash i, \{i\}) \mid (\alpha_0 \alpha_1, \text{var}(\alpha_0) \cup \text{var}(\alpha_1)) \mid (\alpha_0 + \alpha_1, \text{var}(\alpha_0) \cup \text{var}(\alpha_1)) \\ & \mid (\alpha_0^*, \text{var}(\alpha_0)) \mid ((\alpha_0)_j, \text{var}(\alpha_0) \uplus \{j\}) \text{ where } j \in [k] \setminus \text{var}(\alpha_0). \end{aligned}$$

Let 0-rewb denote the set of all regular expressions over  $\Sigma$ . Then, the set of all rewbs is  $\bigcup_{k \geq 0} k\text{-rewb}$ .

**Example 3.** For example,  $\varepsilon, a, \backslash 1, a^* \backslash 1, (1a^*)_1, ((1a^*)_1)^*, (2a^*)_2 \backslash 2, (1a^*)_1 (2b^*)_2 (\backslash 1 + \backslash 2), (2(1(a+b)^*)_1 \backslash 1)_2 \backslash 2 (2 \backslash 1)_2^*, ((1 \backslash 4 a)_1 (2 \backslash 3)_2 (3 \backslash 2 a)_3 (4 \backslash 1 \backslash 3)_4)^*$  are rewbs. On the other hand,  $(1(1a^*)_1)_1, (1a^* \backslash 1)_1, (1(2(1a^*)_1)_2)_1$  are not rewbs.

Note that this syntax allows multiple occurrences of captures with the same label, that is, we adopt the “may repeat labels” convention. Next, we define the semantics.

**Definition 4.** Let  $B_k = \{[i, ]_i \mid i \in [k]\}$ . The mapping  $\mathcal{R}_k : k\text{-rewb} \rightarrow \mathcal{P}((\Sigma \uplus B_k \uplus [k])^*)$  is defined as follows, where  $a \in \Sigma_\varepsilon$  and  $i \in [k]$ :

$$\begin{aligned} \mathcal{R}_k(a) &= \{a\}, \mathcal{R}_k(\backslash i) = \{i\}, \mathcal{R}_k(\alpha_0 \alpha_1) = \mathcal{R}_k(\alpha_0) \mathcal{R}_k(\alpha_1), \\ \mathcal{R}_k(\alpha_0 + \alpha_1) &= \mathcal{R}_k(\alpha_0) \cup \mathcal{R}_k(\alpha_1), \mathcal{R}_k(\alpha^*) = \mathcal{R}_k(\alpha)^*, \mathcal{R}_k((\alpha)_i) = \{[i] \mathcal{R}_k(\alpha) \{ ]_i \} \}. \end{aligned}$$

We let  $\Sigma_k^{[*]}$  denote  $\bigcup_{\alpha \in k\text{-rewb}} \mathcal{R}_k(\alpha)$ .

**Example 5.**  $\mathcal{R}_k((\{a+b\}^*)_1 \setminus 1) = \{[\{a+b\}^* \{1\}]_1\} = \{[w]_1 \mid w \in \{a,b\}^*\}.$

That is, we first regard a rew  $\alpha$  over  $\Sigma$  as a regular expression over  $\Sigma \uplus B_k \uplus [k]$ , deducing the language  $\mathcal{R}_k(\alpha)$ . The second step, described next, is to apply the *dereferencing (partial) function*  $\mathcal{D}_k : (\Sigma \uplus B_k \uplus [k])^* \rightarrow \Sigma^*$  to each of its element.

We give an intuitive description of  $\mathcal{D}_k$ . First,  $\mathcal{D}_k$  scans its input string from the beginning toward the end, seeking  $i \in [k]$ . If such  $i$  is found,  $\mathcal{D}_k$  replaces this  $i$  with the substring obtained by removing the brackets in  $v$  that comes from the preceding  $[i v]_i$  if  $[i$  exists (if this  $[i$  has no corresponding  $]_i$ ,  $\mathcal{D}_k$  becomes undefined). Otherwise,  $\mathcal{D}_k$  replaces this  $i$  with  $\varepsilon$ . The dereferencing function  $\mathcal{D}_k$  repeats this procedure until all elements of  $[k]$  appearing in the string are exhausted, then removes all remaining brackets. We let  $v_{[r]}$  denote the string which  $\mathcal{D}_k$  scans at the  $r^{\text{th}}$  number  $n_r \in [k]$  at the  $r^{\text{th}}$  loop.

1.  $[1a[2b]_2 2]_1 1$ . In this example,  $\mathcal{D}_k$  encounters  $n_1 = 2$  first, and this 2 corresponds the preceding  $[2b]_2$ , therefore this 2 is replaced with  $v_{[1]} = b$ . As a result, the input string becomes  $[1a[2b]_2 b]_1 1$ .  $\mathcal{D}_k$  repeats this process again. Now,  $\mathcal{D}_k$  locates  $n_2 = 1$  corresponding the preceding  $[1a[2b]_2 b]_1$ , so this 1 is replaced with  $v_{[2]} = a[2b]_2 b$  but with the brackets erased. Therefore, we gain  $[1a[2b]_2 b]_1 abb$ . Finally,  $\mathcal{D}_k$  removes all remaining brackets and produces  $abbabb$ . Here is the diagram:  $[1a[2b]_2 2]_1 1 \rightarrow [1a[2b]_2 b]_1 1 \rightarrow [1a[2b]_2 b]_1 abb \rightarrow abbabb$ .
2.  $[1a]_1 1 [1bb]_1 1$ . In this example,  $n_1 = n_2 = 1$ ,  $v_{[1]} = a$ ,  $v_{[2]} = bb$ , and

$$[1a]_1 1 [1bb]_1 1 \rightarrow [1a]_1 a [1bb]_1 1 \rightarrow [1a]_1 a [1bb]_1 bb \rightarrow aabbabb.$$

3.  $abc 1 2$ . In this example,  $n_1 = n_2 = 1$ ,  $v_{[1]} = v_{[2]} = \varepsilon$ , and  $abc 1 2 \rightarrow abc 2 \rightarrow abc$ .

Note that an unbound reference is replaced with the empty string  $\varepsilon$ , that is, we adopt the  $\varepsilon$ -semantics. However, as mentioned in Section 2, this semantics' expressive power is equivalent to that of the  $\emptyset$ -semantics under the “may repeat labels” convention (see [4] for the proof). We define the language  $L(\alpha)$  denoted by a  $k$ -rew  $\alpha$  to be  $\mathcal{D}_k(\mathcal{R}_k(\alpha)) = \{\mathcal{D}_k(v) \mid v \in \mathcal{R}_k(\alpha)\}$  (Lemmas 6 and 8 ensure that  $L(\alpha)$  is well-defined), and define the language class of rews as the union of the language classes of  $k$ -rews for all  $k$ .

Let  $g : (\Sigma \uplus B_k)^* \rightarrow \Sigma^*$  denote the free monoid homomorphism where  $g(x)$  is  $x$  for each  $x \in \Sigma$ , and  $\varepsilon$  for each  $x \in B_k$ . Every  $v \in (\Sigma \uplus B_k \uplus [k])^*$  can be written uniquely in the form  $v = v_0 n_1 v_1 \cdots n_m v_m$ , where  $m \geq 0$  (denoted by  $\text{cnt } v$ ), and  $v_r \in (\Sigma \uplus B_k)^*$  and  $n_r \in [k]$  for each  $r \in \{0, \dots, m\}$ . Here, let  $y_0 \triangleq v_0$  and for each  $r \in \{1, \dots, m\}$ ,  $y_r \triangleq v_0 n_1 v_1 \cdots n_r v_r$ . A string  $v = v_0 n_1 v_1 \cdots n_m v_m$  over  $\Sigma \uplus B_k \uplus [k]$  is said to be *matching* if

$$\forall r \in \{1, \dots, m\}. \forall x_1, x_2. y_{r-1} = x_1 [n_r x_2]_{n_r} \implies (\exists x'_2. x_3. x_2 = x'_2]_{n_r} x_3 \wedge x'_2 \notin [n_r, ]_{n_r})$$

holds. Intuitively, a string  $v$  being matching means that for all  $n_r \in [k]$  in  $v$ , if there exists a left bracket  $[n_r$  in the string immediately before  $n_r$ , then there is a right bracket  $]_{n_r}$  in between this  $[n_r$  and  $n_r$ . The following three lemmas follow.

**Lemma 6.** *Given a matching string  $v$ ,  $\mathcal{D}_k(v) = g(v_0) g(v_{[1]}) g(v_1) \cdots g(v_{[m]}) g(v_m)$ .*

**Lemma 7.** *A prefix of a matching string is matching. That is, if we decompose a string  $v$  into  $v = xy$ ,  $x$  is matching. Moreover,  $x_{[r]} = v_{[r]}$  holds for each  $r = 1, \dots, \text{cnt } x (\leq \text{cnt } v)$ .*

**Lemma 8.** *Every  $v \in \Sigma_k^{[*]}$  is matching.*

In the rest of this subsection, we formally define the dereferencing function  $\mathcal{D}_k$  and the notation  $v_{[r]}$ , and prove the three lemmas above by observing some basic features of  $\mathcal{D}_k$ . The reader may skip and go to the next subsection.

**Definition 9.** Suppose that  $\perp \notin \Sigma \uplus B_k \uplus [k]$ , and we define a Turing machine with one tape  $\mathcal{D}_k$  as follows:

$\mathcal{D}_k =$  “On input string  $v \in (\Sigma \uplus B_k \uplus [k])^*$ ,

- Step 1. Move the head to the right until it reads an  $i \in [k]$ , then go to Step 2. If such  $i$  does not exist, go to Step 6.
- Step 2. The assertion  $P_2 \triangleq$  ‘The symbol the head points out now is the leftmost natural number  $i \in [k]$  on the tape’ holds. Move the head to the left until it reads a  $[i$ , then go to Step 3. If such  $[i$  does not exist, go to Step 5.

- Step 3. Let  $P_3$  be ‘There is a right bracket  $]_i$  between the current head position and this  $i$ ’. If  $P_3$  holds, go to Step 4. Otherwise, go to Step 7.
- Step 4. Move the head to the right one by one, seeking a bracket  $]_i$ . Note that no  $j \in [k]$  can appear in this scan since  $P_2$  and  $P_3$  holds. Now, if the symbol written on the tape cell that the head points to is  $a \in \Sigma$ , then insert  $a$  into the position immediately preceding this  $i$ , and go right; if it is  $b \in B_k \setminus \{ ]_i \}$ , simply go right; if it is  $]_i$ , go to Step 5.
- Step 5. Go back to Step 1 and remove this  $i$ .
- Step 6.  $P_6 \triangleq$  ‘No  $j \in [k]$  is written on the tape’ holds. Again scan the tape from the beginning, and remove all brackets.
- Step 7. Erase all symbols written on the tape, and write the symbol  $\perp$ .”

Henceforth, we refer to these step numbers as encircled numbers ①, ②, etc. The order of execution of  $\mathcal{D}_k$  is  $(\textcircled{1}\textcircled{2}(\textcircled{3}\textcircled{4}\textcircled{5} + \textcircled{5}))^*\textcircled{1}(\textcircled{6} + \textcircled{2}\textcircled{3}\textcircled{7})$ . Observe that  $\mathcal{D}_k$  halts for any input because ⑤ cannot be taken arbitrarily many times. Thus, we think of  $\mathcal{D}_k$  as a computable function  $(\Sigma \uplus B_k \uplus [k])^* \rightarrow \Sigma^* \uplus \{\perp\}$ .

**Lemma 10.** *Suppose that the loop unit (namely  $\textcircled{1}\textcircled{2}\textcircled{3}\textcircled{4}\textcircled{5}$  or  $\textcircled{1}\textcircled{2}\textcircled{5}$ ) is executed exactly  $m'$  times when an input string  $v = v_0 n_1 v_1 \cdots n_m v_m$  ( $m = \text{cnt } v$ ) is given to  $\mathcal{D}_k$ . Then, for each  $r \in \{0, 1, \dots, m'\}$ , let  $v^{(r)}$  be the string written on the tape immediately after the  $r^{\text{th}}$  loop, and let  $v_{[r]}$  be the string over  $\Sigma \uplus B_k$  defined as follows:*

$$v_{[r]} \triangleq \begin{cases} \text{the string bracketed in } [{}_i \text{ and } ]_i \text{ scanned at } \textcircled{4}, & (\text{if } \textcircled{1}\textcircled{2}\textcircled{3}\textcircled{4}\textcircled{5} \text{ is executed}) \\ \varepsilon. & (\text{if } \textcircled{1}\textcircled{2}\textcircled{5} \text{ is executed}) \end{cases}$$

Under the assumptions above, for each  $r \in \{0, 1, \dots, m'\}$  the following equality holds:

$$v^{(r)} = v_0 g(v_{[1]}) v_1 \cdots g(v_{[r]}) v_r \underline{n_{r+1}} v_{r+1} \cdots n_m v_m.$$

*Proof.* When  $r = 0$ , the left side is  $v^{(0)}$  and the right side is  $v_0 n_1 v_1 \cdots n_m v_m = v$ , as required. Recall that immediately before the  $(r + 1)^{\text{st}}$  loop, the string written on the tape is  $v^{(r)}$ . It continues as follows:

- At ①, the head of  $\mathcal{D}_k$  is placed at immediately after  $g(v_{[r]})$  in

$$v^{(r)} = v_0 g(v_{[1]}) v_1 \cdots g(v_{[r]}) v_r \underline{n_{r+1}} v_{r+1} \cdots n_m v_m,$$

and moves to  $n_{r+1}$ . Henceforth, we write simply  $s_r$  for  $v_0 g(v_{[1]}) v_1 \cdots g(v_{[r]}) v_r$ , which appears before  $n_{r+1}$ .

- At ②, there are two cases:
  - When  $\mathcal{D}_k$  follows  $\textcircled{3}\textcircled{4}\textcircled{5}$ , since it moves from ② to ③ and  $P_3$  holds,  $s_r$  is of the form  $x_0 [n_{r+1} v_{[r+1]}]_{n_{r+1}} x_1$ . At ④ and ⑤, the strings written on the tape after each step are

$$\textcircled{4} \ s_r g(v_{[r+1]}) \underline{n_{r+1}} v_{r+1} \cdots n_m v_m, \text{ and } \textcircled{5} \ s_r g(v_{[r+1]}) v_{r+1} \cdots n_m v_m.$$

By definition, this string after ⑤ is  $v^{(r+1)}$ .

- When  $\mathcal{D}_k$  follows ⑤, since  $v_{[r+1]} = \varepsilon$ , the string immediately after the execution of ⑤, namely  $v^{(r+1)}$ , is  $s_r v_{r+1} \cdots n_m v_m = s_r g(v_{[r+1]}) v_{r+1} \cdots n_m v_m$ .

In both cases, the equation holds for  $r + 1$ .

This completes the proof. □

**Lemma 11.** *For a matching string  $v$ , the loop of  $\mathcal{D}_k$  runs exactly  $m = \text{cnt } v$  times. Hence,  $\mathcal{D}_k(v) \in \Sigma^*$ .*

*Proof.* Let  $m'$  be the number of loop iterations. Because obviously  $m' \leq m$ , we suppose  $m' < m$  for contradiction. If so, since it is the case that the execution moves from ①②③ to ⑦ at the  $(m' + 1)^{\text{st}}$  loop,  $P_3$  does not hold for  $v^{(m')}$ . By Lemma 10,

$$v^{(m')} = \underbrace{v_0 g(v_{[1]}) v_1 \cdots g(v_{[m']}) v_{m'}}_{s_{m'}} \underline{n_{m'+1} v_{m'+1} \cdots n_m v_m}$$

follows. Since  $i$  in the  $(m' + 1)^{\text{st}}$  loop is this  $n_{m'+1}$  and  $s_{m'} \ni [n_{m'+1}]$ , there is  $v_j$  among  $v_0, \dots, v_{m'}$  such that  $v_j \ni [n_{m'+1}]$ . This  $v_j$  can be written in the form  $u_0 [n_{m'+1}] u_1$ . Because  $v$  is matching and  $y_{m'} \supseteq v_j \ni [n_{m'+1}]$ , one of  $u_1, v_{j+1}, \dots, v_{m'}$  contains  $]n_{m'+1}$ . This contradicts the fact that  $P_3$  does not hold at ③.  $\square$

*Proof of Lemma 6.* In Lemma 11, the string  $v^{(m)}$ , which is written on the tape immediately after the  $m^{\text{th}}$  loop, becomes  $g(v^{(m)})$  at ⑥, and therefore  $\mathcal{D}_k$  halts, as required.  $\square$

*Proof of Lemma 7.* Immediate from Lemma 6.  $\square$

Finally, we prove that every  $v \in \Sigma_k^{[*]}$  (see Definition 4) is matching, concluding  $L(\alpha) \subseteq \Sigma^*$  with Lemma 11.

**Lemma 12.** *The following facts hold where  $\alpha \in k\text{-rewb}$  and  $i \in [k]$ :*

- (a)  $\mathcal{R}_k(\alpha) \subseteq (\Sigma \cup \{[j, ]_j \mid j \in \text{var}(\alpha)\} \cup \text{var}(\alpha))^*$ ,
- (b)  $\forall v_1, v_2. v_1 [i] v_2 \in \mathcal{R}_k(\alpha) \implies \exists v'_2, v_3. v_2 = v'_2 [i] v_3 \wedge v'_2 \not\ni [i, ]_i, i$ .

*Proof.* Immediate from the definitions of  $\mathcal{R}_k(\alpha)$  and  $\text{var}(\alpha)$ .  $\square$

*Proof of Lemma 8.* There exists  $\alpha \in k\text{-rewb}$  such that  $v \in \mathcal{R}_k(\alpha)$ . Hereafter, we let  $m = \text{cnt } v$  and  $v = v_0 n_1 v_1 \cdots n_m v_m$ . For each  $r \in \{1, \dots, m\}$ , if  $y_{r-1}$  can be written in the form  $x_1 [n_r, x_2]$ , we obtain  $v = x_1 [n_r, x_2 n_r v_r \cdots n_m v_m]$  and by Lemma 12 (b),  $x_2 n_r v_r \cdots n_m v_m$  can be written in the form  $x'_2 ]n_r, x_3$  with  $x'_2 \not\ni [n_r, ]_{n_r}, n_r$ . Here  $x'_2 ]n_r, \not\ni n_r$  holds. Therefore,  $x'_2 ]n_r$  is a prefix of  $x_2$  and of course  $x'_2 \not\ni [n_r, ]_{n_r}$ . Hence,  $v$  is matching.  $\square$

### 3.2. Classical automata models

Next, we recall the notions of SA, NESAs, and Nested-SA. In this paper, we unify their definitions based on [2, 12] to clarify the different capabilities of these models. First, we review NFA. Here is the definition in the textbook by Hopcroft et al. [15]:

**Definition 13 ([15], p.57).** A *nondeterministic finite automaton*  $N$  is a 5-tuple  $(Q, \Sigma, \delta, q_0, F)$ , where  $Q$  is a finite set of states,  $\Sigma$  a finite set of input symbols (also called alphabet),  $q_0 \in Q$  a start state,  $F \subseteq Q$  a set of final states, and  $\delta : Q \times \Sigma \rightarrow \mathcal{P}(Q)$  a transition function.

As well known, the transition function  $\delta$  can be extended to  $\hat{\delta} : Q \times \Sigma^* \rightarrow \mathcal{P}(Q)$  where  $\hat{\delta}(q, w)$  represents the set of all states reachable from  $q$  via  $w$ . Let  $q \xrightarrow{a}_N q'$  denote  $q' \in \delta(q, a)$ , and  $q \xrightarrow{w}_N q'$  denote  $q' \in \hat{\delta}(q, w)$ . With this notation,

the language of an NFA  $N$  can be written as follows:  $L(N) = \left\{ w \in \Sigma^* \mid \exists q_f \in F. q_0 \xrightarrow{w}_N q_f \right\}$ .

A pushdown automaton (PDA) is an NFA equipped with a *stack* such that the PDA may write and read its *stack top* with a transition. A *stack automaton* (SA) is “an extended PDA”, which can reference not only the top but inner content of the stack. That is, while the *stack pointer* of a PDA is fixed to the top, an SA allows its pointer to move left and right and read a stack symbol pointed to by the pointer. However, the only place on the stack that can be rewritten is the top, as in PDA. Formally, a (1N) SA  $A$  is a 10-tuple  $(Q, \Sigma, \Gamma, \delta, q_0, Z_0, \vdash, \#, \$, F)$  satisfying the following conditions: the components  $Q, \Sigma, q_0$  and  $F$  are the same as those of NFA.  $\Gamma (\neq \emptyset)$  is a finite set of stack symbols, and  $Z_0 \in \Gamma$  is an initial stack symbol. The symbol  $\vdash \notin \Sigma$  is the endmarker of the input and the symbol  $\# \notin \Sigma \cup \Gamma$  (resp.  $\$ \notin \Sigma \cup \Gamma$ ) is always and only written at the leftmost (bottom) (resp. the right most (top)) of the stack.<sup>4</sup>

<sup>4</sup>These special symbols  $\#, \$$  representing “bottom” and “top” of the stack respectively do not appear in [12] and are introduced anew in this paper to define NESAs and Nested-SA, which will be defined later, in the style of [2]. In fact, SA defined in [12] is not capable of directly discerning whether the stack pointer is at the top or not. Although it is not difficult to see that directly adding the ability does not increase the expressive power of SA, the ability is directly in NESAs as seen in [16, 20]. Therefore, to make it easy to see that NESAs are a restriction of SA, we define SA to also directly have the ability.

The transition function  $\delta$  has the following two modes, where  $L, S, R \notin (\Sigma \cup \Gamma) \cup \{-1, \#, \$\}$ ,  $\Delta_i \triangleq \{S, R\}$ ,  $\Delta_s \triangleq \{L, S, R\}$  and  $\Sigma' \triangleq \Sigma \cup \{-1\}$ :

- (i) (pushdown mode)  $Q \times \Sigma' \times \Gamma\$ \rightarrow \mathcal{P}(Q \times \Delta_i \times \Gamma^*\$)$ ,
- (ii) (stack reading mode) (a)  $Q \times \Sigma' \times \Gamma\$ \rightarrow \mathcal{P}(Q \times \Delta_i \times \{L, S\})$ , (b)  $Q \times \Sigma' \times \Gamma \rightarrow \mathcal{P}(Q \times \Delta_i \times \Delta_s)$ , (c)  $Q \times \Sigma' \times \{\#\} \rightarrow \mathcal{P}(Q \times \Delta_i \times \{S, R\})$ .

Intuitively,  $\delta$  works as follows (Definition 14 provides the formal semantics). (i) The statement  $(q', d, w\$) \in \delta(q, a, Z\$)$  says that whenever the current state is  $q$ , the input symbol is  $a$ , and the pointer references the top symbol  $Z$ , the machine can move to the state  $q'$ , move the input cursor along  $d$ , and replace  $Z$  with the string  $w$ . (ii) The statement  $(q', d, e) \in \delta(q, a, Z)$  says that whenever the current state is  $q$ , the input symbol is  $a$ , and the pointer references the symbol  $Z$ , the machine can move to the state  $q'$ , move the input cursor along  $d$ , and move the pointer along  $e$ . The statements (a) and (c) are similar to (b) except that the direction in which the pointer can move is restricted lest the pointer go out of the stack. In particular, an SA that cannot erase a symbol once written on the stack is called a *nonerasing stack automaton* (NESA). That is, a (1N) nonerasing stack automaton is an SA whose transition function  $\delta$  satisfies the condition that, in (i) (pushdown mode),  $(q', d, y\$) \in \delta(q, a, Z\$)$  implies  $y \in Z\Gamma^*$ . To formally describe how SA works, we define a tuple called *instantaneous description* (ID), which consists of a state, an input string, and a string representation of the stack, and define the binary relation  $\vdash_A$  over the set of these tuples. Let  $\bar{L} = -1$ ,  $\bar{S} = 0$ , and  $\bar{R} = 1$ .

**Definition 14.** Let  $A$  be an SA  $(Q, \Sigma, \Gamma, \delta, q_0, Z_0, \bar{L}, \bar{S}, \bar{R}, F)$ . An element of the set  $I = Q \times \Sigma^* \{-1\} \times \{\#\} (\Gamma \cup \{\bar{L}\})^* \{\bar{S}\}$  is called *instantaneous description*, where the stack symbol  $\bar{L} \notin \Gamma$  stands for the position of stack pointer. Moreover, let  $\vdash_A$  (or  $\vdash$  when  $A$  is clear) be the smallest binary relation over  $I$  satisfying the following conditions for all  $q, q' \in Q$ ,  $k \geq 0$ ,  $a_0, \dots, a_k \in \Sigma \cup \{-1\}$ ,  $Z, Z_1, \dots, Z_n \in \Gamma$ ,  $\gamma, y \in \Gamma^*$ ,  $d \in \Delta_i$  and  $e \in \Delta_s$  with  $k = 0 \rightarrow d \neq R$ :

- (i)  $(q, a_0 \dots a_k, \# \gamma Z \bar{L} \$) \vdash_A (q', a_{\bar{d}} \dots a_k, \# \gamma y \bar{L} \$)$  if  $(q', d, y\$) \in \delta(q, a_0, Z\$)$ .
- (ii) (a)  $(q, a_0 \dots a_k, \# \gamma Z \bar{L} \$) \vdash_A (q', a_{\bar{d}} \dots a_k, \# \gamma \bar{L} Z \$)$  if  $(q', d, L) \in \delta(q, a_0, Z\$)$ , and  
 $(q, a_0 \dots a_k, \# \gamma Z \bar{L} \$) \vdash_A (q', a_{\bar{d}} \dots a_k, \# \gamma Z \bar{L} \$)$  if  $(q', d, S) \in \delta(q, a_0, Z\$)$ .  
(b) if  $(q', d, e) \in \delta(q, a_0, Z)$ ,  $Z = Z_j$  and  $1 \leq j < n$  with  $j = 1 \rightarrow e \neq L$  and  $j = n \rightarrow e \neq R$ , then  
 $(q, a_0 \dots a_k, \# Z_1 \dots Z_j \bar{L} \dots Z_n \$) \vdash_A (q', a_{\bar{d}} \dots a_k, \# Z_1 \dots Z_{j+\bar{e}} \bar{L} \dots Z_n \$)$ .  
(c)  $(q, a_0 \dots a_k, \# \bar{L} Z \gamma \$) \vdash_A (q', a_{\bar{d}} \dots a_k, \# Z \bar{L} \gamma \$)$  if  $(q', d, R) \in \delta(q, a_0, \#)$ , and  
 $(q, a_0 \dots a_k, \# \bar{L} Z \gamma \$) \vdash_A (q', a_{\bar{d}} \dots a_k, \# \bar{L} Z \gamma \$)$  if  $(q', d, S) \in \delta(q, a_0, \#)$ .

Note that  $L \notin \Delta_i$ , which means the input cursor will not move back to left. We say that  $A$  accepts  $w \in \Sigma^*$  if there exist  $\gamma_1, \gamma_2 \in \Gamma^*$ , and  $q_f \in F$  such that  $(q_0, w\bar{L}, \# Z_0 \bar{L} \$) \vdash_A^* (q_f, \bar{L}, \# \gamma_1 \bar{L} \gamma_2 \$)$ . Let  $L(A)$  denote the set of all strings accepted by  $A$ .

**Notation 15.** We introduce some arrow notations to draw SA diagrams. We denote the rule of (i)  $(q', d, y\$) \in \delta(q, a, Z\$)$  as  $q \xrightarrow{a, d/Z\$ \rightarrow y\$} q'$ ; (ii)  $(q', d, e) \in \delta(q, a, \square)$  as  $q \xrightarrow{a, d/\square, e} q'$  for  $\square = Z\$, Z$  or  $\#$ . If  $d = R$ , we simply write  $a, d/\dots$  as  $a/\dots$ . For convenience, we also introduce intuitive syntax sugars to represent multiple transitions at once. For instance, an  $\varepsilon$ -transition  $\varepsilon/\dots$  consists of putting  $c, S/\dots$  for all  $c \in \Sigma$ , and we may omit the upper part in front of  $/$  and write the part  $\dots$  only. In (i),  $\$ \rightarrow Z\$$  is a shorthand for putting  $Z' \$ \rightarrow Z' Z \$$  for all  $Z' \in \Gamma$ . In (ii), for a finite subset  $U \subseteq \Gamma$ , the notation  $\dots/U, e$  (resp.  $\dots/U \$, e$ ) means putting  $\dots/u, e$  (resp.  $\dots/u \$, e$ ) for all  $u \in U$ . In particular,  $U = \Gamma$ , we omit  $U$  and merely type  $\dots/e$ . The negation  $\neg Z$  is a shorthand for  $\Gamma \setminus \{Z\}$ .

We next define *nested stack automaton* (Nested-SA) which is SA extended with the capability to create and remove substacks. For instance, suppose that the stack is  $\# a_1 a_2 \bar{L} a_3 \$$  and we are to create a new substack containing  $b_1 b_2$ :

$$\# a_1 \underline{c} b_1 b_2 \bar{L} \$ a_2 a_3 \$ \tag{1}$$

Note that the new substack  $\underline{c} b_1 b_2 \$$  is embedded below the symbol  $a_2$  indicated by the stack pointer, and the pointer moves to the top of the created substack. The creation of the inner substack narrows the range within which the stack pointer can move as indicated by the underlined part  $\# a_1 \underline{c} b_1 b_2 \bar{L} \$$ . While the bottom of the entire stack is always fixed by the leftmost symbol  $\#$ , the top of the embedded substack is regarded as the top of the entire stack. The

inner substacks are allowed to be embedded endlessly and everywhere, whereas the writing in the pushdown mode is still restricted to the top of the stack:

$$\frac{\#a_1\epsilon b_1b_2\mid \$a_2a_3\$}{\#a_1\epsilon b_1\mid b_2\$a_2a_3\$} \xrightarrow{L} \frac{\#a_1\epsilon b_1\mid b_2\$a_2a_3\$}{\#a_1\epsilon\epsilon c_1c_2\mid \$b_1b_2\$a_2a_3\$} \xrightarrow{\text{create}} \quad (2)$$

$$\frac{\#a_1\epsilon\mid b_1b_2\$a_2a_3\$}{\#a_1\mid \epsilon b_1b_2\$a_2a_3\$} \xrightarrow{L} \frac{\#a_1\mid \epsilon b_1b_2\$a_2a_3\$}{\#\epsilon c_1c_2\mid \$a_1\epsilon b_1b_2\$a_2a_3\$} \xrightarrow{\text{create}} \quad (3)$$

We must empty the inner substack and then remove itself in advance whenever we want to reference the right side of the inner substack such as  $a_2, a_3$ . For example, let us empty the inner substack by popping twice from (1) and then removing it:

$$\frac{\#a_1\epsilon b_1b_2\mid \$a_2a_3\$}{\#a_1\epsilon b_1\mid \$a_2a_3\$} \xrightarrow{\text{pop}} \frac{\#a_1\epsilon b_1\mid \$a_2a_3\$}{\#a_1\epsilon\mid \$a_2a_3\$} \xrightarrow{\text{pop}} \frac{\#a_1\epsilon\mid \$a_2a_3\$}{\#a_1a_2\mid a_3\$} \xrightarrow{\text{destruct}} \quad (4)$$

Notice that the stack pointer moves to the right after removing the inner substack. We now define Nested-SA formally. A (1N) nested stack automaton  $A$  is a 10-tuple  $(Q, \Sigma, \Gamma, \delta, q_0, Z_0, \#, \epsilon, \$, F)$  satisfying the following conditions: the components  $Q, \Sigma, \Gamma, q_0, Z_0, \#, \$$  and  $F$  are the same as those of SA. The stack symbol  $\epsilon \notin \Sigma \cup \Gamma$  represents the bottom of a substack.<sup>5</sup> The transition function  $\delta$  has the following four modes, where  $\Sigma' \triangleq \Sigma \cup \{\cdot\}$  and  $\Gamma' \triangleq \Gamma \cup \{\epsilon\}$ :

- (i) (pushdown mode)  $Q \times \Sigma' \times \Gamma\$ \rightarrow \mathcal{P}(Q \times \Delta_i \times \Gamma^*\$)$ .
- (ii) (stack reading mode) (a)  $Q \times \Sigma' \times \Gamma^*\$ \rightarrow \mathcal{P}(Q \times \Delta_i \times \{L, S\})$ , (b)  $Q \times \Sigma' \times \Gamma' \rightarrow \mathcal{P}(Q \times \Delta_i \times \Delta_s)$ , (c)  $Q \times \Sigma' \times \{\#\} \rightarrow \mathcal{P}(Q \times \Delta_i \times \{S, R\})$ .
- (iii) (stack creation mode)  $Q \times \Sigma' \times (\Gamma' \cup \Gamma^*\$) \rightarrow \mathcal{P}(Q \times \Delta_i \times \{\epsilon\} \Gamma^*\$)$ .
- (iv) (stack destruction mode)  $Q \times \Sigma' \times \{\epsilon\} \$ \rightarrow \mathcal{P}(Q \times \Delta_i)$ .

Moreover, we define how Nested-SA works with ID and  $\vdash$  in the same manner as SA. Given a Nested-SA  $A = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, \cdot, \#, \epsilon, \$, F)$ , we define ID,  $\vdash_A$ , and  $L(A)$  in the same way as Definition 14 (however, we let  $I$  be  $Q \times \Sigma^* \{\cdot\} \times \{\#\} (\Gamma \cup \{\epsilon, \$, \cdot\})^* \{\#\}$ ). Here, we only give the rules corresponding to (iii) and (iv) in the definition of  $\delta$  (the others are essentially the same as those of SA):

- (iii) if  $(q', d, \epsilon y\$) \in \delta(q, a_0, Z)$  and  $Z = Z_j, 1 \leq j < n$ , then
  - $(q, a_0 \cdots a_k, \#Z_1 \cdots Z_j \mid \cdots Z_n\$) \vdash_A (q', a_{i+\bar{d}} \cdots a_k, \#Z_1 \cdots \epsilon y \mid \$Z_j \cdots Z_n\$)$ ,
  - and  $(q, a_0 \cdots a_k, \# \gamma Z \mid \$) \vdash_A (q', a_{\bar{d}} \cdots a_k, \# \gamma \epsilon y \mid \$Z\$)$  if  $(q', d, \epsilon y\$) \in \delta(q, a_0, Z\$)$ .
- (iv)  $(q, a_0 \cdots a_k, \# \gamma_1 \epsilon \mid \$Z \gamma_2\$) \vdash_A (q', a_{\bar{d}} \cdots a_k, \# \gamma_1 Z \mid \gamma_2\$)$  if  $(q', d) \in \delta(q, a_0, \epsilon\$)$ .

## 4. Every rewb describes an indexed language

### 4.1. Main theorem

As described above, to obtain the language  $L(\alpha)$  described by a  $k$ -rewb  $\alpha$ , we derive the regular language  $\mathcal{R}_k(\alpha)$  over the alphabet  $\Sigma \cup B_k \cup [k]$  first, then apply the dereferencing function  $\mathcal{D}_k$  to every element of  $\mathcal{R}_k(\alpha)$ . Using this observation, we construct a Nested-SA  $A_\alpha$  recognizing the language  $L(\alpha)$  as follows.

The Nested-SA  $A_\alpha$  is based on an NFA  $N$  recognizing the language  $\mathcal{R}_k(\alpha)$ , in the sense that each transition in  $A_\alpha$  comes from a corresponding transition of  $N$ . The NFA  $N$  has the alphabet  $\Sigma \cup B_k \cup [k]$ , and so handles three types of characters. For each transition  $q \xrightarrow[N]{a} q'$  with  $a \in \Sigma$ , i.e., moving from  $q$  to  $q'$  by an input symbol  $a$ ,  $A_\alpha$  also has the same transition except pushing  $a$  to the stack, denoted by  $q \xrightarrow{a/\$ \rightarrow a\$} q'$ . For each transition  $q \xrightarrow[N]{b} q'$  with  $b \in B_k$ , i.e., moving by a bracket  $b$ ,  $A_\alpha$  has the transition pushing  $b$  without consuming input symbols, denoted by  $q \xrightarrow{\epsilon/\$ \rightarrow b\$} q'$ .<sup>6</sup> For each transition  $q \xrightarrow[N]{i} q'$  with  $i \in [k]$ ,  $A_\alpha$  has a large “transition” that consists of several transitions. In this “transition,”  $A_\alpha$  first seeks the left bracket  $[i$  of the bracketed string  $[i \ v]_i$  within the stack, and checks if the input from the cursor position matches  $v$  character by character while consuming the input, and finally moves to  $q'$  if all characters of  $v$  matched.

<sup>5</sup>Note that the bottom of the entire stack is always represented by  $\#$  and not  $\epsilon$ , as mentioned above.

<sup>6</sup>See also Notation 15 for the  $\epsilon$ -transition notation.



A difficult yet interesting point is that Nested-SA cannot check  $v$  against the stack and push  $v$  onto the stack at the same time, that is, after checking a character  $c$  of  $v$ , if  $A_\alpha$  wants to push  $c$  to the stack,  $A_\alpha$  must leave from  $v$ , climb up the stack toward the top, and write  $c$ . However, after the push,  $A_\alpha$  becomes lost by not knowing where to go back to. How about marking the place where  $A_\alpha$  should return in advance? Unfortunately, that does not work; Nested-SA can insert such marks anywhere by creating substacks, but due to the restriction of Nested-SA, it cannot go above the position of the mark, much less climb up to the top. Therefore, Nested-SA cannot directly push the result of a dereferencing onto the stack.

We cope with this problem as follows. We allow  $j \in [k]$  to appear in  $v$ , and for each appearance of  $j$  in the checking of  $v$ ,  $A_\alpha$  pauses the checking and puts a substack containing the current state as a marker at the stack pointer position. Then,  $A_\alpha$  searches down the stack for the corresponding bracketed string  $[_j v']_j$ , and begins checking  $v'$  if it is found. By repeating this process,  $A_\alpha$  eventually reaches a string  $v'' \in (\Sigma \uplus B_k)^*$  containing no characters of  $[k]$ . Once done with the check of  $v''$ ,  $A_\alpha$  climbs up toward the stack top, finds a marker  $p$  denoting the state to return to, and resumes from  $p$  after deleting the substack containing the marker. By repeating this, if  $A_\alpha$  returns to the position where it initially found  $j$ , it has successfully consumed the substring of the input string corresponding to the dereferencing of  $j$ . The following lemma is immediate.

**Lemma 16.** *Let  $k \geq 1$  and  $\alpha \in k\text{-rewb}$ . There exists an NFA  $(Q, \Sigma \uplus B_k \uplus [k], \delta, q_0, F)$  over  $\Sigma \uplus B_k \uplus [k]$  recognizing  $\mathcal{R}_k(\alpha)$  all of whose states can reach some final state, that is,  $\forall q \in Q. \exists w \in (\Sigma \uplus B_k \uplus [k])^*. \exists q_f \in F. q \xRightarrow[N]{w} q_f$ .*

**Corollary 17.** *Let  $N$  be the NFA in Lemma 16. For all  $q \in Q$  and for all  $w \in (\Sigma \uplus B_k \uplus [k])^*$ , if  $q_0 \xRightarrow[N]{w} q$  then  $w$  is matching.*

*Proof.* There are a string  $w' \in (\Sigma \uplus B_k \uplus [k])^*$  and a final state  $q_f \in F$  such that  $q_0 \xRightarrow[N]{w} q \xRightarrow[N]{w'} q_f$ . Hence,  $ww' \in L(N) = \mathcal{R}_k(\alpha)$  follows. By Lemma 8 the string  $ww'$  is matching, therefore by Lemma 7 its prefix  $w$  is matching.  $\square$

We show the main theorem (the proof sketch is coming later):

**Theorem 18.** *For every rew  $\alpha$ , there exists a Nested-SA that recognizes  $L(\alpha)$ .*

The claim obviously holds when  $\alpha$  is a pure regular expression (i.e.,  $\alpha \in 0\text{-rewb}$ ). Suppose that  $\alpha \in k\text{-rewb}$  with  $k \geq 1$ . By Lemma 16, there is an NFA  $N = (Q_N, \Sigma \uplus B_k \uplus [k], \delta_N, q_0, F)$  that recognizes  $\mathcal{R}_k(\alpha)$  and satisfies Corollary 17. We construct a Nested-SA  $A_\alpha = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, \#, \$, F)$  as follows. Let  $Q \triangleq Q_N \uplus \{c_i, e_i, r_i \mid i \in [k]\} \uplus \{W_q \mid q \in Q_N\} \uplus \{E_{p,i}, L_{p,i} \mid p \in Q_N \uplus \{e_i \mid i \in [k]\}, i \in [k]\}$ ,  $\Gamma \triangleq \Sigma \uplus B_k \uplus [k] \uplus Q \uplus \{Z_0\}$ , and let  $\delta$  be the smallest relation that, for all  $a \in \Sigma$ ,  $b \in B_k$ ,  $c \in \Sigma \uplus \{\cdot\}$ ,  $i, j \in [k]$ ,  $q, q' \in Q_N$ ,  $Z \in \Gamma$  and  $p \in Q_N \uplus \{e_i \mid i \in [k]\}$ , satisfies the following conditions:

- |  |  |
|--|--|
| (1) $\delta_N(q, a) \ni q' \implies \delta(q, a, Z\$) \ni (q', R, Za\$)$     | (10) $\delta(e_i, c, [j]) = \{(e_i, S, R)\}$ where $i \neq j$  |
| (2) $\delta_N(q, b) \ni q' \implies \delta(q, c, Z\$) \ni (q', S, Zb\$)$     | (11) $\delta(e_i, c, [j]) = \begin{cases} \{(r_i, S, R)\} & (i = j) \\ \{(e_i, S, R)\} & (i \neq j) \end{cases}$ |
| (3) $\delta_N(q, i) \ni q' \implies \delta(q, c, Z\$) \ni (W_{q'}, S, Zi\$)$ | (12) $\delta(e_i, c, j) = \{(c_j, S, \$e_i\$)\}$ where $i \neq j$  |
| (4) $\delta(W_q, c, i\$) = \{(c_i, S, \$q\$)\}$                              | (13) $\delta(r_i, c, Z) = \{(r_i, S, R)\}$   |
| (5) $\delta(c_i, c, p\$) = \{(c_i, S, L)\}$                                  | (14) $\delta(r_i, c, p\$) = \{(E_{p,i}, S, \$)\}$  |
| (6) $\delta(c_i, c, Z) = \{(c_i, S, L)\}$ where $Z \neq [i, Z_0$             | (15) $\delta(E_{p,i}, c, \$) = \{(L_{p,i}, S)\}$   |
| (7) $\delta(c_i, c, Z_0) = \{(r_i, S, R)\}$                                  | (16) $\delta(L_{e_j,i}, c, i) = \{(e_j, S, R)\}$   |
| (8) $\delta(c_i, c, [i]) = \{(e_i, S, R)\}$                                  | (17) $\delta(L_{q,i}, c, i\$) = \{(q, S, S)\}$   |
| (9) $\delta(e_i, a, a) = \{(e_i, R, R)\}$                                    |  |

Rule (1) translates  $q \xrightarrow[N]{a} q'$  into  $q \xrightarrow{a/\$ \rightarrow a\$} q'$ , (2) translates  $q \xrightarrow[N]{b} q'$  into  $q \xrightarrow{\$/\$ \rightarrow b\$} q'$ , and rules (3)–(17) translate  $q \xrightarrow[N]{i} q'$  into a large “transition” to consume the string that corresponds to the dereferencing of  $i$ . The details of the “transition” are as follows. By looking at the underlying  $N$  with rule (3),  $A_\alpha$  finds a state  $q'$  that it should go back

to after going throughout the “transition,” and goes to the state  $W_{q'}$  by pushing  $i$  to the stack. At  $W_{q'}$ , by rule (4),  $A_\alpha$  inserts  $\$q'\$$  just below  $i$ , and goes to the state  $c_i$ . The state  $c_i$  represents the *call mode* in which  $A_\alpha$  looks for the left-nearest  $[i$  by rules (5) and (6) and proceeds to the state  $e_i$  (*execution mode*) by (8) if it finds  $[i$ . Otherwise (i.e., the case when  $A_\alpha$  arrives at the bottom of the stack), it proceeds to the state  $r_i$  (*return mode*) by rule (7). At  $e_i$ ,  $A_\alpha$  consumes input symbols by checking them against the symbols on the stack (rules (9)–(12)). In particular, rule (9) handles the case when the symbols match. Rules (10) and (11) handle the cases when brackets are read from the stack. The first case of (11) handles the case when the right bracket  $]i$  is read, and the rules handle the other brackets (i.e.,  $[j$  or  $]j$  with  $i \neq j$ ) by simply skipping them (note that  $]i = [i$  cannot happen since we started from the left-nearest  $[i$ ). Reading  $j \in [k]$ , by rule (12),  $A_\alpha$  inserts  $\$e_j\$$  just below  $j$  and goes to  $c_j$  to locate the corresponding  $[j$  (here,  $j \neq i$  holds by the definition of the syntax). At  $r_i$ ,  $A_\alpha$  proceeds to return to the state  $p$  that passed the control to  $c_i$  (rules (13)–(17)). Since this  $p$  was pushed at the stack top,  $A_\alpha$  first climbs up to the stack top by rule (13), transits to the state  $E_{p,i}$  popping  $p$  by (14), then goes to  $L_{p,i}$  removing the embedded substack by (15), and finally goes back to  $p$  by (16) and (17). A subtle point in the last step is that where the stack pointer should be placed depends on whether  $p$  is a state  $e_j$  (for some  $j \in [k]$ ) or in  $Q_N$ . In the former case, after (15) removes the embedded substack  $\$e_j\$$  that was created just below the call to  $i$ , the stack pointer points to  $i$ . However, the stack pointer should shift one more to the right, lest  $A_\alpha$  begins to repeat the call reading  $i$  again by (12). Therefore, (16) correctly handles the case by doing the shift. In the latter case, as stipulated by (17), the stack pointer should point to the stack top symbol  $i$  since  $p$  is the state stored at (3).

*Proof sketch of Theorem 18.* For proving  $L(\alpha) \subseteq L(A_\alpha)$ , we take  $w \in L(\alpha)$  and  $v \in \mathcal{R}_k(\alpha)$  such that  $w = \mathcal{D}_k(v)$ . Decomposing  $v$  into  $v_0 n_1 v_1 \cdots n_m v_m$  (where  $m = \text{cnt } v$ ), we obtain a transition sequence in the underlying NFA  $N$ , denoted by  $q_0 \xRightarrow[N]{v_0} q_{(0)} \xRightarrow[N]{n_1 v_1} q_{(1)} \xRightarrow[N]{n_2 v_2} \cdots \xRightarrow[N]{n_m v_m} q_{(m)} \in F$ . We prove by induction on  $r = 0, \dots, m$  that  $A_\alpha$  can reach  $q_{(r)}$  while consuming  $z_r = g(v_0)g(v_{[1]})g(v_1) \cdots g(v_{[r]})g(v_r)$  from the input and pushing  $y_r = v_0 n_1 v_1 \cdots n_r v_r$  to the stack. Conversely, we suppose a calculation in  $A_\alpha$ , denoted by  $C_{(1)} = (q_0, w \uparrow, \#Z_0 \uparrow \$) \vdash \cdots \vdash C_{(r)} \vdash \cdots \vdash C_{(m)} = (p_m, \uparrow, \#\beta_m \$)$ , where  $p_m \in F$  and  $C_{(r)} = (p_r, w_r \uparrow, \#\beta_r \$)$  for each  $r \in \{1, \dots, m\}$ . By induction on  $r = 1, \dots, m$ , we extract an underlying transition  $q_0 \xRightarrow[N]{\gamma_r} p_r$  step by step while maintaining the invariants  $\gamma_r \in (\Sigma \uplus B_k \uplus [k])^*$  and  $w = \mathcal{D}_k(\gamma_r) w_r$ , as long as  $p_r \in Q_N$ .  $\square$

Hereafter, we shall refine this proof sketch. We first state two lemmas used to write our full proof of Theorem 18. Let  $\vdash_{(n)}$  denote the subrelation of  $\vdash$  derived from the rule numbered (n). The following lemma is immediate from the definition of  $\vdash_{(n)}$ .

**Lemma 19.** For all  $q, q' \in Q_N$ ,  $w, w' \in \Sigma^*$  and  $\gamma, \gamma' \in \Gamma^*$ ,

- (i) (a) for each  $a \in \Sigma$ ,  $(q, aw \uparrow, \#Z_0 \gamma \uparrow \$) \vdash_{(1)} (q', w \uparrow, \#Z_0 \gamma a \uparrow \$)$  if  $q \xrightarrow[N]{a} q'$ ,
- (b)  $\exists a \in \Sigma. q \xrightarrow[N]{a} q' \wedge w = aw' \wedge \beta = Z_0 \gamma a \uparrow$  if  $(q, w \uparrow, \#Z_0 \gamma \uparrow \$) \vdash_{(1)} (q', w' \uparrow, \#\beta \$)$ ,
- (ii) (a) for each  $b \in B_k$ ,  $(q, w \uparrow, \#Z_0 \gamma \uparrow \$) \vdash_{(2)} (q', w \uparrow, \#Z_0 \gamma b \uparrow \$)$  if  $q \xrightarrow[N]{b} q'$ ,
- (b)  $\exists b \in B_k. q \xrightarrow[N]{b} q' \wedge w = w' \wedge \beta = Z_0 \gamma b \uparrow$  if  $(q, w \uparrow, \#Z_0 \gamma \uparrow \$) \vdash_{(2)} (q', w' \uparrow, \#\beta \$)$ .

In particular, letting  $\vdash_{(1),(2)} = \vdash_{(1)} \uplus \vdash_{(2)}$ , we obtain the following statement by repeating (i)(a) and (ii)(a) zero or more times: For all  $v \in (\Sigma \uplus B_k)^*$ ,  $(q, g(v) w \uparrow, \#Z_0 \gamma \uparrow \$) \vdash_{(1),(2)}^* (q', w \uparrow, \#Z_0 \gamma v \uparrow \$)$  if  $q \xrightarrow[N]{v} q'$ .

**Lemma 20.** Suppose that  $q \xrightarrow[N]{i} q'$ , and  $\gamma i$  is matching. Let  $m = \text{cnt}(\gamma i)$ . For all  $p \in Q_N$ ,  $w, w' \in \Sigma^*$  and  $\beta \in (\Gamma \uplus \{\$, \$, \uparrow\})^*$ , the following (a) and (b) are equivalent:

- (a)  $p = q'$ ,  $w = g((\gamma i)_{[m]}) w'$ , and  $\beta = Z_0 \gamma i \uparrow$ .
- (b)  $(q, w \uparrow, \#Z_0 \gamma \uparrow \$) \vdash_{(3)} (W_{q'}, w \uparrow, \#Z_0 \gamma i \uparrow \$) \vdash \cdots \vdash (p, w' \uparrow, \#\beta \$)$ , where no ID with a state in  $Q_N$  appears in the calculation  $\cdots$ .

For the proof of Lemma 20, we recall the notations  $v_{[r]}$  and  $v^{(r)}$  informally (see Section 3 for the formal definitions of  $v_{[r]}$  and  $v^{(r)}$ ). Let  $k$  be a positive integer and  $v = v_0 n_1 v_1 \dots n_m v_m$  ( $m = \text{cnt } v$ ) a matching string over  $\Sigma \uplus B_k \uplus [k]$ . For each  $r = 1, 2, \dots$ , the notation  $v_{[r]}$  denotes the string which  $\mathcal{D}_k$  scans at the  $r^{\text{th}}$  number  $n_r$  and  $v^{(r)}$  the string immediately after the  $r^{\text{th}}$  replacement (also we let  $v^{(0)} = v$ ). For example, in the case of  $v = [1a[2b]_2 2]_1 1$ ,  $\mathcal{D}_k$  processes  $v$  as  $v^{(0)} = [1a[2b]_2 2]_1 1 \rightarrow v^{(1)} = [1a[2b]_2 b]_1 1 \rightarrow v^{(2)} = [1a[2b]_2 b]_1 abb \rightarrow abbabbb$ , and therefore  $v_{[1]} = b$  and  $v_{[2]} = a[2b]_2 b$ . We can easily prove the following relationship between two notations  $v_{[r]}$  and  $v^{(r)}$  (see Lemma 10 for the formal proof): For each  $r \in \{0, 1, \dots, m\}$ ,

$$v^{(r)} = v_0 g(v_{[1]}) v_1 \dots g(v_{[r]}) v_r n_{r+1} v_{r+1} \dots n_m v_m.$$

We continue preparing some more notations for the proof. Let  $\Sigma_{\perp}^* \triangleq \Sigma^* \uplus \{\perp\}$  and for every  $w \in \Sigma_{\perp}^*$  and  $s \in \Sigma^*$ , let  $w/s$  denote the string  $w$  but with the suffix  $s$  erased if  $w$  ends with  $s$ , and otherwise  $\perp$ . To use this notation, we expand the set of all IDs  $I = Q \times \Sigma^* \{ \vdash \} \times \{ \# \} (\Gamma \uplus \{ \mathfrak{c}, \$, \uparrow \})^* \{ \$ \}$  to  $I_{\perp} \triangleq Q \times \Sigma^* \{ \vdash \} \uplus \{ \perp \} \times \{ \# \} (\Gamma \uplus \{ \mathfrak{c}, \$, \uparrow \})^* \{ \$ \}$ , and we let  $C(u)$  denote an ID  $C = (\cdot, u, \dots)$  and  $\vdash'_{(n)}$  denote the following binary relation over  $I_{\perp}$ :

$$\vdash'_{(n)} \triangleq \left\{ (C(u), C'(u/(a_0 \dots a_{d-1}))) \mid C(a_0 \dots a_k) \vdash_{(n)} C'(a_{\bar{d}} \dots a_k), u \in \Sigma^* \{ \vdash \} \uplus \{ \perp \} \right\}^7.$$

We define  $\vdash' \triangleq \bigcup_n \vdash'_{(n)}$ . Then,  $\vdash \subseteq \vdash'$  is immediate and we show that the converse partially holds, in the sense that:

**Lemma 21.** *For every string  $w \in \Sigma^*$  and  $w' \in \Sigma^*$ ,  $C(w \vdash) \vdash' C'(w' \vdash)$  implies  $C(w \vdash) \vdash C'(w' \vdash)$ .*

*Proof.* By the definition of  $\vdash'$ , there is  $(C(a_0 \dots a_k), C'(a_{\bar{d}} \dots a_k)) \in \vdash$  such that  $w' \vdash = w \vdash / (a_0 \dots a_{d-1})$ . By the definition of  $\vdash$ ,  $C(w \vdash) = C(a_0 \dots a_{d-1} w' \vdash) \vdash C'(w' \vdash)$  holds.  $\square$

**Definition 22.** Given  $C, C' \in I_{\perp}$ , we write  $C \models_{(n)} C'$  if  $\forall j, C''. C \vdash'_{(j)} C'' \iff j = n \wedge C'' = C'$ . We often omit the subscript  $(n)$  and simply write  $C \models C'$ . Note that  $C \models C'$  implies not only  $C \vdash' C'$  but also determinism:  $\forall C'' \in I_{\perp}. C \vdash' C'' \implies C' = C''$ .

**Lemma 23.** *Suppose that  $\gamma \in (\Sigma \uplus B_k \uplus [k])^*$ ,  $i \in [k]$ ,  $w \in \Sigma^*$ ,  $\beta \in (\Gamma \uplus \{ \mathfrak{c}, \$ \})^*$  and  $p \in Q_N \uplus \{ e_i \mid i \in [k] \}$ . Let  $m = \text{cnt}(\gamma i) (\geq 1)$ . If  $\gamma i$  is matching,*

$$(c_i, w \vdash, \#Z_0 \gamma \mathfrak{c} p \uparrow \$i\beta\$) \models \dots \models (r_i, w \vdash / g((\gamma i)_{[m]}), \#Z_0 \gamma \mathfrak{c} p \uparrow \$i\beta\$)$$

*holds, where no ID with a state in  $Q_N$  appears in the calculation  $\dots$ .*

*Proof.* In this proof, we sometimes write the stack representation  $\# \dots Z \uparrow \dots \$$  as  $\# \dots \uparrow Z \dots \$$  with the head-reversed arrow  $\uparrow$ . First, if  $\gamma \not\vdash [i]$ , it holds that

$$\begin{aligned} (c_i, w \vdash, \#Z_0 \gamma \mathfrak{c} p \uparrow \$i\beta\$) &\models^* (c_i, w \vdash, \#Z_0 \uparrow \gamma \mathfrak{c} p \$i\beta\$) \\ &\models (r_i, w \vdash, \#Z_0 \uparrow \gamma \mathfrak{c} p \$i\beta\$) \models^* (r_i, w \vdash, \#Z_0 \gamma \mathfrak{c} p \uparrow \$i\beta\$), \end{aligned}$$

and by  $(\gamma i)_{[m]} = \varepsilon$ , we have  $w = w/g((\gamma i)_{[m]})$ , as required. Henceforth, we assume that  $\gamma \ni [i]$  and the decomposition  $\gamma = \gamma_0 [i \gamma_1]$  ( $\gamma_1 \not\vdash [i]$ ). Moreover, we can further decompose  $\gamma_1 = \gamma_2 [i \gamma_3]$  ( $\gamma_2 \not\vdash [i]$ ,  $\gamma_3 \not\vdash [i]$ ) because  $\gamma i$  is matching. We prove by induction on  $m$ .

Case  $m = 1$ : By  $\text{cnt } \gamma = 0$ ,  $\gamma_2 \in (\Sigma \uplus B_k)^*$  follows. Letting  $w' \triangleq w/g(\gamma_2)$ , we have

$$\begin{aligned} (c_i, w \vdash, \#Z_0 \gamma_0 [i \gamma_1 \mathfrak{c} p \uparrow \$i\beta\$) &\models^* (c_i, w \vdash, \#Z_0 \gamma_0 [i \uparrow \gamma_1 \mathfrak{c} p \$i\beta\$) \models (e_i, w \vdash, \#Z_0 \gamma_0 [i \uparrow \gamma_1 \mathfrak{c} p \$i\beta\$) \\ &\models^* (e_i, w' \vdash, \#Z_0 \gamma_0 [i \gamma_2]_i \uparrow \gamma_3 \mathfrak{c} p \$i\beta\$) \models (r_i, w' \vdash, \#Z_0 \gamma_0 [i \gamma_2]_i \uparrow \gamma_3 \mathfrak{c} p \$i\beta\$) \\ &\models^* (r_i, w' \vdash, \#Z_0 \gamma_0 [i \gamma_2]_i \gamma_3 \mathfrak{c} p \uparrow \$i\beta\$). \end{aligned}$$

Therefore, the claim holds since no ID with a state in  $Q_N$  appears in this calculation and  $\gamma_2 = (\gamma i)_{[m]}$  follows from  $(\gamma i)^{(0)} = \gamma i = \gamma_0 [i \gamma_2]_i \gamma_3 i$ ,  $\gamma_3 \not\vdash [i]$ .

<sup>7</sup>We regard  $a_0 \dots a_{-1}$  as the empty string  $\varepsilon$ .

Case  $\{1, \dots, m\} \implies m+1$ : Let  $m_0 \triangleq \text{cnt } \gamma_0$  and  $l \triangleq \text{cnt } \gamma_2 (\geq 0)$ . Now,  $m_0 + l \leq m = \text{cnt } \gamma$  holds and we write  $\gamma_2 = \lambda_0 n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \lambda_l$ . We also define  $\eta_r \triangleq \gamma_0[i \lambda_0 n_{m_0+1} \cdots \lambda_{r-1} n_{m_0+r}]$  for each  $r \in \{1, \dots, l\}$ . By  $\eta_r$  being a prefix of  $\gamma i$  and Lemma 7,  $\eta_r$  is matching and  $(\eta_r)_{[m_0+r]} = (\gamma i)_{[m_0+r]}$ ,  $r \in \{1, \dots, l\}$  holds. In particular, it follows that  $n_{m_0+r} \neq i$  for every  $r$  (if there is  $r$  such that  $n_{m_0+r} = i$ ,  $\gamma_2 \supseteq \lambda_0 n_{m_0+1} \cdots \lambda_{r-1} \ni i$  holds but this contradicts  $\gamma_2 \not\supseteq i$ ). Thus, letting  $u_0 \triangleq w \cdot i$ ,  $u'_r \triangleq u_{r-1} / g(\lambda_{r-1})$ ,  $u_r \triangleq u'_r / g((\eta_r)_{[m_0+r]})$  and  $u' = u_l / g(\lambda_l)$ , we have

$$\begin{aligned}
& (c_i, w \cdot i, \#Z_0 \gamma \wp p \upharpoonright \$i\beta\$) \\
& \models^* (c_i, w \cdot i, \#Z_0 \gamma_0[i \upharpoonright \lambda_0 n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \gamma_3 \wp p \$i\beta\$) \\
& \models (e_i, u_0, \#Z_0 \gamma_0[i \upharpoonright \lambda_0 n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \gamma_3 \wp p \$i\beta\$) \\
& \models^* (e_i, u'_1, \#Z_0 \gamma_0[i \lambda_0 n_{m_0+1} \upharpoonright \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \gamma_3 \wp p \$i\beta\$) \\
& \models (c_{n_{m_0+1}}, u'_1, \#Z_0 \gamma_0[i \lambda_0 \wp r_i \upharpoonright \$n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \gamma_3 \wp p \$i\beta\$) \\
& \models^* (r_{n_{m_0+1}}, u_1, \#Z_0 \gamma_0[i \lambda_0 \wp r_i \upharpoonright \$n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \gamma_3 \wp p \$i\beta\$) \\
& \quad \text{(by } \eta_1 \text{ being matching and induction hypothesis)} \\
& \models (E_{e_i, n_{m_0+1}}, u_1, \#Z_0 \gamma_0[i \lambda_0 \wp \upharpoonright \$n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \gamma_3 \wp p \$i\beta\$) \\
& \models (L_{e_i, n_{m_0+1}}, u_1, \#Z_0 \gamma_0[i \lambda_0 n_{m_0+1} \upharpoonright \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \gamma_3 \wp p \$i\beta\$) \\
& \models (e_i, u_1, \#Z_0 \gamma_0[i \lambda_0 n_{m_0+1} \upharpoonright \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \gamma_3 \wp p \$i\beta\$) \\
& \models^* \cdots \models^* (e_i, u_l, \#Z_0 \gamma_0[i \lambda_0 n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \upharpoonright \lambda_l]_i \gamma_3 \wp p \$i\beta\$) \\
& \quad \text{(by similar calculation and induction hypothesis)} \\
& \models^* (e_i, u', \#Z_0 \gamma_0[i \lambda_0 n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \upharpoonright \gamma_3 \wp p \$i\beta\$) \\
& \models (r_i, u', \#Z_0 \gamma_0[i \lambda_0 n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \upharpoonright \gamma_3 \wp p \$i\beta\$) \\
& \models^* (r_i, u', \#Z_0 \gamma_0[i \lambda_0 n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \gamma_3 \wp p \upharpoonright \$i\beta\$),
\end{aligned}$$

and

$$\begin{aligned}
u' &= w \cdot i / g(\lambda_0) g((\eta_1)_{[m_0+1]}) g(\lambda_1) \cdots g((\eta_l)_{[m_0+l]}) g(\lambda_l) \\
&= w \cdot i / g(\lambda_0) g((\gamma i)_{[m_0+1]}) g(\lambda_1) \cdots g((\gamma i)_{[m_0+l]}) g(\lambda_l),
\end{aligned}$$

where no ID with a state in  $Q_N$  appears in this calculation. Here, we write

$$\gamma = \gamma_0[i \lambda_0 n_{m_0+1} \lambda_1 \cdots n_{m_0+l} \lambda_l]_i \gamma_3 = v_0 n_1 v_1 \cdots n_m v_m$$

and decompose its substrings as

$$v_{m_0} = \chi_0[i \lambda_0, \quad v_{m_0+l} = \lambda_l]_i \chi_1, \text{ and } \gamma_3 = \chi_1 n_{m_0+l+1} v_{m_0+l+1} \cdots n_m v_m.$$

Then, by Lemma 10 (or the relationship between two notations  $v_{[r]}$  and  $v^{(r)}$  mentioned above), we can write  $(\gamma i)^{(m)}$  as

$$v_0 \cdots \underbrace{\chi_0[i \lambda_0 g((\gamma i)_{[m_0+1]}) v_{m_0+1} \cdots g((\gamma i)_{[m_0+l]}) \lambda_l]_i \chi_1}_{v_{m_0}} \underbrace{g((\gamma i)_{[m_0+l+1]}) v_{m_0+l+1} \cdots g((\gamma i)_{[m]}) v_m}_{\gamma'_3 \triangleq} i.$$

That is, it holds that  $\gamma'_3 \triangleq \chi_1 g((\gamma i)_{[m_0+l+1]}) v_{m_0+l+1} \cdots g((\gamma i)_{[m]}) v_m \not\supseteq i$  by  $\gamma_3 \not\supseteq i$ , and we obtain  $(\gamma i)_{[m+1]} = \lambda_0 g((\gamma i)_{[m_0+1]}) \lambda_1 \cdots g((\gamma i)_{[m_0+l]}) \lambda_l$ , as shown above. Therefore, the claim holds for  $m+1$  since  $u' = w \cdot i / g((\gamma i)_{[m+1]})$ .  $\square$

*Proof of Lemma 20.* For an arbitrary  $w \in \Sigma^*$ , by Lemma 23,

$$\begin{aligned}
& (q, w \cdot i, \#Z_0 \gamma \upharpoonright \$) \vdash_{(3)} (W_{q'}, w \cdot i, \#Z_0 \gamma i \upharpoonright \$) \models_{(4)} (c_i, w \cdot i, \#Z_0 \gamma \wp q' \upharpoonright \$i\$) \\
& \models^* (r_i, w \cdot i / g((\gamma i)_{[m]}), \#Z_0 \gamma \wp q' \upharpoonright \$i\$) \models_{(14)} (E_{q', i}, w \cdot i / g((\gamma i)_{[m]}), \#Z_0 \gamma \wp \upharpoonright \$i\$)
\end{aligned}$$

$$\models_{(15)} (L_{q',i}, w \dashv / g((\gamma i)_{[m]}), \#Z_0 \gamma i \upharpoonright \$) \models_{(17)} (q', w \dashv / g((\gamma i)_{[m]}), \#Z_0 \gamma i \upharpoonright \$) \quad (*)$$

holds. Assuming (a), we can replace  $\models$  in equation  $(*)$  with  $\vdash$  by Lemma 21 because  $w/g((\gamma i)_{[m]}) = w' \in \Sigma^*$  holds, and therefore, (b) follows. Supposing (b) conversely, we have  $(q, w \dashv, \#Z_0 \gamma \upharpoonright \$) \vdash_{(3)} (W_{q'}, w \dashv, \#Z_0 \gamma i \upharpoonright \$) \vdash' \dots \vdash' (p, w' \dashv, \# \beta \$)$ , where no ID with a state in  $Q_N$  appears in either this calculation or  $(*)$  except in their leftmost and rightmost IDs. Therefore, their two calculations coincide by the determinism of  $\models$ . In particular, we obtain  $p = q'$ ,  $w' \dashv = w \dashv / g((\gamma i)_{[m]})$  and  $\beta = Z_0 \gamma i \upharpoonright$  by the equality of their rightmost IDs, and thus, (a) follows because  $w' \in \Sigma^*$ .  $\square$

*Proof of Theorem 18.* Taking any  $w \in L(\alpha)$ , we have  $v = v_0 n_1 v_1 \dots n_m v_m \in \mathcal{R}_k(\alpha) = L(N)$  (here,  $m = \text{cnt } v$ ) such that  $w = \mathcal{D}_k(v)$ . Now, suppose that  $q_0 \xrightarrow[N]{v_0} q_{(0)} \xrightarrow[N]{n_1 v_1} q_{(1)} \xrightarrow[N]{n_2 v_2} \dots \xrightarrow[N]{n_m v_m} q_{(m)} \in F$ . Letting  $y_r \triangleq v_0 n_1 v_1 \dots n_r v_r$  and  $z_r \triangleq g(v_0) g(v_{[1]}) g(v_1) \dots g(v_{[r]}) g(v_r)$  for each  $r \in \{0, 1, \dots, m\}$ , we show by induction on  $r$  that for any  $w' \in \Sigma^*$ ,  $(q_0, z_r w' \dashv, \#Z_0 \upharpoonright \$) \vdash^* (q_{(r)}, w' \dashv, \#Z_0 y_r \upharpoonright \$)$  holds.

Case  $r = 0$ : It holds that  $q_0 \xrightarrow[N]{v_0} q_{(0)}$  and  $y_0 = v_0 \in (\Sigma \uplus B_k)^*$ . Letting  $\gamma = \varepsilon$  in Lemma 19, we obtain  $(q_0, g(v_0) w' \dashv, \#Z_0 \upharpoonright \$) \vdash^* (q_{(0)}, w' \dashv, \#Z_0 y_0 \upharpoonright \$)$ , as required.

Case  $r \Rightarrow r+1$ : By the induction hypothesis, it holds that  $(q_0, z_{r+1} w' \dashv, \#Z_0 \upharpoonright \$) \vdash^* (q_{(r)}, g(v_{[r+1]}) g(v_{r+1}) w' \dashv, \#Z_0 y_r \upharpoonright \$)$ . By Lemma 8,  $v$  is matching. Hence,  $v$ 's prefix  $y_r n_{r+1}$  is also matching by Lemma 7. Let  $q_{(r)} \xrightarrow[N]{n_{r+1} v_{r+1}} q_{(r+1)}$  be  $q_{(r)} \xrightarrow[N]{n_{r+1}} q' \xrightarrow[N]{v_{r+1}} q_{(r+1)}$ . Because  $\text{cnt}(y_r n_{r+1}) = r+1$ , the following calculation holds by Lemma 20 (a)  $\Rightarrow$  (b):

$$\begin{aligned} & (q_{(r)}, g((y_r n_{r+1})_{[r+1]}) g(v_{r+1}) w' \dashv, \#Z_0 y_r \upharpoonright \$) \\ & \vdash_{(3)} (W_{q'}, g((y_r n_{r+1})_{[r+1]}) g(v_{r+1}) w' \dashv, \#Z_0 y_r n_{r+1} \upharpoonright \$) \\ & \vdash^* (q', g(v_{r+1}) w' \dashv, \#Z_0 y_r n_{r+1} \upharpoonright \$) \vdash^* (q_{(r+1)}, w' \dashv, \#Z_0 y_r n_{r+1} v_{r+1} \upharpoonright \$). \end{aligned}$$

By Lemma 7, it holds that  $(y_r n_{r+1})_{[r+1]} = v_{[r+1]}$  and  $y_r n_{r+1} v_{r+1} = y_{r+1}$ , as required. In particular, letting  $r = m$  and  $w' = \varepsilon$  in the claim, we obtain  $(q_0, w \dashv, \#Z_0 \upharpoonright \$) \vdash^* (q_{(m)}, \dashv, \#Z_0 y_m \upharpoonright \$)$  because  $w = \mathcal{D}_k(v) = z_m$  by Lemma 6. Therefore,  $w \in L(A_\alpha)$  holds since  $q_{(m)} \in F$ .

Conversely, take any  $w \in L(A_\alpha)$ . There exist  $m$  and an ID  $C_{(r)} = (p_r, w_r \dashv, \# \beta_r \$)$  for each  $r \in \{1, \dots, m\}$  such that  $C_{(1)} = (q_0, w \dashv, \#Z_0 \upharpoonright \$) \vdash \dots \vdash C_{(r)} \vdash \dots \vdash C_{(m)} = (p_m, \dashv, \# \beta_m \$)$ , where  $p_m \in F$ . We show by induction on  $r$  that for each  $r = 1, \dots, m$ , the following claim holds: if  $p_r \in Q_N$ , there is  $\gamma_r$  such that  $C_{(r)} = (p_r, w_r \dashv, \#Z_0 \gamma_r \upharpoonright \$)$ , (a)  $\gamma_r \in (\Sigma \uplus B_k \uplus [k])^*$ , (b)  $w = \mathcal{D}_k(\gamma_r) w_r$  and (c)  $q_0 \xrightarrow[N]{\gamma_r} p_r$ .

Case  $r = 1$ : It holds that  $p_1 = q_0 \in Q_N$  and  $C_{(1)} = (q_0, w \dashv, \#Z_0 \upharpoonright \$)$ . Letting  $p_1 = q_0$ ,  $w_1 = w$  and  $\gamma_1 = \varepsilon$ , we have (a)  $\gamma_1 = \varepsilon \in (\Sigma \uplus B_k \uplus [k])^*$ , (b)  $w = \mathcal{D}_k(\gamma_1) w_1$  and (c)  $q_0 \xrightarrow[N]{\gamma_1} p_1 = q_0$ , as required.

Case  $\{1, \dots, r\} \Rightarrow r+1$ : Suppose that  $p_{r+1} \in Q_N$ . We can define  $j$  as the maximum of the set  $\{1 \leq j \leq r \mid p_j \in Q_N\}$  since  $p_1 \in Q_N$ . Rules that can be applied to  $C_{(j)}$  are limited to (1), (2) and (3) because  $p_j \in Q_N$ .

In the case of (1),  $j = r$  holds because  $p_{j+1} \in Q_N$ . By  $p_r \in Q_N$  and the induction hypothesis, we have  $C_{(r)} = (p_r, w_r \dashv, \#Z_0 \gamma_r \upharpoonright \$)$ , and  $p_r$ ,  $w_r$  and  $\gamma_r$  satisfy conditions (a), (b) and (c). Hence, by Lemma 19 (i)(b), there is  $a \in \Sigma$  such that  $p_r \xrightarrow[N]{a} p_{r+1}$ ,  $w_r = a w_{r+1}$  and  $\beta_{r+1} = Z_0 \gamma_r a \upharpoonright$ . Now, we let  $\gamma_{r+1} = \gamma_r a$ . Since  $\beta_{r+1} = Z_0 \gamma_{r+1} \upharpoonright$ , (a)  $\gamma_{r+1} = \gamma_r a \in (\Sigma \uplus B_k \uplus [k])^*$  because  $\gamma_r \in (\Sigma \uplus B_k \uplus [k])^*$ , (b)  $\mathcal{D}_k(\gamma_{r+1}) w_{r+1} = \mathcal{D}_k(\gamma_r) a w_{r+1} = \mathcal{D}_k(\gamma_r) w_r = w$ , and (c)  $q_0 \xrightarrow[N]{\gamma_{r+1}} p_{r+1}$  because  $q_0 \xrightarrow[N]{\gamma_r} p_r \xrightarrow[N]{a} p_{r+1}$ . The case of (2) follows similarly as the case of (1) with Lemma 19 (ii)(b). In the case of (3),

there is  $q' \in Q_N$  such that  $p_j \xrightarrow[N]{i} q'$ . By  $p_j \in Q_N$  and the induction hypothesis, we have  $C_{(j)} = (p_j, w_j \dashv, \#Z_0 \gamma_j \upharpoonright \$)$ , and  $p_j$ ,  $w_j$  and  $\gamma_j$  satisfy the conditions (a), (b) and (c). Because  $q_0 \xrightarrow[N]{\gamma_j} p_j \xrightarrow[N]{i} q'$ ,  $\gamma_j i$  is matching by Corollary 17. Besides, it holds that  $p_{r+1} \in Q_N$  and  $(p_j, w_j \dashv, \#Z_0 \gamma_j \upharpoonright \$) \vdash_{(3)} (W_{q'}, w_j \dashv, \#Z_0 \gamma_j i \upharpoonright \$) \vdash \dots \vdash C_{(r+1)}$ , where no ID with a state in  $Q_N$  appears in the calculation  $\dots$ . Hence, by Lemma 20 (b)  $\Rightarrow$  (a), we have  $p_{r+1} = q'$ ,  $w_j = g((\gamma_j i)_{[m]}) w_{r+1}$ , and  $\beta_{r+1} = Z_0 \gamma_j i \upharpoonright$ . Now, we let  $\gamma_{r+1} = \gamma_j i$ . Since  $\beta_{r+1} = Z_0 \gamma_{r+1} \upharpoonright$ , (a)  $\gamma_{r+1} = \gamma_j i \in (\Sigma \uplus B_k \uplus [k])^*$  because  $\gamma_j \in (\Sigma \uplus B_k \uplus [k])^*$ , (b)  $\mathcal{D}_k(\gamma_{r+1}) w_{r+1} = \mathcal{D}_k(\gamma_j) g((\gamma_j i)_{[m]}) w_{r+1} = \mathcal{D}_k(\gamma_j) w_j = w$ , and (c)  $q_0 \xrightarrow[N]{\gamma_{r+1}} p_{r+1}$  because  $q_0 \xrightarrow[N]{\gamma_j} p_j \xrightarrow[N]{i} p_{r+1}$ .

Therefore, the claim holds for the case of  $r + 1$ . In particular, letting  $r = m$  in the claim, we have  $C_{(m)} = (p_m, w_m^{-1}, \#Z_0\gamma_m\uparrow\$)$ , and  $p_m, w_m$  and  $\gamma_m$  satisfy (a), (b) and (c) (note that  $p_m \in F \subseteq Q_N$ ). Because  $w_m = \varepsilon$ , it holds that  $w = \mathcal{D}_k(\gamma_m)$  and that  $q_0 \xRightarrow[N]{\gamma_m} p_m \in F$ , or  $\gamma_m \in L(N) = \mathcal{R}_k(\alpha)$ . Therefore, we have  $w \in \mathcal{D}_k(\mathcal{R}_k(\alpha)) = L(\alpha)$ .  $\square$

**Corollary 24.** *Every rew b describes an indexed language, but not vice versa.*

*Proof.* The first half follows by Theorem 18 since 1N Nested-SA and indexed grammars are equivalent [2]. The second half also follows since the class of CSLs is a subclass of indexed languages [1], and the language class of rewbs and that of CFLs are incomparable [4].  $\square$

#### 4.2. The case without a captured reference

In the case of a rew b  $\alpha$  without a captured reference (that is, one in which no reference  $\backslash i$  appears as a subexpression of an expression of the form  $(j \dots)_j$ ), we can transform  $A_\alpha$  into an NESAs  $A''_\alpha$  recognizing  $L(\alpha)$ , i.e., one that neither uses substacks nor pops its stack. First, we transform  $A_\alpha$  to a Nested-SA without substacks (i.e., SA)  $A'_\alpha$ . Inspecting how substacks are used in  $A_\alpha$ , we can drop rules (12) and (16) in  $A'_\alpha$  because there is no captured reference in  $\alpha$ . We also remove the uses of substacks from rules (3) and (4), which correspond to calling, and rules (14), (15) and (17), which correspond to returning. Namely, while  $A_\alpha$ , upon a call, stores the substack  $\phi q'\$$  that consists of just the state  $q'$  where the control should return,  $A'_\alpha$  simply pushes  $q'$  to the stack top. That is, we remove (4), (15) and (17), and change (3) and (14) to the following (3') and (14'), respectively:

$$(3') \delta_N(q, i) \ni q' \implies \delta(q, c, Z\$) \ni (c_i, S, Zi q' \$), \quad (14') \delta(r_i, c, q\$) = \{(q, S, \$)\}.$$

Furthermore, we transform  $A'_\alpha$  to an SA without stack popping (i.e., NESAs)  $A''_\alpha$ . Observe that  $A'_\alpha$  pops only when returning via (14') and popping a state that was pushed in a preceding call. Thus,  $A''_\alpha$ , rather than popping  $q'$ , leaves it on the stack, and has the modes  $c_i, e_i$  and  $r_i$  skip all state symbols on the stack except the ones at the top. Here, we only need to modify  $e_i$  since  $A_\alpha$  already skips them at  $c_i$  and  $r_i$  (rules (6) and (13)). In short, we add the new rule (9\*) and change (14') to (14''), as follows:

$$(9^*) \delta(e_i, c, q) = \{(e_i, S, R)\}, \quad (14'') \delta(r_i, c, q\$) = \{(q, S, q\$)\}.$$

This NESAs  $A''_\alpha$  whose transition function consists of the rules (1),(2),(3'),(5)–(9),(9\*), (10), (11), (13) and (14'') recognizes  $L(\alpha)$ . Therefore,

**Corollary 25.** *Every rew b without a captured reference describes a nonerasing stack language, but not vice versa.*

For the latter part of Corollary 25, we can take the language  $T \triangleq \{a^n b^n \mid n \in \mathbb{N}\}$  that can be described by an NESAs but not by any rew b [4]. Here, we show that  $T$  is described by the NESAs  $A$  given in Figure 1.

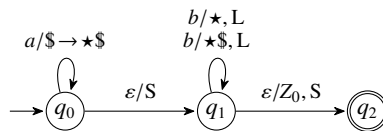


Figure 1: An NESAs that recognizes  $T = \{a^n b^n \mid n \in \mathbb{N}\}$  (see Notation 15 for the notations)

The set of stack symbols  $\Gamma$  is  $\{Z_0, \star\}$ , where  $\star$  is a distinguished character. Intuitively,  $A$  first pushes  $\star$  while consuming the input  $a$  to count the occurrences of  $a$  and leaves  $q_0$  for  $q_1$  nondeterministically. At  $q_1$ ,  $A$  moves its stack pointer leftward while consuming the input  $b$ , and leaves there for the accepting state  $q_2$  if and only if the stack pointer reaches the bottom of the stack. Finally,  $A$  halts at  $q_2$ . Here is the proof of  $L(A) = T$ :

*Proof.* The inclusion  $T \subseteq L(A)$  follows from the following calculation for any  $n \in \mathbb{N}$ :

$$\begin{aligned} (q_0, a^n b^n \uparrow, \#Z_0 \uparrow \$) &\vdash^n (q_0, b^n \uparrow, \#Z_0 \star^n \uparrow \$) \\ &\vdash (q_1, b^n \uparrow, \#Z_0 \star^n \uparrow \$) \vdash^n (q_1, \uparrow, \#Z_0 \uparrow \star^n \$) \vdash (q_2, \uparrow, \#Z_0 \uparrow \star^n \$). \end{aligned}$$

Conversely, take any  $w \in L(A)$ . By the construction of  $A$ , we can assume that  $w = a^n b^m$  (where  $n, m \in \mathbb{N}$ ) and its accepting calculation is in the following form for some  $l \in \mathbb{N}$ :

$$\begin{aligned} (q_0, a^n b^m \dashv, \#Z_0 \upharpoonright \$) \vdash^l (q_0, a^{n-l} b^m \dashv, \#Z_0 \star^l \upharpoonright \$) \\ \vdash (q_1, a^{n-l} b^m \dashv, \#Z_0 \star^l \upharpoonright \$) \vdash^* (q_1, \dashv, \#Z_0 \upharpoonright \star^l \$) \vdash (q_2, \dashv, \#Z_0 \upharpoonright \star^l \$). \end{aligned}$$

In this calculation, the steps  $\vdash^*$  indicate  $n = l = m$ . □

Notice that there exists a rew *with* a captured reference that describes a nonerasing stack language. The rew  $(1a)_1(2 \setminus 1)_2 \setminus 2$  is a simple counterexample. In addition, as shown later in Section 7, NESA can recognize nontrivial language (hierarchy) with a captured reference such as Larsen's hierarchy [17].

## 5. A rew that describes a non-stack language

We just showed that every rew describes an indexed language and in particular every rew without a captured reference describes a nonerasing stack language. So, a natural question is whether every rew describes a (nonerasing) stack language. We show that the answer is *no*. That is, there is a rew that describes a non-stack language.

Ogden has proposed a pumping lemma for stack languages and shown that the language  $\{a^{n^3} \mid n \in \mathbb{N}\}$  is a non-stack language as an application (see [20], Theorem 2). A key point in the proof is that the exponential  $n^3$  of  $a$  is a cubic polynomial, and we can show that for every cubic polynomial  $f : \mathbb{N} \rightarrow \mathbb{N}$ , the language  $\{a^{f(n)} \mid n \in \mathbb{N}\}$  is also a non-stack language by the same proof. Thus, a rew that describes a language in this form is a counterexample. We borrow the technique in [10] (Example 1) which shows that the rew  $\alpha = ((1 \setminus 2)_1(2 \setminus 1a)_2)^*$  describes  $L(\alpha) = \{a^{n^2} \mid n \in \mathbb{N}\}$ . This follows since  $\mathcal{D}_k((1 \setminus 2)_1 [2 \setminus 1a]_2)^n = a^{n^2}$  holds by recording the iteration count of the Kleene star,  $n$ , in the capture  $(2)_2$  as  $a^n$ , and extending the length by  $2n + 1$ , as shown below:

$$\begin{aligned} \mathcal{D}_k((1 \setminus 2)_1 [2 \setminus 1a]_2)^{n+1} &= \mathcal{D}_k((1 \setminus 2)_1 [2 \setminus 1a]_2)^n [1 \setminus 2]_1 [2 \setminus 1a]_2 = \mathcal{D}_k(\cdots [2a^n]_2 [1 \setminus 2]_1 [2 \setminus 1a]_2) \\ &= \mathcal{D}_k(\cdots [2a^n]_2 [1a^n]_1 [2 \setminus 1a]_2) = \mathcal{D}_k(\cdots [2a^n]_2 [1a^n]_1 [2a^{n+1}]_2) = \underline{\underline{a^{n^2} a^{2n+1}}} = a^{(n+1)^2}. \end{aligned}$$

**Theorem 26.** *There exists a rew that describes a non-stack language.*

*Proof.* The rew  $\alpha_0 \triangleq ((1 \setminus 4a)_1 (2 \setminus 3)_2 (3 \setminus 2a)_3 (4 \setminus 1 \setminus 3)_4)^*$  describes  $L_0 \triangleq \{a^{n(n+7)(2n+1)/6} \mid n \in \mathbb{N}\}$  and extends the length by a quadratic in  $n$  instead. Let  $k = 4$ . It is easy to see that  $\mathcal{R}_k(\alpha_0) = \{([1 \setminus 4a]_1 [2 \setminus 3]_2 [3 \setminus 2a]_3 [4 \setminus 1 \setminus 3]_4)^n \mid n \in \mathbb{N}\}$ . Let  $b_n$  denote  $([1 \setminus 4a]_1 [2 \setminus 3]_2 [3 \setminus 2a]_3 [4 \setminus 1 \setminus 3]_4)^n$  (therefore  $L(\alpha_0) = \{\mathcal{D}_k(b_n) \mid n \in \mathbb{N}\}$ ). We show by induction that  $\mathcal{D}_k(b_n) = \mathcal{D}_k(c_1 \cdots c_n)$  where  $c_n \triangleq [1a^{n(n+1)/2}]_1 [2a^{n-1}]_2 [3a^n]_3 [4a^{n(n+3)/2}]_4$  for every  $n \geq 1$ . First, when  $n = 1$ ,

$$\begin{aligned} \mathcal{D}_k(b_1) &= \mathcal{D}_k([1 \setminus 4a]_1 [2 \setminus 3]_2 [3 \setminus 2a]_3 [4 \setminus 1 \setminus 3]_4) = \mathcal{D}_k([1a]_1 [2 \setminus 3]_2 [3 \setminus 2a]_3 [4 \setminus 1 \setminus 3]_4) \\ &= \mathcal{D}_k([1a]_1 [2]_2 [3 \setminus 2a]_3 [4 \setminus 1 \setminus 3]_4) = \mathcal{D}_k([1a]_1 [2]_2 [3a]_3 [4 \setminus 1 \setminus 3]_4) \\ &= \mathcal{D}_k([1a]_1 [2]_2 [3a]_3 [4a]_4) = \mathcal{D}_k([1a]_1 [2]_2 [3a]_3 [4a]_4) = \mathcal{D}_k(c_1). \end{aligned}$$

Next, in the case of  $n + 1$ ,

$$\begin{aligned} \mathcal{D}_k(b_{n+1}) &= \mathcal{D}_k(b_n [1 \setminus 4a]_1 [2 \setminus 3]_2 [3 \setminus 2a]_3 [4 \setminus 1 \setminus 3]_4) \\ &= \mathcal{D}_k(c_1 \cdots c_n [1 \setminus 4a]_1 [2 \setminus 3]_2 [3 \setminus 2a]_3 [4 \setminus 1 \setminus 3]_4) \\ &= \mathcal{D}_k(c_1 \cdots c_n [1a^{n(n+3)/2}]_1 [2a^n]_2 [3a^n]_3 [4a^{n(n+3)/2}]_4) \\ &= \mathcal{D}_k(c_1 \cdots c_n [1a^{(n+1)(n+2)/2}]_1 [2a^n]_2 [3a^{n+1}]_3 [4a^{(n+1)(n+4)/2}]_4) \\ &= \mathcal{D}_k(c_1 \cdots c_n c_{n+1}). \end{aligned}$$

$$\therefore |\mathcal{D}_k(b_n)| = \sum_{i=1}^n |g(c_i)| = \sum_{i=1}^n (i^2 + 4i - 1) = \frac{n(n+7)(2n+1)}{6}. \quad (\text{also for } n = 0)$$

Therefore, we have  $L(\alpha_0) = L_0$ . □

From this theorem and Corollary 25, this rew  $\alpha_0$  needs a captured reference, in the sense that:

**Corollary 27.** *There exists a rew that describes a language that no rew without a captured reference can describe.*

The rew  $\alpha_0$  consists of mutual references by a Kleene star and the fact that it describes  $L_0$  deeply depends on the formalization of rews that we have adopted.<sup>8</sup> However, we can show that, even with only more mundane uses of backreferences that do not involve mutual references, rews can go beyond the class of stack languages. In addition, we illustrate that backreferences can be used to easily cross over the gaps between the nonerasing stack, stack, and nested stack (i.e., indexed) languages. Berglund and van der Merwe formalized mutual references on a rew with the name *circular* [4]:

**Definition 28.** Let  $\alpha$  be a rew. We define the relationship  $\rightarrow_\alpha$  as  $j \rightarrow_\alpha i$  means the rew  $\alpha$  has a subexpression in the form  $(\dots j \dots)_i$  and  $\rightarrow_\alpha^+$  as the transitive closure. The rew  $\alpha$  is called *circular* if there is a number  $i$  with  $i \rightarrow_\alpha^+ i$ .

**Definition 29.** Let  $\alpha_1, \alpha_2$  and  $\alpha_3$  be the following non-circular rews, and  $L_1, L_2$  and  $L_3$  denote their languages:

$$\begin{aligned} \alpha_1 &\triangleq (1a^*)_1 c (\backslash 1\#)^*, & L_1 &\triangleq L(\alpha_1) = \{wc(w\#)^n \mid w \in \{a\}^*, n \in \mathbb{N}\}, \\ \alpha_2 &\triangleq (1a^*)_1 c (2(\backslash 1\#)^*)_2 c \backslash 2, & L_2 &\triangleq L(\alpha_2) = \{wc(w\#)^n c(w\#)^n \mid w \in \{a\}^*, n \in \mathbb{N}\}, \\ \alpha_3 &\triangleq (1a^*)_1 c (2(\backslash 1\#)^*)_2 c \backslash 2 c \backslash 2, & L_3 &\triangleq L(\alpha_3) = \{wc(w\#)^n c(w\#)^n c(w\#)^n \mid w \in \{a\}^*, n \in \mathbb{N}\}. \end{aligned}$$

**Lemma 30.**  $L_1$  is nonerasing stack,  $L_2$  is stack but not nonerasing stack, and  $L_3$  is non-stack.

The proof is coming later. Observe that the lemma portrays a surprising power of backreferences: a slight modification of the expressions by using a simple backreference to copy a substring crosses the borders between the nonerasing stack, stack, and nested stack (i.e., indexed) languages. We obtain the following three corollaries. The first and second corollaries are stronger versions of Theorem 26 and Corollary 27, and the third corollary claims the existence of a rew language that is located between the stack and nonerasing stack languages:

**Corollary 31.** *There exists a non-circular rew that describes a non-stack language.*

**Corollary 32.** *There exists a non-circular rew that describes a language that no rew without a captured reference can describe.*

**Corollary 33.** *There exists a non-circular rew that describes a stack language but not a nonerasing stack language.*

In the rest of the section, we shall prove Lemma 30. The first claim follows from Corollary 25 because  $\alpha_1$  is free of captured references. For the other claims, we need Ogden's pumping lemmas for the stack languages and the nonerasing stack languages [20]:<sup>9</sup>

**Theorem 34.** *For an arbitrary (one-way) SA  $A$ , there exists some integer  $k$  satisfying the following conditions: for an arbitrary word  $\xi_0 \in L(A)$  with the length no less than  $k$ , there exist strings  $\mu, \nu$  and  $\rho_i, \sigma_i, \tau_i$  ( $i \geq 0$ ) that satisfy the following conditions:*

- (i)  $\xi_0 = \mu \rho_0 \sigma_0 \tau_0 \nu$ ,
- (ii) for each  $j > 0$ , it holds that  $\xi_j = \mu \rho_0 \dots \rho_j \sigma_j \tau_j \dots \tau_0 \nu \in L(A)$ ,
- (iii) there are integers  $0 < m < n, p > 0$  such that  $\rho_i, \tau_i, \sigma_i$  are of the following forms:<sup>10</sup>

$$\begin{aligned} \rho_0 &= \theta_0 & \delta_1 & \theta_1 & \dots & \delta_{m-1} & \theta_{m-1} \\ \rho_{i+1} &= \alpha_0 \beta_1^i \gamma_1 & \delta_1 & \gamma_2 \beta_2^i \alpha_1 \beta_3^i \gamma_3 & \dots & \delta_{m-1} & \gamma_{2m-2} \beta_{2m-2}^i \alpha_{m-1}, \\ \\ \tau_0 &= \theta_m & \delta_m & \theta_{m+1} & \dots & \delta_{n-2} & \theta_{n-1} \\ \tau_{i+1} &= \alpha_m \beta_{2m-1}^i \gamma_{2m-1} & \delta_m & \gamma_{2m} \beta_{2m}^i \alpha_{m+1} \beta_{2m+1}^i \gamma_{2m+1} & \dots & \delta_{n-2} & \gamma_{2n-4} \beta_{2n-4}^i \alpha_{n-1}, \\ \\ \sigma_0 &= \chi_0 & & \chi_1 & \dots & & \chi_{p-1} \\ \sigma_{i+1} &= \chi_0 & \psi_1^{i+1} & \chi_1 & \dots & \psi_{p-1}^{i+1} & \chi_{p-1}, \end{aligned}$$

<sup>8</sup>For instance, if we adopt the semantics of rews given in the ECMAScript 2023 specification [9],  $\alpha_0$  would describe a different language because that semantics processes a Kleene star expression  $\alpha^*$  by iterating the matching of the input string against  $\alpha$  and “resets” the strings captured within  $\alpha$  to  $\epsilon$  before each iteration.

<sup>9</sup>Our excerpts (Theorems 34 and 35) work as limited versions of the original pumping lemmas, which are more general statements.

<sup>10</sup>The ellipsis in  $\rho_{i+1}$  continues as follows:  $\delta_2 \gamma_4 \beta_4^i \alpha_2 \beta_5^i \gamma_5 \delta_3 \dots \delta_{m-2} \gamma_{2m-4} \beta_{2m-4}^i \alpha_{m-2} \beta_{2m-3}^i \gamma_{2m-3}$ . As well for  $\tau_{i+1}$ .



where  $\alpha_j, \beta_j, \gamma_j, \delta_j, \theta_j, \psi_j, \chi_j \in \Sigma^*$ ,

(iv)  $|\rho_0|$  (resp.  $|\tau_0|$ ) = 0  $\iff |\rho_i|$  (resp.  $|\tau_i|$ ) = 0 ( $i > 0$ ),

(v)  $|\psi_j| > 0$  and  $|\delta_j| > 0$ ,

(vi) At least one of the following conditions holds:

(a)  $|\rho_0| > 0$ ,

(b)  $|\tau_0| > 0$ ,

(c)  $p \geq 2$  and at least two  $|\chi_j| > 0$ ,

(vii) The following inequality holds:

$$\sum_j |\beta_j| + \sum_j |\psi_j| < k.$$

We obtain the pumping lemma for the nonerasing stack languages by trimming  $\tau_j$  and  $\nu$  from this theorem:

**Theorem 35.** *For an arbitrary (one-way) NESLA  $A$ , there exists some integer  $k$  satisfying the following conditions: for an arbitrary word  $\xi_0 \in L(A)$  with the length no less than  $k$ , there exist strings  $\mu$  and  $\rho_i, \sigma_i$  ( $i \geq 0$ ) that satisfy the following conditions:*

(i)  $\xi_0 = \mu\rho_0\sigma_0$ ,

(ii) for each  $j > 0$ , it holds that  $\xi_j = \mu\rho_0 \cdots \rho_j \sigma_j \in L(A)$ ,

(iii) there are integers  $m, p > 0$  such that  $\rho_i, \sigma_i$  are of the following forms:

$$\begin{aligned} \rho_0 &= \theta_0 & \delta_1 & \theta_1 & \cdots & \delta_{m-1} & \theta_{m-1} \\ \rho_{i+1} &= \alpha_0 \beta_1^i \gamma_1 & \delta_1 & \gamma_2 \beta_2^i \alpha_1 \beta_3^i \gamma_3 & \cdots & \delta_{m-1} & \gamma_{2m-2} \beta_{2m-2}^i \alpha_{m-1}, \\ \\ \sigma_0 &= \chi_0 & & \chi_1 & \cdots & & \chi_{p-1} \\ \sigma_{i+1} &= \chi_0 & \psi_1^{i+1} & \chi_1 & \cdots & \psi_{p-1}^{i+1} & \chi_{p-1}, \end{aligned}$$

where  $\alpha_j, \beta_j, \gamma_j, \delta_j, \theta_j, \psi_j, \chi_j \in \Sigma^*$ ,

(iv)  $|\rho_0|$  (resp.  $|\tau_0|$ ) = 0  $\iff |\rho_i|$  (resp.  $|\tau_i|$ ) = 0 ( $i > 0$ ),

(v)  $|\psi_j| > 0$  and  $|\delta_j| > 0$ ,

(vi) At least one of the following conditions holds:

(a)  $|\rho_0| > 0$ ,

(b)  $p \geq 2$  and at least two  $|\chi_j| > 0$ ,

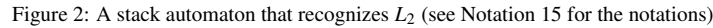
(vii) The following inequality holds:

$$\sum_j |\beta_j| + \sum_j |\psi_j| < k.$$

Let us explain the proof idea for  $L_2$  and  $L_3$ . Assume that  $L_2$  is a nonerasing stack language and fix an integer  $k$  satisfying the conditions of Theorem 35. Let  $\xi_0 \triangleq wc(w\#)^k c(w\#)^k \in L_2$  where  $w$  is a string of length no less than  $k$ , for example  $w = a^k$ . The first key point is considering the position of the second  $c$  of  $\xi_0$ . A word in  $L_2$  has the property that the prefix up to the second  $c$  determines the entire string. Therefore, the second  $c$  of  $\xi_0$  is not in  $\mu$  or  $\rho_0$ ; if so,  $\xi_1 = \xi_0$  would hold because both  $\xi_1$  and  $\xi_0$  would have a common prefix up to the second  $c$ , but  $\xi_1$  must be strictly longer than  $\xi_0$  because of the conditions of Theorem 35. Thus, it must be the case that  $c \in \sigma_0$ . In essence, in the case of  $L_2$ , the substring  $(w\#)^n$  has all information of the word it accepts, hence there is a prefix property in the sense that any two words that have this substring in common must be equal. Moreover, because  $\xi_j$ 's have a common prefix  $\mu_0\rho_0$ , the second  $c$  must not be in this prefix. For proving the claim that  $L_3$  is not a stack language, we also use a suffix property of  $L_3$  and derive that both the second and third  $c$  must be in  $\sigma_0$ . The second key point is that the informative substring  $(w\#)^n$  will be included in  $\sigma_0$ . From this, we can show that  $\sigma_1$  must be longer than  $\sigma_0$  by at least  $k$ . However, this contradicts the last condition of Theorem 35. A similar argument will also be used in the proof for  $L_3$ .

Let  $s, t$  be strings. We write  $s \leq t$  if there is a partition  $s = s_1 \cdots s_r$  and strings  $u_0, \dots, u_r$  such that  $t = u_0 s_1 u_1 \cdots s_r u_r$ .

However, this contradicts the last condition of the pumping lemma. Therefore,  $L_2$  must not be nonerasing stack. The fact that  $L_2$  is stack follows by the SA shown in Figure 2.  $\square$



## 18

**Claim 37.** Let  $\alpha_1$  and  $\alpha_2$  denote rewbs. Without loss of generality, we assume  $\alpha_1$  and  $\alpha_2$  have no common capturing group indexes. Fix a fresh symbol  $\$$  not appearing in  $\alpha_1$  and  $\alpha_2$ . The  $\text{rewbl}_p$   $(?= \alpha_1 \$) \alpha_2 \$$  (resp.  $\text{rewbl}_n$   $(?! (?! \alpha_1 \$)) \alpha_2 \$$ ) describes the language  $L(\alpha_1 \$) \cap L(\alpha_2 \$)$ .

Note that whether the language class of  $\text{rewbl}_p$ s is closed under intersections remains open but those of  $\text{rewbl}_n$ s and  $\text{rewbl}$ s are known to be closed under intersections [8].

**Lemma 38.** For an arbitrary recursively enumerable language  $E$ , there exist a  $\text{rewbl}_p$   $\alpha$  and a GSM mapping<sup>12</sup>  $g$  such that  $E = g(L(\alpha))$ .

*Proof.* This proof is quite similar to the proofs in [12, 7]. Fix a Turing machine  $M = (Q, \Sigma, \Gamma, \delta, q_0, \sqcup, F)$  that recognizes  $E$ , where  $\sqcup \in \Gamma$  is the blank symbol and  $F$  is a set of accepting states (see [15] for the definition of the other symbols and semantics). Without loss of generality, we can assume that the tape of  $M$  is semi-infinite and  $M$  never attempts to move left at the leftmost tape position. Let  $c, \phi, \$$  be fresh symbols and  $\bigvee_{i \in [l]} e_i$  denote  $e_1 + e_2 + \dots + e_l$ . Define rewbs

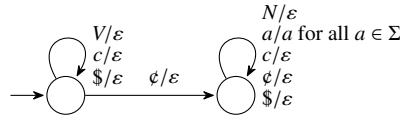
$$\alpha_1 \triangleq \left( \bigvee_{\substack{q \in Q, a \in \Gamma \\ \delta(q, a) = (p, b, L)}} ({}_1\Gamma^*)_1 ({}_2\Gamma)_2 q a ({}_3\Gamma^*)_3 c \setminus 1p \setminus 2b \setminus 3c + \bigvee_{\substack{q \in Q, a \in \Gamma \\ \delta(q, a) = (p, b, R)}} ({}_1\Gamma^*)_1 q a ({}_2\Gamma^*)_2 c \setminus 1bp \setminus 2 \sqcup^* c \right) \phi \Sigma^* \$,$$

$$\alpha_2 \triangleq q_0 ({}_1\Sigma^*)_1 \sqcup c ({}_2\Gamma \uplus Q)_2 c \setminus 2c)^* \Gamma^* F \Gamma^* c \phi \setminus 1\$.$$

Notice that a string that  $\alpha_1$  matches stands for a valid computation of  $M$  that is separated by the delimiter  $c$  and ends with  $\phi \dots \$$ , and  $\alpha_2$  imposes three constraints on the strings it matches: (1) the first component must be equal to the initial configuration of  $M$ , (2) for each  $m \geq 1$  both  $2m^{\text{th}}$  and  $(2m + 1)^{\text{st}}$  components must coincide, and (3) the penultimate component must be an accepting configuration. A subtle trick is the extra copy of the first component surrounded by  $\phi$  and  $\$$ , which lies in the last component. It is used to indicate by  $\phi$  the cutting position for a GSM and by  $\$$  where the string ends for a positive lookahead. Thus, we can easily show that  $L(\alpha_1) \cap L(\alpha_2)$  is

$$\left\{ q_0 w \sqcup c t_1 c t_1 c \dots t_{m-1} c t_{m-1} c t_m c \phi w \$ \mid q_0 w \sqcup \vdash t_1 \vdash \dots \vdash t_m \text{ is an accepting computation, } m \geq 1, w \in \Sigma^* \right\}.$$

Let  $g$  be the GSM mapping defined as follows:



By Claim 37,  $E = g(L((?= \alpha_1) \alpha_2))$  holds. □

**Corollary 39.** None of the classes of  $\text{rewbl}_p$ ,  $\text{rewbl}_n$ , and  $\text{rewbl}$  is a subclass of indexed languages.

*Proof.* Since the class of indexed languages is closed under GSM mappings [1], the class of  $\text{rewbl}_p$  cannot be a subclass of indexed languages by Lemma 38. As for  $\text{rewbl}_n$ , replace  $(?= \alpha_1)$  with  $(?! (?! \alpha_1))$ . The last is immediate. □

<sup>12</sup>A generalized sequential machine (GSM) [11] is a 6-tuple  $(Q, \Sigma, \Delta, \delta, \lambda, q_0)$ , where  $Q$  is a finite set of states,  $\Sigma$  a finite set of input symbols,  $\Delta$  a finite set of output symbols,  $\delta : Q \times \Sigma \rightarrow Q$  a transition function,  $\lambda : Q \times \Sigma \rightarrow \Delta^*$  an output function, and  $q_0 \in Q$  a start state. The functions  $\delta$  and  $\lambda$  can be extended to  $\hat{\delta} : Q \times \Sigma^* \rightarrow Q$  and  $\hat{\lambda} : Q \times \Sigma^* \rightarrow \Delta^*$  with  $\hat{\delta}(q, \varepsilon) = q$ ;  $\hat{\delta}(q, wa) = \delta(\hat{\delta}(q, w), a)$  and  $\hat{\lambda}(q, \varepsilon) = \varepsilon$ ;  $\hat{\lambda}(q, wa) = \hat{\lambda}(q, w)\lambda(\hat{\delta}(q, w), a)$ , respectively. We write  $q \xrightarrow{a/x} q'$  for  $\delta(q, a) = q'$  and  $\lambda(q, a) = x$ . A GSM mapping is a mapping  $\Sigma^* \rightarrow \Delta^*$ ,  $w \mapsto \hat{\lambda}(q_0, w)$  for some GSM  $(Q, \Sigma, \Delta, \delta, \lambda, q_0)$ .

## 7. Larsen's hierarchy is within the class of nonerasing stack languages

In this section, we construct an NESA  $A_i$  that describes  $L(x_i)$ , where the rewb  $x_i$  is given by Larsen [17] and defined over the alphabet  $\Sigma = \{a'_0, a''_0, a'_r, a'_l, a''_l, a'_r, \dots\}$  as follows:  $x_0 \triangleq (a'_0 a''_0 a'_r)^*$ ,  $x_{i+1} \triangleq (a'_{i+1} (i \ x_i); a''_{i+1} \setminus i a'_r)^*$  ( $i \geq 0$ ). Our result implies that Larsen's hierarchy is within the class of nonerasing stack languages. Since Larsen showed that no rewb with its nested level less than  $i$  can describe  $L(x_i)$  [17], it also implies that for every  $i \in \mathbb{N}$ , there is a nonerasing stack language that needs a rewb of nested level at least  $i$ .<sup>13</sup>

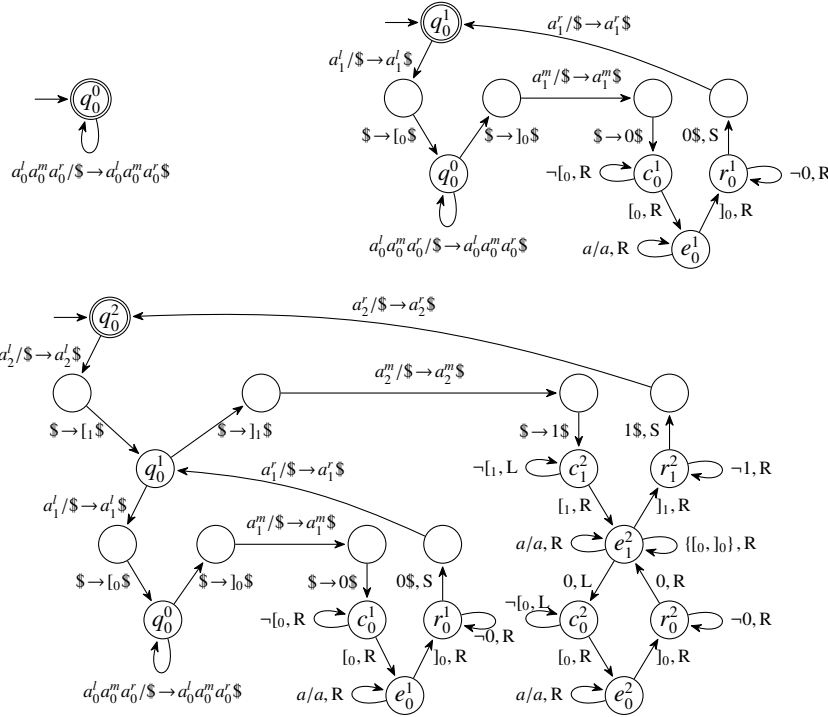


Figure 3:  $A_0$  (upper left),  $A_1$  (upper right),  $A_2$  (lower) (see Notation 15 for the notations)

The NESA  $A_i$  has the start state  $q_i^0$  which is also its only final state. Figure 3 depicts  $A_0$ ,  $A_1$ , and  $A_2$ .  $A_0$  is easy.  $A_1$  is obtained by connecting the eight states to  $q_0^0$  and making  $q_0^0$  the start/final state, as shown in the figure. The five states on the right handle the dereferencing of  $\backslash 0$  in  $x_1$ . That is, at  $c_0^1$ ,  $A_1$  first seeks the left-nearest  $[0$ , passes the control to  $e_0^1$ , checks the input string against the stack at  $e_0^1$ , passes the control to  $r_0^1$ , and at  $r_0^1$ , finally goes back to the right-nearest  $0$  which must be written on the stack top. In much the same way,  $A_2$  is obtained from  $A_1$  but we must be sensitive to the handling of the dereferencing of  $\backslash 1$  because  $A_2$  must handle the dereferencing of not only  $\backslash 1$  but also  $\backslash 0$  that appears in a string captured by  $[1]_1$  whereas no backreference appears in a string captured by  $[0]_0$  in the case of  $A_1$ . To deal with this issue, we connect the three new states  $c_0^2$ ,  $e_0^2$  and  $r_0^2$  to  $e_1^1$ . At  $e_1^1$ , if  $A_2$  encounters  $0$  in a checking,  $A_2$  suspends the checking and first goes to  $c_0^2$  to seek  $[0$ , goes to  $e_0^2$  to check the input against the stack by reading out a  $]_0$  (no number appears in this checking), and finally goes to  $r_0^2$  to go back to  $0$  which passed the control to  $c_0^2$ . We repeat this modification until  $A_i$  is obtained. (Thus,  $A_i$  has such states  $c_j^i, e_j^i, r_j^i$  for each  $j \in \{0, \dots, i-1\}$ .) Therefore,

**Theorem 40.** *There exists an NESAs  $A_i$  that recognizes  $L(x_i)$ .*

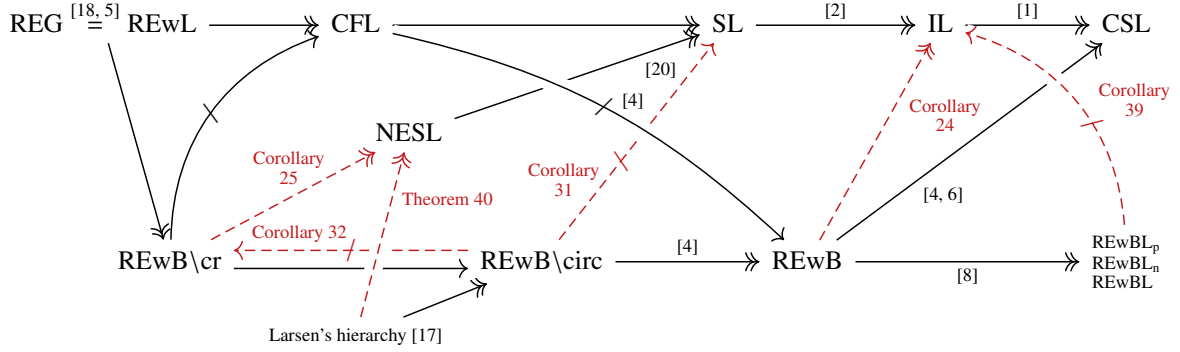


Figure 4: The inclusion relations between the language classes

## 8. Conclusions

In this paper, we have shown the following seven results: (1) that every rew describes an indexed language (Corollary 24), (2) in particular that every rew without a captured reference describes a nonerasing stack language (Corollary 25), however that there exists (3) a (non-circular) rew that describes a non-stack language (Theorem 26 and Corollary 31) and (4) a (non-circular) rew that describes a stack language but not a nonerasing stack language (Corollary 33), (5) therefore that there exists a (non-circular) rew that needs a captured reference (Corollaries 27 and 32), (6) in contrast with (1) that a rew with lookaheads can describe a non-indexed language (Corollary 39), and (7) finally that Larsen’s hierarchy  $\{L(x_i) \mid i \in \mathbb{N}\}$  given in [17] is within the class of nonerasing stack languages (Theorem 40). We have obtained the results by using three automata models, namely NESAs, SAs, and Nested-SAs, and using the semantics of rews given in [21, 10] that treats a rew as a regular expression allowing us to obtain the underlying NFA.

Figure 4 depicts the inclusion relations between the classes mentioned in the paper. Here, REG stands for the class of regular languages, CFL context free languages, NESL nonerasing stack languages, SL stack languages, IL indexed languages, CSL context sensitive languages, REwL the language class of regular expressions with lookaheads, REwB the language class of rews, REwB\cr the language class of rews without a captured reference, REwB\circ the language class of rews without a circular, REwBL<sub>p</sub> (resp. REwBL<sub>n</sub>, REwBL) the language class of rews with positive (resp. negative, both positive and negative) lookaheads. For the arrows,  $A \rightarrow B$  stands for  $A \subseteq B$ ,  $A \twoheadrightarrow B$  for  $A \subsetneq B$ , and  $A \nrightarrow B$  for  $A \not\subseteq B$ , respectively. A label on an arrow refers to the evidence. A red dashed arrow indicates a novel result proved in this paper, where for a strict inclusion, we show for the first time the inclusion itself in addition to the fact that it is strict.

As future work, we would like to investigate the use of the pumping lemma for rews without a captured reference that can be derived from the contraposition of our Corollary 25 and a pumping lemma for NESAs [20]. We expect it to be a useful tool for discerning which rews need captured references. Additionally, we suspect that our construction of NESAs in Theorem 40 is useful for not just  $x_i$  of [17] but also for more general rews that have only one  $\backslash i$  for each  $(i)_i$ , and we would like to investigate further uses of the construction.

## References

- [1] Alfred V Aho. Indexed grammars—an extension of context-free grammars. *Journal of the ACM (JACM)*, 15(4):647–671, 1968.
- [2] Alfred V Aho. Nested stack automata. *Journal of the ACM (JACM)*, 16(3):383–406, 1969.
- [3] Alfred V. Aho. *Algorithms for finding patterns in strings*, page 255–300. MIT Press, Cambridge, MA, USA, 1991.
- [4] Martin Berglund and Brink van der Merwe. Re-examining regular expressions with backreferences. *Theoretical Computer Science*, 940:66–80, 2023.

<sup>13</sup>Technically, Larsen [17] adopts a syntax that excludes unbound references, and so this implied result applies only to rews with no unbound references.

- [5] Martin Berglund, Brink van Der Merwe, and Steyn van Litsenborgh. Regular expressions with lookahead. *Journal of universal computer science (Online)*, 27(4):324–340, 2021.
- [6] 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.
- [7] Benjamin Carle and Paliath Narendran. On extended regular expressions. In *Language and Automata Theory and Applications: Third International Conference, LATA 2009, Tarragona, Spain, April 2-8, 2009. Proceedings 3*, pages 279–289. Springer, 2009.
- [8] Nariyoshi Chida and Tachio Terauchi. On lookaheads in regular expressions with backreferences. In *7th International Conference on Formal Structures for Computation and Deduction (FSCD 2022)*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2022.
- [9] ECMA International. EcmaScript® 2023 language specification, 2022. <https://tc39.es/ecma262/multipage/\#sec-intro>.
- [10] Dominik D Freydenberger and Markus L Schmid. Deterministic regular expressions with back-references. *Journal of Computer and System Sciences*, 105:1–39, 2019.
- [11] Seymour Ginsburg. *The mathematical theory of context-free languages*. McGraw-Hill Book, 1966.
- [12] Seymour Ginsburg, Sheila A Greibach, and Michael A Harrison. One-way stack automata. *Journal of the ACM (JACM)*, 14(2):389–418, 1967.
- [13] Seymour Ginsburg, Sheila A Greibach, and Michael A Harrison. Stack automata and compiling. *Journal of the ACM (JACM)*, 14(1):172–201, 1967.
- [14] Takeshi Hayashi. On derivation trees of indexed grammars—an extension of the uvwxy-theorem—. *Publications of the Research Institute for Mathematical Sciences*, 9(1):61–92, 1973.
- [15] John E Hopcroft, Rajeev Motwani, and Jeffrey D Ullman. *Introduction to automata theory, languages, and computation, 2nd Edition*. Addison-Wesley, 2001.
- [16] John E. Hopcroft and Jeffrey D. Ullman. Nonerasing stack automata. *Journal of Computer and System Sciences*, 1(2):166–186, 1967.
- [17] Kim S Larsen. Regular expressions with nested levels of back referencing form a hierarchy. *Information Processing Letters*, 65(4):169–172, 1998.
- [18] Akimasa Morihata. Translation of regular expression with lookahead into finite state automaton. *Computer Software*, 29(1):147–158, 2012.
- [19] Taisei Nogami and Tachio Terauchi. On the expressive power of regular expressions with backreferences. In *48th International Symposium on Mathematical Foundations of Computer Science (MFCS 2023)*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2023.
- [20] William F Ogden. Intercalation theorems for stack languages. In *Proceedings of the first annual ACM symposium on Theory of computing*, pages 31–42, 1969.
- [21] Markus L Schmid. Characterising regex languages by regular languages equipped with factor-referencing. *Information and Computation*, 249:1–17, 2016.