

Measuring the Expressive Power of Practical Regular Expressions by Classical Stacking Automata Models[★]

Taisei Nogami^{a,*}, Tachio Terauchi^{a,*}

^aWaseda 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 Σ 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, α first *captures* a preceding string $w \in L((a+b)^*)$ by $({}_1)_1$, and second *references* that w by following $\setminus 1$. Therefore, α 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.

*This work was supported by JSPS KAKENHI Grant Numbers JP20K20625 and JP23K24826.

*Corresponding authors.

¹One-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 b 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 b describe a (nonerasing) stack language? We show that the answer is no. Namely, we show a rew b that describes a non-stack language. Simultaneously, we also show a rew b that describes a stack language but not a nonerasing stack language.²

Additionally, we investigate the relationship between stacking automata models and rew b extended with another popular practical extension called *lookaheads*.³ Namely, does every rew b 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 b and showed that they give rise to a concrete increasing hierarchy of expressive powers of rew bs by exhibiting, for each nested level $i \in \mathbb{N}$, a language L_i that is expressible by a rew b 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 b describes an indexed language. (Section 4, Corollary 24)
- (b) Every rew b without a captured reference describes a nonerasing stack language. (Section 4, Corollary 25)
- (c) There exists a rew b that describes a non-stack language. (Section 5, Theorem 26 and Corollary 31)
- (d) There exists a rew b that describes a stack language but not a nonerasing stack language. (Section 5, Corollary 33)
- (e) A rew b 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 b language that no rew b 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 b, 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 b that describes a non-stack language. The non-stack rew b $\alpha_0 = ((_1 \backslash 4 a)_1 (_2 \backslash 3)_2 (_3 \backslash 2 a)_3 (_4 \backslash 1 \backslash 3)_4)^*$ given in the conference version deeply relied on mutual backreferences and was somewhat artificial. We present a new non-stack rew b that does not depend on such a complex structure. Furthermore, we show that a slight modification of the rew b yields the *first* rew b 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 rew bs extended with lookaheads are not always indexed languages.

2. Related work

First, we discuss related work on rew bs. There are several variants of the syntax and semantics of rew bs 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

²Note 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].

³See 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 (ε -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 ε -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 ε -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 ε -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 Cămpeanu 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, NESAs, and Nested-SAs. Ginsburg et al. introduced SA as a mathematical model that is more powerful than pushdown automaton (PDA), and NESAs as a restricted version of SA [13]. Hopcroft and Ullman discovered a type of Turing machine corresponding to the class of two-way NESAs [16]. Ogden proposed a pumping lemma for stack languages and non-erasing stack languages [20]. Aho proposed Nested-SAs with a proof of the fact that (1N)Nested-SAs and indexed grammars given by himself in [1] are equivalent in their expressive powers, and recognized PDA and SA as special cases of Nested-SAs [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 Σ , 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 $i \in [k]$:

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

Let 0-rewb denote the set of all regular expressions over Σ . 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 $i \in [k]$:

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

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

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

That is, we first regard a rew α over Σ 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 $\lceil v \rceil_i$ if $\lceil v \rceil_i$ exists (if this $\lceil v \rceil_i$ has no corresponding $\lceil v \rceil_i$, \mathcal{D}_k becomes undefined). Otherwise, \mathcal{D}_k replaces this i with ε . 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 $v_{\lceil r \rceil}$ denote the string which \mathcal{D}_k scans at the r^{th} number $n_r \in [k]$ at the r^{th} loop.

1. $\lceil \lceil a \rceil \lceil \lceil 2b \rceil \rceil \lceil \lceil 2 \rceil \rceil \lceil \lceil 1 \rceil \rceil$. In this example, \mathcal{D}_k encounters $n_1 = 2$ first, and this 2 corresponds the preceding $\lceil 2b \rceil_2$, therefore this 2 is replaced with $v_{\lceil 1 \rceil} = b$. As a result, the input string becomes $\lceil \lceil a \rceil \lceil \lceil 2b \rceil \rceil b \lceil \lceil 1 \rceil \rceil$. \mathcal{D}_k repeats this process again. Now, \mathcal{D}_k locates $n_2 = 1$ corresponding the preceding $\lceil \lceil a \rceil \lceil \lceil 2b \rceil \rceil b \lceil \lceil 1 \rceil \rceil$, so this 1 is replaced with $v_{\lceil 2 \rceil} = a \lceil 2b \rceil_2 b$ but with the brackets erased. Therefore, we gain $\lceil \lceil a \rceil \lceil \lceil 2b \rceil \rceil b \lceil \lceil 1 \rceil \rceil abb$. Finally, \mathcal{D}_k removes all remaining brackets and produces $abbabb$. Here is the diagram: $\lceil \lceil a \rceil \lceil \lceil 2b \rceil \rceil \lceil \lceil 2 \rceil \rceil \lceil \lceil 1 \rceil \rceil \rightarrow \lceil \lceil a \rceil \lceil \lceil 2b \rceil \rceil b \lceil \lceil 1 \rceil \rceil \rightarrow \lceil \lceil a \rceil \lceil \lceil 2b \rceil \rceil b \lceil \lceil 1 \rceil \rceil abb \rightarrow abbabb$.
2. $\lceil \lceil 1a \rceil \rceil \lceil \lceil 1 \rceil \rceil \lceil \lceil 1bb \rceil \rceil \lceil \lceil 1 \rceil \rceil$. In this example, $n_1 = n_2 = 1$, $v_{\lceil 1 \rceil} = a$, $v_{\lceil 2 \rceil} = bb$, and

$$\lceil \lceil 1a \rceil \rceil \lceil \lceil 1 \rceil \rceil \lceil \lceil 1bb \rceil \rceil \lceil \lceil 1 \rceil \rceil \rightarrow \lceil \lceil 1a \rceil \rceil a \lceil \lceil 1bb \rceil \rceil \lceil \lceil 1 \rceil \rceil \rightarrow \lceil \lceil 1a \rceil \rceil a \lceil \lceil 1bb \rceil \rceil bb \rightarrow aabbabb$$

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

Note that an unbound reference is replaced with the empty string ε , that is, we adopt the ε -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 α 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 ε 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

$$\forall r \in \{1, \dots, m\}. \forall x_1, x_2. y_{r-1} = x_1 \lceil n_r x_2 \rceil \implies (\exists x'_2, x_3. x_2 = x'_2 \lceil n_r x_3 \rceil \wedge x'_2 \not\in \lceil n_r \rceil, \lceil n_r \rceil)$$

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

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

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, $x_{\lceil r \rceil} = v_{\lceil r \rceil}$ holds for each $r = 1, \dots, \text{cnt } x (\leq \text{cnt } v)$.*

Lemma 8. *Every $v \in \Sigma_k^{\lceil * \rceil}$ is matching.*

In the rest of this subsection, we formally define the dereferencing function \mathcal{D}_k and the notation $v_{\lceil r \rceil}$, 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:

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

- 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 $\lceil i$, then go to Step 3. If such $\lceil i$ 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 .”

The order of execution of \mathcal{D}_k is $(12(345 + 5))^*(6 + 237)$. Observe that \mathcal{D}_k halts for any input because Step 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 Steps 1,2,3,4,5 or Steps 1,2,5) 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^{th} loop, and let $v_{[r]}$ be the string over $\Sigma \uplus B_k$ defined as follows:*

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

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

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

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 Step 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 Step 2, there are two cases:
 - When \mathcal{D}_k follows Steps 3,4,5, since it moves from Step 2 to Step 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 Step 4 and Step 5, the strings written on the tape after each step are

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

By definition, this string after Step 5 is $v^{(r+1)}$.
 - When \mathcal{D}_k follows Step 5, since $v_{[r+1]} = \varepsilon$, the string immediately after the execution of Step 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 Steps 1,2,3 to Step 7 at the $(m' + 1)^{\text{st}}$ loop, P_3 does not hold for $v^{(m')}$. By Lemma 10,

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

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 Step 3. □

Proof of Lemma 6. In Lemma 11, the string $v^{(m)}$, which is written on the tape immediately after the m^{th} loop, becomes $g(v^{(m)})$ at Step 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) $\mathcal{R}_k(\alpha) \subseteq (\Sigma \uplus \{[j,]_j \mid j \in \text{var}(\alpha)\} \uplus \text{var}(\alpha))^*$,
- (b) $\forall v_1, v_2. v_1[i]v_2 \in \mathcal{R}_k(\alpha) \implies \exists v'_2, v_3. v_2 = v'_2[i]v_3 \wedge v'_2 \not\# [i,]_i, i$.

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 $x'_2[n_r,]_{n_r} x_3$ with $x'_2 \not\# [n_r,]_{n_r}, n_r$. Here $x'_2[n_r,]_{n_r} \not\# n_r$ holds. Therefore, $x'_2[n_r,]_{n_r}$ is a prefix of x_2 and of course $x'_2 \not\# [n_r,]_{n_r}$. Hence, v is matching. \square

3.2. Classical automata models

Next, we recall the notions of SA, NESAs, and Nested-SAs. 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, Σ 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 δ 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, Σ, 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.⁴ The transition function δ has the following two modes, where $L, S, R \notin (\Sigma \cup \Gamma) \uplus \{\vdash, \#, \$\}$, $\Delta_i \triangleq \{S, R\}$, $\Delta_s \triangleq \{L, S, R\}$ and $\Sigma' \triangleq \Sigma \uplus \{\vdash\}$:

- (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) $Q \times \Sigma' \times \{\#\} \rightarrow \mathcal{P}(Q \times \Delta_i \times \{S, R\})$.

⁴These 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-SAs, 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 SAs, we define SA to also directly have the ability.

Intuitively, δ 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 δ 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 $(Q, \Sigma, \Gamma, \delta, q_0, Z_0, \bar{\cdot}, \#, \$, F)$. An element of the set $I = Q \times \Sigma^*\{\bar{\cdot}\} \times \{\#\}(\Gamma \cup \{\bar{\cdot}\})^*\{\$\}$ is called *instantaneous description*, where the stack symbol $\bar{\cdot} \notin \Gamma$ 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 $q, q' \in Q$, $k \geq 0$, $a_0, \dots, a_k \in \Sigma \cup \{\bar{\cdot}\}$, $Z, Z_1, \dots, Z_n \in \Gamma$, $\gamma, \gamma' \in \Gamma^*$, $d \in \Delta_i$ and $e \in \Delta_s$ with $k = 0 \rightarrow d \neq R$:

- (i) $(q, a_0 \cdots a_k, \#\gamma Z \bar{\cdot} \$) \vdash_A (q', a_{\bar{d}} \cdots a_k, \#\gamma' y \bar{\cdot} \$)$ if $(q', d, y\$) \in \delta(q, a_0, Z\$)$.
- (ii) (a) $(q, a_0 \cdots a_k, \#\gamma Z \bar{\cdot} \$) \vdash_A (q', a_{\bar{d}} \cdots a_k, \#\gamma \bar{\cdot} Z \$)$ if $(q', d, L) \in \delta(q, a_0, Z\$)$, and $(q, a_0 \cdots a_k, \#\gamma Z \bar{\cdot} \$) \vdash_A (q', a_{\bar{d}} \cdots a_k, \#\gamma Z \bar{\cdot} \$)$ 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 $(q, a_0 \cdots a_k, \#Z_1 \cdots Z_j \bar{\cdot} \cdots Z_n \$) \vdash_A (q', a_{\bar{d}} \cdots a_k, \#Z_1 \cdots Z_{j+\bar{e}} \bar{\cdot} \cdots Z_n \$)$.
- (c) $(q, a_0 \cdots a_k, \#\bar{\cdot} Z \gamma \$) \vdash_A (q', a_{\bar{d}} \cdots a_k, \#\bar{\cdot} Z \bar{\cdot} \gamma \$)$ if $(q', d, R) \in \delta(q, a_0, \#)$, and $(q, a_0 \cdots a_k, \#\bar{\cdot} Z \gamma \$) \vdash_A (q', a_{\bar{d}} \cdots a_k, \#\bar{\cdot} Z \gamma \$)$ 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 $(q_0, w\bar{\cdot}, \#Z_0 \bar{\cdot} \$) \vdash_A^* (q_f, \bar{\cdot}, \#\gamma_1 \bar{\cdot} \gamma_2 \$)$. 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 $a, d / \cdots$ as a / \cdots . For convenience, we also introduce intuitive syntax sugars to represent multiple transitions at once. For instance, an ε -transition ε / \cdots consists of putting $c, S / \cdots$ for all $c \in \Sigma$, and we may omit the upper part in front of $/$ and write the part \cdots 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 $\cdots / U, e$ (resp. $\cdots / U \$, e$) means putting $\cdots / u, e$ (resp. $\cdots / u \$, e$) for all $u \in U$. In particular, $U = \Gamma$, we omit U and merely type \cdots / 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 $\#a_1 a_2 \bar{\cdot} a_3 \$$ and we are to create a new substack containing $b_1 b_2$:

$$\#a_1 \underline{\phi} b_1 b_2 \bar{\cdot} a_2 a_3 \$ \tag{1}$$

Note that the new substack $\phi b_1 b_2 \$$ 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 $\#a_1 \phi b_1 b_2 \bar{\cdot} \$$. 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:

$$\#a_1 \underline{\phi} b_1 b_2 \bar{\cdot} a_2 a_3 \$ \xrightarrow{L} \#a_1 \underline{\phi} b_1 \bar{\cdot} b_2 a_2 a_3 \$ \xrightarrow{\text{create}} \#a_1 \underline{\phi} \phi c_1 c_2 \bar{\cdot} b_1 b_2 a_2 a_3 \$, \tag{2}$$

$$\#a_1 \underline{\phi} \bar{\cdot} b_1 b_2 a_2 a_3 \$ \xrightarrow{L} \#a_1 \bar{\cdot} \phi b_1 b_2 a_2 a_3 \$ \xrightarrow{\text{create}} \# \phi c_1 c_2 \bar{\cdot} a_1 \phi b_1 b_2 a_2 a_3 \$ \tag{3}$$

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:

$$\#a_1\#b_1b_2\uparrow\$a_2a_3\$ \xrightarrow{\text{pop}} \#a_1\#b_1\uparrow\$a_2a_3\$ \xrightarrow{\text{pop}} \#a_1\#\uparrow\$a_2a_3\$ \xrightarrow{\text{destruct}} \#a_1a_2\uparrow a_3\$. \quad (4)$$

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 $(Q, \Sigma, \Gamma, \delta, q_0, Z_0, \#, \$, F)$ 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 $\# \notin \Sigma \cup \Gamma$ represents the bottom of a substack.⁵ The transition function δ has the following four modes, where $\Sigma' \triangleq \Sigma \uplus \{-\}$ and $\Gamma' \triangleq \Gamma \uplus \{\#\}$:

- (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) $Q \times \Sigma' \times (\Gamma' \uplus \Gamma'\$) \rightarrow \mathcal{P}(Q \times \Delta_i \times \{\#\} \Gamma^*\$)$.
- (iv) (stack destruction mode) $Q \times \Sigma' \times \{\#\} \$ \rightarrow \mathcal{P}(Q \times \Delta_i)$.

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

- (iii) if $(q', d, \#y\$) \in \delta(q, a_0, Z)$ and $Z = Z_j, 1 \leq j < n$, then
 - $(q, a_0 \cdots a_k, \#Z_1 \cdots Z_j \uparrow \cdots Z_n\$) \vdash_A (q', a_{i+d} \cdots a_k, \#Z_1 \cdots \#y \uparrow \$Z_j \cdots Z_n\$)$,
 - and $(q, a_0 \cdots a_k, \#yZ \uparrow \$) \vdash_A (q', a_{\bar{d}} \cdots a_k, \#y \#y \uparrow \$Z\$)$ if $(q', d, \#y\$) \in \delta(q, a_0, Z\$)$.
- (iv) $(q, a_0 \cdots a_k, \#y_1\#\uparrow \$Z\gamma_2\$) \vdash_A (q', a_{\bar{d}} \cdots a_k, \#y_1Z \uparrow \gamma_2\$)$ if $(q', d) \in \delta(q, a_0, \#\$)$.

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 α , we derive the regular language $\mathcal{R}_k(\alpha)$ over the alphabet $\Sigma \uplus B_k \uplus [k]$ 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_α recognizing the language $L(\alpha)$ as follows.

The Nested-SA A_α is based on an NFA N recognizing the language $\mathcal{R}_k(\alpha)$, in the sense that each transition in A_α comes from a corresponding transition of N . The NFA N has the alphabet $\Sigma \uplus B_k \uplus [k]$, 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_α 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_α has the transition pushing b without consuming input symbols, denoted by $q \xrightarrow{\varepsilon/\$ \rightarrow b\$} q'$.⁶ For each transition $q \xrightarrow[N]{i} q'$ with $i \in [k]$, A_α has a large “transition” that consists of several transitions. In this “transition,” A_α 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.

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_α wants to push c to the stack, A_α must leave from v , climb up the stack toward the top, and write c . However, after the push, A_α becomes lost by not knowing where to go back to. How about marking the place where A_α 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

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

⁶See also Notation 15 for the ε -transition notation.

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_α pauses the checking and puts a substack containing the current state as a marker at the stack pointer position. Then, A_α 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_α eventually reaches a string $v'' \in (\Sigma \uplus B_k)^*$ containing no characters of $[k]$. Once done with the check of v'' , A_α 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_α 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 $\alpha \in k$ -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 α , there exists a Nested-SA that recognizes $L(\alpha)$.*

The claim obviously holds when α is a pure regular expression (i.e., $\alpha \in 0$ -rewb). Suppose that $\alpha \in k$ -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 δ be the smallest relation that, for all $a \in \Sigma, b \in B_k, c \in \Sigma \uplus \{\#\}$, $i, j \in [k], q, q' \in Q_N, Z \in \Gamma$ 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 $q \xrightarrow[N]{a/\$ \rightarrow a\$} q'$, (2) translates $q \xrightarrow[N]{b} q'$ into $q \xrightarrow[N]{\epsilon/\$ \rightarrow b\$} q'$, 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_α 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_α inserts $\#q'\$$ just below i , and goes to the state c_i . The state c_i represents the *call mode* in which A_α 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_α arrives at the bottom of the stack), it proceeds to the state r_i (*return mode*) by rule (7). At e_i , A_α 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 $]_j$ is read, and the rules handle the other brackets (i.e., $[_j$ or $]_j$ with $i \neq j$) by simply skipping them (note that $]_j = [_i$ cannot happen since we started from the left-nearest $[_i$). Reading $j \in [k]$, by rule (12), A_α 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_α 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_α 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_α 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_α 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_α , 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 γi 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 $v = v_0 n_1 v_1 \cdots n_m v_m$ ($m = \text{cnt } v$) 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^{th} number n_r and $v^{(r)}$ the string immediately after the r^{th} replacement (also we let $v^{(0)} = v$). For example, in the case of $v = [1a [2b]_2]_1 1$, \mathcal{D}_k processes v as $v^{(0)} = [1a [2b]_2]_1 1 \rightarrow v^{(1)} = [1a [2b]_2 b]_1 1 \rightarrow v^{(2)} = [1a [2b]_2 b]_1 abb \rightarrow abbabb$, and therefore

$v_{[1]} = b$ and $v_{[2]} = a[2b]_2b$. We can easily prove the following relationship between two notations $v_{[r]}$ and $v^{(r)}$ (see Lemma 10 for the formal proof): For each $r \in \{0, 1, \dots, m\}$,

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

We continue preparing some more notations for the proof. Let $\Sigma_{\perp} \triangleq \Sigma^* \cup \{\perp\}$ 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 $I = Q \times \Sigma^*\{-\} \times \{\#\} (\Gamma \cup \{\mathcal{C}, \$, \uparrow\})^* \{\$\}$ to $I_{\perp} \triangleq Q \times \Sigma^*\{-\} \cup \{\perp\} \times \{\#\} (\Gamma \cup \{\mathcal{C}, \$, \uparrow\})^* \{\$\}$, and we let $C(u)$ denote an ID $C = (\cdot, u, \cdots)$ and $\vdash'_{(n)}$ denote the following binary relation over I_{\perp} (we regard $a_0 \cdots a_{-1}$ as the empty string ε):

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

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 $w' \in \Sigma^*$, $C(w \dashv) \vdash' C'(w' \dashv)$ implies $C(w \dashv) \vdash C'(w' \dashv)$.*

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

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

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

$$(c_i, w \dashv, \#Z_0 \gamma \mathcal{C} p \uparrow \$i\beta\$) \vDash \cdots \vDash (r_i, w \dashv / g((\gamma i)_{[m]}), \#Z_0 \gamma \mathcal{C} p \uparrow \$i\beta\$)$$

holds, where no ID with a state in Q_N appears in the calculation \cdots .

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

$$\begin{aligned} (c_i, w \dashv, \#Z_0 \gamma \mathcal{C} p \uparrow \$i\beta\$) &\vDash^* (c_i, w \dashv, \#Z_0 \uparrow \gamma \mathcal{C} p \$i\beta\$) \\ &\vDash (r_i, w \dashv, \#Z_0 \uparrow \gamma \mathcal{C} p \$i\beta\$) \vDash^* (r_i, w \dashv, \#Z_0 \gamma \mathcal{C} p \uparrow \$i\beta\$), \end{aligned}$$

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 $\gamma = \gamma_0 [i] \gamma_1$ ($\gamma_1 \not\# [i]$). Moreover, we can further decompose $\gamma_1 = \gamma_2 [i] \gamma_3$ ($\gamma_2 \not\# [i]$, $\gamma_3 \not\# [i]$) because γi is matching. We prove by induction on m .

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

$$\begin{aligned} (c_i, w \dashv, \#Z_0 \gamma_0 [i] \gamma_1 \mathcal{C} p \uparrow \$i\beta\$) &\vDash^* (c_i, w \dashv, \#Z_0 \gamma_0 [i] \uparrow \gamma_1 \mathcal{C} p \$i\beta\$) \vDash (e_i, w \dashv, \#Z_0 \gamma_0 [i] \uparrow \gamma_1 \mathcal{C} p \$i\beta\$) \\ &\vDash^* (e_i, w' \dashv, \#Z_0 \gamma_0 [i] \gamma_2 [i] \uparrow \gamma_3 \mathcal{C} p \$i\beta\$) \vDash (r_i, w' \dashv, \#Z_0 \gamma_0 [i] \gamma_2 [i] \uparrow \gamma_3 \mathcal{C} p \$i\beta\$) \\ &\vDash^* (r_i, w' \dashv, \#Z_0 \gamma_0 [i] \gamma_2 [i] \gamma_3 \mathcal{C} p \uparrow \$i\beta\$). \end{aligned}$$

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 $(\gamma i)^{(0)} = \gamma i = \gamma_0 [i] \gamma_2 [i] \gamma_3 i$, $\gamma_3 \not\# [i]$.

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 η_r being a prefix of γi and Lemma 7, η_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 $\gamma_2 \not\# [i]$). Thus, letting $u_0 \triangleq w \dashv$, $u'_r \triangleq u_{r-1} / g(\lambda_{r-1})$, $u_r \triangleq u'_r / g((\eta_r)_{[m_0+r]})$ and $u' = u_l / g(\lambda_l)$, we have

$$(c_i, w \dashv, \#Z_0 \gamma \mathcal{C} p \uparrow \$i\beta\$)$$

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

and

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

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

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

That is, it holds that $\gamma'_3 \triangleq \chi_1 g((\gamma i)_{[m_0+l+1]}) v_{m_0+l+1} \cdots g((\gamma i)_{[m]}) v_m \not\equiv [i$ by $\gamma_3 \not\equiv [i$, 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 $u' = w\uparrow/g((\gamma i)_{[m+1]})$. \square

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

$$\begin{aligned}
& (q, w\uparrow, \#Z_0\gamma \uparrow \$) \vdash_{(3)} (W_{q'}, w\uparrow, \#Z_0\gamma i \uparrow \$) \vdash_{(4)} (c_i, w\uparrow, \#Z_0\gamma \wp q' \uparrow \$i\$) \\
& \models^* (r_i, w\uparrow/g((\gamma i)_{[m]}), \#Z_0\gamma \wp q' \uparrow \$i\$) \vdash_{(14)} (E_{q', i}, w\uparrow/g((\gamma i)_{[m]}), \#Z_0\gamma \wp \uparrow \$i\$) \\
& \vdash_{(15)} (L_{q', i}, w\uparrow/g((\gamma i)_{[m]}), \#Z_0\gamma i \uparrow \$) \vdash_{(17)} (q', w\uparrow/g((\gamma i)_{[m]}), \#Z_0\gamma i \uparrow \$) \quad (*)
\end{aligned}$$

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 $(q, w\uparrow, \#Z_0\gamma \uparrow \$) \vdash_{(3)} (W_{q'}, w\uparrow, \#Z_0\gamma i \uparrow \$) \vdash \cdots \vdash' (p, w'\uparrow, \#\beta\$)$, where no ID with a state in Q_N appears in 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 $p = q'$, $w'\uparrow = w\uparrow/g((\gamma i)_{[m]})$ and $\beta = Z_0\gamma i \uparrow$ 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 $v = v_0 n_1 v_1 \cdots n_m v_m \in \mathcal{R}_k(\alpha) = L(N)$ (here, $m = \text{cnt } v$) such that $w = \mathcal{D}_k(v)$. Now, suppose that $q_0 \xrightarrow{v_0/N} q_{(0)} \xrightarrow{n_1 v_1/N} q_{(1)} \xrightarrow{n_2 v_2/N} \cdots \xrightarrow{n_m v_m/N} q_{(m)} \in F$. Letting $y_r \triangleq v_0 n_1 v_1 \cdots n_r v_r$ and $z_r \triangleq g(v_0) g(v_{[1]}) g(v_1) \cdots g(v_{[r]}) g(v_r)$ for each $r \in \{0, 1, \dots, m\}$, we show by induction on r that for any $w' \in \Sigma^*$, $(q_0, z_r w'^{-1}, \#Z_0 \uparrow \$) \vdash^* (q_{(r)}, w'^{-1}, \#Z_0 y_r \uparrow \$)$ holds.

Case $r = 0$: It holds that $q_0 \xrightarrow{v_0/N} q_{(0)}$ and $y_0 = v_0 \in (\Sigma \uplus B_k)^*$. Letting $\gamma = \varepsilon$ in Lemma 19, we obtain $(q_0, g(v_0) w'^{-1}, \#Z_0 \uparrow \$) \vdash^* (q_{(0)}, w'^{-1}, \#Z_0 y_0 \uparrow \$)$, as required.

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

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

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 $(q_0, w^{-1}, \#Z_0 \uparrow \$) \vdash^* (q_{(m)}, \varepsilon^{-1}, \#Z_0 y_m \uparrow \$)$ 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 $C_{(r)} = (p_r, w_r^{-1}, \#\beta_r, \$)$ for each $r \in \{1, \dots, m\}$ such that $C_{(1)} = (q_0, w^{-1}, \#Z_0 \uparrow \$) \vdash \cdots \vdash C_{(r)} \vdash \cdots \vdash C_{(m)} = (p_m, \varepsilon^{-1}, \#\beta_m, \$)$, 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 γ_r such that $C_{(r)} = (p_r, w_r^{-1}, \#Z_0 \gamma_r \uparrow \$)$, (a) $\gamma_r \in (\Sigma \uplus B_k \uplus [k])^*$, (b) $w = \mathcal{D}_k(\gamma_r) w_r$ and (c) $q_0 \xrightarrow{\gamma_r/N} p_r$.

Case $r = 1$: It holds that $p_1 = q_0 \in Q_N$ and $C_{(1)} = (q_0, w^{-1}, \#Z_0 \uparrow \$)$. 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) $q_0 \xrightarrow{\gamma_1/N} p_1 = q_0$, 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 $C_{(r)} = (p_r, w_r^{-1}, \#Z_0 \gamma_r \uparrow \$)$, and p_r, w_r and γ_r satisfy conditions (a), (b) and (c). Hence, by Lemma 19 (i)(b), there is $a \in \Sigma$ such that $p_r \xrightarrow{a/N} p_{r+1}$, $w_r = a w_{r+1}$ and $\beta_{r+1} = Z_0 \gamma_r a \uparrow$. Now, we let $\gamma_{r+1} = \gamma_r a$. Since $\beta_{r+1} = Z_0 \gamma_{r+1} \uparrow$, (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) $q_0 \xrightarrow{\gamma_{r+1}/N} p_{r+1}$ because $q_0 \xrightarrow{\gamma_r/N} p_r \xrightarrow{a/N} p_{r+1}$. 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 $p_j \xrightarrow{i/N} q'$. By $p_j \in Q_N$ and the induction hypothesis, we have $C_{(j)} = (p_j, w_j^{-1}, \#Z_0 \gamma_j \uparrow \$)$, and p_j, w_j and γ_j satisfy the conditions (a), (b) and (c). Because $q_0 \xrightarrow{\gamma_j/N} p_j \xrightarrow{i/N} q'$, $\gamma_j i$ is matching by Corollary 17. Besides, it holds that $p_{r+1} \in Q_N$ and $(p_j, w_j^{-1}, \#Z_0 \gamma_j \uparrow \$) \vdash_{(3)} (W_{q'}, w_j^{-1}, \#Z_0 \gamma_j i \uparrow \$) \vdash \cdots \vdash C_{(r+1)}$, where no ID with a state in Q_N appears in the calculation \cdots . Hence, by Lemma 20 (b) \Rightarrow (a), we have $p_{r+1} = q'$, $w_j = g((\gamma_j i)_{[m]}) w_{r+1}$, and $\beta_{r+1} = Z_0 \gamma_j i \uparrow$. Now, we let $\gamma_{r+1} = \gamma_j i$. Since $\beta_{r+1} = Z_0 \gamma_{r+1} \uparrow$, (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) $q_0 \xrightarrow{\gamma_{r+1}/N} p_{r+1}$ because $q_0 \xrightarrow{\gamma_j/N} p_j \xrightarrow{i/N} p_{r+1}$.

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 γ_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 $q_0 \xrightarrow{\gamma_m/N} p_m \in F$, 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 α without a captured reference (that is, one in which no reference $\setminus i$ appears as a subexpression of an expression of the form $(j \dots)_j$), we can transform A_α into an NESAs A'_α recognizing $L(\alpha)$, i.e., one that neither uses substacks nor pops its stack. First, we transform A_α to a Nested-SA without substacks (i.e., SA) A'_α . Inspecting how substacks are used in A_α , we can drop rules (12) and (16) in A'_α because there is no captured reference in α . 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_α , upon a call, stores the substack $\notin q' \$$ that consists of just the state q' where the control should return, A'_α 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'_α to an SA without stack popping (i.e., NESAs) A''_α . Observe that A'_α pops only when returning via (14') and popping a state that was pushed in a preceding call. Thus, A''_α , 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_α 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''_α 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 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 [4]. Here, we show that T is described by the NESAs A given in Figure 1.

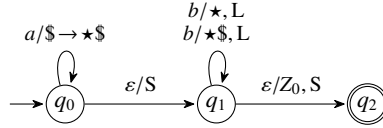


Figure 1: An NESAs that recognizes $T = \{a^n b^n \mid n \in \mathbb{N}\}$ (see Notation 15 for the notations)

The set of stack symbols Γ 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 $n \in \mathbb{N}$:

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

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, NESAs can recognize nontrivial language (hierarchy) with a captured reference such as Larsen's hierarchy [17].

5. A rewb that describes a non-stack language

We just showed that every rewb describes an indexed language and in particular every rewb without a captured reference describes a nonerasing stack language. So, a natural question is whether every rewb describes a (nonerasing) stack language. We show that the answer is *no*. That is, there is a rewb 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 rewb that describes a language in this form is a counterexample. We borrow the technique in [10] (Example 1) which shows that the rewb $\alpha = ((1 \setminus 2)_1 (2 \setminus 1 a)_2)^*$ describes $L(\alpha) = \{a^{n^2} \mid n \in \mathbb{N}\}$. This follows since $\mathcal{D}_k((\llbracket 1 \ 2 \rrbracket_1 \llbracket 2 \ 1 \ a \rrbracket_2)^n) = a^{n^2}$ 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:

$$\begin{aligned} \mathcal{D}_k((\llbracket 1 \ 2 \rrbracket_1 \llbracket 2 \ 1 \ a \rrbracket_2)^{n+1}) &= \mathcal{D}_k((\llbracket 1 \ 2 \rrbracket_1 \llbracket 2 \ 1 \ a \rrbracket_2)^n \llbracket 1 \ 2 \rrbracket_1 \llbracket 2 \ 1 \ a \rrbracket_2) = \mathcal{D}_k(\cdots \llbracket 2 \ a^n \rrbracket_2 \llbracket 1 \ 2 \rrbracket_1 \llbracket 2 \ 1 \ a \rrbracket_2) \\ &= \mathcal{D}_k(\cdots \llbracket 2 \ a^n \rrbracket_2 \llbracket 1 \ a^n \rrbracket_1 \llbracket 2 \ 1 \ a \rrbracket_2) = \mathcal{D}_k(\cdots \llbracket 2 \ a^n \rrbracket_2 \llbracket 1 \ a^n \rrbracket_1 \llbracket 2 \ a^{n+1} \rrbracket_2) = \underline{\underline{a^{n^2} a^{2n+1}}} = a^{(n+1)^2}. \end{aligned}$$

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

Proof. The rewb $\alpha_0 \triangleq ((1 \setminus 4 a)_1 (2 \setminus 3)_2 (3 \setminus 2 a)_3 (4 \setminus 1 \setminus 3)_4)^*$ 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 $\mathcal{R}_k(\alpha_0) = \{(\llbracket 1 \ 4 \ a \rrbracket_1 \llbracket 2 \ 3 \rrbracket_2 \llbracket 3 \ 2 \ a \rrbracket_3 \llbracket 4 \ 1 \ 3 \rrbracket_4)^n \mid n \in \mathbb{N}\}$. Let b_n denote $(\llbracket 1 \ 4 \ a \rrbracket_1 \llbracket 2 \ 3 \rrbracket_2 \llbracket 3 \ 2 \ a \rrbracket_3 \llbracket 4 \ 1 \ 3 \rrbracket_4)^n$ (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 $c_n \triangleq \llbracket 1 \ a^{n(n+1)/2} \rrbracket_1 \llbracket 2 \ a^{n-1} \rrbracket_2 \llbracket 3 \ a^n \rrbracket_3 \llbracket 4 \ a^{n(n+3)/2} \rrbracket_4$ for every $n \geq 1$. First, when $n = 1$,

$$\begin{aligned} \mathcal{D}_k(b_1) &= \mathcal{D}_k(\llbracket 1 \ 4 \ a \rrbracket_1 \llbracket 2 \ 3 \rrbracket_2 \llbracket 3 \ 2 \ a \rrbracket_3 \llbracket 4 \ 1 \ 3 \rrbracket_4) = \mathcal{D}_k(\llbracket 1 \ a \rrbracket_1 \llbracket 2 \ 3 \rrbracket_2 \llbracket 3 \ 2 \ a \rrbracket_3 \llbracket 4 \ 1 \ 3 \rrbracket_4) \\ &= \mathcal{D}_k(\llbracket 1 \ a \rrbracket_1 \llbracket 2 \rrbracket_2 \llbracket 3 \ 2 \ a \rrbracket_3 \llbracket 4 \ 1 \ 3 \rrbracket_4) = \mathcal{D}_k(\llbracket 1 \ a \rrbracket_1 \llbracket 2 \rrbracket_2 \llbracket 3 \ a \rrbracket_3 \llbracket 4 \ 1 \ 3 \rrbracket_4) \\ &= \mathcal{D}_k(\llbracket 1 \ a \rrbracket_1 \llbracket 2 \rrbracket_2 \llbracket 3 \ a \rrbracket_3 \llbracket 4 \ a \ 3 \rrbracket_4) = \mathcal{D}_k(\llbracket 1 \ a \rrbracket_1 \llbracket 2 \rrbracket_2 \llbracket 3 \ a \rrbracket_3 \llbracket 4 \ a \ a \rrbracket_4) = \mathcal{D}_k(c_1). \end{aligned}$$

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

$$\begin{aligned} \mathcal{D}_k(b_{n+1}) &= \mathcal{D}_k(b_n \llbracket 1 \ 4 \ a \rrbracket_1 \llbracket 2 \ 3 \rrbracket_2 \llbracket 3 \ 2 \ a \rrbracket_3 \llbracket 4 \ 1 \ 3 \rrbracket_4) \\ &= \mathcal{D}_k(c_1 \cdots c_n \llbracket 1 \ 4 \ a \rrbracket_1 \llbracket 2 \ 3 \rrbracket_2 \llbracket 3 \ 2 \ a \rrbracket_3 \llbracket 4 \ 1 \ 3 \rrbracket_4) \\ &= \mathcal{D}_k(c_1 \cdots c_n \llbracket 1 \ a^{n(n+3)/2} \rrbracket_1 \llbracket 2 \ a^n \rrbracket_2 \llbracket 3 \ a^n \rrbracket_3 \llbracket 4 \ a^{n(n+3)/2} \rrbracket_4) \\ &= \mathcal{D}_k(c_1 \cdots c_n \llbracket 1 \ a^{(n+1)(n+2)/2} \rrbracket_1 \llbracket 2 \ a^n \rrbracket_2 \llbracket 3 \ a^{n+1} \rrbracket_3 \llbracket 4 \ a^{(n+1)(n+4)/2} \rrbracket_4) \\ &= \mathcal{D}_k(c_1 \cdots c_n c_{n+1}). \end{aligned}$$

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

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

From this theorem and Corollary 25, this rewb α_0 needs a captured reference, in the sense that:

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

The rewb α_0 consists of mutual references by a Kleene star and the fact that it describes L_0 deeply depends on the formalization of rewb that we have adopted.⁷ However, we can show that, even with only more mundane uses of

⁷For instance, if we adopt the semantics of rewb given in the ECMAScript 2023 specification [9], α_0 would describe a different language because that semantics processes a Kleene star expression α^* by iterating the matching of the input string against α and “resets” the strings captured within α to ϵ before each iteration.

backreferences that do not involve mutual references, rewbs 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 rewbs with the name *circular* [4]:

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

Definition 29. Let α_1, α_2 and α_3 be the following non-circular rewbs, 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 rewbs language that is located between the stack and nonerasing stack languages:

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

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

Corollary 33. *There exists a non-circular rewbs 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 α_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]:⁸

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 μ, ν and ρ_i, σ_i, τ_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 $\xi_j = \mu\rho_0 \cdots \rho_j \sigma_j \tau_j \cdots \tau_0 \nu \in L(A)$,
- (iii) there are integers $0 < m < n, p > 0$ such that ρ_i, τ_i, σ_i are of the following forms:⁹

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

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 ($i > 0$),

⁸Our excerpts (Theorems 34 and 35) work as limited versions of the original pumping lemmas, which are more general statements.

⁹The ellipsis in ρ_{i+1} continues as follows: $\delta_2 \gamma_4 \beta_4^i \alpha_2 \beta_5^i \gamma_5 \delta_3 \cdots \delta_{m-2} \gamma_{2m-4} \beta_{2m-4}^i \alpha_{m-2} \beta_{2m-3}^i \gamma_{2m-3}$. As well for τ_{i+1} .

- (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 τ_j and ν 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 μ and ρ_i, σ_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 ρ_i, σ_i are of the following forms:

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

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 ξ_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 ξ_0 is not in μ or ρ_0 ; if so, $\xi_1 = \xi_0$ would hold because both ξ_1 and ξ_0 would have a common prefix up to the second c , but ξ_1 must be strictly longer than ξ_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 ξ_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 σ_0 . The second key point is that the informative substring $(w\#)^n$ will be included in σ_0 . From this, we can show that σ_1 must be longer than σ_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 $t = u_0 s_1 u_1 \cdots s_r u_r$.

Proof of Lemma 30. We only prove that L_2 is a stack but not nonerasing stack language (L_3 being non-stack is quite similar). Assume to the contrary that L_2 is a nonerasing stack language and fix k given by Theorem 35. Let $\xi_0 = wc(w\#)^k c(w\#)^k \in L_2$ where $w = a^k$. As $|\xi_0| \geq k$, there exist μ, ρ_i and σ_i ($i \geq 0$) that satisfy the conditions of Theorem 35. By conditions (iii), (iv) and (v), the inequality (*) $|\xi_0| < |\xi_1|$ easily follows.

Henceforth, we write the first and second c of $\xi_0 = \mu\rho_0\sigma_0 = wc(w\#)^k c(w\#)^k$ as c_1, c_2 respectively. Let us prove $c_2 \in \sigma_0$. Assume $c_2 \in \mu$. Now, μ is of the form $wc(w\#)^k c \dots$. However, $\xi_1 = \mu\rho_0\rho_1\sigma_1 = wc(w\#)^k c \dots \in L_2$ and so $\xi_1 = \xi_0$, contrary to (*). Similarly, $c_2 \in \rho_0$ also does not hold. Thus, σ_0 is of the form $\eta_0 c(w\#)^k$. Because $\sigma_0 \leq \sigma_1$, $c \in \sigma_1$. Now, we can derive $|\sigma_1| - |\sigma_0| \geq k$ as follows.

If $c_1 \in \mu\rho_0$, both ξ_0 and ξ_1 have the prefix wc in common. Because $\xi_1 \in L_2$, σ_1 can be written as $\sigma_1 = \eta_1 c(w\#)^m$ where $\eta_0 \leq \eta_1$ and $c \notin \eta_0, \eta_1$. Inequality (*) requires $m > k$. Then, $|\sigma_1| - |\sigma_0| = |\eta_1| - |\eta_0| + |(w\#)^m| - |(w\#)^k| \geq k$. Otherwise, $c_1 \in \sigma_0$ and we can write $\mu\rho_0 = u_0$, $\sigma_0 = v_0 c(w\#)^k c(w\#)^k$ and $w = u_0 v_0$. Because $\sigma_0 \leq \sigma_1$ and $\xi_1 \in L_2$, σ_1 is of the form $v_1 c(x\#)^m c(x\#)^m$, $v_0 \leq v_1$. Because $\xi_1 = \mu\rho_0\rho_1\sigma_1 = u_0 \dots v_1 c(x\#)^m c(x\#)^m = xc(x\#)^m c(x\#)^m$, $|x| \geq |w|$ follows. On the other hand, since $\sigma_0 \leq \sigma_1$, $(w\#)^k \leq (x\#)^m$ holds and $m \geq k$ by counting #. Inequality (*) ensures that either $|x| > |w|$ or $m > k$ holds. In either case, $|\sigma_1| - |\sigma_0| \geq k$.¹⁰

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

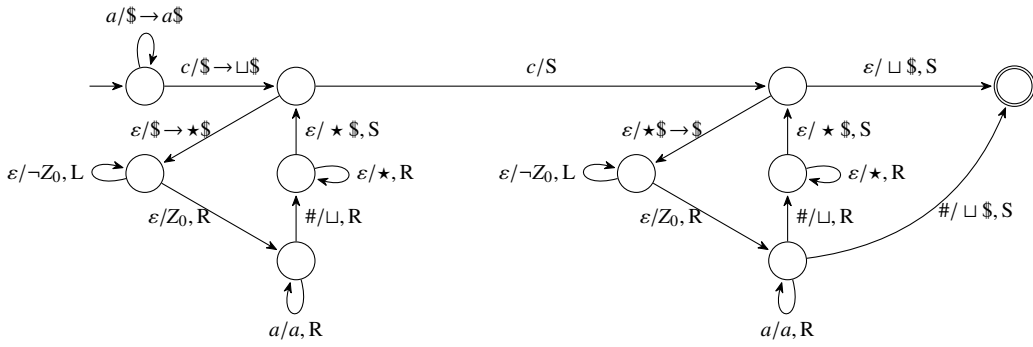


Figure 2: A stack automaton that recognizes L_2 (see Notation 15 for the notations)

6. Backreferences and lookaheads go beyond indexed languages

Chida and Terauchi formalized *rewbl*, a regular expression extended with backreferences and lookaheads, and showed that its expressive power strictly exceeds that of *rewbs* [8]. There are two kinds of lookaheads: positive lookaheads $?=$ and negative lookaheads $?!$. In particular, a *rewbl_p* (resp. *rewbl_n*) is a *rewb* with positive (resp. negative) lookaheads. We write *rewbl* for a *rewb* with both positive and negative lookaheads. In what follows, for brevity, we use the *character set* notation, which regards a nonempty finite set of characters $C = \{a_1, \dots, a_n\} \subseteq \Sigma$ as a shorthand for the expression $a_1 + \dots + a_n$.

Example 36. Let Σ denote the alphabet $\{a, b\}$. The *rewbl_p* $\alpha = (?=ab)\Sigma\Sigma\Sigma$ describes the language $L(\alpha) = \{aba, abb\}$ (i.e., the set of all three length words starting with ab). This is because the subexpression $(?=ab)$ first “filters” a *prefix* of the input string by matching with the regular expression $ab\Sigma^*$ without consuming any input. Subsequently, the expression $\Sigma\Sigma\Sigma$ will match any unfiltered three length strings. In short, $(?= \alpha)$ tests for free if the matching between $\alpha\Sigma^*$ and the rest of the input string succeeds or not. Similarly, $(?! \alpha)$ tests if the matching fails or not.

In this section, in contrast with the fact that every *rewb* language is indexed (Section 4, Corollary 24), we show that it is not the case for *rewbs* extended with lookaheads. We do not give a formal semantics of *rewbs* with lookaheads (see [8] for such a semantics) and only assume the following claim which is sufficient for our purpose:

Claim 37. Let α_1 and α_2 denote *rewbs*. Without loss of generality, we assume α_1 and α_2 have no common capturing group indexes. Fix a fresh symbol $\$$ not appearing in α_1 and α_2 . The *rewbl_p* $(?= \alpha_1 \$)\alpha_2 \$$ (resp. *rewbl_n* $(?! \alpha_1 \$)\alpha_2 \$$) describes the language $L(\alpha_1 \$) \cap L(\alpha_2 \$)$.

¹⁰One may think that the proof can be simplified by only comparing the suffix $(w\#)^k \leq (x\#)^m$ to derive both $|x| \geq |w|$ and $m \geq k$ at once. However, that does not derive the former inequality $|x| \geq |w|$ because, for example, $aaa\# \leq aa(\#a)a\# = (aa\#)^2$ but not $|aaa| \leq |aa|$.

Note that whether the language class of rewbl_p s is closed under intersections remains open but those of rewbl_n s and rewbl s are known to be closed under intersections [8].

Lemma 38. *For an arbitrary recursively enumerable language E , there exist a rewbl_p α and a GSM mapping¹¹ 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

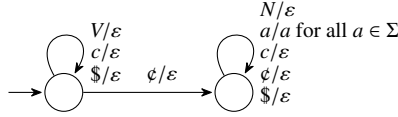
$$\alpha_1 \triangleq \left(\bigvee_{\substack{q \in Q, a \in \Gamma \\ \delta(q,a) = (p,b,L)}} (1\Gamma^*)_1 (2\Gamma)_2 q a (3\Gamma^*)_3 c \setminus 1p \setminus 2b \setminus 3c + \bigvee_{\substack{q \in Q, a \in \Gamma \\ \delta(q,a) = (p,b,R)}} (1\Gamma^*)_1 q a (2\Gamma^*)_2 c \setminus 1bp \setminus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

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

Notice that a string that α_1 matches stands for a valid computation of M that is separated by the delimiter c and ends with $\phi \dots \$$, and α_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 ϕ and $\$$, which lies in the last component. It is used to indicate by ϕ 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 rewbl_p , rewbl_n , and rewbl is a subclass of indexed languages.*

Proof. Since the class of indexed languages is closed under GSM mappings [1], the class of rewbl_p cannot be a subclass of indexed languages by Lemma 38. As for rewbl_n , replace (α_1) with $(?! (\alpha_1))$. The last is immediate. □

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

In this section, we construct an NESAs A_i that describes $L(x_i)$, where the rew x_i is given by Larsen [17] and defined over the alphabet $\Sigma = \{a_0^l, a_0^m, a_0^r, a_1^l, a_1^m, a_1^r, \dots\}$ as follows: $x_0 \triangleq (a_0^l a_0^m a_0^r)^*$, $x_{i+1} \triangleq (a_{i+1}^l (x_i)_i a_{i+1}^m \setminus i a_{i+1}^r)^*$ ($i \geq 0$). Our result implies that Larsen's hierarchy is within the class of nonerasing stack languages. Since Larsen showed that no rew 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 rew of nested level at least i .¹²

¹¹A *generalized sequential machine* (GSM) [11] is a 6-tuple $(Q, \Sigma, \Delta, \delta, \lambda, q_0)$, where Q is a finite set of states, Σ a finite set of input symbols, Δ 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 δ and λ 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 $\hat{\delta}(q, a) = q'$ and $\hat{\lambda}(q, a) = x$. A *GSM mapping* is a mapping $\Sigma^* \rightarrow \Delta^*$, $w \mapsto \hat{\lambda}(q_0, w)$ for some GSM $(Q, \Sigma, \Delta, \delta, \lambda, q_0)$.

¹²Technically, Larsen [17] adopts a syntax that excludes unbound references, and so this implied result applies only to rews with no unbound references.

- [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.