

Complexity Results and the Growths of Hairpin Completions of Regular Languages

Volker Diekert and Steffen Kopecki

Universität Stuttgart,
Institut für Formale Methoden der Informatik
Universitätsstr. 38, D-70569 Stuttgart, Germany
{diekert,kopecki}@fmi.uni-stuttgart.de

June 28, 2010

Abstract

The hairpin completion is a natural operation on formal languages which has been inspired by molecular phenomena in biology and by DNA-computing. In 2009 we presented in [6] a (polynomial time) decision algorithm to decide regularity of the hairpin completion. In this paper we provide four new results: 1.) We show that the decision problem is NL-complete. 2.) There is a polynomial time decision algorithm which runs in time $\mathcal{O}(n^8)$, this improves [6], which provided $\mathcal{O}(n^{20})$. 3.) For the one-sided case (which is closer to DNA computing) the time is $\mathcal{O}(n^2)$, only. 4.) The hairpin completion is unambiguous linear context-free. This result allows to compute the growth (generating function) of the hairpin completion and to compare it with the growth of the underlying regular language.

1 Introduction

The hairpin completion is a natural operation of formal languages which has been inspired by molecular phenomena in biology and by DNA-computing. An intramolecular base pairing, known as a *hairpin*, is a pattern that can occur in single-stranded DNA and, more commonly, in RNA. Hairpin or hairpin-free structures have numerous applications to DNA computing and molecular genetics, see [5, 7, 8, 12, 13] and the references within. For example, an instance of 3-SAT has been solved with a DNA-algorithm and one of the main concepts was to eliminate all molecules with a hairpin structure, see [17].

In this paper we study the hairpin completion from a purely formal language viewpoint. The hairpin completion of a formal language was first defined in [4]; here we use a slightly more general definition which was introduced in [6]. The formal operation of the hairpin completion on words is best explained in Figure 1. In that picture as in the rest of the paper we mean by putting a *bar* on a word (like \bar{a}) to read it from right-to-left in addition to replacing a with the Watson-Crick complement \bar{a} for letters. The hairpin completion of a regular language is linear context-free [4]. For some time it was not known

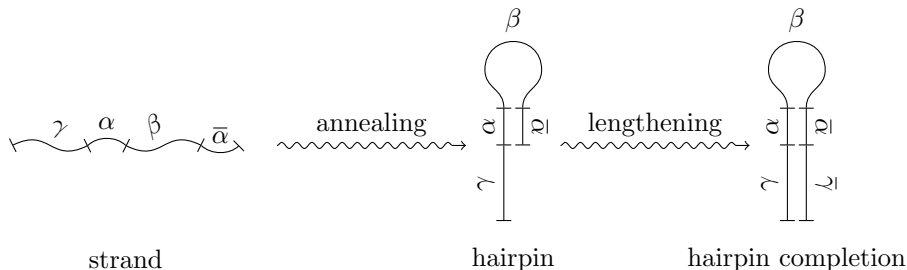


Figure 1: Hairpin completion of a DNA-strand (or a word).

whether regularity of the hairpin completion is decidable. It was only in 2009 when we presented in [6] a decision algorithm. The runtime of that algorithm is in $\mathcal{O}(n^{20})$, hence polynomial.

Here we present a modified approach to solve the decision problem. The new approach leads to improved complexity results and a new structure theorem. We show that the decision problem is **NL**-complete (Theorem 3.1). We show that there is a polynomial time decision algorithm which runs in time $\mathcal{O}(n^8)$ (Theorem 3.2, i.). So, the improvement is from $\mathcal{O}(n^{20})$ down to $\mathcal{O}(n^8)$. Moreover, in the biological model the one-sided hairpin completion is of particular interest, and in that special case we need quadratic time, only (Theorem 3.2, iii.). We also argue why the time bounds might be optimal in the worst case.

A byproduct of the method yields that the hairpin completion of a regular language is unambiguous linear context-free (Theorem 3.7). The result about unambiguity allows to compute the growth (generating function) of the hairpin completion and to compare it with the growth of the underlying regular language (Theorem 3.3 and Corrolary 3.8).

This takes us back to a challenging open problem in formal languages. Regularity of linear context-free languages is undecidable in general [1, 10]. But the situation for unambiguous context-free languages is open for more than 40 years. Hence, we have new a positive result within the classical context of deciding regularity within a class of unambiguous (linear) context-free languages.

2 Preliminaries and Notation

We assume the reader to be familiar with the fundamental concepts of formal language theory, automata theory, and complexity theory, see [11, 16]. By **NL** we mean the complexity class **NLOGSPACE**, which contains the problems which can be decided with a non-deterministic algorithm using $\mathcal{O}(\log n)$ space. We heavily rely on the well-known result that **NL** is closed under complementation. We also use the fact that if L can be reduced to L' via some single-valued non-deterministic transduction in $\mathcal{O}(\log n)$ space and $L' \in \mathbf{NL}$, then we have $L \in \mathbf{NL}$, too. This reduction is performed by a non-deterministic log-space Turing machine. In case the machine stops on input w , the output is always the same, independently of non-deterministic moves during the computation. So, we can call the output $r(w)$. The reduction property tells us $w \in L$ if and only if both, the machine sometimes stops on input w and $r(w) \in L'$.

By Σ we denote a finite alphabet with at least two letters which is equipped with an *involution* $\bar{\cdot} : \Sigma \rightarrow \Sigma$. An involution for a set is a bijection such that $\bar{\bar{a}} = a$. We extend the involution to words $a_1 \cdots a_n$ by $\overline{a_1 \cdots a_n} = \bar{a}_n \cdots \bar{a}_1$. (Just like taking inverses in groups.) For languages \bar{L} denotes the set $\{\bar{w} \mid w \in L\}$. The set of words over Σ is denoted Σ^* ; and the *empty word* is denoted by 1. Given a word w , we denote by $|w|$ its length and $w(m) \in \Sigma$ its m -th letter. By $\Sigma^{\leq m}$ we mean the set of all words with length at most m . If $w = xyz$ for some $x, y, z \in \Sigma^*$, then x and z are called *prefix* and *suffix*, respectively. The prefix relation between words x and y is denoted by $x \leq y$.

Throughout the paper L_1, L_2 mean two regular languages in Σ^* and by k we mean a (small) constant, say $k = 10$. We define the *hairpin completion* $\mathcal{H}_k(L_1, L_2)$ by

$$\mathcal{H}_k(L_1, L_2) = \{\gamma\alpha\beta\bar{\alpha}\bar{\gamma} \mid (\gamma\alpha\beta\bar{\alpha} \in L_1 \vee \alpha\beta\bar{\alpha}\bar{\gamma} \in L_2) \wedge |\alpha| = k\}.$$

Three cases are of main interest: 1.) $L_1 = L_2$, 2.) $L_1 = \bar{L}_2$, and 3.) $L_1 = \emptyset$ or $L_2 = \emptyset$. Compared to the definition of the hairpin completion in [4, 15], case 1.) corresponds to the two-sided hairpin completion and case 3.) to the one-sided hairpin completion. Since we have better time complexities for 2.) and 3.) than for 1.) or in the general case we make the time bounds rather precise.

A regular languages can be specified by a non-deterministic finite automaton (NFA) $\mathcal{A} = (\mathcal{Q}, \Sigma, E, \mathcal{I}, \mathcal{F})$, where \mathcal{Q} is the finite set of *states*, $\mathcal{I} \subseteq \mathcal{Q}$ is the set of *initial states*, and $\mathcal{F} \subseteq \mathcal{Q}$ is the set of *final states*. The set E contains labeled *edges* (or *arcs*), it is a subset of $\mathcal{Q} \times \Sigma \times \mathcal{Q}$. For a word $u \in \Sigma^*$ we write $p \xrightarrow{u} q$, if there is a path from state p to q which is labeled by the word u . Thus, the accepted language becomes

$$L(\mathcal{A}) = \left\{ u \in \Sigma^* \mid \exists p \in \mathcal{I}, \exists q \in \mathcal{F} : p \xrightarrow{u} q \right\}.$$

Later it will be crucial to use also paths which avoid final states. For this we introduce a special notation. First remove all arcs (p, a, q) where $q \in \mathcal{F}$ is a final state. Thus, final states do not have incoming arcs anymore in this reduced automaton. Let us write $p \xRightarrow{u} q$, if there is a path in this reduced automaton from state p to q which is labeled by the word u . Note that for such a path $p \xRightarrow{u} q$ we allow $p \in \mathcal{F}$, but on the path we never meet any final state again.

An NFA is called a deterministic finite automaton (DFA), if it has one initial state and for every state $p \in \mathcal{Q}$ and every letter $a \in \Sigma$ there is exactly one arc $(p, a, q) \in E$. In particular, a DFA in this paper is always complete, thus we can read every word to its end. We also write $p \cdot u = q$, if $p \xrightarrow{u} q$. This yields a (totally defined) function $\mathcal{Q} \times \Sigma^* \rightarrow \mathcal{Q}$, which defines an action of Σ^* on \mathcal{Q} on the right.

In the following we need a DFA accepting L_1 as well as a DFA accepting L_2 , but the DFA for L_2 has to work from right-to-left. Instead of introducing this concept we use a DFA (working as usual from left-to-right), which accepts \bar{L}_2 . This automaton has the same number of states (and is structurally isomorphic to) as a DFA accepting the *reversal language* of L_2 .

As input we assume that the regular languages L_1 and \bar{L}_2 are specified by DFAs with state set \mathcal{Q}_i , state $q_{0i} \in \mathcal{Q}_i$ as initial state, and $\mathcal{F}_i \subseteq \mathcal{Q}_i$ as final states. We fix $n_i = |\mathcal{Q}_i|$ to be the number of states, $i = 1, 2$. By n we mean $\max\{n_1, n_2\}$. The input size is therefore the number n .

We also need the usual product DFA with state space

$$\mathcal{Q} = \{(p_1, p_2) \in \mathcal{Q}_1 \times \mathcal{Q}_2 \mid \exists w \in \Sigma^* : (p_1, p_2) = (q_{01} \cdot w, q_{02} \cdot w)\}.$$

The action is given by $(p_1, p_2) \cdot a = (p_1 \cdot a, p_2 \cdot a)$. We let $n_{12} = |\mathcal{Q}|$. Hence, $n \leq n_{12} \leq n_1 \cdot n_2 \leq n^2$, and $n = n_1 = n_{12}$ if $L_2 = \emptyset$ or $L_1 = \overline{L_2}$. In the following we work simultaneously in all three automata defined so far. Moreover, in \mathcal{Q}_1 and \mathcal{Q}_2 we have to work backwards. This leads to nondeterminism. Our first new construction concerns a special NFA in Section 3.1.

3 Main results

The complexity results of this paper are the following:

Theorem 3.1. *The problem whether the hairpin completion $\mathcal{H}_k(L_1, L_2)$ is regular is NL-complete.*

Theorem 3.2. *i.) The problem whether the hairpin completion $\mathcal{H}_k(L_1, L_2)$ is regular can be decided in time $\mathcal{O}(n^8)$.*

ii.) For $L_1 = \overline{L_2}$ it can be decided in time $\mathcal{O}(n^6)$.

iii.) For $L_2 = \emptyset$ it can be decided in time $\mathcal{O}(n^2)$.

An algorithm solving this problem is sketched in Section 3.3 and the time complexity is proved in Section 3.4.

The *growth* or *generating function* g_L of a formal language L is defined as:

$$g_L(z) = \sum_{m \geq 0} |L \cap \Sigma^{\leq m}| z^m.$$

We can view g_L as a formal power series or as an analytic function in one complex variable where the radius of convergence is strictly positive. The radius of convergence is at least $1/|\Sigma|$.

It is well-known that the growth of a regular language L is effectively rational, i.e., a quotient of two polynomials. The same is true for unambiguous linear context-free languages. In particular, the growth is either polynomial or exponential. If the growth is exponential, then we find an algebraic number $\rho \in \mathbb{R}$ such that $|L \cap \Sigma^{\leq m}|$ behaves essentially as ρ^m , see [2, 3, 9].

It was shown in [4] that $\mathcal{H}_k(L_1, L_2)$ is an linear context-free language. As a byproduct to our techniques to prove the complexity results above we find that $\mathcal{H}_k(L_1, L_2)$ is unambiguous, and hence its growth (i.e., generating function) is a rational function, see e.g. [14] for this well-known fact. We obtain:

Theorem 3.3. *The hairpin completion $\mathcal{H}_k(L_1, L_2)$ is an unambiguous linear context-free language with an effectively computable rational growth function.*

This result is proved of Section 3.2.

3.1 The NFA \mathcal{A}

In this section we define a certain NFA which is called simply \mathcal{A} . Almost all further results are done by exploring properties of this NFA. The NFA is a sort of product automaton over $\mathcal{Q} \times \mathcal{Q}_1 \times \mathcal{Q}_2 \subseteq \mathcal{Q}_1 \times \mathcal{Q}_2 \times \mathcal{Q}_1 \times \mathcal{Q}_2$ where $\mathcal{Q}_1, \mathcal{Q}_2$ and \mathcal{Q} are defined as in Section 2. The size of this automaton is $\mathcal{O}(n^4)$ in the worst case, and our decision algorithm will take into account all pairs of states in this NFA. Hence, $\mathcal{O}(n^8)$ might be an optimal time bound and the decision algorithm is not worse than quadratic in the size of the NFA \mathcal{A} .

For every quadruple $(p_1, p_2, q_1, q_2) \in \mathcal{Q}_1 \times \mathcal{Q}_2 \times \mathcal{Q}_1 \times \mathcal{Q}_2$ we define a regular language $B(p_1, p_2, q_1, q_2)$ as follows:

$$B(p_1, p_2, q_1, q_2) = \{w \in \Sigma^* \mid p_1 \cdot w = q_1 \wedge p_2 \cdot \bar{w} = q_2\}.$$

We say that (p_1, p_2, q_1, q_2) is a *basic bridge* if $B(p_1, p_2, q_1, q_2) \neq \emptyset$. The idea behind of this notation is that $B(p_1, p_2, q_1, q_2)$ closes a gap between pairs (p_1, p_2) and (q_1, q_2) (which are on different sides). For a letter $a \in \Sigma$ we call (p_1, p_2, q_1, q_2) an *a-bridge* if $B(p_1, p_2, q_1, q_2) \cap a\Sigma^* \neq \emptyset$.

Lemma 3.4. *The number of basic bridges and a-bridges is bounded by $\mathcal{O}(n_1^2 n_2^2)$. A table containing all these bridges can be computed in time $\mathcal{O}(n_1^2 n_2^2) \subseteq \mathcal{O}(n^4)$, and there is a single-valued non-deterministic transduction working in $\mathcal{O}(\log n)$ space which outputs this table.*

Proof. To compute the basic bridges amounts to compute the transitive closure in some graph where the number of nodes and edges is in $\mathcal{O}(n_1 n_2)$. This gives the time bound. Once we have the bridges we can compute the *a*-bridges in time $\mathcal{O}(n_1^2 n_2^2)$.

If (p_1, p_2, q_1, q_2) is a basic bridge, we can verify this property in **NL**. Since **NL** is closed under complementation, we can output the whole table by a single-valued non-deterministic transduction in $\mathcal{O}(\log n)$ space. \square

We also need *levels* for $0 \leq \ell \leq k$, hence there are $k + 1$ levels. By $[k]$ we denote in this paper the set $\{0, \dots, k\}$. Define

$$\{((p_1, p_2), q_1, q_2, \ell) \in \mathcal{Q} \times \mathcal{Q}_1 \times \mathcal{Q}_2 \times [k] \mid (p_1, p_2, q_1, q_2) \text{ is a basic bridge}\}$$

as the state space of an NFA called \mathcal{A} . Its size is bounded by $N \cdot (k + 1) \in \mathcal{O}(N) \subseteq \mathcal{O}(n^4)$, where $N = n_{12} n_1 n_2$. We have $N = n^2$ for $L_2 = \emptyset$, and $N = n^3$ for $L_2 = \bar{L}_1$.

By a (slight) abuse of languages we call a state $((p_1, p_2), q_1, q_2, \ell)$ a *bridge*, and we keep in mind that there exists a word w such that $p_1 \cdot w = q_1$ and $p_2 \cdot \bar{w} = q_2$. Bridges are frequently denoted by (P, q_1, q_2, ℓ) with $P = (p_1, p_2) \in \mathcal{Q}$, $q_i \in \mathcal{Q}_i$, $i = 1, 2$, and $\ell \in [k]$. Bridges are a central concept in the following.

The *a*-transitions in the NFA for $a \in \Sigma$ are given by the following arcs:

$$\begin{aligned} (P, q_1 \cdot \bar{a}, q_2 \cdot \bar{a}, 0) &\xrightarrow{a} (P \cdot a, q_1, q_2, 0) && \text{for } q_i \cdot \bar{a} \notin \mathcal{F}_i, i = 1, 2, \\ (P, q_1 \cdot \bar{a}, q_2 \cdot \bar{a}, 0) &\xrightarrow{a} (P \cdot a, q_1, q_2, 1) && \text{for } q_1 \cdot \bar{a} \in \mathcal{F}_1 \text{ or } q_2 \cdot \bar{a} \in \mathcal{F}_2, \\ (P, q_1 \cdot \bar{a}, q_2 \cdot \bar{a}, \ell) &\xrightarrow{a} (P \cdot a, q_1, q_2, \ell + 1) && \text{for } 1 \leq \ell < k. \end{aligned}$$

Observe that no state of the form $(P, q_1, q_2, 0)$ with $q_1 \in \mathcal{F}_1$ or $q_2 \in \mathcal{F}_2$ has an outgoing arc to level zero; we must switch to level one. There are no

outgoing arcs on level k , and for each $(a, P, q_1, q_2, \ell) \in \Sigma \times \mathcal{Q} \times \mathcal{Q}_1 \times \mathcal{Q}_2 \times [k-1]$ there exists at most one arc $(P, q'_1, q'_2, \ell) \xrightarrow{a} (P \cdot a, q_1, q_2, \ell')$. Indeed, the triple (q'_1, q'_2, ℓ') is determined by (q_1, q_2, ℓ) and the letter a . Not all arcs exist because (P, q'_1, q'_2, ℓ) can be a bridge whereas $(P \cdot a, q_1, q_2, \ell')$ is not. Thus, there are at most $|\Sigma| \cdot N \cdot k \in \mathcal{O}(N)$ arcs in the NFA.

The set of initial states \mathcal{I} contains all bridges of the form $(Q_0, q'_1, q'_2, 0)$ with $Q_0 = (q_{01}, q_{02})$. The set of final states \mathcal{F} is given by all bridges (P, q_1, q_2, k) on level k .

For an example and a graphical presentation of the NFA, see Figure 2.

Remark 3.5. *The NFA \mathcal{A} can be computed by Lemma 3.4 in time $\mathcal{O}(n_1^2 n_2^2)$ and by a single-valued non-deterministic transduction in $\mathcal{O}(\log n)$ space. Thus for both the polynomial time and the NL algorithm we can have direct access to \mathcal{A} and we can assume that \mathcal{A} is written on the input tape.*

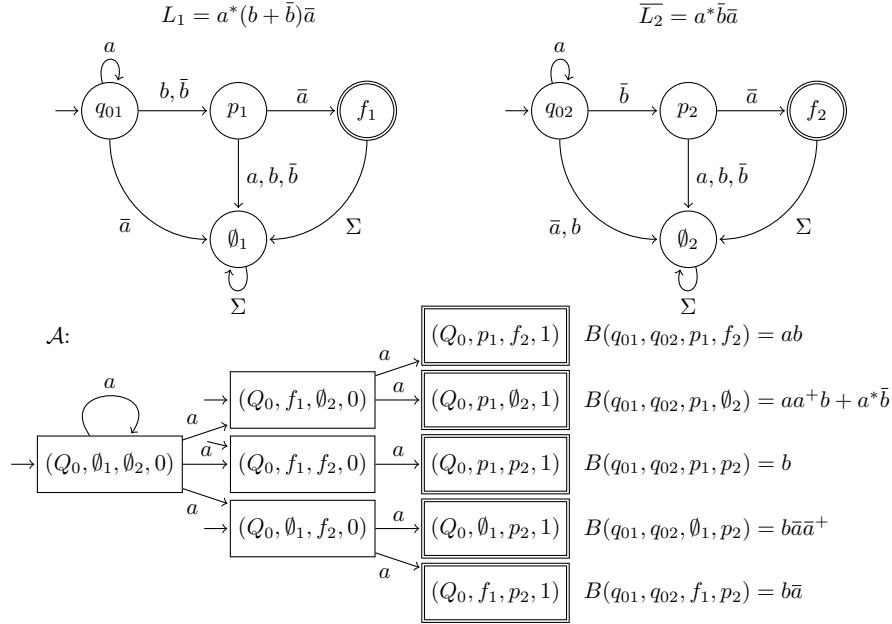


Figure 2: DFAs for L_1 and \bar{L}_2 and the resulting NFA \mathcal{A} with 4 initial states and 5 final states associated to the (linear context-free) hairpin completion $\mathcal{H}_k(L_1, L_2) = a^+b\bar{a}^+ \cup \{a^s\bar{b}\bar{a}^t \mid s \geq t \geq 1\}$ with $k = 1$.

The next result shows the unambiguity of paths in the automaton \mathcal{A} .

Lemma 3.6. *Let $w \in \Sigma^*$ be the label of a path in \mathcal{A} from a bridge $A = (P, p_1, p_2, \ell)$ to $A' = (P', p'_1, p'_2, \ell')$, then the path is unique. This means that $B = B'$ whenever $w = uv$ and*

$$A \xrightarrow{u} B \xrightarrow{v} A', \quad A \xrightarrow{u} B' \xrightarrow{v} A'.$$

Proof. It is enough to consider $u = a \in \Sigma$. Let $B = (Q, q_1, q_2, m)$. Then we have $Q = P \cdot a$ and $q_i = p'_i \cdot \bar{v}$. If $\ell = 0$ and $p_i \notin \mathcal{F}_i$ for $i = 1, 2$, then $m = 0$, too; otherwise $m = \ell + 1$. Thus, B is defined by A, A' , and u, v . We conclude $B = B'$. \square

3.2 Structure theorem and rational growth

For languages U and V we define the language V^U as follows:

$$V^U = \{uv\bar{u} \mid u \in U, v \in V\}.$$

Clearly, if U and V are regular, then V^U is linear context-free. We are interested in a disjoint union of languages V^U where for $w \in V^U$ the factorization $w = uv\bar{u}$ with $u \in U$ and $v \in V$ is unambiguous.

Theorem 3.7. *Let $T = \mathcal{I} \times \mathcal{F}$. For each pair $\tau = (I, F) \in T$ with $F = ((d_1, d_2), e_1, e_2, k)$ let R_τ be the (regular) set of words which label a path from the initial bridge I to the final bridge F and let $B_\tau = B(d_1, d_2, e_1, e_2)$. The hairpin completion is a disjoint union*

$$\mathcal{H}_k(L_1, L_2) = \bigcup_{\tau \in T} B_\tau^{R_\tau}.$$

Moreover, for each word in some $w \in B_\tau^{R_\tau}$ there is a unique factorization $w = \rho\beta\bar{\rho}$ with $\rho \in R_\tau$ and $\beta \in B_\tau$.

Proof. Let $w \in \mathcal{H}_k(L_1, L_2)$. There exists exactly one factorization $w = \gamma\alpha\beta\bar{\alpha}\bar{\gamma}$ such that $|\alpha| = k$ and there are runs in the DFAs

$$\begin{aligned} L_1 : \quad q_{01} &\xrightarrow{\gamma} c_1 \xrightarrow{\alpha} d_1 \xrightarrow{\beta} e_1 \xrightarrow{\bar{\alpha}} f_1 \xrightarrow{\bar{\gamma}} q'_1, \\ \bar{L}_2 : \quad q_{02} &\xrightarrow{\gamma} c_2 \xrightarrow{\alpha} d_2 \xrightarrow{\bar{\beta}} e_2 \xrightarrow{\bar{\alpha}} f_2 \xrightarrow{\bar{\gamma}} q'_2 \end{aligned}$$

where $f_1 \in \mathcal{F}_1$ or $f_2 \in \mathcal{F}_2$ (or both). In other words, $\gamma\alpha\beta\bar{\alpha}$ is the longest prefix of w belonging to L_1 or $\alpha\beta\bar{\alpha}\bar{\gamma}$ is the longest suffix of w belonging to L_2 .

In the NFA \mathcal{A} we find the run

$$(Q_0, q'_1, q'_2, 0) \xrightarrow{\gamma} ((c_1, c_2), f_1, f_2, 0) \xrightarrow{\alpha} ((d_1, d_2), e_1, e_2, k).$$

Now let $I = (Q_0, q'_1, q'_2, 0)$, $F = ((d_1, d_2), e_1, e_2, k)$, and $\tau = (I, F)$. Obviously, we have $\gamma\alpha \in R_\tau$ and $\beta \in B_\tau$.

Conversely, for every $\tau \in T$, $\rho = \gamma\alpha \in R_\tau$ with $|\alpha| = k$, and $\beta \in B_\tau$, we find runs in the DFAs as above, hence $\rho\beta\bar{\rho} \in \mathcal{H}_k(L_1, L_2)$. Since the states d_i , e_i , and q'_i ($i = 1, 2$) on the runs are uniquely defined by the word $\rho\beta\bar{\rho}$, we cannot have that $\tau \neq \tau' \in T$ exists and $\rho\beta\bar{\rho} \in B_{\tau'}^{R_{\tau'}}$. \square

Corollary 3.8. *The hairpin completion $\mathcal{H}_k(L_1, L_2)$ is an unambiguous linear context-free language and it has a rational growth function. The growth can be directly calculated by the growth of the regular languages R_τ and B_τ .*

Corollary 3.8 allows to compare the growth of L_1 and L_2 with the growth of their hairpin completion $\mathcal{H}_k(L_1, L_2)$. It is also a slightly more precise version of Theorem 3.3.

3.3 Complexity for testing the regularity of $\mathcal{H}_k(L_1, L_2)$

3.3.1 First Test

The automaton \mathcal{A} accepts the union of the languages R_τ as defined in Theorem 3.7. If the accepted language is finite then all R_τ are finite and hence all $B_\tau^{R_\tau}$ are regular. This leads to the following result:

Proposition 3.9. *i.) If the accepted language of the NFA \mathcal{A} is finite, then the hairpin completion $\mathcal{H}_k(L_1, L_2)$ is regular.*

ii.) If L_1 or L_2 is finite, but the accepted language of \mathcal{A} is infinite, then the hairpin completion $\mathcal{H}_k(L_1, L_2)$ is not regular.

Proof. The first statement follows directly from Theorem 3.7.

For the second statement, assume $L(\mathcal{A})$ is infinite. We find a long word $uvw\alpha$ with $|\alpha| = k$ and $v \neq 1$ such that all $uv^i w\alpha$ label some path in \mathcal{A} from an initial bridge $((q_{01}, q_{02}), q'_1, q'_2, 0)$ to a final bridge $((d_1, d_2), e_1, e_2, k)$. We obtain for $\beta \in B(d_1, d_2, e_1, e_2)$ that $uv^i w\alpha\beta\bar{\alpha}\bar{w}\bar{v}^i\bar{u} \in \mathcal{H}_k(L_1, L_2)$ for all $i \geq 0$. If $\mathcal{H}_k(L_1, L_2)$ is regular, by pumping, there exists $s \in \mathbb{N}$ and infinitely many $t \in \mathbb{N}$ such that $\pi_t = uv^s w\alpha\beta\bar{\alpha}\bar{w}\bar{v}^t\bar{u} \in \mathcal{H}_k(L_1, L_2)$. By the construction of the NFA \mathcal{A} , for π_t we find the run

$$q_{01} \xrightarrow{uv^s w\alpha\beta} d_1 \xrightarrow{\alpha} f_1 \xrightarrow{\bar{w}\bar{v}^t\bar{u}} q'_1$$

in the DFA for L_1 . Hence, the longest prefix of π_t belonging to L_1 is always a prefix of $uv^s w\alpha\beta\bar{\alpha}$. This is far too short to create the hairpin completion if t becomes huge. Thus we must use a suffix belonging to L_2 and which has at least half the length of π_t to do the job. Hence, L_2 is infinite and, by a symmetric argument, L_1 is infinite, too. \square

Test 1: Check either by some **NL**-algorithm or in time $\mathcal{O}(N) \subseteq \mathcal{O}(n^4)$ whether the accepted language of the NFA \mathcal{A} is finite. If “yes” ($=L(\mathcal{A})$ is finite), then output that $\mathcal{H}_k(L_1, L_2)$ is regular. If “no”, but L_1 or L_2 is finite, then output that $\mathcal{H}_k(L_1, L_2)$ is not regular.

Remark 3.10. *If L_1 or L_2 is empty, we have $\mathcal{O}(N) \subseteq \mathcal{O}(n^2)$. In that case it suffices to perform Test 1 in order to decide the regularity of $\mathcal{H}_k(L_1, L_2)$. This proves iii.) of Theorem 3.2.*

3.3.2 Second Test

From now on we may assume that the automaton \mathcal{A} accepts an infinite language and both L_1 and L_2 are infinite as well. We assume that all states are reachable from initial bridges and lead to some final bridges. (Recall that graph reachability can be checked in **NL**.)

Let K be the set of non-trivial strongly connected components of the automaton \mathcal{A} (read as a directed graph). For $\kappa \in K$ let $N_\kappa = |\kappa|$ the number of states in the component κ . Let us choose some $A_\kappa \in \kappa$ and some shortest non-empty word $v_\kappa \in \Sigma^+$ such that there is a path in \mathcal{A} labeled by v_κ from A_κ to A_κ .

The next lemma tells us that for a regular hairpin completion $\mathcal{H}_k(L_1, L_2)$ the word v_κ is uniquely defined by A_κ , its length is N_κ , and its conjugacy class depends only on κ .

Lemma 3.11. *Assume that the hairpin completion $\mathcal{H}_k(L_1, L_2)$ is regular.*

1.) *Let $A_\kappa \xrightarrow{v_\kappa} A_\kappa$ as above and $A_\kappa \xrightarrow{w} C$ be a path in \mathcal{A} to some final bridge. Then the word w is a prefix of some word in v_κ^+ .*

- 2.) The word v_κ and the loop $A_\kappa \xrightarrow{v_\kappa} A_\kappa$ are uniquely defined by the state A_κ and we have $|v_\kappa| = N_\kappa$.
- 3.) The loop $A_\kappa \xrightarrow{v_\kappa} A_\kappa$ visits every other state $B \in \kappa$ exactly once. Thus, the loop defines an Hamiltonian cycle of κ .

Proof. Let $A = A_\kappa$ and $v = v_\kappa$. As $A \xrightarrow{v} A$ is a non-trivial loop, we see that A is on level zero. Consider a path labeled by w from A to a final bridge $((p_1, p_2), q_1, q_2, k)$. By assumption, all states in \mathcal{A} are reachable from some initial state. Thus, we find a word u such that the automaton \mathcal{A} accepts $uv^i w$ for all $i \geq 0$. We see next that $uv^i w \beta \bar{w} \bar{v}^i \bar{u} \in \mathcal{H}_k(L_1, L_2)$ for all $i \geq 0$ and all $\beta \in B(p_1, p_2, q_1, q_2)$. As $\mathcal{H}_k(L_1, L_2)$ is regular there are s, t with $uv^s w \beta \bar{w} \bar{v}^{s+t} \bar{u} \in \mathcal{H}_k(L_1, L_2)$ and $t > |w\beta|$. This means that the hairpin completion is forced to use a suffix in L_2 , and due to the definition of \mathcal{A} we conclude that $uv^s w$ must be a prefix of $uv^{s+t} w$. This implies that w is a prefix of v^t , and it proves the first assertion.

For the second one, observe first that $|v| \leq N_\kappa$ is trivial. Now, let $A \neq B \in \kappa$ and $A \xrightarrow{v'} B \xrightarrow{v''} A$. For some $i, j > 0$ we have $|v^i| = |(v'v'')^j|$. Thus, $v^i = (v'v'')^j$ by the first property. By the unique-path-property stated in Lemma 3.6 we obtain that the loop $A \xrightarrow{(v'v'')^j} A$ just uses the shortest loop $A \xrightarrow{v} A$ several times. If $|v| = |v'v''|$, then we conclude $v = v'v''$. Moreover, for each $A \neq B \in \kappa$ we find $A \xrightarrow{v'} B \xrightarrow{v''} A$. In particular, B is on the shortest loop around A . This yields $|v| \geq N_\kappa$ and hence the second and third assertion. \square

Example 3.12. In the example given in Figure 2 the state $(Q_0, \emptyset_1, \emptyset_2, 0)$ forms the only strongly connected component and the corresponding path is labeled with a . As one can easily observe the automaton \mathcal{A} satisfies condition 1.) to 3.) of Lemma 3.11 even though the hairpin completion is not regular.

Remark 3.13. We decompose the automaton \mathcal{A} in its strongly connected components by the algorithm of Tarjan in time $\mathcal{O}(N)$. (Note that we have $K \neq \emptyset$ since $|L(\mathcal{A})|$ is infinite.) This is also possible by some single-valued non-deterministic transduction. Putting some linear order on the set of bridges, we can assign to each $\kappa \in K$ the least $A_\kappa \in \kappa$. If $\mathcal{H}_k(L_1, L_2)$ is regular, then (by Lemma 3.11) we can output the uniquely defined words v_κ for all $\kappa \in K$. We observe that

$$\sum_{\kappa \in K} |v_\kappa| = \sum_{\kappa \in K} N_\kappa \leq N.$$

So, the list of all v_κ is computable in time $\mathcal{O}(N)$ and also by some single-valued non-deterministic transduction, in case $\mathcal{H}_k(L_1, L_2)$ is regular.

Test 2: It has two parts. Part I: For each strongly connected component $\kappa \in K$ compute a shortest word v with $0 < |v| \leq N_\kappa$ such that $A_\kappa \xrightarrow{v} A_\kappa$ is a loop in the automaton \mathcal{A} . If $|v| \neq N_\kappa$, then **stop** and output that $\mathcal{H}_k(L_1, L_2)$ is not regular. Part II: If $|v| = N_\kappa$ for all κ , then let L_κ be the accepted language of \mathcal{A} when the bridge A_κ is used as initial state. Let $\text{Pref}(v^+)$ be the language of prefixes of words in v^+ . (Note that a DFA for the complement of $\text{Pref}(v^+)$ has $N_\kappa + 1$ states.) If we do not find $L_\kappa \subseteq \text{Pref}(v^+)$, then **stop** and output that $\mathcal{H}_k(L_1, L_2)$ is not regular.

Part I can be done in time $\mathcal{O}(\sum_{\kappa \in K} N_\kappa) \subseteq \mathcal{O}(N)$. Part II can be done in time $\mathcal{O}(\sum_{\kappa \in K} N_\kappa \cdot N) \subseteq \mathcal{O}(N^2) \subseteq \mathcal{O}(n^8)$. The **NL**-algorithm is based on the fact that we can guess a position m where the m -th letter of $w \in L_\kappa$ differs from the $(m \bmod N_\kappa)$ -th letter of a word v which labels a path $A_\kappa \xrightarrow{v} A_\kappa$.

Henceforth we may assume that Test 2 was successful and following Remark 3.13 we assume that the list of all words v_κ is available. Thus, we can think that the list $(v_\kappa; \kappa \in K)$ is written on the input tape. For the **NL**-algorithm we perform another single-valued non-deterministic transduction to achieve this.

3.3.3 Third and Fourth Test

We fix a strongly connected component $\kappa \in K$ of \mathcal{A} . We let $A = A_\kappa = ((p_1, p_2), q_1, q_2, 0)$ and $v = v_\kappa$ as above. By u we denote some word leading from an initial bridge $((q_{01}, q_{02}), q'_1, q'_2, 0)$ to A . (The following tests do not rely on the choice of u .) The main idea is to investigate runs through the DFAs for L_1 and $\overline{L_2}$ where $s, t \geq n$.

$$\begin{aligned} L_1 : & \quad q_{01} \xrightarrow{u} p_1 \xrightarrow{v^s} p_1 \xrightarrow{xy} c_1 \xrightarrow{z} d_1 \xrightarrow{\bar{x}} e_1 \xrightarrow{\bar{v}^{n_1}} q_1 \xrightarrow{\bar{v}^*} q_1 \xrightarrow{\bar{u}} q'_1 \\ \overline{L_2} : & \quad q_{02} \xrightarrow{u} p_2 \xrightarrow{v^t} p_2 \xrightarrow{x} c_2 \xrightarrow{\bar{z}} d_2 \xrightarrow{\bar{y}\bar{x}} e_2 \xrightarrow{\bar{v}^{n_2}} q_2 \xrightarrow{\bar{v}^*} q_2 \xrightarrow{\bar{u}} q'_2 \end{aligned}$$

We investigate the case where $uv^sxyz\bar{x}\bar{v}^t\bar{u} \in \mathcal{H}_k(L_1, L_2)$ for all $s \geq t$ and where (by symmetry) this property is due to the longest prefix belonging to L_1 .

The following lemma is the most technical one in our paper.

Lemma 3.14. *Let $x, y, z \in \Sigma^*$ be words and $(d_1, d_2) \in \mathcal{Q}_1 \times \mathcal{Q}_2$ with the following properties:*

- 1.) $k \leq |x| < |v| + k$ and x is a prefix of some word in v^+ .
- 2.) $0 \leq |y| < |v|$ and xy is the longest common prefix of xyz and some word in v^+ .
- 3.) $z \in B(c_1, c_2, d_1, d_2)$, where $c_1 = p_1 \cdot xy$ and $c_2 = p_2 \cdot x$.
- 4.) $q_1 = d_1 \cdot \bar{x}\bar{v}^{n_1}$ and during the computation of $d_1 \cdot \bar{x}\bar{v}^{n_1}$ we see after exactly k steps a final state in \mathcal{F}_1 and then never again.
- 5.) $q_2 = d_2 \cdot \bar{y}\bar{x}\bar{v}^{n_2}$ and, let $e_2 = d_2 \cdot \bar{y}\bar{x}$, during the computation of $e_2 \cdot \bar{v}^{n_2}$ we do not see a final state in \mathcal{F}_2 .

If $\mathcal{H}_k(L_1, L_2)$ is regular, then $xyz\bar{x}\bar{v} = \mu\delta\bar{\beta}\bar{\delta}\bar{\mu}$ where $|\delta| = k$ and $\delta\bar{\beta}\bar{\delta}\bar{\mu} \in L_2$.

Proof. The conditions say that $uv^sxyz\bar{x}\bar{v}^t\bar{u} \in \mathcal{H}_k(L_1, L_2)$ for all $s \geq t \geq n$. Moreover, by 4.) the hairpin completion can be achieved with a prefix in L_1 , and the longest prefix of $uv^sxyz\bar{x}\bar{v}^t\bar{u}$ belonging to L_1 is a prefix of $uv^sxyz\bar{x}$.

If $\mathcal{H}_k(L_1, L_2)$ is regular, then we have $uv^sxyz\bar{x}\bar{v}^{s+1}\bar{u} \in \mathcal{H}_k(L_1, L_2)$, too, as soon as s is large enough, by a simple pumping argument. For this hairpin completion we must use a suffix belonging to L_2 . This follows from $|y| < |v|$ and a case distinction whether or not z is empty. For $z \neq 1$ we need condition 2.) to see this.

By 5.) the longest suffix of $uv^sxyz\bar{x}\bar{v}^{s+1}\bar{u}$ belonging to L_2 is a suffix of $xyz\bar{x}\bar{v}^{s+1}\bar{u}$. Thus, we can write

$$uv^sxyz\bar{x}\bar{v}^{s+1}\bar{u} = uv^sxyz\bar{x}\bar{v}\bar{v}^s\bar{u} = uv^s\mu\delta\bar{\beta}\bar{\delta}\bar{\mu}\bar{v}^s\bar{u}$$

where $\delta\bar{\beta}\bar{\delta}\bar{\mu}\bar{v}^s\bar{u} \in L_2$ as soon as s is large enough.

(Recall that our second DFA accepts \bar{L}_2 .) Hence, as $p_2 = q_{02} \cdot u$ and $p_2 = p_2 \cdot v$, we see that $\delta\bar{\beta}\bar{\delta}\bar{\mu}\bar{v}^s\bar{u} \in L_2$ if and only if $\delta\bar{\beta}\bar{\delta}\bar{\mu}\bar{u} \in L_2$. Thus, if $\mathcal{H}_k(L_1, L_2)$ is regular, then $\delta\bar{\beta}\bar{\delta}\bar{\mu}\bar{u} \in L_2$. \square

Example 3.15. *Let us take a look at Figure 2 again. Let $A_\kappa = (Q_0, \emptyset_1, \emptyset_2, 0)$, $v_\kappa = a$ and $u = 1$. If we choose $x = a$, $y = 1$, $z = \bar{b}$ and $(d_1, d_2) = (p_1, p_2)$ we can see, that conditions 1.) to 5.) of Lemma 3.14 are satisfied but there is no factorization $ab\bar{a}\bar{a} = \mu\delta\bar{\beta}\bar{\delta}\bar{\mu}$ with $|\delta| = k$ such that $\delta\bar{\beta}\bar{\delta}\bar{\mu}\bar{u} \in L_2$. Hence, the hairpin completion is not regular.*

Lemma 3.16. *The existence of words $x, y, z \in \Sigma^*$ and states $(d_1, d_2) \in \mathcal{Q}_1 \times \mathcal{Q}_2$ satisfying 1.) to 5.) of Lemma 3.14, but where for all factorizations $xyz\bar{x}\bar{v} = \mu\delta\bar{\beta}\bar{\delta}\bar{\mu}$ we have $p_2 \cdot \mu\delta\bar{\beta}\bar{\delta} \notin \mathcal{F}_2$ (and accordingly $\delta\bar{\beta}\bar{\delta}\bar{\mu}\bar{u} \notin L_2$), can be decided in time $\mathcal{O}(n_1^2 n_2^2) \subseteq \mathcal{O}(n^8)$ and in **NL**.*

Proof. It is enough to perform Test 3 and 4 below and to prove the complexity. The tests distinguish whether the word z is empty or non-empty.

Test 3: Decide the existence of words $x, y, z \in \Sigma^*$ with $z \neq 1$ and states $(d_1, d_2) \in \mathcal{Q}_1 \times \mathcal{Q}_2$ satisfying 1.) to 5.) of Lemma 3.14, but where for all factorizations $xyz\bar{x}\bar{v} = \mu\delta\bar{\beta}\bar{\delta}\bar{\mu}$ we have $p_2 \cdot \mu\delta\bar{\beta}\bar{\delta} \notin \mathcal{F}_2$. If we find such a situation, then **stop** and output that $\mathcal{H}_k(L_1, L_2)$ is not regular.

Test 4: Decide the existence of words $x, y \in \Sigma^*$ and states $(d_1, d_2) \in \mathcal{Q}_1 \times \mathcal{Q}_2$ satisfying 1.) to 5.) of Lemma 3.14 with $z = 1$, but where for all factorizations $xy\bar{x}\bar{v} = \mu\delta\bar{\beta}\bar{\delta}\bar{\mu}$ we have $p_2 \cdot \mu\delta\bar{\beta}\bar{\delta} \notin \mathcal{F}_2$. If we find such a situation, then **stop** and output that $\mathcal{H}_k(L_1, L_2)$ is not regular.

The correctness of both tests follows by Lemma 3.14, but even termination of Test 3 is not completely obvious. Termination is due to condition that xy is the longest common prefix of xyz and some word in v^+ . This means, if $z \neq 1$, then there exists a letter a such that $z \in a\Sigma^*$ and xya is no prefix of any word in v^+ . Now $|y| < |v|$, hence we see that $xyz\bar{x}\bar{v} = \mu\delta\bar{\beta}\bar{\delta}\bar{\mu}$ implies $\mu\delta \leq xy$.

Thus it is enough to check the computation starting in state $d_2 \in \mathcal{Q}_2$ when reading the word $\bar{y}\bar{x}$. Test 3 yields “not regular” if we find such a computation which after more than $k-1$ steps does not meet any final state in \mathcal{F}_2 . We do not need the word z , we just have to know that (c_1, c_2, d_1, d_2) is in the precomputed table of a -bridges (cf. Lemma 3.4) where a is a letter such that xya is no prefix of any word in v^+ . It is obvious that Test 3 can be performed in polynomial time as well as in **NL**.

Test 4 is for $z = 1$, so in any case the number of factorizations $xy\bar{x}\bar{v} = \mu\delta\bar{\beta}\bar{\delta}\bar{\mu}$ is polynomial. It is again obvious that Test 4 can be performed polynomial time as well as in **NL**.

For the exact time complexity we refer to Section 3.4. \square

The following lemmas complete the proof of Theorem 3.1 and 3.2.

Lemma 3.17. *Suppose no outcome of Tests 1, 2, 3, and 4 is “not regular”. Then the hairpin completion $\mathcal{H}_k(L_1, L_2)$ is regular.*

Proof. Let $\pi \in \mathcal{H}_k(L_1, L_2)$. Write $\pi = \gamma\alpha\beta\bar{\alpha}\bar{\gamma}$ with $|\gamma|$ minimal such that either $\gamma\alpha\beta\bar{\alpha} \in L_1$ or $\alpha\beta\bar{\alpha}\bar{\gamma} \in L_2$. By symmetry we assume $\gamma\alpha\beta\bar{\alpha} \in L_1$. We may assume that $|\gamma| > n^4$. We can factorize $\gamma = uvw$ with $|uv| \leq n^4$ and $|v| \geq 1$ such that there are runs as follows:

$$1.) \quad q_{01} \xrightarrow{u} p_1 \xrightarrow{v} p_1 \xrightarrow{w\alpha\beta\bar{\alpha}} f_1 \xrightarrow{\bar{w}} q_1 \xrightarrow{\bar{v}} q_1 \xrightarrow{\bar{u}} q'_1,$$

$$2.) \quad q_{02} \xrightarrow{u} p_2 \xrightarrow{v} p_2 \xrightarrow{w\alpha\beta\bar{\alpha}} f_2 \xrightarrow{\bar{w}} q_2 \xrightarrow{\bar{v}} q_2 \xrightarrow{\bar{u}} q'_2,$$

$$3.) \quad f_1 \in \mathcal{F}_1.$$

We infer from Test 1/2 that $w\alpha$ is a prefix of some word in v^+ . Hence, we can write $w\alpha\beta = v^mxyz$ with $m \geq 0$ such that v^mxy is the maximal common prefix of $w\alpha\beta$ and some word in v^+ , $w\alpha \in v^*x$ with $k \leq |x| < |v| + k$, and $|y| < |v|$.

We see that for some $s \geq t \geq 0$ we can write

$$\pi = uv^sxyz\bar{x}\bar{v}^t\bar{u}.$$

Moreover, $uv^sxyz\bar{x}\bar{v}^t\bar{u} \in \mathcal{H}_k(L_1, L_2)$ for all $s \geq t \geq 0$. There are only finitely many choices for u, v, x, y (due to the lengths bounds) and for each of them there is a regular set R_z in a finite collection of regular sets such that

$$\pi \in \{uv^sxyz\bar{x}\bar{v}^t\bar{u} \mid s \geq t \geq 0\} \subseteq \mathcal{H}_k(L_1, L_2).$$

Note that the sets $\{uv^sxyz\bar{x}\bar{v}^t\bar{u} \mid s \geq t \geq 0\}$ need not to be regular, in general. If we bound however t by n then the finite union

$$\bigcup_{0 \leq t \leq n} \{uv^sxyz\bar{x}\bar{v}^t\bar{u} \mid s \geq t\}$$

is regular. Thus, we may assume that $t > n$. Let $e_2 = p_2 \cdot x\bar{z}\bar{y}\bar{x}$. We have $e_2 \cdot \bar{v}^n = q_2$ and, if we see a final state during the computation of $e_2 \cdot \bar{v}^n$, then for all $t \geq s \geq n$ and $z \in R_z$ we see that $uv^sxyz\bar{x}\bar{v}^t\bar{u} \in \mathcal{H}_k(L_1, L_2)$, due to a suffix in L_2 and, $uv^n v^+ xyR_z \bar{x}\bar{v}^+ \bar{v}^n \bar{u} \subseteq \mathcal{H}_k(L_1, L_2)$.

Otherwise Test 3/4 tells us that for all $z \in R_z$ the word $xyz\bar{x}\bar{v}$ has a factorization $\mu\delta\nu\bar{\delta}\bar{\mu}$ such that $|\delta| = k$ and $\delta\nu\bar{\delta}\bar{\mu}\bar{u} \in L_2$. The paths $q_{02} \cdot u = p_2$ and $p_2 \cdot v = p_2$ yield $\delta\nu\bar{\delta}\bar{\mu}\bar{v}^+ \bar{u} \subseteq L_2$ and, again, $uv^n v^+ xyR_z \bar{x}\bar{v}^+ \bar{v}^n \bar{u} \subseteq \mathcal{H}_k(L_1, L_2)$.

Hence, the hairpin completion $\mathcal{H}_k(L_1, L_2)$ is a finite union of regular languages. \square

Lemma 3.18. *It is NL-hard to decide whether the hairpin completion $\mathcal{H}_k(L_1, \emptyset)$ is regular.*

The well-known **NL**-complete *Graph-Accessibility-Problem* can easily be reduced to the following decision problem for DFAs.

- Input: A DFA where the accepted language L satisfies $L \subseteq b(ab + ba)^*$ with $a \neq b$.
- Problem: $L = \emptyset$?

Now, for L as above and $k > 2$ consider $L' = a^+L\bar{a}^k$. Then

$$\mathcal{H}_k(L', \emptyset) = \{a^\ell w \bar{a}^m \mid w \in L \wedge \ell \geq m \geq k\}.$$

Hence, $\mathcal{H}_k(L', \emptyset)$ is regular if and only if $L = \emptyset$.

3.4 Time Complexity

Let us recall that the construction of the NFA \mathcal{A} , Test 1, and Test 2 can be performed in time $\mathcal{O}(N^2)$. This is $\mathcal{O}(n^8)$ in the general case and $\mathcal{O}(n^6)$ for $L_1 = \overline{L_2}$. In this section we will sketch how Test 3 and 4 can be implemented in order to meet the same time bounds.

3.4.1 Test 3

Let $\kappa \in K$ be fixed, let $v = v_\kappa$, $A = A_\kappa = ((p_1, p_2), q_1, q_2, 0)$ and u be some word leading from an initial bridge to A .

In order to perform Test 3 we create two tables T_1 and T_2 . The table T_1 holds all pairs $(c_2, d_1) \in \mathcal{Q}_2 \times \mathcal{Q}_1$ such that a word x exists with

- 1.) $k \leq |x| < |v| + k$ and x is a prefix of a word in v^+ ,
- 2.) $p_2 \cdot x = c_2$,
- 3.) $d_1 \cdot \overline{xv}^{n_1} = q_1$, and during the computation of $d_1 \cdot \overline{xv}^{n_1}$ we see a final state after exactly k steps and then never again.

We call x a witness for $(c_2, d_1) \in T_1$.

The table T_2 holds all triples $(c_1, d_2, a) \in \mathcal{Q}_1 \times \mathcal{Q}_2 \times \Sigma$ such that a prefix $y' < v$ exists with

- 1.) $y'a$ is no prefix of v ,
- 2.) $p_1 \cdot y' = c_1$,
- 3.) $d_2 \cdot \overline{y'v}^{n_2} = q_2$, and during the computation of $d_2 \cdot \overline{y'v}^{n_2}$ we do not see a final state after more than $k - 1$ steps.

We call y' a witness for $(c_1, d_2, a) \in T_2$.

By backwards computing in the second component, the tables T_1 and T_2 can be created in $\mathcal{O}(N_\kappa n_1^2)$ and $\mathcal{O}(N_\kappa n_2^2)$, respectively.

Lemma 3.19. *The outcome of Test 3 is “not regular” if and only if there exists a pair $(c_2, d_1) \in T_1$ and a triple $(c_1, d_2, a) \in T_2$ such that (c_1, c_2, d_1, d_2) is an a -bridge.*

Proof. Assume $(c_2, d_1) \in T_1$, $(c_1, d_2, a) \in T_2$, and (c_1, c_2, d_1, d_2) is an a -bridge. Let x and y' be the witnesses for $(c_2, d_1) \in T_1$ and $(c_1, d_2, a) \in T_2$, respectively. Choose $z \in B(c_1, c_2, d_1, d_2) \cap a\Sigma^*$ and y such that xy is a prefix of some word in v^+ , $|xy| \equiv |y'| \pmod{|v|}$, and $|y| < |v|$. Verify that x, y, z and (d_1, d_2) satisfy the conditions 1.) to 5.) of Lemma 3.14.

For any factorization $xyz\overline{xv} = \mu\delta\beta\overline{\delta\mu}$ with $|\delta| = k$, the word $\mu\delta$ has to be a prefix of xy , since xya is no prefix of vx . During the computation of $d_2 \cdot \overline{y'v}^{n_2}$ we do not see a final state after more than $k - 1$ steps. The same holds for the computation of $d_2 \cdot \overline{y\overline{xv}}^{n_2}$. Hence, $\delta\beta\overline{\delta\mu}$ is not included in L_2 .

Now assume that $x, y, z \in \Sigma^*$, $z \neq 1$ and $(d_1, d_2) \in \mathcal{Q}_1 \times \mathcal{Q}_2$ exist, which satisfy the conditions 1.) to 5.) of Lemma 3.14 but where for all factorizations $xyz\overline{xv} = \mu\delta\beta\overline{\delta\mu}$ we have $\delta\beta\overline{\delta\mu} \notin L_2$. Choose $y' < v$ such that $|xy| \equiv |y'| \pmod{|v|}$. Let $c_2 = p_2 \cdot x$, $c_1 = p_1 \cdot y'$ and $a \in \Sigma$ be the first letter of z . (c_1, c_2, d_1, d_2) is an a -bridge. If we saw a final state after more than $k - 1$

steps during the computation of $d_2 \cdot \bar{y}\bar{v}^{n_2}$, then a factorization $xyz\bar{x}\bar{v} = \mu\delta\beta\bar{\delta}\bar{\mu}$ with $\delta\beta\bar{\delta}\bar{\mu}\bar{u} \in L_2$ would exist. Hence, y' is a witness for $(c_1, d_2, a) \in T_2$ and, obviously, x is a witness for $(c_2, d_1) \in T_1$. \square

The set of all first components of T_1 (T_2) is bounded by both, the size of N_κ and n_2 (resp. n_1). Therefore it is of size $\mathcal{O}(n_1 \cdot \min(N_\kappa, n_2))$ (resp. $\mathcal{O}(n_2 \cdot \min(N_\kappa, n_1))$). By symmetry, assume $n_2 \leq n_1$. Since the table of a -bridges is precomputed, we can perform Test 3 in

$$\begin{aligned} & \mathcal{O} \left(\sum_{\kappa \in K} (N_\kappa n_1^2 + N_\kappa n_2^2 + n_1 n_2 \cdot \min(N_\kappa, n_1) \cdot \min(N_\kappa, n_2)) \right) \subseteq \\ & \mathcal{O} \left(n_{12} n_1^3 n_2 + n_{12} n_1 n_2^3 + \sum_{\kappa \in K, N_\kappa \geq n_2} n_1^2 n_2^2 + \sum_{\kappa \in K, N_\kappa < n_2} N_\kappa^2 n_1 n_2 \right) \end{aligned}$$

(Recall that $n_1 \leq n \leq n_{12} \leq n_1 n_2 \leq n^2$ and $\sum_{\kappa \in K} N_\kappa \leq N = n_{12} n_1 n_2$.)

Since there are at most $n_{12} n_1$ strongly connected components with a size of more than n_2 states, we have

$$\sum_{\kappa \in K, N_\kappa \geq n_2} n_1^2 n_2^2 \leq n_{12} n_1^3 n_2^2.$$

For the last term we can use the approximation

$$\sum_{\kappa \in K, N_\kappa < n_2} N_\kappa^2 n_1 n_2 \leq \sum_{\kappa \in K, N_\kappa < n_2} N_\kappa n_1 n_2^2 \leq n_{12} n_1^2 n_2^3.$$

Test 3 can be performed in time $\mathcal{O}(n_{12} n_1^3 n_2^2) \subseteq \mathcal{O}(n^7)$ in the general case and in time $\mathcal{O}(n^6)$ for $L_1 = \bar{L}_2$.

3.4.2 Test 4

Let $\kappa \in K$ be fixed, let $v = v_\kappa$, $A = A_\kappa = ((p_1, p_2), q_1, q_2, 0)$ and u be some word leading from an initial bridge to A . For the final test we have to compute all words x and y such that there are runs

$$p_1 \xrightarrow{xy} d_1 \xrightarrow{\bar{x}\bar{v}^{n_1}} q_1 \quad \text{and} \quad p_2 \xrightarrow{x} d_2 \xrightarrow{\bar{y}\bar{x}\bar{v}^{n_2}} q_2$$

and together with $z = 1$ the conditions 1.) to 5.) of Lemma 3.14 are satisfied. Moreover, in addition and similar as in Test 3 we demand that in the computation of $d_2 \cdot \bar{y}\bar{x}\bar{v}^{n_2}$ we do not meet any final state after more than $k - 1$ steps. (In case such a final state exists, either condition 5.) is breached or a factorization $xy\bar{x}\bar{v} = \mu\delta\beta\bar{\delta}\bar{\mu}$ with $|\delta| = k$ and $\delta\beta\bar{\delta}\bar{\mu}\bar{u} \in L_2$ exists.) In time $\mathcal{O}(N_\kappa^2)$ we compute all pairs (x, y) satisfying these conditions.

We also compute at this stage a number $0 \leq \ell(x) \leq N_\kappa + |x| + k$ as follows. Let xx' be the prefix of some word in v^+ of length $|x| + k$. (Thus, $|x'| = k$.) On the run $p_2 \cdot vxx'$ we let $\ell(x) = |x''|$ such that $x'' \leq vxx'$ is the maximal prefix where $p_2 \cdot x'' \in \mathcal{F}_2$ is a final state. If there is no final state on this run, we let $\ell(x) = 0$. If a factorization $xy\bar{x}\bar{v} = \mu\delta\beta\bar{\delta}\bar{\mu}$ with $|\delta| = k$ and $\delta\beta\bar{\delta}\bar{\mu}\bar{u} \in L_2$ exists, then the value $\ell(x)$ gives us the lower bound $|xy\bar{x}\bar{v}| - \ell(x)$ for the length of μ . (Note that $|xy\bar{x}\bar{v}| \geq \ell(x)$.)

Let $m(x, xy)$ be the length of the longest $\mu \leq vx$ such that a factorization $xy\bar{x}\bar{v} = \mu\delta\beta\bar{\delta}\bar{\mu}$ with $|\delta| = k$ exists (without the condition $\delta\beta\bar{\delta}\bar{\mu}\bar{u} \in L_2$). The outcome of Test 4 is negative if and only if $m(x, xy) \geq |xy\bar{x}\bar{v}| - \ell(x)$ and $\ell(x) - k \geq |xy\bar{x}\bar{v}|/2$.

We need to precompute the values $m(x, xy)$ efficiently, which turns out to be a little bit tricky.

First let us fix x and write $vx = xv'$. We wish to match positions in $v'v'$ with positions in \bar{v}^2 . let us *mark* those $1 \leq j \leq |v| + k$ where the j th letter $v'(j)$ in v' is equal to the j -th letter $\bar{v}(j)$ in \bar{v}^2 . For each x one scan through v' and \bar{v}^2 is enough. Having stored these marked positions in a table of size $N_\kappa + k$ we can compute for each $1 \leq j \leq |v|$ the maximal value m such that all positions $j, \dots, j + m$ are marked. This is possible in $\mathcal{O}(N_\kappa^2)$.

All in all Test 4 can be performed in

$$\mathcal{O}\left(\sum_{\kappa \in K} N_\kappa^2\right) \subseteq \mathcal{O}(N^2).$$

References

- [1] B. S. Baker and R. V. Book. Reversal-bounded multi-pushdown machines. *Foundations of Computer Science, Annual IEEE Symposium on*, 0:207–211, 1972.
- [2] J. Berstel and C. Reutenauer. *Rational series and their languages*. Springer-Verlag New York, Inc., New York, NY, USA, 1988.
- [3] T. Ceccherini-Silberstein. On the growth of linear languages. *Advances in Applied Mathematics*, 35(3):243 – 253, 2005.
- [4] D. Cheptea, C. Martin-Vide, and V. Mitrana. A new operation on words suggested by DNA biochemistry: Hairpin completion. *Transgressive Computing*, pages 216–228, 2006.
- [5] R. Deaton, R. Murphy, M. Garzon, D. Franceschetti, and S. Stevens. Good encodings for DNA-based solutions to combinatorial problems. *Proc. of DNA-based computers DIMACS Series*, 44:247–258, 1998.
- [6] V. Diekert, S. Kopecki, and V. Mitrana. On the hairpin completion of regular languages. In M. Leucker and C. Morgan, editors, *ICTAC*, volume 5684 of *Lecture Notes in Computer Science*, pages 170–184. Springer, 2009.
- [7] M. Garzon, R. Deaton, P. Neathery, R. Murphy, D. Franceschetti, and E. Stevens. On the encoding problem for DNA computing. *The Third DIMACS Workshop on DNA-Based Computing*, pages 230–237, 1997.
- [8] M. Garzon, R. Deaton, L. Nino, S. Stevens Jr., and M. Wittner. Genome encoding for DNA computing. *Proc. Third Genetic Programming Conference*, pages 684–690, 1998.
- [9] P. Gawrychowski, D. Krieger, N. Rampersad, and J. Shallit. Finding the growth rate of a regular or context-free language in polynomial time. In *Developments in Language Theory*, pages 339–358, 2008.

- [10] S. A. Greibach. A note on undecidable properties of formal languages. *Mathematical Systems Theory*, 2(1):1–6, 1968.
- [11] J. E. Hopcroft and J. D. Ulman. *Introduction to Automata Theory, Languages and Computation*. Addison-Wesley, 1979.
- [12] L. Kari, S. Konstantinidis, E. Losseva, P. Sosík, and G. Thierrin. Hairpin structures in DNA words. In A. Carbone and N. A. Pierce, editors, *DNA*, volume 3892 of *Lecture Notes in Computer Science*, pages 158–170. Springer, 2005.
- [13] L. Kari, K. Mahalingam, and G. Thierrin. The syntactic monoid of hairpin-free languages. *Acta Inf.*, 44(3-4):153–166, 2007.
- [14] W. Kuich. On the entropy of context-free languages. *Information and Control*, 16:173–200, 1970.
- [15] F. Manea, V. Mitrana, and T. Yokomori. Two complementary operations inspired by the DNA hairpin formation: Completion and reduction. *Theor. Comput. Sci.*, 410(4-5):417–425, 2009.
- [16] C. H. Papadimitriou. *Computational Complexity*. Addison Wesley, 1994.
- [17] K. Sakamoto, H. Gouzu, K. Komiya, D. Kiga, S. Yokoyama, T. Yokomori, and M. Hagiya. Molecular Computation by DNA Hairpin Formation. *Science*, 288(5469):1223–1226, 2000.