

Iterated TGR Languages: Membership Problem and Effective Closure Properties (Extended Abstract)

Ian McQuillan¹, Kai Salomaa², and Mark Daley³

¹ Department of Computer Science, University of Saskatchewan, Saskatoon,
Saskatchewan, Canada S7N 5A9

mcquillan@cs.usask.ca

² School of Computing, Queen's University, Kingston, Ontario, Canada K7L 3N6

ksalomaa@cs.queensu.ca

³ Department of Computer Science and Department of Biology, University of Western
Ontario, London, Ontario, Canada N6A 5B7

daley@csd.uwo.ca

Abstract. We show that membership is decidable for languages defined by iterated template-guided recombination systems when the set of templates is regular and the initial language is context-free. Using this result we show that when the set of templates is regular and the initial language is context-free (respectively, regular) we can effectively construct a push-down automaton (respectively, finite automaton) for the corresponding iterated template-guided recombination language.

1 Introduction

The spirotrichous ciliates are a type of unicellular protozoa which possess a unique and fascinating genetic behaviour. Each ciliate cell contains two types of nuclei, *macronuclei* which are responsible for the day-to-day “genetic house-keeping” of the cell, and *miconuclei* which are functionally inert, but used in reproduction. This is in contrast to, e.g., mammalian cells which have only one micronucleus. Although they reproduce asexually, ciliates are also capable of sexual activity in which they exchange haploid micronuclear genomes. This results in each ciliate getting a “genetic facelift” by combining its own genes with those of a mate. After creating a new, hybrid, micronucleus, each ciliate will then regenerate its macronucleus. It is this process of macronuclear regeneration that is of principle interest to us here.

In the spirotrichous ciliates in particular, this macronuclear regeneration involves an intricate process of genetic gymnastics. Suppose that a functional gene in the macronucleus can be divided into 5 sections and written as follows: 1-2-3-4-5. In many cases, the micronuclear form of the same gene may have the segments in a completely different order and include additional segments not found in the macronucleus. For the example given above, a micronuclear gene

may appear as: 3-x-5-y-1-z-4-2. For the ciliate to produce a functional macronucleus and continue living, it must *descramble* these micronuclear genes. (See, e.g., [10] for further detail).

A biological model for this descrambling process, based on template-guided DNA recombination, was proposed in [11]. This model was formalized as an operation on words and languages in [3] which also introduced the notion of a template-guided recombination system (TGR system). It was then shown in [4] that a TGR system with a regular set of templates preserves regularity, that is, for a regular initial language, the language resulting from iterated application of the TGR system is always regular. This is in striking contrast to splicing systems since the splicing language generated by a regular set of rules and a finite initial language need not be recursive [8]. In fact, [4] shows much more generally that the operation defined by a TGR system with a regular set of templates preserves any language family that is a full AFL [7,12].

However, the above results are non-constructive and, in particular, do not give an algorithm to decide the membership problem for the language defined by a TGR system, even in the case where the initial language is finite and the set of templates is regular. Here we show that the uniform membership problem for the language defined by a TGR system is decidable when the initial language is context-free and the set of templates is regular. The nonuniform membership problem (where the TGR system is fixed) can be decided in polynomial time. The decidability result is extended for languages that are extensions of the context-free languages, such as the indexed languages, or, more generally, for languages that belong to a full AFL satisfying certain natural effectiveness conditions.

Moreover, we use this result to positively solve the main open problem from [4]. That is, given a context-free (respectively, regular) initial language and a regular set of templates, we can effectively construct a pushdown automaton (respectively, a finite automaton) for the language defined by the TGR system. Using a variant of the decision algorithm for the membership problem, we effectively find a deterministic finite automaton (DFA) for the subset of templates that can be used in some recombination operation and this, together with the results of [4], enables us to construct the pushdown automaton (respectively, the finite automaton) for the language defined by the TGR system. This result also holds for regular sets of templates and initial languages from an arbitrary full AFL that satisfies certain effectiveness conditions.

Both the algorithm for the membership problem and the method for finding the set of useful templates use expensive brute-force techniques. It remains an open question, whether it is possible to find a more efficient algorithm, at least in the case where both the initial language and the set of templates are regular.

2 Preliminaries

Here we recall some basic definitions needed in the next section. For all unexplained notions related to formal languages we refer the reader e.g. to [12]. Recent work on language classes and bio-operations can be found e.g. in [2].

In the following Σ is a finite alphabet and the set of all words over Σ is Σ^* . The length of a word $w \in \Sigma^*$ is $|w|$. The i th symbol of a word $w \in \Sigma^*$ is denoted $w[i]$, $i = 1, \dots, |w|$. A language is a subset of Σ^* . The sets of all prefixes, all suffixes and all subwords of words in L are denoted, respectively, $\text{pref}(L)$, $\text{suf}(L)$, $\text{subw}(L)$.

A family of languages is said to be a *full abstract family of languages* (full AFL) [7,12] if it contains a nonempty language and is closed under the following operations: union, Kleene plus, homomorphism, inverse homomorphism, and intersection with regular languages.

Definition 2.1. [3,4] A template-guided recombination system (*TGR system*) is a tuple $\varrho = (T, \Sigma, n_1, n_2)$, where Σ is a finite alphabet, $T \subseteq \Sigma^*$ is the template language, and $n_1, n_2 \in \mathbb{N}$.

Let $x, y \in \Sigma^*$ and $t \in T$. The recombination operation defined by ϱ is given by: $(x, y) \vdash_t^\varrho w$ if and only if we can write

$$x = u\alpha\beta d, \quad y = e\beta\gamma v, \quad t = \alpha\beta\gamma \quad \text{and} \quad w = u\alpha\beta\gamma v$$

for some $u, v, d, e \in \Sigma^*$, $\alpha, \gamma \in \Sigma^{\geq n_1}$ and $\beta \in \Sigma^{n_2}$. For $L \subseteq \Sigma^*$ we define $\varrho(L) = \{w \in \Sigma^* \mid (x, y) \vdash_t^\varrho w \text{ for some } x, y \in L, t \in T\}$.

Let $\varrho = (T, \Sigma, n_1, n_2)$ be a TGR system and let $L \subseteq \Sigma^*$. We define the iteration $\varrho^{(*)}$ of the operation ϱ by setting $\varrho^{(0)}(L) = L$, and defining

$$\varrho^{(i+1)}(L) = \varrho^{(i)}(L) \cup \varrho(\varrho^{(i)}(L)) \quad \text{for all } i \geq 0. \quad (1)$$

Denote $\varrho^{(*)}(L) = \bigcup_{i=0}^{\infty} \varrho^{(i)}(L)$.

Let $\varrho = (T, \Sigma, n_1, n_2)$ be a TGR system and let $L \subseteq \Sigma^*$. A word $t \in T$ is said to be *useful* on (L, ϱ) if t can be used in iterated application of ϱ on the initial language L . It is shown in [4] that $t \in T$ is useful on (L, ϱ) if and only if $|t| \geq 2n_1 + n_2$ and t is a subword of some word in $\varrho^{(*)}(L)$. The TGR system ϱ is said to be *useful on L* if every word of T is useful on (L, ϱ) . The *useful subset of ϱ on L* is the set of all words in T which are useful on (L, ϱ) .

3 Membership Problem

Here we show that for a context-free language L and a TGR system $\varrho = (T, \Sigma, n_1, n_2)$ where T is regular, the uniform membership problem for the language $\varrho^{(*)}(L)$ is decidable.

We want to establish properties concerning how many recombination operations are required to produce some subword of a word w when it is known that w requires a given number of recombination operations. For this purpose it turns out to be useful to consider “marked variants” of words over Σ . The marked variants associate states of a DFA recognizing the set of templates T and length information with certain positions in the word. This additional control information is used to keep track of the templates (or strictly speaking equivalence classes of templates) that can be used in the recombination operations.

For the above purpose we next introduce some technical notation. Let $\varrho = (T, \Sigma, n_1, n_2)$ be a TGR system and

$$A = (\Sigma, Q, q_0, F, \delta) \quad (2)$$

be a DFA that recognizes T . Denote $\vec{Q} = \{\vec{q} \mid q \in Q\}$, $\overleftarrow{Q} = \{\overleftarrow{q} \mid q \in Q\}$. For $n \in \mathbb{N}$ let $[n] = \{0, 1, \dots, n\}$. We define the extended alphabet $\Sigma[\varrho]$ as

$$\Sigma[\varrho] = \Sigma \times \mathcal{P}((\vec{Q} \cup \overleftarrow{Q}) \times [n_1]). \quad (3)$$

The first component of elements of $\Sigma[\varrho]$ is an element of Σ and the second component consists of a set of states of Q each marked with a “right arrow” or a “left arrow”. Additionally, each state is associated with an index from $\{0, 1, \dots, n_1\}$.

The projections from $\Sigma[\varrho]$ to Σ and to $\mathcal{P}((\vec{Q} \cup \overleftarrow{Q}) \times [n_1])$ are denoted, respectively, π_1^ϱ and π_2^ϱ . When ϱ is clear from the context, we denote the projections simply as π_1 and π_2 . The projection π_1 is in the natural way extended to a morphism $\Sigma[\varrho]^* \rightarrow \Sigma^*$.

Let $L \subseteq \Sigma^*$. The *T-controlled marked variant* of L is the largest language $C_T(L) \subseteq \Sigma[\varrho]^*$ such that the below conditions (i) and (ii) hold¹. The notations refer to (2) that gives a DFA for the language T .

- (i) For every $w \in C_T(L)$, $\pi_1(w) \in L$.
- (ii) Assume that $w \in C_T(L)$ and $(p, j) \in \pi_2(w[i])$, $1 \leq i \leq |w|$, $p \in \vec{Q} \cup \overleftarrow{Q}$, $j \in [n_1]$.
 - (a) If $p \in \vec{Q}$, then $\pi_1(w)$ has a subword u starting at the $(i + 1)$ th position such that $|u| \geq j$ and $\delta(p, u) \in F$.
 - (b) If $p \in \overleftarrow{Q}$, then $\pi_1(w)$ has a subword u ending at the $(i - 1)$ th position such that $|u| \geq j$ and $\delta(q_0, u) = p$.

Note that for any $w \in L$, the word w' is in $C_T(L)$ where w' is obtained from w by replacing each symbol $c \in \Sigma$ by $(c, \emptyset) \in \Sigma[\varrho]$. We identify words w and w' and in this way we can view L to be a subset of $C_T(L)$.

According to (i) and (ii) above, the elements (p, j) , $p \in \vec{Q} \cup \overleftarrow{Q}$ occurring in symbols of a word $w \in C_T(L)$ place conditions on what kind of subwords w must have starting directly after or ending directly before that position. If $p \in \vec{Q}$, this means that $\pi_1(w)$ must have a subword u starting from the next position that is a suffix of a word in T , u is of length at least j , and the state p corresponds to this suffix (that is, $\delta(p, u) \in F$). If $p \in \overleftarrow{Q}$, this means that $\pi_1(w)$ must have a subword u ending at the previous position that is a prefix of a word in T , u has length at least j , and the state p corresponds to this prefix.

¹ Note that the union of languages satisfying this property also satisfies this property, and so the largest language must exist.

We still need the following notation to manipulate words over the alphabet $\Sigma[\varrho]$. Let $w \in \Sigma[\varrho]^*$, $1 \leq i \leq |w|$, $p \in (\overrightarrow{Q} \cup \overleftarrow{Q})$ and $j \in [n_1]$. Then $w[i \leftarrow (p, j)]$ denotes the word obtained from w by adding (p, j) to the second component of the i th symbol, that is, the second component of the i th symbol is changed to be $\pi_2(w[i]) \cup \{(p, j)\}$.

We say that a word $w \in \Sigma[\varrho]^*$ is *well formed* if $|w| \geq 2$ and the following three conditions hold: (i) $\pi_2(w[1]) \subseteq \overleftarrow{Q} \times [n_1]$, (ii) $\pi_2(w[|w|]) \subseteq \overrightarrow{Q} \times [n_1]$, and (iii) $\pi_2(w[j]) = \emptyset$ when $1 < j < |w|$.

In a well formed marked word the first symbol contains only elements of the type (\overleftarrow{p}, j) as markers, and the last symbol contains only elements of the type (\overrightarrow{p}, j) as markers, $p \in Q$, $j \in [n_1]$. Symbols of w other than the first or the last symbol have \emptyset as the second component.

The set of all well formed words over $\Sigma[\varrho]$ is denoted by $\mathcal{WF}(\Sigma[\varrho])$

The following lemma says, very roughly speaking, that if w is a subword of $\varrho^{(k+1)}(L)$ but w is not a subword of $\varrho^{(k)}(L)$, then w has a proper subword that is a subword of $\varrho^{(k)}(L)$ but not a subword of $\varrho^{(k-1)}(L)$. The statement in the previous sentence is oversimplified and does not hold as such. To be precise, in order to be able to establish the required property we need to add to the subwords information on the states of the DFA for T associated with the templates used in the recombination operations, that is, we need to consider subwords of the T -controlled marked variant of $\varrho^{(k)}(L)$, $k \geq 1$.

For $m, n \in \mathbb{N}$ we define the non-negative difference of m and n , $m \ominus n$, as $m - n$ if $m \geq n$ and $m \ominus n = 0$ otherwise.

Lemma 3.1. *Let $\varrho = (T, \Sigma, n_1, n_2)$ where T is regular and let A as in (2) be a DFA that recognizes T . Let $k \geq 1$ and $L \subseteq \Sigma^*$.*

We claim that if $w \in \mathcal{WF}(\Sigma[\varrho])$ and

$$w \in \text{subw}(C_T(\varrho^{(k+1)}(L))) - \text{subw}(C_T(\varrho^{(k)}(L))) \quad (4)$$

then one of the below cases (P1)–(P4) holds:

(P1) $w = u\alpha\beta\gamma v$, $\pi_1(\alpha\beta\gamma) \in T$, $|\beta| = n_2$, $|\alpha|, |\gamma| \geq n_1$, $u\alpha\beta \in \text{subw}(C_T(\varrho^{(k)}(L))) \cap \mathcal{WF}(\Sigma[\varrho])$, $\beta\gamma v \in \text{subw}(C_T(\varrho^{(k)}(L))) \cap \mathcal{WF}(\Sigma[\varrho])$,

(P2) $w = u\alpha\beta\gamma'$, $|\beta| = n_2$, $|\alpha| \geq n_1$, $|\gamma'| \geq 1$, $u\alpha\beta \in \text{subw}(C_T(\varrho^{(k)}(L))) \cap \mathcal{WF}(\Sigma[\varrho])$, $\beta\gamma'[\beta\gamma'] \leftarrow (\overleftarrow{p}, n_1 \ominus |\gamma'|) \in \text{subw}(C_T(\varrho^{(k)}(L))) \cap \mathcal{WF}(\Sigma[\varrho])$, where $p = \delta(q_0, \alpha\beta\gamma')$.

(P3) $w = \alpha'\beta\gamma v$, $|\beta| = n_2$, $|\gamma| \geq n_1$, $|\alpha'| \geq 1$, $\alpha'\beta[1 \leftarrow (\overleftarrow{p}, n_1 \ominus |\alpha'|)] \in \text{subw}(C_T(\varrho^{(k)}(L))) \cap \mathcal{WF}(\Sigma[\varrho])$, $p \in Q$, $\beta\gamma v \in \text{subw}(C_T(\varrho^{(k)}(L))) \cap \mathcal{WF}(\Sigma[\varrho])$, where $\delta(p, \alpha'\beta\gamma) \in F$.

(P4) $w = \alpha'\beta\gamma'$, $|\beta| = n_2$, $|\alpha'|, |\gamma'| \geq 1$, $\alpha'\beta[1 \leftarrow (\overleftarrow{p}, n_1 \ominus |\alpha'|)] \in \text{subw}(C_T(\varrho^{(k)}(L))) \cap \mathcal{WF}(\Sigma[\varrho])$, $p \in Q$, $\beta\gamma'[\beta\gamma'] \leftarrow (\overrightarrow{p_1}, n_1 \ominus |\gamma'|) \in \text{subw}(C_T(\varrho^{(k)}(L))) \cap \mathcal{WF}(\Sigma[\varrho])$, where $\delta(p, \alpha'\beta\gamma') = p_1$.

Furthermore, in any decomposition of w as in (P1)–(P4) at most one of the two mentioned marked words of $\text{subw}(C_T(\varrho^{(k)}(L)))$ can be in $\text{subw}(C_T(\varrho^{(k-1)}(L)))$.

We should note that in (P2), (P3) and (P4) in Lemma 3.1 it is essential that we add the new marker states to the resulting subwords. For example, using the notations of (P4), it is quite possible that $\pi_1(\alpha'\beta) \in \text{subw}(\varrho^{(k-1)}(L))$ and $\pi_1(\beta\gamma') \in \text{subw}(\varrho^{(k-1)}(L))$ because $\alpha'\beta$ could be part of a word that does not allow recombination using any template of T with the words where $\beta\gamma'$ occurs as a subword. The marked variants of the words prevent this possibility by storing the appropriate states and length information in the first symbol of α' and in the last symbol of γ' . The marker information forces that $\alpha'\beta$ (respectively, $\beta\gamma'$) must occur in a position where the immediately preceding (respectively, immediately following) subword contains a suffix (respectively, a prefix) that allows us to complete $\alpha'\beta\gamma'$ into a template of T .

Due to length restrictions the technical proof of Lemma 3.1 is omitted. We refer the reader to [9] for the proof of Lemma 3.1.

Using Lemma 3.1 we get the following property that will be essential for deciding the membership problem. Also we note that Lemma 3.2 (i) is not a special case of (ii) (although their proofs are similar) and hence we include both statements. The proof of Lemma 3.2 is available in [9].

Lemma 3.2. *Let $\varrho = (T, \Sigma, n_1, n_2)$ be a TGR-system where T is regular and $L \subseteq \Sigma^*$.*

- (i) *If $w \in \varrho^{(k)}(L) - \varrho^{(k-1)}(L)$, $k \geq 1$, then $|w| - n_2 - 1 \geq k$.*
- (ii) *If $w \in \text{subw}(\varrho^{(k)}(L)) - \text{subw}(\varrho^{(k-1)}(L))$, then $|w| - n_2 - 1 \geq k$.*

Theorem 3.1. *Given a TGR system $\varrho = (T, \Sigma, n_1, n_2)$ with T regular, a context-free language L and a word $w \in \Sigma^*$, it is decidable whether or not $w \in \varrho^{(*)}(L)$.*

Furthermore, it is decidable whether or not $w \in \text{subw}(\varrho^{()}(L))$.*

Proof. Let $A = (\Sigma, Q, q_0, F, \delta)$ be a DFA that recognizes T . Given a pushdown automaton B_i for $\varrho^{(i)}(L)$, $i \geq 0$, we can construct a pushdown automaton B_{i+1} for $\varrho^{(i+1)}(L)$ as follows. Let $\beta \in \Sigma^{n_2}$ and $q \in Q$. We define $L_1(B_i, \beta, q) = \{ w \in \text{pref}(L(B_i)) \mid w = u\alpha\beta, |\alpha| \geq n_1, \delta(q_0, \alpha\beta) = q \}$, $L_2(B_i, \beta, q) = \{ w \in \beta^{-1}\text{suf}(L(B_i)) \mid w = \gamma v, |\gamma| \geq n_1, \delta(q, \gamma) \in F \}$. Now it is clear that

$$\varrho^{(i+1)}(L) = \varrho^{(i)}(L) \cup \bigcup_{\beta \in \Sigma^{n_2}, q \in Q} L_1(B_i, \beta, q) \cdot L_2(B_i, \beta, q). \quad (5)$$

Since context-free languages are effectively closed under prefix, suffix, union, and quotient and intersection with a regular language, using (5) we can construct a pushdown automaton B_{i+1} for $\varrho^{(i+1)}(L)$.

By Lemma 3.2, it is sufficient to construct the pushdown automaton $B_{|w|-n_2-1}$ and decide whether or not $B_{|w|-n_2-1}$ accepts w . The latter can be done effectively since membership is decidable for context-free languages.

Also, context-free languages are effectively closed under subword. Thus, we can test whether $w \in \text{subw}(\varrho^{(|w|-n_2-1)}(L))$ and, by Lemma 3.2 (ii), this holds if and only if $w \in \text{subw}(\varrho^{(*)}(L))$. ■

The operation (5) uses union indexed over all words of length n_2 and consequently the algorithm given by Theorem 3.1 for the uniform membership problem requires exponential time. However, if ϱ is fixed, i.e., if we consider the non-uniform membership problem then the algorithm given by Theorem 3.1 uses polynomial time. The same is true even if only the value of n_2 is fixed. Note that the number of iterations of (5) is upper bounded by the length of w , i.e., the number of iterations is given in unary notation.

Corollary 3.1. *Let n_2 be fixed. Given a TGR system $\varrho = (T, \Sigma, n_1, n_2)$ with T regular, a context-free language L and a word $w \in \Sigma^*$, it is decidable in polynomial time whether or not $w \in \varrho^{(*)}(L)$.*

Lemma 3.2 does not make any assumptions on the initial language. The proof of Theorem 3.1 uses certain closure and decidability properties of context-free languages. A full AFL satisfies the required conditions, assuming that membership is decidable and closure under the AFL operations is effective, and a corresponding extended result is stated below in Corollary 3.2. Before that we introduce some terminology dealing with AFL's consisting of recursive languages. The terminology will be useful also in the next section in order to be able to rely in a uniform way on results from [4] that are formulated in terms of AFL's.

Definition 3.1. *We say that a property P of Turing machines is syntactic if given a Turing machine M it is decidable whether or not M has property P . The class of Turing machines satisfying a property P is denoted $\text{TM}[P]$.*

A language family \mathcal{L} is said to be a constructive full AFL if \mathcal{L} contains a nonempty language and there exists a syntactic property of Turing machines $P_{\mathcal{L}}$ such that

- (i) *a language L is in \mathcal{L} if and only if L is recognized by some Turing machine in $\text{TM}[P_{\mathcal{L}}]$,*
- (ii) *given $M \in \text{TM}[P_{\mathcal{L}}]$ and an input word w , it is decidable whether or not $w \in L(M)$, and*
- (iii) *languages recognized by machines in $\text{TM}[P_{\mathcal{L}}]$ are effectively closed under the AFL operations. That is, there is an algorithm that for given $M_1, M_2 \in \text{TM}[P_{\mathcal{L}}]$ constructs $M_{\text{union}} \in \text{TM}[P_{\mathcal{L}}]$ such that $L(M_{\text{union}}) = L(M_1) \cup L(M_2)$, and for any AFL operation σ other than union there is an algorithm to construct $M \in \text{TM}[P_{\mathcal{L}}]$ such that $L(M) = \sigma(L(M_1))$.*

Well known examples of constructive full AFL's are the regular and the context-free languages. An example of a more general constructive full AFL is the family of languages recognized by (one-way, single head) k -iterated pushdown automata, $k \geq 1$, [6]. It is easy to verify that any (k -iterated) pushdown automaton can be simulated by a Turing machine where the transition relation satisfies a suitably defined syntactic property that forces the work tape to simulate a (k -iterated) pushdown store. It seems that any full AFL consisting only of recursive languages that is defined by a "reasonable" machine model could be characterized in the above way. The family of recursively enumerable languages is a full AFL that is not a constructive full AFL.

Corollary 3.2. *Let \mathcal{L} be a constructive full AFL. Given a TGR system $\varrho = (T, \Sigma, n_1, n_2)$ where T is regular and $L \in \mathcal{L}$, the membership problem for $\varrho^{(*)}(L)$ is decidable.*

The set of useful templates of a TGR system $\varrho = (T, \Sigma, n_1, n_2)$ with an initial language L is the set $T \cap \text{subw}(\varrho^{(*)}(L)) \cap \Sigma^{\geq 2n_1+n_2}$ [4]. Thus by Theorem 3.1:

Corollary 3.3. *Given a TGR system $\varrho = (T, \Sigma, n_1, n_2)$ where T is regular and a context-free initial language L , we can effectively decide whether or not a given template is useful on (L, ϱ) .*

Corollary 3.4. *Let \mathcal{L} be a constructive full AFL. Given a TGR system $\varrho = (T, \Sigma, n_1, n_2)$ where T is regular and $L \in \mathcal{L}$, we can effectively decide whether or not a given template is useful on (L, ϱ) .*

To conclude this section we make a couple of remarks on limitations in attempting to extend the previous results. The 2-iterated pushdown automata recognize the indexed languages [1] and, thus, from Corollary 3.2 we get a decidability result for the membership problem when the initial language is an indexed language. However, there is no known polynomial time parsing algorithm for general indexed languages and Corollary 3.1 cannot be extended for the case where the initial language is indexed.

4 Effective Closure Properties

We would now like to attack the question of, given $\varrho = (T, \Sigma, n_1, n_2)$, with T regular, and L recognized by a pushdown automaton (respectively, a finite automaton), can we effectively construct a pushdown automaton (respectively, a finite automaton) which recognizes $\varrho^{(*)}(L)$? Note that in the former case it is known that $\varrho^{(*)}(L)$ is context-free (and in the latter case regular) [4] but the results are non-constructive.

We first need to provide some details from [4]. The main non-constructive proof from this paper shows that, for an arbitrary TGR system $\varrho = (T, \Sigma, n_1, n_2)$ with T regular, and an arbitrary full AFL \mathcal{L} the following holds: If $L \in \mathcal{L}$, then $\varrho^{(*)}(L) \in \mathcal{L}$. The proof of this result relies on two auxiliary results, the first one of which is the following:

Proposition 4.1. (Theorem 4.2 of [4]) *Let $\varrho = (T, \Sigma, n_1, n_2)$ be a TGR system and let $L \subseteq \Sigma^*$. Let T_u be the useful subset of ϱ on L . If T is a regular language, then T_u is also regular.*

The proof of the above result [4] is not constructive, even in the case where we have some effective representation for L . However, the proof does give some information as to the structure of the DFA which accepts T_u . If Q is the state set of a DFA which accepts T , then the proof creates a finite set of automata $\mathcal{X}_{T,L}$, each automaton with a state set of size $q_{T,L} = (|Q| + 1)^n \cdot (|\Sigma| + 1)^{n-1}$ where $n = 2n_1 + n_2 - 1$. Moreover, the proof establishes that one of these automata accepts T_u , but does not tell us which one is the correct automaton.

Indeed, let $\varrho = (T, \Sigma, n_1, n_2)$ be a TGR system where T is regular and let \mathcal{L} be a constructive full AFL, and let $L \in \mathcal{L}$. Then, by Corollary 3.4, we can decide whether or not a given template is useful on (L, ϱ) . Consider $T_u \cap \Sigma^{\leq 2 \cdot q_{T,L}}$, the finite set of all words which are useful on (L, ϱ) and which are of length less than or equal to $2 \cdot q_{T,L}$. Using Corollary 3.4 we can now effectively determine this set. In addition, for each automaton $M = (Q, \Sigma, q_0, F, \delta) \in \mathcal{X}_{T,L}$, we can check whether or not $T_u \cap \Sigma^{\leq 2 \cdot q_{T,L}} = L(M) \cap \Sigma^{\leq 2 \cdot q_{T,L}}$.

Claim. $T_u \cap \Sigma^{\leq 2 \cdot q_{T,L}} = L(M) \cap \Sigma^{\leq 2 \cdot q_{T,L}}$ if and only if $T_u = L(M)$.

Proof of the claim. It is sufficient to show the implication from left to right. According to Proposition 6.3 of [5], the following is true: Let M_1, M_2 be two DFAs with state sets Q_1, Q_2 respectively. Then $L(M_1) = L(M_2)$ whenever for all $s \in \Sigma^*$ such that $|s| < |Q_1| + |Q_2|$ we have $s \in L(M_1)$ if and only if $s \in L(M_2)$.

Assume by contradiction that $T_u \neq L(M)$. But there exists $M' \in \mathcal{X}_{T,L}$ (also with a state set of size $q_{T,L}$) such that $L(M') = T_u$, and hence

$$L(M') \cap \Sigma^{\leq 2q_{T,L}} = T_u \cap \Sigma^{\leq 2q_{T,L}} = L(M) \cap \Sigma^{\leq 2q_{T,L}}.$$

However, according to the proposition from [5], this implies $L(M') = L(M)$, a contradiction. This concludes the proof of the claim.

By the above claim, we can find from $\mathcal{X}_{T,L}$ the correct automaton which accepts T_u . Hence, we can effectively construct a deterministic finite automaton which accepts T_u . Thus we have shown that the following holds:

Lemma 4.1. *Let \mathcal{L} be a constructive full AFL. Given $\varrho = (T, \Sigma, n_1, n_2)$ with T regular and $L \in \mathcal{L}$, we can construct a DFA for the useful subset of ϱ on L .*

Corollary 4.1. *Let \mathcal{L} be a constructive full AFL. Given $\varrho = (T, \Sigma, n_1, n_2)$ with T regular and $L \in \mathcal{L}$, we can effectively find a regular set of templates T_1 such that if $\varrho_1 = (T_1, \Sigma, n_1, n_2)$ then $\varrho_1^{(*)}(L) = \varrho^{(*)}(L)$ and ϱ_1 is useful on L .*

The second result from [4] that turns out to be useful is the following:

Proposition 4.2. (Theorem 4.1 of [4]) *If \mathcal{L} is a full AFL, $\varrho = (T, \Sigma, n_1, n_2)$ is a TGR system and $L, T \in \mathcal{L}$, $L \subseteq \Sigma^*$, are such that ϱ is useful on L , then $\varrho^{(*)}(L) \in \mathcal{L}$.*

The proof of Proposition 4.2 in [4] establishes that $\varrho^{(*)}(L)$ is in \mathcal{L} by showing that $\varrho^{(*)}(L)$ is obtained from L using a finite number of operations that can be expressed as compositions of AFL operations. This gives the following:

Corollary 4.2. *Let \mathcal{L} be a constructive full AFL. Given a TGR system $\varrho = (T, \Sigma, n_1, n_2)$ where $T \in \mathcal{L}$, an initial language $L \in \mathcal{L}$, $L \subseteq \Sigma^*$, such that ϱ is useful on L , we can effectively construct (a Turing machine in $\text{TM}[P_{\mathcal{L}}]$ for) $\varrho^{(*)}(L) \in \mathcal{L}$.*

Now we are ready to prove the main result of this section.

Theorem 4.1. *Let \mathcal{L} be a constructive full AFL. Given $L \in \mathcal{L}$ and a TGR system $\varrho = (T, \Sigma, n_1, n_2)$ where T is regular, we can effectively construct (a Turing machine in $\text{TM}[P_{\mathcal{L}}]$ for) the language $\varrho^{(*)}(L)$ (which is always in \mathcal{L}).*

Proof. By Corollary 4.1 we can effectively find a regular set of templates $T_1 (\subseteq T)$ such that if $\varrho_1 = (T_1, \Sigma, n_1, n_2)$ then ϱ_1 is useful on L and $\varrho_1^{(*)}(L) = \varrho^{(*)}(L)$.

Since any full AFL contains all regular languages, we have $T_1 \in \mathcal{L}$. Now, by Corollary 4.2, given L and T_1 we can effectively construct a Turing machine for $\varrho_1^{(*)}(L)$ and we are done. ■

Since the regular and the context-free languages are examples of constructive full AFL's, as particular cases Theorem 4.1 implies that if ϱ is a TGR system with a regular set of templates, given a finite automaton (respectively, a pushdown automaton) for a language L , we can effectively construct a finite automaton (respectively, a pushdown automaton) for the language $\varrho^{(*)}(L)$.

Finally, it can be noted that Theorem 4.1 relies on Corollary 4.1 and Corollary 3.4 (that in turn relies on Corollary 3.2), and these results use brute-force constructions that basically enumerate all words up to a given length. It would be interesting to know whether for a regular initial language L and a regular set of templates there is some reasonably efficient algorithm to construct a (not necessarily deterministic) finite automaton for $\varrho^{(*)}(L)$.

References

1. Aho, A.V.: Indexed grammars – an extension of context-free grammars. *Journal of the ACM* **15** (1968) 647–671
2. Daley, M., Ibarra, O., Kari, L., McQuillan, I., Nakano, K.: Closure and decision properties of some language classes under ld and dlad bio-operations. *J. Automata, Languages and Combinatorics* **8** (2003) 477–498
3. Daley, M., McQuillan, I.: Template-guided DNA recombination. *Theoret. Comput. Sci.* **330** (2005) 237–250
4. Daley, M., McQuillan, I.: Useful templates and iterated template-guided DNA recombination in ciliates. *Theory of Computing Systems*, in press. E-print available doi:10.1007/s00224-005-1206-6
5. Eilenberg, S.: *Automata, Languages, and Machines*, volume A. Academic Press, Inc., New York, NY, (1974)
6. Engelfriet, J.: Iterated stack automata and complexity classes. *Information and Computation* **95** (1991) 21–75
7. Ginsburg, S.: *Algebraic and Automata-Theoretic Properties of Formal Languages*. North-Holland, Amsterdam (1975)
8. Head, T., Pixton, D.: Splicing and regularity, to appear. Available at www.math.binghamton.edu/dennis/Papers
9. McQuillan, I., Salomaa, K., Daley, M.: Iterated TGR languages: Membership problem and effective closure properties. Queen's School of Computing Technical Report No. 2006-513. Available at www.cs.queensu.ca/TechReports
10. Prescott, D.M.: Genome Gymnastics: Unique modes of DNA evolution and processing in ciliates. *Nature Reviews Genetics* **1** (2000) 191–198
11. Prescott, D.M., Ehrenfeucht, A., Rozenberg, G.: Template guided recombination for IES elimination and unscrambling of genes in stichotrichous ciliates. *J. Theoretical Biology* **222** (2003) 323-330
12. Rozenberg, G., Salomaa, A. (eds.): *Handbook of Formal Languages*, Vols. 1–3. Springer Verlag, (1997)