

Three models for gene assembly in ciliates: a comparison

Miika Langille, Ion Petre, Vladimir Rogojin

Abstract

We survey in this paper the main differences among three variants of an intramolecular model for gene assembly: the general, the simple, and the elementary models. We present all of them in terms of sorting signed permutations and compare their behavior with respect to: (i) completeness, (ii) confluence (with the notion defined in three different setups), (iii) decidability, (iv) characterization of the sortable permutations in each model, (v) sequential complexity, and (vi) experimental validation.

Keywords: Ciliate, simple gene assembly, simple model, elementary model, confluence, completeness, characterization, sequential complexity, model validation, signed permutations, sorting.

1 Introduction

Gene assembly in ciliates has been subject of intense research in the last few years, both regarding the molecular details driving it, as well as the theoretical implications of some mathematical models proposed for it, see [7, 10, 25, 17, 16, 26, 1, 20].

Ciliates form an ancient and rich group of eukaryotes. There are about 8000 species of ciliates currently known. Two characteristics which are common for all ciliates distinguish them from other groups of unicellular eukaryotes. First, they all have “*cilia*”, organs used for motility and for feeding. Second, they all have two types of nuclei presented in each organism. Almost all RNA-transcriptions happen in

macronuclei (somatic nuclei) during the life of a ciliate. The DNA-molecules in the *micronuclei* (germline nuclei) seem to remain silent until the sexual reproduction begins (see [24]).

The genetical information is stored in different ways on micro- and macronuclear molecules. The macronuclear genes are contiguous sequences of nucleotides. The micronuclear genes however, are split into coding blocks (called MDSs), shuffled and separated by noncoding blocks (called IESs). This shuffling and inversion of MDSs is especially visible in a species of ciliates called *stichotrichs*. Macronuclear molecules are known to be the shortest DNA in Nature, ranging in the *Sterkiella nova* organisms between 200bp and 3700bp with an average of 2200 bp in length (see [11, 5, 6, 23, 24, 27]). Macronuclear molecules consist mainly of coding sequences. On the other hand, coding sequences occupy as little as 2 – 5 % of the micronuclear molecules of the length about 10^7 bp (in *Sterkiella nova*, see [5, 23]).

At some point during sexual reproduction, ciliates destroy all macronuclei and develop new ones from the micronuclei. In the process they must excise non-coding sequences and assemble correctly all coding blocks of the micronuclear genes. This process is called gene assembly. For a brief introduction to the biology of ciliates, especially to the gene assembly process we refer to [7].

Two molecular models have been proposed for gene assembly in ciliates. The intermolecular model [17, 16] and the intramolecular model [10, 25] suggest splicing of gene fragments via short nucleotide sequences called pointers. Each pointer at the end of an MDS repeats at the beginning of the MDS which follows it in the assembled gene. Recent results [1, 20] suggest that some template molecules may assist the correct alignment of the recombining molecules. The intermolecular model suggests that two molecules may participate in the recombination, while the intramolecular model considers folding and recombination within a single molecule.

We focus in this paper on the intramolecular model (called in the sequel the general model) and on two of its variants: the simple model, introduced in [15] and the elementary model, introduced in [14].

The general model consists of three molecular operations, ld, hi,

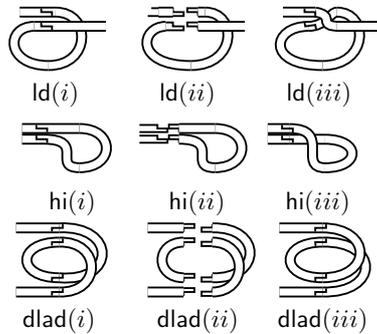


Figure 1. Illustration of the *ld*, *hi*, *dlad* molecular operation showing in each case: (i) the folding, (ii) the recombination, (iii) the result. Courtesy of Tero Harju.

dlad, see [10, 25]. The three operations are illustrated in Figure 1 where in each case we show the folding of the molecule on itself, the recombination that takes place and the subsequent result. A characteristic of this model is that all three operations operate on a single molecule that folds on itself in a specific way. One thus says that the model is *intramolecular*.

Note that the three intramolecular operations allow in their general formulation that the MDSs participating in an operation may be located anywhere along the molecule. Arguing on the principle of parsimony, a simplified model was introduced in [15], asking that all operations are applied ‘locally’. This simple model consists of the same three molecular operations as the general model, requiring however that there is at most one coding block involved in each of the three operations. This idea was then further developed into two separate models, both using the terminology of *simple gene assembly*. In the first one, that we will refer to in here as the *elementary model*, introduced in [13, 14], the model was further restricted so that only *micronuclear*, but not *composite*, MDSs could be manipulated by the molecular operations. Consequently, once two or more micronuclear

MDSs are combined into a larger composite MDS, they can no longer be moved along the sequence. The second model, that we will refer to as the *simple model* [18], allowed that both micronuclear, as well as composite MDSs may be manipulated in each of the three molecular operations.

However minor the difference between the frameworks of the simple and the elementary models may seem, it does have a great impact on the characteristics of each model. We survey in this paper the main known results on the simple and elementary gene assembly, comparing them also with the corresponding properties of the general model with respect to: (i) completeness, (ii) confluence (with the notion defined in three different setups), (iii) decidability, (iv) characterization of the sortable permutations in each model, (v) sequential complexity, and (vi) experimental validation. For this, we introduce in this paper a permutation-based presentation of the general model. We discuss in particular the question of model validation and consider the assembly of all currently known ciliate gene patterns, see [4]. We also present several open problems in this area.

The results in this paper have been previously published in [8, 14, 18, 19] using non-uniform (and even conflicting) terminology and notation. In here we give the topic a uniform presentation, fix the terminology and discuss in some details differences among the three models of interest.

2 Mathematical preliminaries

For a finite alphabet $A = \{a_1, \dots, a_n\}$, we denote by A^* the free monoid generated by A and call any element of A^* a *word*. For any $v \in A^*$, we denote $\text{dom}(v) = \{a \in A \mid a \text{ occurs in } v\}$.

Let $\bar{A} = \{\bar{a}_1, \dots, \bar{a}_n\}$, where $A \cap \bar{A} = \emptyset$. For $p, q \in A \cup \bar{A}$, we say that p, q have the same *signature* if either $p, q \in A$, or $p, q \in \bar{A}$ and we say that they have *different signatures* otherwise. For any $u \in (A \cup \bar{A})^*$, $u = x_1 \dots x_k$, with $x_i \in A \cup \bar{A}$, for all $1 \leq i \leq k$, we denote $\|u\| = \|x_1\| \dots \|x_k\|$, where $\|a\| = \|\bar{a}\| = a$, for all $a \in A$. We also denote $\bar{u} = \bar{x}_k \dots \bar{x}_1$, where $\bar{\bar{a}} = a$, for all $a \in A$. We say, that u

is *uniformly signed*, if either $x_i \in A$ for all $1 \leq i \leq k$, or $x_i \in \bar{A}$ for all $1 \leq i \leq k$.

For strings u, v over Σ , we say that u is a *substring* of v , denoted by $u \leq v$, if $v = xuy$, for some strings x, y . We say that u is a *subsequence* of v , denoted by $u \leq_s v$, if $u = a_1 a_2 \dots a_m$, $a_i \in \Sigma \cup \bar{\Sigma}$ and $v = v_0 a_1 v_1 a_2 v_2 \dots a_m v_m$, for some strings $v_i, 0 \leq i \leq m$, over Σ .

A *permutation* π over A is a bijection $\pi : A \rightarrow A$. Fixing the order relation (a_1, a_2, \dots, a_m) over A , we often denote π as the word $\pi(a_1) \dots \pi(a_m) \in A^*$. A *signed permutation* over A is a string $\psi \in (A \cup \bar{A})^*$, where $\|\psi\|$ is a permutation over A . We say that a signed permutation π is (*circularly*) *sorted* if it is of either of the following forms:

- (i) $\pi = a_k a_{k+1} \dots a_n a_1 \dots a_{k-1}$, for some $k \geq 1$. In this case, we say that π is an *orthodox sorted permutation*.
- (ii) $\pi = \bar{a}_{k-1} \dots \bar{a}_1 \bar{a}_n \dots \bar{a}_{k+1} \bar{a}_k$, for some $k \geq 1$. In this case, we say that π is an *inverted sorted permutation*.

In both cases, if $k = 1$, then we say that π is a *linear sorted permutation*; otherwise, we say that it is *circular*.

A *sorted block* in the signed permutation π is a substring of π either of the form $a_i a_{i+1} \dots a_j$, or of the form $\bar{a}_j \dots \bar{a}_{i+1} \bar{a}_i$, $1 \leq i \leq j \leq n$, where $a_{i-1} a_i, \bar{a}_i \bar{a}_{i-1}, a_j a_{j+1}, \bar{a}_{j+1} \bar{a}_j$ are not substrings of π . By $S(\pi)$ we denote the total number of sorted blocks in π . Clearly, the permutation is cyclically sorted if we have $S(\pi) \leq 2$.

The notion of structure of a permutation will be useful in the paper. To define it, we first introduce the morphism $\xi_i : (A \cup \bar{A})^* \rightarrow (A \cup \bar{A})^*$, for any $1 \leq i \leq |A|$:

$$\xi_i(a_j) = \begin{cases} \lambda & \text{if } j = i; \\ a_j & \text{if } j < i; \\ a_{j-1} & \text{if } j > i; \end{cases}$$

where $a_j \in A \cup \bar{A}$.

Consider the mapping $\sigma_i : (A \cup \bar{A})^* \rightarrow (A \cup \bar{A})^*$, where for any string $u \in (A \cup \bar{A})^*$, $\sigma_i(u)$ is defined as follows:

- (a) $\sigma_i(u) = u$, if $a_i a_{i+1} \not\prec u$, with $a_i, a_{i+1} \in A$, or $a_{i+1} a_i \not\prec u$, with $a_i, a_{i+1} \in \bar{A}$, and
- (b) $\sigma_i(u) = \xi_i(u)$ otherwise.

Then, the structure of a string is the mapping $\sigma : (A \cup \bar{A})^* \rightarrow (A \cup \bar{A})^*$, such that $\sigma(u) = (\sigma_1 \circ \sigma_2 \circ \dots \circ \sigma_{|A|-1} \circ \sigma_{|A|})(u)$. Note that the structure of a sorted permutation π is either $\sigma(\pi) = a_1$, or $\sigma(\pi) = a_2 a_1$, where $a_1, a_2 \in A$, or $\sigma(\pi) = a_1 a_2$, where $a_1, a_2 \in \bar{A}$.

Example 1. Consider a sorted permutation $\pi = 34512$. We find its structure $\sigma(\pi)$ as follows:

$$\begin{aligned} \pi_5 &= \sigma_5(\pi) = \pi & \pi_2 &= \sigma_2(\pi_3) = \pi_3 \\ \pi_4 &= \sigma_4(\pi_5) = \xi_4(\pi_5) = 3412 & \pi_1 &= \sigma_1(\pi_2) = \\ \pi_3 &= \sigma_3(\pi_4) = \xi_3(\pi_4) = 312 & &= \xi_1(\pi_2) = 21 \\ & & & \sigma(\pi) = \pi_1 = 21 \end{aligned}$$

3 Gene assembly as a sorting of signed permutations

As discussed in [18, 13, 14], a natural formalization of the simple and elementary operations is through rewriting rules for signed permutations. A given gene is represented as a signed permutation by denoting the sequence and the orientation of its MDSs and assembling the gene is modeled through the sorting of the associated permutation. As shown in Definitions 1, 2, 3, the formalization of the molecular models in terms of sorting permutations is somewhat intricate: a high number of cases needs to be considered. For the general model, a more concise formalization can be done in terms of *signed double occurrence* (also called *legal*) strings, see [7]. The two main advantages of the legal string framework are: (i) it abstracts from denoting the sequence of gene blocks to denoting only the sequence of pointers (and in the process it ignores the two markers); (ii) it models gene assembly as a process of consecutive pointer removals, based on the observation that the assembled gene contains no pointers.

The simple model makes crucial use of the two markers. Consequently, this model can only be formalized through *extended* legal strings that denote the pointers, as well as the markers of the gene, as done in [3]. The resulting model is equivalent with the permutation-based model for simple operations but more concise for the same reasons (i)-(ii) discussed above.

In the case of the elementary model however, it is crucial that all pointers and markers are indicated throughout the gene assembly, rather than being removed as in the (extended) legal string framework. The main reason is that the elementary model distinguishes between the original (micronuclear) gene blocks and the larger (composite) blocks that are being formed throughout the process of assembly. It is an open problem whether a more concise formalizations may be introduced also for the elementary model.

In the following we consider a presentation based on signed permutations for all three models. This presentation in the case of the general model appears to be given here for the first time, although an equivalent presentation in terms of *MDS descriptors* was reported before, see [7]. As observed also in the case of simple and elementary operations, it is a characteristic of permutation-based models for gene assembly that the *ld* operation is not explicitly modeled. Instead, it is just assumed that two consecutive blocks are going to be spliced together in a bigger composite block at some arbitrary point, independently of the other operations applied to the permutation.

3.1 Modeling of the general operations

Consider a gene pattern formalized as a signed permutation over alphabet $\Pi_n = \{1, 2, \dots, n\}$. We formalize the general operations over signed permutations as follows:

Definition 1. *i. For each $1 \leq p < n$, hi_p is defined as follows:*

$$\begin{aligned}\text{hi}_p(xpy(\overline{p+1})z) &= xp(p+1)\bar{y}z, \\ \text{hi}_p(x\bar{p}y(p+1)z) &= x\bar{y}p(p+1)z, \\ \text{hi}_p(x(p+1)y\bar{p}z) &= x\bar{y}(\overline{p+1})\bar{p}z, \\ \text{hi}_p(x(\overline{p+1})ypz) &= x(\overline{p+1})\bar{p}\bar{y}z,\end{aligned}$$

where x, y, z are signed strings over Π_n . We denote $\text{Hi} = \{\text{hi}_i \mid 1 \leq i < n\}$.

ii. For each $1 \leq p, q < n$, where $|p - q| > 1$, $\text{dlad}_{p,q}$ is defined as follows:

$$\begin{aligned}\text{dlad}_{p,q}(xp''uq''vp'wq'z) &= xwq'q''vp'p''uz, \\ \text{dlad}_{p,q}(xp''uq'vp'wq''z) &= xwvp'p''uq'q''z, \\ \text{dlad}_{p,q}(xp'uq''vp''wq'z) &= xp'p''wq'q''vuz, \\ \text{dlad}_{p,q}(xp'uq'vp''wq''z) &= xp'p''wv uq'q''z,\end{aligned}$$

where $p' = p, p'' = p+1$, or $p' = (\overline{p+1}), p'' = \bar{p}$, and $q' = q, q'' = q+1$, or $q' = (\overline{q+1}), q'' = \bar{q}$, and x, u, v, w, z are signed strings over Π_n . In all these cases, we also denote $\text{dlad}_{q,p} = \text{dlad}_{p,q}$.

For each $1 < p < n$, we define $\text{dlad}_{p-1,p}$ and $\text{dlad}_{p,p-1}$ as follows:

$$\begin{aligned}\text{dlad}_{p-1,p}(xp'''up''wp'z) &= xwp'p''p'''uz, \\ \text{dlad}_{p-1,p}(xp''vp'wp'''z) &= xwvp'p''p'''z, \\ \text{dlad}_{p-1,p}(xp'up'''vp''z) &= xp'p''p'''vuz,\end{aligned}$$

where $p' = p-1, p'' = p, p''' = p+1$, or $p''' = (\overline{p+1}), p'' = \bar{p}, p' = (\overline{p-1})$, x, u, v, w, z are signed strings over Π_n . We denote $\text{Dlad} = \{\text{dlad}_{i,j} \mid 1 \leq i, j < n, i \neq j\}$.

Example 2. Consider the permutation $\pi_1 = 2\bar{5}1\bar{4}376$. We sort it by hi and dlad as follows:

$$\begin{aligned}\text{hi}_5(2\bar{5}1\bar{4}376) &= 2\bar{7}\bar{3}4\bar{1}56 & \text{hi}_4(\bar{4}\bar{7}\bar{3}\bar{2}\bar{1}56) &= 1237456 \\ \text{hi}_2(2\bar{7}\bar{3}4\bar{1}56) &= 2374\bar{1}56 & \text{dlad}_{3,6}(1237456) &= 1234567 \\ \text{hi}_1(2374\bar{1}56) &= \bar{4}\bar{7}\bar{3}\bar{2}\bar{1}56\end{aligned}$$

3.2 Modeling of the simple operations

Simple operations are a restriction of the general operations [9, 7]: they rearrange pieces of DNA containing at most one MDS, be that micronuclear, or composite.

Definition 2. *The molecular model of simple hi and simple dlad can be formalized as follows.*

i. For each $1 \leq p < n$, sh_p is defined as follows:

$$\begin{aligned} \text{sh}_p(xp \dots (p+i)(\overline{p+k}) \dots (\overline{p+i+1})y) &= \\ &= xp \dots (p+i)(p+i+1) \dots (p+k)y, \\ \text{sh}_p(x(\overline{p+i}) \dots \overline{p}(p+i+1) \dots (p+k)y) &= \\ &= xp \dots (p+i)(p+i+1) \dots (p+k)y, \\ \text{sh}_p(x(p+i+1) \dots (p+k)(\overline{p+i}) \dots \overline{p}y) &= \\ &= x(\overline{p+k}) \dots (\overline{p+i+1})(\overline{p+i}) \dots \overline{p}y, \\ \text{sh}_p(x(\overline{p+k}) \dots (\overline{p+i+1})p \dots (p+i)y) &= \\ &= x(\overline{p+k}) \dots (\overline{p+i+1})(\overline{p+i}) \dots \overline{p}y, \end{aligned}$$

where $k > i \geq 0$ and x, y are signed strings over Π_n . We denote $\text{Sh} = \{\text{sh}_i \mid 1 \leq i \leq n\}$.

ii. For each p , $2 \leq p \leq n-1$, sd_p is defined as follows:

$$\begin{aligned} \text{sd}_p(xp \dots (p+i)y(p-1)(p+i+1)z) &= \\ &= xy(p-1)p \dots (p+i)(p+i+1)z, \\ \text{sd}_p(x(p-1)(p+i+1)yp \dots (p+i)z) &= \\ &= x(p-1)p \dots (p+i)(p+i+1)yz, \\ \text{sd}_{\overline{p}}(x(\overline{p+i+1})(\overline{p-1})y(\overline{p+i}) \dots \overline{p}z) &= \\ &= x(\overline{p+i+1})(\overline{p+i}) \dots \overline{p}(\overline{p-1})yz, \\ \text{sd}_{\overline{p}}(x(\overline{p+i}) \dots \overline{p}y(\overline{p+i+1})(\overline{p-1})z) &= \\ &= xy(\overline{p+i+1})(\overline{p+i}) \dots \overline{p}(\overline{p-1})z, \end{aligned}$$

where $i \geq 0$ and x, y, z are signed strings over Π_n . We denote $\text{Sd} = \{\text{sd}_i, \text{sd}_{\overline{i}} \mid 1 \leq i \leq n\}$.

Example 3. Consider the following signed permutation $\pi_1 = 5\bar{4}\bar{7}\bar{6}\bar{3}\bar{1}\bar{2}$. It can be sorted by the following composition of simple operations

$$\begin{aligned} \text{sh}_6(\pi) &= 5\bar{4}\bar{7}\bar{6}\bar{3}\bar{1}\bar{2}, & \text{sh}_4 \circ \text{sd}_{\bar{2}} \circ \text{sh}_6(\pi) &= \bar{5}\bar{4}\bar{7}\bar{6}\bar{3}\bar{2}\bar{1}, \\ \text{sd}_{\bar{2}} \circ \text{sh}_6(\pi) &= 5\bar{4}\bar{7}\bar{6}\bar{3}\bar{2}\bar{1}, & \text{sd}_{\bar{4}} \circ \text{sh}_4 \circ \text{sd}_{\bar{2}} \circ \text{sh}_6(\pi) &= \\ & & &= \bar{7}\bar{6}\bar{5}\bar{4}\bar{3}\bar{2}\bar{1}. \end{aligned}$$

3.3 Modeling of the elementary operations

The elementary model is a restriction of the simple model: elementary intramolecular operations rearrange only micronuclear MDSs. This leads to the following formalization for elementary operations.

Definition 3. *i.* For each $p \geq 1$, eh_p is defined as follows:

$$\begin{aligned} \text{eh}_p(x\overline{p(p+1)}z) &= xp(p+1)z, \\ \text{eh}_p(x\bar{p}(p+1)z) &= xp(p+1)z, \\ \text{eh}_p(x(p+1)\bar{p}z) &= x(\overline{p+1})\bar{p}z, \\ \text{eh}_p(x(\overline{p+1})pz) &= x(\overline{p+1})\bar{p}z, \end{aligned}$$

where x, z are signed strings over Π_n . We denote $\text{Eh} = \{\text{eh}_p \mid 1 \leq p \leq n\}$.

ii. For each p , $2 \leq p \leq n-1$, ed_p is defined as follows:

$$\begin{aligned} \text{ed}_p(xpy(p-1)(p+1)z) &= xy(p-1)p(p+1)z, \\ \text{ed}_p(x(p-1)(p+1)ypz) &= x(p-1)p(p+1)yz, \\ \text{ed}_p(x\bar{p}y(\overline{p+1})(\overline{p-1})z) &= xy(\overline{p+1})\bar{p}(\overline{p-1})z, \\ \text{ed}_p(x(\overline{p+1})(\overline{p-1})y\bar{p}z) &= x(\overline{p+1})\bar{p}(\overline{p-1})yz, \end{aligned}$$

where x, y, z are signed strings over Π_n . We denote $\text{Ed} = \{\text{ed}_p \mid 1 < p < n\}$.

Note that $\text{Eh} \subset \text{Sh} \subset \text{Hi}$ and $\text{Ed} \subset \text{Sd} \subset \text{Dlad}$.

Example 4. Assume the signed permutation $\pi = 315\bar{2}46$. It can be sorted by a composition of elementary operations as follows

$$\begin{aligned} \text{ed}_5(\pi) &= 31\bar{2}456, & \text{ed}_3 \circ \text{eh}_1 \circ \text{ed}_5(\pi) &= 123456, \\ \text{eh}_1 \circ \text{ed}_5(\pi) &= 312456, \end{aligned}$$

3.4 Sorting strategies: terminology

A composition of operations $\Phi = \phi_k \circ \phi_{k-1} \circ \dots \circ \phi_2 \circ \phi_1$, where all operations are from either $\text{Hi} \cup \text{Dlad}$, or $\text{Sh} \cup \text{Sd}$, or $\text{Eh} \cup \text{Ed}$ is called a *strategy*. A composition $\Phi = \phi_k \circ \phi_{k-1} \circ \dots \circ \phi_2 \circ \phi_1$ of operations is called a *sorting strategy* for π , if $\Phi(\pi)$ is a (circularly) sorted permutation. If $\phi \in (\text{Hi} \cup \text{Dlad})$ for all $1 \leq i \leq k$, we say that Φ is a *general sorting strategy*. If $\phi \in (\text{Sh} \cup \text{Sd})$ for all $1 \leq i \leq k$, we say that Φ is a *simple sorting strategy*. If $\phi \in (\text{Eh} \cup \text{Ed})$ for all $1 \leq i \leq k$, we say that Φ is an *elementary sorting strategy*. We say that an unsigned permutation π is *blocked* if no (simple, elementary) operation is applicable to it. We say that Φ is an *unsuccessful strategy* for π , if $\Phi(\pi)$ is blocked. If there are no sorting strategies for π , then we say that π is an *unsortable permutation*.

4 Comparison of the three models

In this section we compare the general, simple and elementary intramolecular models for gene assembly by different criteria:

- completeness: whether any gene pattern may be assembled or not;
- confluence, defined in three different ways:
 - (i) whether there are permutations having both successful and unsuccessful strategies,
 - (ii) whether different assembly strategies starting from the same gene pattern lead to assembled genes with the same structure,
 - (iii) if different assembly strategies starting from the same gene pattern lead to the same assembled gene;
- decidability of assembly: whether it is possible to decide effectively if a given gene pattern can be assembled or not;

- characterization of gene patterns that can be assembled (starting from certain characteristics of a given gene pattern we can conclude whether the gene pattern can be assembled);
- sequential complexity is constant: whether all assembly strategies apply the same number of intramolecular operations;
- model validation: whether it is consistent with biological data.

4.1 Completeness

It was shown in [9, 7] that the general model is complete, i.e., it assembles any gene pattern. The result was proved in terms of MDS-descriptors. To prove it for signed permutations, one may take two different approaches.

On one hand, one may observe that the set of signed permutations and that of MDS descriptors are in an one-to-one correspondence. Moreover, for a signed permutation π , if $\psi(\pi)$ is its corresponding MDS descriptor, then for any operation $f \in \text{Hi} \cup \text{Dlad}$, $\psi(f(\pi)) = f(\psi(\pi))$. The completeness result for signed permutations then follows easily from the corresponding result for MDS descriptors.

On the other hand, one may give a direct proof of the completeness, by essentially mimicking the proof in the case of MDS descriptors. The essential observation in this case is that for any $\phi \in \text{Hi} \cup \text{Dlad}$ and any signed permutation π , the number of sorted blocks of $\phi(\pi)$ is smaller than that of π (i.e., $S(\phi(\pi)) < S(\pi)$). One needs to observe then that a signed permutation π is sorted if and only if $S(\pi) \leq 2$ and π is uniformly signed.

Theorem 1. *All signed permutations are sortable over $\text{Hi} \cup \text{Dlad}$.*

Note however that the simple and the elementary models are not complete, as shown by the following example.

Example 5. *Consider the permutation $\pi = 321$. We cannot apply either eh or sh operations as all pointers have the same signature, and there is no applicable ed or sd operation either. On the other hand, π is successful in the general model: $\text{dlad}_{1,2}(\pi) = 123$.*

4.2 Confluence

We consider the notion of *confluence* in three different setups, so as to reflect the success of different assembly strategies, the resulting gene structure, or the resulting gene pattern. These aspects are discussed below stressing the differences between the three models for gene assembly.

Consider first the most common notion of confluence, requiring that the result of all assemblies of a given input is the same. Equivalently, all strategies for a given signed permutation are confluent. It is easy to see that neither of the three models for gene assembly is confluent in this sense. For this, consider the permutation $\pi = 2413$. Then $\text{dlad}_{2,1}(\pi) = \text{sd}_2(\pi) = \text{ed}_2(\pi) = 4123$, while $\text{dlad}_{2,3}(\pi) = \text{sd}_3(\pi) = \text{ed}_3(\pi) = 2341$. Note that this observation does not contradict earlier invariant results of [8, 21], see also [7], where it was proved that the result of all assembly strategies of a given gene/string is always the same. The difference comes from considering a so-called *boundary* ld operation that would be applied as a last step in both strategies above to yield a circular string 1234 (that may also be denoted as 2341, 3412, or 4123, or even their inverses). In the permutation-based presentation, we have chosen to consider only standard linear permutations, rather than circular ones. The non-confluence result above is a direct consequence of this choice. We discuss more aspects of this matter in Section 5.

The example above shows that all three models are nondeterministic in the sense that different sorting strategies may lead to different results. A natural question is then whether a given signed permutation may have both successful, as well as unsuccessful strategies in any of the three models. Consider then the following notion of confluence. We say that the general (simple, elementary, resp.) model is confluent if there are no signed permutations having both successful and unsuccessful strategies.

It follows from Theorem 1 that the general model is indeed confluent in the sense above. As shown in [18], the simple model is also confluent. However, the elementary model is not confluent. To see it, consider the permutation $\pi = 24135$. Then $\text{ed}_3(\pi) = 23415$ is a blocked

permutation, while $\text{ed}_2 \circ \text{ed}_4(\pi) = 12345$, a sorted permutation.

It was proved in [8, 21], see also [2], that for any gene pattern, either all general assembly strategies assemble it to a linear molecule, or all of them assemble it to a circular one. Consequently, even though if the assembly process is non-deterministic, the results of all possible assemblies of a given gene pattern have the same structure. I.e., the results of all *sorting* strategies applicable to a permutation have the same structure. As such, the same result holds also for all *sorting* strategies in the simple and in the elementary models. The question may however be asked also for the unsuccessful strategies. In this context, we say that a model for gene assembly is confluent if, for any signed permutation, all its sorting strategies lead to permutations having the same structure. Based on the considerations above, it follows easily that the general model is confluent in this sense, while the elementary model is not (since a permutation may have both successful and unsuccessful elementary strategies). Interestingly, it was proved in [18] that the simple model is in fact confluent in this sense.

Example 6. *Consider permutation $\pi = 623514$. There are only two simple strategies applicable to π : $\pi_1 = \text{sd}_2(\pi) = 651234$ and $\pi_2 = \text{sd}_4(\pi) = 623451$. These strategies are unsuccessful, and there are no other simple strategies applicable to π . Permutation π cannot be sorted by simple operations. Note however, that permutations π_1 and π_2 have the same structure $\sigma(\pi_1) = 321 = \sigma(\pi_2)$.*

The following table captures the behavior of the three models for gene assembly with respect to the three notions of confluence above. Interestingly, none of these notions distinguishes the simple and the general model. One property that does distinguish between the two is the completeness, valid only for the general model.

4.3 Deciding the sortability problem

For the simple and elementary models, which are not complete, deciding the sortability of a given signed permutation is an interesting problem. Based on the confluence results in the previous section, it

	Success	Same result	Same structure
General	confluent	not confluent	confluent
Simple	confluent	not confluent	confluent
Elementary	not confluent	not confluent	not confluent

Table 1. The results of considering confluence with regard to the three aspects are summarized here.

turns out that the problem is easy for the simple model: for any signed permutation, either all its sorting strategies are successful, or they are all unsuccessful. As such, to decide the sortability problem, it is enough to find an arbitrary strategy (e.g., using a straightforward procedure having quadratic time complexity) and answer ‘yes’/‘no’, depending on whether or not that strategy is successful.

For the elementary model the problem of the eh-sortability of a signed permutation is easy.

Theorem 2 ([14]). *The signed permutation π is eh-sortable if and only if either*

$$(i) \|\pi\| = k(k+1)\dots n12\dots(k-1) \text{ and for some } 1 \leq i \leq k-1, \\ k \leq j \leq n \text{ we have } i, j \text{ unsigned, or}$$

$$(ii) \|\pi\| = (k-1)\dots 21n\dots(k+1)k, \text{ and for some } 1 \leq i \leq k-1, \\ k \leq j \leq n \text{ we have } i, j \text{ signed.}$$

The problem of the ed-sortability turns out to be technically more involved, since a signed permutation may have both successful, and unsuccessful strategies. A complete characterization of the ed-sortable signed permutation has been given in [13, 14, 22]. The main notions used in the result are those of dependency graphs and forbidden elements. We only present here these notions for unsigned permutation; in the case of signed permutation, the setup is technically more complex, see [14]. Note also that an efficient decision procedure for the sortability problem is only known for unsigned permutation, see [22].

Dependency graphs in the elementary model

Dependency graphs suggest in which order elementary operations should be used to assemble a given gene pattern. Let π be an unsigned permutation with $\text{dom}(\pi) = \{1, 2, \dots, n\}$. We associate to it a *dependency graph* $\Gamma_\pi = (V_\pi, E_\pi)$, where $V_\pi = \text{dom}(\pi)$, and

$$E_\pi = \{(1, 1), (n, n)\} \cup \{(i, i) \mid (i+1)(i-1) \leq_s \pi\} \cup \{(j, i) \mid (i-1)j(i+1) \leq_s \pi\}.$$

Intuitively, an edge (j, i) in Γ_π shows that in any sorting strategy for π , the operation ed_j should be used first, in order for ed_i to become applicable. If there is a loop (i, i) in Γ_π , then ed_i cannot be applied in any strategy applicable to π . We refer to [14] for a proof of these observations.

Example 7. Consider the unsigned permutation $\pi = 628\ 41071359$. Its associated dependency graph $\Gamma_\pi = (V_\pi, E_\pi)$ is shown in Figure 2.

We have loops $(1, 1)$, $(5, 5)$, $(6, 6)$, $(10, 10)$ in the dependency graph, and so, the operations ed_1 , ed_5 , ed_6 and ed_{10} cannot be applied in any strategy applicable to G . We have cycle 838 in Γ_π and so, neither operation ed_3 , nor operation ed_8 can be applied in any strategy applicable to π . The dependency graph Γ_π suggests the following order of operations to be applied in any sorting strategy of π : ed_2 should be applied before ed_7 , and ed_4 should be applied before ed_9 . Indeed, for instance, strategy $\text{ed}_9 \circ \text{ed}_4 \circ \text{ed}_7 \circ \text{ed}_2(\pi)$ sorts π : $\text{ed}_9 \circ \text{ed}_4 \circ \text{ed}_7 \circ \text{ed}_2(\pi) = 67891012345$.

Forbidden elements and ed– sortability of unsigned permutations

For a signed permutation π , we say that $p \in \text{dom}(\pi)$ is *forbidden* in π if and only if there exists no composition of eh and ed operations applicable to π with p in the domain of one of them. We denote U_π the set of all forbidden elements of π . It was proved in [14] that $p \in U(\pi)$ if and only if

- (i) p is on a cycle of Γ_π or

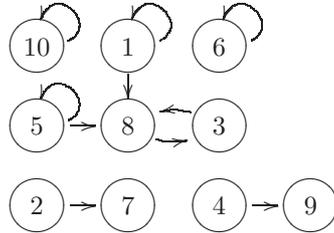


Figure 2. The dependency graph associated to $\pi = 62841071359$.

- (ii) there is a path from q to p in Γ_π , for some q on a cycle of Γ_π or
- (iii) there exists $r > 1$ such that there are paths from $r - 1$ to p and from r to p in Γ_π .

The following result gives the ed-sortability of unsigned permutations.

Theorem 3 ([14]). *The unsigned permutation π is ed-sortable if and only if $\pi|_{U_\pi}$ is sorted.*

Finding an efficient method for the eh, ed-sortability of a signed permutation remains an open problem.

4.4 Characterization of sortable permutations

The following theorem characterizes ed-sortable unsigned permutations. A similar, albeit technically more involved, characterization exists also for signed permutations, see [14].

Theorem 4 ([14]). *Let π be a unsigned permutation. Then π is Ed-sortable if and only if there exists a partition $\{1, 2, \dots, n\} = D \cup U$, such that the following conditions are satisfied:*

- (i) $\pi|_U$ is sorted;
- (ii) The subgraph induced by D in G_π is acyclic;

- (iii) If $(p, q) \in G_\pi$ with $q \in D$, then $p \in D$;
- (iv) For any $p \in D$, $(p-1)(p+1) \leq_s \pi$;
- (v) For any $p \in D$, $(p-1), (p+1) \in U$.

For simple operations we do not have a characterization of sortable permutations for the moment. For general operations the question is moot since all signed permutations are sortable.

4.5 Sequential complexity

We focus now on the length of various sorting strategies of a given signed permutation, where the length is defined as the number of operations in the strategy. Consider first the general model and let $\pi_1 = 15\bar{2}\bar{4}36$. One can sort it by applying $\text{dlad}_{1,5} \circ \text{hi}_2$, or by applying $\text{hi}_2 \circ \text{hi}_3 \circ \text{hi}_1$. These two sorting strategies are of different length, and use a different combination of operations.

Somewhat surprisingly, the situation is different in the simple model and by consequence, also in the elementary model. It was established in [19] (using a string-based formalism) that any two sorting strategies for a given signed permutation have the same assembly length.

Theorem 5 ([19]). *Let π be a signed permutation and ϕ, ψ be two simple sorting strategies for π . Then ϕ and ψ have the same sequential assembly length. Moreover, they have the same number of **sh** and the same number of **sd** operations.*

The differences between the general model and the two restricted models go beyond Theorem 5. E.g., when choosing operations in the simple model, we may always just choose the first available operation as the number of operations required in the end remains the same. If the operations were given different weights or costs, then the general model may have optimal and sub-optimal sorting strategies. We refer to [12] for a detailed discussion on various measures of complexity for gene assembly.

	General	Simple	Elementary
Completeness	complete	not complete	not complete
Confluence (Success)	confluent	confluent	not confluent
Confluence (Structure)	confluent	confluent	not confluent
Confluence (Result)	not confluent	not confluent	not confluent
Deciding Sortability	yes: trivial	yes: confluence	open for eh + ed
Characterizing sortable permutations	trivial	open	yes
Sequential Complexity constant	no	yes	yes
Model Validation	valid	valid	not valid

Table 2. Summary for general, simple and elementary intramolecular models

4.6 Model validation

A database of known sequences of micronuclear and macronuclear ciliate genes can be found in [4]. Based on the completeness result for the general model, it is clear that all the gene patterns have an assembly strategy in the general model. As it turns out however, the elementary model cannot account for the assembly of some of the gene patterns in [4].

Example 8. *Actin I gene in it Sterkiella nova is represented by the permutation $\pi = 346579\bar{2}18$. It is easy to check that there is no elementary sorting strategy applicable to π . However, we can sort π by applying the simple sorting strategy*

$$sh_1 \circ sh_2 \circ sd_8 \circ sd_5(\pi) = \bar{9}\bar{8}\bar{7}\bar{6}\bar{5}\bar{4}\bar{3}\bar{2}\bar{1}.$$

Below we will outline all the available scrambled gene patterns in [4], together with one simple sorting strategy. Genes that are not scrambled in their micronuclear form or the ones that have missing MDSs will not be included.

Actin I, Sterkiella n. : $\pi = 346579\bar{2}18$;

$$sh_1 \circ sh_2 \circ sd_8 \circ sd_6(\pi) = \overline{987654321}.$$

Actin I, Sterkiella h. : $\pi = 34657910\bar{2}18$;

$$\text{sh}_1 \circ \text{sh}_2 \circ \text{sd}_8 \circ \text{sd}_6(\pi) = \overline{10987654321}.$$

Actin I, Stylonychia p. : $\pi = 346578\bar{2}1$;

$$\text{sh}_1 \circ \text{sh}_2 \circ \text{sd}_6(\pi) = \overline{87654321}.$$

α **Telomere Binding Protein, Sterkiella n.** :

$$\pi = 1357911246810121314;$$

$$\text{sd}_{10} \circ \text{sd}_8 \circ \text{sd}_6 \circ \text{sd}_4 \circ \text{sd}_2(\pi) =$$

$$= 1234567891011121314.$$

DNA Polymerase α , Paraurostyla weissei:

$$\begin{aligned} \pi = & \overline{44\ 42\ 40\ 38\ 36\ 34\ 32\ 30\ 28\ 26\ 24\ 22\ 20\ 18\ 16\ 14\ 12\ 10\ 8} \\ & \overline{6\ 1\ 2\ 3\ 4\ 5\ 7\ 9\ 11\ 13\ 15\ 17\ 19\ 21\ 23\ 25\ 27\ 29\ 31\ 33\ 35\ 37\ 39\ 41} \\ & 43\ 45\ 46\ 47\ 48 \end{aligned}$$

The signed permutation sorting strategy for this gene is just sh_1 repeated 40 times.

4.7 Summary

Table 2 summarizes properties of general, simple and elementary models considered in this paper.

5 Discussion and open problems

There has been significant interest in the last few years in the so-called simple operations for gene assembly, both for their biological appeal as a minimal, parsimonious model, but also for the properties of their mathematical formalization. The term simple has been used in connection with two different versions of the model. In this survey we review

these two models and fix the proper terminology. We also compare the mathematical properties of these two models with those of the general model.

For reasons detailed already in Section 3 we chose in this paper to follow a permutation-based presentation, rather than a string-based one. Indeed, a string-based presentation that would be more concise than the permutation-based one is still missing for the elementary model. Our choice of using permutations rather than strings has one direct consequence that we mentioned already in Section 4.2. Rather than eliminating all pointers as in the legal string and ending up with linear, or circular strings, we always end up with sorted linear permutations, where the term ‘sorted’ is extended to cover also permutations such as 3412. We call such a permutation circularly sorted, see Section 2. For this reason, a permutation such as 2413 may be sorted to two seemingly different results: either 2341, or 4123. Clearly, the two results correspond to the same circular string in the framework of legal strings. This ambiguity leads nevertheless to some open problems of independent interest. E.g., given a permutation that may be sorted circularly, enumerate efficiently all the circularly sorted permutations it can be sorted to. Similarly, the permutation $2\bar{1}3$ may be sorted to either 231 or $\bar{2}1\bar{3}$. One may also ask about the properties of those permutations that have sortings both to an unsigned permutation, as well as to a signed one. The properties of the three models may even be different in this respect.

There are two currently open problems related to the simple model: the linear decidability of the sortability problem and computing the number of sortable permutations of length n . It is however possible that these two problems are intertwined and an answer to one may at least partly solve the other.

Decidability. It was shown in [18] that it is possible to decide whether a permutation is sortable or unsortable in the simple model by applying available operations in an arbitrary order until the permutation is blocked or sorted. This gives us a quadratic method for deciding. Our first open problem is related to the optimality of this

method: is there a procedure to decide in linear time the sortability problem in the simple model?

For the elementary model, finding an efficient decision procedure for $\{eh, ed\}$ -sortability problem is also open.

Sortable permutations of length n . As we pointed out also in this paper, not all permutations may be sorted using the simple operations. This differs from the general model which has been shown to be complete. Thus, an interesting problem is computing how many permutations of length n are sortable in the simple/elementary models. As a related problem, it should even be interesting to see whether the ratio of sortable signed permutations tends to 0 when n tends to infinity. Both problems are open also in the case of unsigned permutations.

6 Acknowledgments.

Ion Petre and Vladimir Rogojin are supported by Academy of Finland, project 108421. Vladimir Rogojin is on leave of absence from Institute of Mathematics and Computer Science of Academy of Sciences of Moldova, Chisinau MD-2028 Moldova. Vladimir Rogojin is supported by Science and Technology Center in Ukraine, project 4032.

References

- [1] A. Angeleska, N. Jonoska, M. Saito, and L.F.Landweber. Rna-template guided dna assembly. *Journal of Theoretical Biology*, 248:706–720, 2007.
- [2] R. Brijder, H. Hoogeboom, and G. Rozenberg. Reducibility of gene patterns in ciliates using the breakpoint graph. *Theoretical Computer Science*, 356:26–45, 2006.
- [3] R. Brijder, M. Langille, and I. Petre. A string-based model for simple gene assembly. In E. Csuhaj-Varju and Z. Esik, editors,

- FCT 2007, Proceedings*, volume 4639 of *Lecture Notes in Computer Science*, pages 161–172. Springer-Verlag Berlin Heidelberg, 2007.
- [4] A. Cavalcanti, T. Clarke, and L. Landweber. Mds.ies_db: a database of macronuclear and micronuclear genes in spirotrichous ciliates. *Nucleic Acids Research*, 33:396–398, 2005.
- [5] W. Chang, P. Bryson, H. Liang, M. Shin, and L. Landweber. The evolutionary origin of a complex scrambled gene. In *Proceedings of the National Academy of Sciences of the US*, volume 102, pages 15149–15154, 2005.
- [6] W. Chang, S. Kuo, and L. Landweber. A new scrambled gene in the ciliate *uroleptus*. *Gene*, 368:72–77, 2006.
- [7] A. Ehrenfeucht, T. Harju, I. Petre, D. M. Prescott, and G. Rozenberg. *Computation in Living Cells: Gene Assembly in Ciliates*. Springer, 2003.
- [8] A. Ehrenfeucht, I. Petre, D. M. Prescott, and G. Rozenberg. Circularity and other invariants of gene assembly in ciliates. In M. Ito, G. Paun, and S. Yu, editors, *Words, semigroups, and transductions*, pages 81–97. World Scientific, Singapore, 2001.
- [9] A. Ehrenfeucht, I. Petre, D. M. Prescott, and G. Rozenberg. *Universal and simple operations for gene assembly in ciliates*, pages 329–342. Kluwer Academic, Dordrecht, 2001.
- [10] A. Ehrenfeucht, D. M. Prescott, and G. Rozenberg. Computational aspects of gene (un)scrambling in ciliates. In L. F. Landweber and E. Winfree, editors, *Evolution as Computation*, pages 216–256. Springer, Berlin, Heidelberg, New York, 2001.
- [11] W. Foissner and H. Berger. Identification and ontogenesis of the *nomen nudum* ypotrichs (protozoa: Ciliophora) *oxytricha nova* (= *sterkiella nova* sp. n.) and *o. trifallax* (= *s. histriomuscorum*). *Acta Protozool.*, 38:215–248, 1999.

- [12] T. Harju, C. Li, I. Petre, and G. Rozenberg. Complexity measures for gene assembly. In K. Tuyls, editor, *Proceedings of the Knowledge Discovery and Emergent Complexity in Bioinformatics workshop*, volume 4366 of *Lecture Notes in Bioinformatics*, pages 42–60. Springer, 2007.
- [13] T. Harju, I. Petre, V. Rogojin, and G. Rozenberg. Simple operations for gene assembly. In A. Carbone and N. A. Pierci, editors, *Proceedings of DNA-based computers 11*, volume 3892 of *Lecture Notes in Computer Science*, pages 96–111. Springer, 2006.
- [14] T. Harju, I. Petre, V. Rogojin, and G. Rozenberg. Patterns of simple gene assembly in ciliates. *Discrete Applied Mathematics*, 2007. to appear.
- [15] T. Harju, I. Petre, and G. Rozenberg. Modelling simple operations for gene assembly. In J.Chen, N.Jonoska, and G.Rozenberg, editors, *Nanotechnology: Science and Computation*, pages 361–376, 2006.
- [16] L. F. Landweber and L. Kari. *Evolution as computation*, chapter Universal molecular computation in ciliates, pages 257–274. Natural computing series. Springer Verlag, Berlin, Heidelberg, New York, 1999.
- [17] L. F. Landweber and L. Kari. The evolution of cellular computing: Nature’s solution to a computational problem. In L. Kari, H. Rubin, and D. H. Wood, editors, *Proceedings of the 4th DIMACS Meeting on DNA-Based Computers*, volume 52, (1–3), pages 3–13. Elsevier, 1999.
- [18] M. Langille and I. Petre. Simple gene assembly is deterministic. *Fundamenta Informaticae*, 73(1-2):179–190, 2006.
- [19] M. Langille and I. Petre. Sequential vs. parallel complexity in simple gene assembly. *Theoretical Computer Science*, 395(1):24–30, 2008.

- [20] M. Nowacki, V. Vijayan, Y. Zhou, K. Schotanus, T. Doak, and L. Landweber. Rna-mediated epigenetic programming of a genome-rearrangement pathway. *Nature*, 451:153–158, Jan. 2008. doi:10.1038/nature06452.
- [21] I. Petre. Invariants of gene assembly in stichotrichous ciliates. *Information Technology*, 48(3):161–167, 2006.
- [22] I. Petre and V. Rogojin. Decision problem for shuffled genes. *Information and Computation*, 2007. to appear.
- [23] D. M. Prescott. The dna of ciliated protozoa. *Microbiology and Molecular Biology Reviews*, 58(2):233–267, 1994.
- [24] D. M. Prescott. DNA manipulations in ciliates. In W. Brauer, H. Ehrig, J. Karhumäki, and A. Salomaa, editors, *Formal and Natural Computing*, volume 2300 of *Lecture Notes in Computer Science*, pages 394–417. Springer, 2002.
- [25] D. M. Prescott, A. Ehrenfeucht, and G. Rozenberg. Molecular operations for dna processing in hypotrichous ciliates. *European Journal of Protistology*, 37:241–260, 2001.
- [26] D. M. Prescott, A. Ehrenfeucht, and G. Rozenberg. Template-guided recombination for ies elimination and unscrambling of genes in stichotrichous ciliates. *Journal of Theoretical Biology*, 222:323–330, 2003.
- [27] M. Swanton, J. Heumann, and D. Prescott. Gene-sized dna molecules of the macronuclei in three species of hypotrichs: size distribution and absence of nicks. *Chromosoma*, 77:217–227, 1980.

Miika Langille, Ion Petre, Vladimir Rogojin,

Received March 22, 2010

Miika Langille,
Department of IT, Åbo Akademi University
ICT-building, Joukahaisenkatu 3-5 A, 5th floor
Turku 20520 Finland
E-mail: miika.langille@abo.fi

Ion Petre
Academy of Finland and
Turku Centre for Computer Science
Department of Computer Science, Åbo Akademi University
Turku 20520 Finland E-mail: ion.petre@abo.fi

Vladimir Rogojin
Turku Centre for Computer Science
Department of Computer Science, Åbo Akademi University
Turku 20520 Finland
E-mail: vrogojin@abo.fi