

# A String-based Model for Simple Gene Assembly

Robert Brijder<sup>1</sup>, Miika Langille<sup>2</sup>, and Ion Petre<sup>2,3</sup>

<sup>1</sup> Leiden Institute of Advanced Computer Science, Universiteit Leiden  
Niels Bohrweg 1, 2333 CA Leiden, The Netherlands

`rbrijder@liacs.nl`

<sup>2</sup> Turku Centre for Computer Science  
Department of IT, Åbo Akademi University

Turku 20520 Finland

`miika.langille@abo.fi`

`ion.petre@abo.fi`

<sup>3</sup> Academy of Finland

**Abstract.** The simple intramolecular model for gene assembly in ciliates is particularly interesting because it can predict the correct assembly of all available experimental data, although it is not universal. The simple model also has a confluence property that is not shared by the general model. A previous formalization of the simple model through sorting of signed permutations is unsatisfactory because it effectively ignores one operation of the model and thus, it cannot be used to answer questions about parallelism in the model, or about measures of complexity. We propose in this paper a string-based model in which a gene is represented through its sequence of pointers and markers and its assembly is represented as a string rewriting process. We prove that this string-based model is equivalent to the permutation-based model as far as gene assembly is concerned, while it tracks all operations of the model.

## 1 Introduction

Gene assembly in ciliates has been subject to extensive combinatorial research in recent years, see [2]. Ciliates are unicellular eukaryotes that organize their genome differently in their two types of nuclei. In micronuclei, genes are split into blocks (called *MDSs*), placed in a shuffled order on the chromosome, separated by non-coding blocks. Moreover, some of the *MDSs* are even presented in an inverted form. In macronuclei however, genes are contiguous sequences of nucleotides, with all blocks sorted in the orthodox order. The assembly of genes from their micronuclear to their macronuclear form has a definite combinatorial and computational flavor: each *MDS*  $M$  ends with a sequence of nucleotides (called a *pointer*) that is repeated identically in the beginning of the *MDS* that should follow  $M$  in the macronuclear gene.

The exact kinetical mechanisms of gene assembly still remain to be clarified through further laboratory experiments. Two models have been proposed for gene assembly: an intermolecular one, see [7, 8] and an intramolecular one, see [3, 10]. The intramolecular model, that we consider in this paper, consists of

three operations: `ld`, `hi`, and `dlad`. For a detailed discussion of these operations, including their formalization on various levels of abstraction, we refer to [2]. We consider in this paper the simple variant of the model, in which the sequences manipulated by each operation are minimal, see [5] for a detailed discussion. It turns out that although not universal, the simple model is capable of correctly predicting the assembly of all currently available data on genes in stichotrichous ciliates, see [6]. Also, the model has an interesting confluence property that does not hold in the general model, see [9].

The simple model for gene assembly has been investigated in [9] as a process of sorting signed permutations. A major shortcoming of this formalization is that the `ld` operations are not explicitly modeled and it is only assumed that they take place eventually at some arbitrary moment of time. Although sufficient for the purpose of [9], this formalization does not allow reasoning about parallelism or complexity of assemblies.

Abstracting the gene to its sequence of pointers has been a solution to this problem in the case of the general model. In this way, the gene assembly process is formalized through a process of string rewriting. Doing the same in the case of simple operations does not work, as shown in Example 9 of this paper: the string-based model is not equivalent with the permutation-based model.

We propose in this paper a simple solution to this problem. Rather than representing a gene only through its sequence of pointers, we preserve also the beginning and ending markers (starting the first MDS and ending the last MDS, respectively). We prove that with this simple addition, the string-based model is equivalent with the permutation-based model.

## 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*. 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.

We denote  $A^{\boxtimes} = (A \cup \bar{A})^*$ . For any  $u \in A^{\boxtimes}$ ,  $u = x_1 \dots x_k$ , 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  $a \in A$ . For two alphabets  $A, B$ , a mapping  $f : A^{\boxtimes} \rightarrow B^{\boxtimes}$  is called a *morphism* if  $f(uv) = f(u)f(v)$  and  $f(\bar{u}) = \overline{f(u)}$ .

A *permutation*  $\pi$  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  $\pi$  as the word  $\pi(a_1) \dots \pi(a_m) \in A^*$ . A *signed permutation* over  $A$  is a string  $\psi \in A^{\boxtimes}$ , where  $\|\psi\|$  is a permutation over  $A$ .

A string  $v \in \Sigma^*$  over an (unsigned) alphabet  $\Sigma$  is a *double occurrence string* if every pointer  $a \in \text{dom}(v)$  occurs exactly twice in  $v$ . A signing of a non-empty double occurrence string is a *legal string*.

Two legal strings,  $u$  and  $v$  over alphabets  $A$  and  $B$ , respectively, are said to have the same *structure* if there is a morphism  $f : A^{\mathfrak{X}} \rightarrow B^{\mathfrak{X}}$  such that  $f(u) = v$ , and is denoted by  $u \equiv v$ .

### 3 Permutations and Strings

Simple operations have previously been defined in terms of signed permutations, see [9]. This formalism was useful for observing simple operations applied sequentially, with the arrangement of explicit MDSs being clearly visible. However, note that with signed permutations the  $\text{ld}$  operation is always assumed to occur at some point in the assembly. This is not sufficient for studying some aspects of simple gene assembly, and so we now task ourselves with providing a formal framework for performing simple gene assembly on *legal strings*. Micronuclear and intermediate genes are represented as a sequence of pointers and the operations defined as a *string pointer reduction system*. Moreover, we will show that simple gene assembly on legal strings is equivalent to that on signed permutations.

Legal strings have been used before to formalise the general operations in [2]. This definition is missing one crucial piece of information, necessary to model simple operations in particular. Namely, the absence of markers indicating the beginning and end of the macronuclear gene make it necessary to extend the definition.

Let us now briefly describe the signed permutations, followed then by the extended definition of legal strings and the string pointer reduction system.

#### Signed Permutations

Genes may be represented as a sequence of MDSs. Using the alphabet  $\Pi_n = \{1, 2, \dots, n\}$  to denote MDSs, where the numbering is given by the order in which MDSs are assembled in the macronuclear gene. Thus, a micronuclear gene will be a signed permutation over  $\Pi_n$  and a macronuclear gene will be a sorted signed permutation over  $\Pi_n$ .

**Definition 1** *We say that a signed permutation  $\pi$  is sorted if either:*

- i.*  $\pi = i(i+1) \dots n1 \dots (i-1)$ , or
- ii.*  $\pi = (\overline{i-1}) \dots \overline{1n} \dots (i+1)i$ ,

*for some  $1 \leq i \leq n$ . If  $i = 1$  we say that  $\pi$  is a linear permutation. We call  $\pi$  circular otherwise. In case *i.* we say that  $\pi$  is sorted in the orthodox order, while in case *ii.* we say that  $\pi$  is sorted in the inverted order.*

The term circular in the above definition refers to a gene that gets assembled, say in the form  $i \dots n1 \dots (i-1)$ , and then gets excised from the chromosome by an  $\text{ld}$  operation applied on the pointer in the beginning of  $i$  and its identical copy at the end of  $(i-1)$ .

Given this definition for genes, we may now consider gene assembly as a sorting of signed permutations. We will define the simple operations as transformation rules for signed permutations in such a way that a simple operation is applicable on a gene pattern if and only if the corresponding rule is applicable on the associated signed permutation.

As mentioned previously, when using signed permutations to model gene assembly, the  $\text{ld}$  operation is ignored. Indeed,  $\text{ld}$  only combines two MDSs  $i$  and  $(i+1)$  already placed next to each other into a bigger composite MDS. To avoid renaming the alphabet after each  $\text{ld}$  operation, we will consider that when  $i$  and  $(i+1)$  are placed next to each other, the operation joining them is already accomplished.

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

i. For each  $p \geq 1$ ,  $\text{sh}_p$  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  $x, y, z$  are signed strings over  $\Pi_n$ . We denote  $\text{Sh} = \{\text{sh}_i \mid 1 \leq i \leq n\}$ .

ii. For each  $p$ ,  $2 \leq p \leq n-1$ ,  $\text{sd}_p$  is defined as follows:

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

where  $i \geq 0$  and  $x, y, z$  are signed strings over  $\Pi_n$ . We also define  $\text{sd}_{\overline{p}}$  as follows:

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

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

We say that a signed permutation  $\pi$  over the set of integers  $\{i, i+1, \dots, i+l\}$  is sortable if there are operations  $\phi_1, \dots, \phi_k \in \text{Sh} \cup \text{Sd}$  such that  $(\phi_k \circ \dots \circ \phi_1)(\pi)$  is a sorted permutation. We say that  $\pi$  is blocked if neither an  $\text{sh}$  operation, nor an  $\text{sd}$  one is applicable to  $\pi$  and  $\pi$  is not sorted.

Let  $\phi = \phi_k \circ \dots \circ \phi_1$ ,  $\phi_i \in \text{Sh} \cup \text{Sd}$ , for all  $1 \leq i \leq k$ . We say that  $\phi$  is a strategy for  $\pi$  if  $\phi(\pi)$  is either sorted, or blocked. In the former case we say that  $\phi$  is a sorting strategy, while in the latter case we say that  $\phi$  is an unsuccessful strategy for  $\pi$ .

If  $\phi = \phi_k \circ \dots \circ \phi_1$  is a sorting strategy for  $\pi$ , we say that  $\pi$  is  $\text{Sh}$ -sortable if  $\phi_1, \dots, \phi_k \in \text{Sh}$  and we say that  $\pi$  is  $\text{Sd}$ -sortable if  $\phi_1, \dots, \phi_k \in \text{Sd}$ .

- Example 1** *i.* The permutation  $\pi_1 = \bar{2}1435$  is sortable. It has the following sorting strategies:  $\text{sh}_1 \circ \text{sd}_4 \circ \text{sh}_1(\pi_1) = 12345$  and  $\text{sh}_1 \circ \text{sh}_1 \circ \text{sd}_4(\pi_1) = 12345$ .
- ii.* The permutation  $\pi_2 = 2\bar{5}\bar{3}14$  is blocked as no operations are applicable and it is not sorted.
- iii.* The permutation  $\pi_3 = 35124$  has two assembly strategies which lead to different results:  $\text{sd}_3(\pi_3) = 51234$  and  $\text{sd}_4(\pi_3) = 34512$ .

## Legal Strings

Using signed permutations, we hold the information of explicit MDSs. It is sufficient, however, to consider the gene as a sequence of pointers. Moving to legal strings, we remove all notation pertaining to explicit MDSs, and represent the gene as a string of pointers, with each pointer occurring twice in the gene. Though some information is lost about the gene, the notation is elegant and the characteristics and rules of the model are maintained.

To represent a micronuclear or intermediate gene as a legal string, we need a set of pointers,  $\Delta_k = \{2, 3, \dots, k\}$ , and the set of markers  $M = \{b, e\}$ . Note that each marker will occur only once in the string, and each pointer will occur twice. As with signed permutations,  $k$  represents the length of the gene, and we shall assume that it will be clear from the context. Let this be our alphabet  $\Sigma = \Delta_k \cup M$ .

The definition for legal strings given in [2] did not include the markers, but as it turns out, they are essential to being able to formalise simple operations on legal strings. Additionally, not including the markers meant that it was no longer possible to uniquely obtain a signed permutation from a legal string, as too much information was lost.

Let  $a \in \Sigma \cup \bar{\Sigma}$  and let  $u \in \Sigma^{\times}$  be a legal string. If  $u$  contains both substrings  $a$  and  $\bar{a}$  then  $a$  is *positive* in  $u$ ; otherwise,  $a$  is *negative* in  $u$ .

- Example 2** *i.* Consider the signed string  $u_1 = b23423\bar{e}\bar{4}$  over  $\Sigma_4$ . Clearly,  $u_1$  is legal. Pointer 4 is positive in  $u_1$ , while 2 and 3 are negative in  $u_1$ .
- ii.* The string  $u_2 = \bar{4}\bar{3}b234$  is not legal, since 2 has only one occurrence in  $u_2$  and the end marker does not occur at all.

Let  $u = a_1a_2 \dots a_n \in \Sigma^{\times}$  be a legal string over  $\Sigma$ , where  $a_i \in \Sigma \cup \bar{\Sigma}$  for each  $i$ . For each letter  $a \in \text{dom}(u)$ , there are indices  $i$  and  $j$  with  $1 \leq i < j \leq n$  such that  $\|a_i\| = a = \|a_j\|$ , excepting of course the two markers present. The substring  $u_{(a)} = a_i a_{i+1} \dots a_j$  is the  $a$ -interval of  $u$ . Two different letters  $a, b \in \Sigma$  are said to *overlap* in  $u$  if the  $a$ -interval and the  $b$ -interval of  $u$  overlap: if  $u_{(a)} = a_{i_1} \dots a_{j_1}$  and  $u_{(b)} = a_{i_2} \dots a_{j_2}$ , then either  $i_1 < i_2 < j_1 < j_2$  or  $i_2 < i_1 < j_2 < j_1$ . Moreover, for each letter  $a$ , we denote by

$$O_u(a) = \{b \in \Sigma \mid b \text{ overlaps with } a \text{ in } u\} \cup \{a\}.$$

- Example 3** Let  $u = b23423\bar{5}\bar{4}\bar{e}\bar{6}56$  be a legal string. The 5-interval of  $u$  is the substring  $u_{(5)} = \bar{5}\bar{4}\bar{6}5$ , which contains only one occurrence of the pointers 4 and

6, but either two or no occurrences of 2 and 3 and so  $O_u(5) = \{4, 5, 6\}$ . Note that the marker is not included in the interval. Similarly,

$$\begin{array}{ll} u_{(2)} = 2342, & \text{and 2 overlaps with 3 and 4,} \\ u_{(3)} = 3423, & \text{and 3 overlaps with 2 and 4,} \\ u_{(4)} = 423\overline{54}, & \text{and 4 overlaps with 2, 3 and 5,} \\ u_{(6)} = \overline{6}56, & \text{and 6 overlaps with 5.} \end{array}$$

We will now consider the transformation from signed permutations to legal strings, carried out by replacing each MDS with the pointers at either end of it. When transforming from signed permutations to legal strings, there is a problem with how to deal with sorted blocks of pointers. Therefore, when doing the transformation, the signed permutation must be assumed to be micronuclear, i.e., no MDSs have been joined together. This allows for the following morphism from a signed permutation  $\pi$  to a legal string  $u$ .  $\zeta : \Pi_n^{\mathfrak{X}} \rightarrow \Sigma_n^{\mathfrak{X}}$  is applied to each letter in  $\pi$ , and is defined as follows

- i.  $i \rightarrow i(i+1)$ ,
- ii.  $1 \rightarrow b2$ ,
- iii.  $n \rightarrow ne$ ,

where  $1 \leq i \leq n$ .

- Example 4** i. The signed permutation  $\pi = 1\overline{2}$ , is equivalent to the following legal string  $\zeta(\pi) = b2\overline{e}2$ .
- ii. The signed permutation  $\pi = 2314\overline{7}6\overline{5}$ , is equivalent to the following legal string  $\zeta(\pi) = 2334b245\overline{e}7\overline{6}7\overline{6}5$ .

We say that a legal string  $u$  is *realistic* if there exists a signed permutation  $\pi$  such that  $u = \zeta(\pi)$ . It is important to note, however, that not all legal strings are realistic. The following example illustrates this point.

**Example 5** The string  $u = be22$  is legal but it is not realistic, since it has no “realistic parsing”, i.e., there is no signed permutation which would transform to this legal string.

Let us first give our full definition of legal strings and the simple operations, and then we will show why the markers are necessary to the abstraction.

### Simple Operations on Signed Strings

Recall that in signed permutations gene assembly was simulated by joining together MDSs to obtain a (circularly) sorted sequence, the macronuclear gene. With the move to legal strings we have now removed all details of individual MDSs, leaving us with a string of pointers. The goal remains the same though, to remove all pointers by matching them together, consequently building the

macronuclear gene. We shall now introduce the *string pointer reduction system* for simple operations to formalise the three molecular operations ld, hi and dlad.

A realistic legal string is considered sorted, once all pointers are removed, leaving us with one of the following possible sorted permutations:  $be$ ,  $\bar{e}\bar{b}$ ,  $eb$  and  $\bar{b}\bar{e}$ . The first two cases represent linear sortings, and the latter two are considered circular.

In the following we shall consider only realistic legal strings, that is, the strings are from  $\Sigma_k^{\mathfrak{X}}$  for some  $k \geq 2$ . Each of the rules below is a function that maps legal strings to legal strings.

- i. The *string negative rule*  $\text{snr}$  for a pointer  $p \in \Delta_k$  ( $k \geq 2$ ) is the equivalent of the ld operation. As with ld, the string negative rule is always simple, and as such, it is the same for the general and simple model. It can only be applied when the pointers are adjacent, i.e., not separated by any other pointers. It can be formalised in the following way

$$\text{snr}_p(u_1ppu_2) = u_1u_2,$$

$$\text{snr}_p(pu_3p) = u_3,$$

where  $u_1, u_2 \in \Sigma^{\mathfrak{X}}$  and  $u_3$  contains only markers (boundary case). Let

$$\text{Snr} = \{\text{snr}_p \mid p \in \Delta_k, k \geq 2\}$$

be the set of all simple string negative rules.

- ii. The *simple string positive rule*  $\text{sspr}$  for a pointer  $p \in \Delta_k$  ( $k \geq 2$ ) is the equivalent of the sh operation. Recall that in the simple hi operation the pointers may only be separated by one MDS, and thus only a single pointer or marker. This gives us the following formalisation

$$\text{sspr}_p(u_1pu_2\bar{p}u_3) = u_1\bar{u}_2u_3,$$

where  $|u_2| = 1$  and  $u_1, u_2, u_3 \in \Sigma^{\mathfrak{X}}$ . Let

$$\text{Sspr} = \{\text{sspr}_p \mid p \in \Delta_k, k \geq 2\}$$

be the set of all simple string positive rules.

- iii. The *simple string double rule*  $\text{ssdr}$  for pointers  $p, q \in \Delta_k$  ( $k \geq 2$ ) is the equivalent of the sd operation. Recall that in the simple dlad operation the first occurrences of  $p$  and  $q$  must be adjacent, with the same condition applied to the second occurrences. This gives us the following formalisation

$$\text{ssdr}_{p,q}(u_1pqu_2pqu_3) = u_1u_2u_3,$$

where  $u_1, u_2, u_3 \in \Sigma^{\mathfrak{X}}$ . Let

$$\text{Ssdr} = \{\text{ssdr}_{p,q} \mid p, q \in \Delta_k, k \geq 2\}$$

be the set of all simple string double rules.

Note that the operations have now become very simple indeed, with only pointers being removed. The  $\text{sspr}$  operation is the only one that affects the remaining permutation, inverting the pointer or marker separating the two occurrences.

**Example 6** Consider the following signed permutation  $\pi = \bar{1}2354$ . It has the following realistic legal string  $u = \zeta(\pi) = \bar{2}b23345e45$ , and these simple operations applicable to it.

- i. The  $\text{snr}$  operation removes two adjacent pointers with the same signature. E.g.,  $\text{snr}_3(u) = \bar{2}b245e45$ .
- ii. The  $\text{sspr}$  operation removes two pointers with different signatures, separated by exactly one pointer or marker, inverting the sequence separating them. E.g.,  $\text{sspr}_2(u) = b3345e45$ .
- iii. The  $\text{ssdr}$  operation removes two overlapping pointers, where the first occurrences are adjacent, as are the second occurrences. E.g.,  $\text{ssdr}_{4,5}(u) = \bar{2}b233e$ .

**Example 7** Let us consider the actin I gene from *Sterkiella nova*, which has the following signed permutation  $\pi = 346579\bar{2}18$ . The corresponding legal string would be

$$u = \zeta(\pi) = 34456756789e\bar{3}\bar{2}b289.$$

It has the following assembly strategy using the simple string pointer reduction system.

$$\begin{aligned} \text{sspr}_2(u) &= 34456756789e\bar{3}\bar{b}89 \\ \text{ssdr}_{5,6} \circ \text{sspr}_2(u) &= 3447789e\bar{3}\bar{b}89 \\ \text{snr}_4 \circ \text{ssdr}_{5,6} \circ \text{sspr}_2(u) &= 37789e\bar{3}\bar{b}89 \\ \text{snr}_7 \circ \text{snr}_4 \circ \text{ssdr}_{5,6} \circ \text{sspr}_2(u) &= 389e\bar{3}\bar{b}89 \\ \text{ssdr}_{8,9} \circ \text{snr}_7 \circ \text{snr}_4 \circ \text{ssdr}_{5,6} \circ \text{sspr}_2(u) &= 3e\bar{3}\bar{b} \\ \text{sspr}_3 \circ \text{ssdr}_{8,9} \circ \text{snr}_7 \circ \text{snr}_4 \circ \text{ssdr}_{5,6} \circ \text{sspr}_2(u) &= \bar{e}\bar{b} \end{aligned}$$

It is important to note that the naming of simple operations on legal strings does not follow the same standard as for signed permutations. On signed permutations an operation was always named by the smallest MDS in the composite involved. However, on legal strings we have lost the information of MDSs, and only remaining pointers are shown. Thus an operation must be named according to the pointer on which the fold is made. This is illustrated in the following example.

**Example 8** Let  $u = b2452334\bar{e}\bar{5}$  be a realistic legal string corresponding to the signed permutation  $\pi = 1423\bar{5}$ . They have the following equivalent assembly strategies,

$$\begin{aligned} \text{snr}_3(u) &= b24524\bar{e}\bar{5}, & (\text{id operation is assumed here}) \\ \text{ssdr}_{2,4} \circ \text{snr}_3(u) &= b\bar{5}\bar{e}\bar{5}, & \text{sd}_2(\pi) = 1234\bar{5}, \\ \text{sspr}_5 \circ \text{ssdr}_{2,4} \circ \text{snr}_3(u) &= be, & \text{sh}_1 \circ \text{sd}_2(\pi) = 12345. \end{aligned}$$



Let us now show that the operations on legal strings are equivalent to those on signed permutations.

**Theorem 1** *For any signed permutations  $\pi$  and  $\pi'$ , with  $\pi' = \phi_m \circ \dots \circ \phi_1(\pi)$ ,  $\phi_i \in \text{Sh} \cup \text{Sd}$ , let  $u = \zeta(\pi)$  and  $u' = \zeta(\pi')$ . Then  $\tau_0(u') = \psi_m \circ \tau_m \circ \dots \circ \psi_1 \circ \tau_1(u)$ , where  $\psi_i$  is the equivalent of  $\phi_i$ , and  $\tau_i$  is a composition of **snr** operations.*

*Proof.* We will only prove the claim for  $m = 1$ . In its full generality, the claim can be proved by iterating the same argument.

**Case 1:**  $\phi_1 \in \text{Sh}$ . Then  $\pi$  is of one of the four forms in Definition 2(i). Assume that  $\pi = xp \dots (p+i)(\overline{p+k}) \dots (\overline{p+i+1})y$ , as the other cases are completely similar. Then  $u = \zeta(x)p(p+1)(p+1) \dots (p+i)(p+i+1)(\overline{p+k+1})(\overline{p+k})(\overline{p+k}) \dots (\overline{p+i+2})(\overline{p+i+2})(\overline{p+i+1})\zeta(y)$  and

$$\begin{aligned} \text{sh}_p(\pi) &= xp \dots (p+k)y, \\ \zeta(\text{sh}_p(u)) &= \zeta(x)p(p+1)(p+1) \dots (p+k)(p+k)(p+k+1) \\ (\text{sspr}_{p+i+1} \circ \text{snr}_{p+1} \circ \text{snr}_{p+i} \text{snr}_{p+i+2} \circ \text{snr}_{p+k})(u) &= \zeta(x)p(p+k+1)\zeta(y) = \\ &= (\text{snr}_{p+1} \circ \dots \circ \text{snr}_{p+k})(\zeta(\text{sh}_p(u))). \end{aligned}$$

**Case 2:**  $\phi_1 \in \text{Sd}$ . Then  $\pi$  is of one of the forms in Definition 2(ii). Assume that  $\pi = xp \dots (p+i)y(p-1)(p+i+1)z$ , as the other cases are completely similar. Then  $u = \zeta(x)p(p+1)(p+1) \dots (p+i)(p+i)(p+i+1)\zeta(y)(p-1)p(p+i+1)(p+i+2)\zeta(z)$  and

$$\begin{aligned} \text{sd}_p(\pi) &= xy(p-1)p \dots (p+i)(p+i+1)z, \\ \zeta(\text{sd}_p(\pi)) &= \zeta(xy)(p-1)pp \dots (p+i+1)(p+i+1)(p+i+2)\zeta(z), \\ (\text{ssdr}_{p,p+i+1} \circ \text{snr}_{p+1} \circ \dots \circ \text{snr}_{p+i})(u) &= \zeta(x)\zeta(y)(p-1)(p+i+2)\zeta(z) = \\ &= (\text{snr}_p \circ \dots \circ \text{snr}_{p+i+1})(\zeta(\text{sd}_p(\pi))). \end{aligned}$$

Thus we have shown that the operations on legal strings are equivalent to those for signed permutations.  $\square$

Now that we have defined legal strings and formalised simple operations on legal strings we will show why it was necessary to extend the definition of legal strings without markers given in [2]. Consider the definition of the **sspr** operation  $\text{sspr}_p(u_1pu_2\bar{p}u_3) = u_1\bar{u}_2u_3$ . The operation  $\text{sspr}_p$  may only be applied if  $|u_2| = 1$ , i.e., it contains a single pointer or a single marker. As no markers are recorded in the definition given in [2],  $u_2$  could in fact contain a pointer and a marker, thus making **sspr** <sub>$p$</sub>  inapplicable. The problem is illustrated in the following example.

**Example 9** *Let  $u = \overline{324234}$  be a legal string with the markers removed. It would seem that **sspr**<sub>2</sub> should be applicable to  $u$ , as they are separated by only a single pointer. However,  $u$  can be obtained from the following signed permutation  $\pi = \overline{2413}$ . It is clear that **sh**<sub>1</sub> is not applicable to  $\pi$  because of MDS 4. Since the permutation-based model and the string-based one should be equivalent, this in turn implies that **sspr**<sub>2</sub> should not be applicable to  $u$ .*

## Confluent Strategies on Legal Strings

It was shown in [9] that assembly strategies for a given signed permutation using simple operations are confluent: they are either all successful, or all unsuccessful and moreover, they lead to final results having the same structure. We will now show that the same applies on legal strings. Note that we have no need to define structure for legal strings, as they are in fact isomorphic and thus share the same characteristics and the same operations may be applied to both. Note also that the results of this section are simpler to prove and more general than the similar results in [9].

We will now show that assembly strategies on legal strings are confluent. Lemma 1 will show the case where one of the operations is an **snr** operation, Lemma 2 considers two **sspr** operations, and Lemma 3 considers two **ssdr** operations.

**Lemma 1** *Let  $u$  be a legal string over  $\Sigma_n$  and  $\phi, \psi \in \text{Snr} \cup \text{Sspr} \cup \text{Ssdr}$  be two operations applicable to  $u$ . If  $\phi \in \text{Snr}$  or  $\psi \in \text{Snr}$ , then  $\phi \circ \psi(u) = \psi \circ \phi(u)$ .*

*Proof.* The proof for this is straightforward as the pointers involved in an **snr** operation must be adjacent, and thus cannot overlap or affect any other pointers.  $\square$

**Lemma 2** *Let  $u$  be a legal string over  $\Sigma_n$  and  $\psi, \phi \in \text{Sspr}$  be two operations applicable to  $u$ . Then either  $\phi \circ \psi(u) = \psi \circ \phi(u)$ , or  $\psi(u) \equiv \phi(u)$ .*

*Proof.* Let  $\psi = \text{sspr}_p$  and  $\phi = \text{sspr}_q$ , for some  $p \neq q$ . If  $pq\bar{p}\bar{q} \not\leq u$  and  $qp\bar{q}\bar{p} \not\leq u$ , then clearly  $\text{sspr}_q \circ \text{sspr}_p(u) = \text{sspr}_p \circ \text{sspr}_q(u)$ .

Now, if  $pq\bar{p}\bar{q} \leq u$  or  $qp\bar{q}\bar{p} \leq u$ , then clearly applying one operation makes the other inapplicable, but if  $\psi(u) = qq$  and  $\phi(u) = pp$  then  $\psi(u) \equiv \phi(u)$ .  $\square$

Note, however, that after applying one of the **sspr** operations, an **snr** operation on the remaining pointer becomes available, giving us the following:

$$\begin{aligned}\text{sspr}_p(u_1pq\bar{p}\bar{q}u_2) &= u_1\bar{q}\bar{q}u_2, \\ \text{sspr}_q(u_1pq\bar{p}\bar{q}u_2) &= u_1ppu_2,\end{aligned}$$

$$\text{snr}_{\bar{q}} \circ \text{sspr}_p(u_1pq\bar{p}\bar{q}u_2) = u_1u_2 = \text{snr}_p \circ \text{sspr}_q(u_1pq\bar{p}\bar{q}u_2),$$

for some legal strings  $u_1, u_2$  over  $\Sigma_n$ .

**Lemma 3** *Let  $u$  be a legal string over  $\Sigma_n$  and  $\psi, \phi \in \text{Ssdr}$  be two operations applicable to  $u$ . Then either  $\phi \circ \psi(u) = \psi \circ \phi(u)$ , or  $\psi(u) \equiv \phi(u)$ .*

*Proof.* Let  $\psi = \text{ssdr}_{p,q}$  and  $\phi = \text{ssdr}_{r,s}$ , for some  $p, q \neq r, s$ . If they have a different signature, then clearly  $\phi \circ \psi(u) = \psi \circ \phi(u)$ . Assume then that  $p, q, r, s \in \Sigma_n$ . The cases when one or more of  $p, q, r, s$  are in  $\bar{\Sigma}_n$  are completely similar.

Also, if  $p, q \neq r, s$  then  $\phi \circ \psi(u) = \psi \circ \phi(u)$ , so let  $q = r$  (the case when  $p = s$  is similar). Thus either  $pqsu_1pqs \leq u$  or  $sqpu_1sqp \leq u$ . Both operations cannot be applied, but  $\text{ssdr}_{p,q}(u)$  and  $\text{ssdr}_{q,s}(u)$  arrive at equivalent results. Indeed,  $\text{ssdr}_{p,q}(u) = u_2su_1su_3$ ,  $\text{ssdr}_{q,s}(u) = u_2pu_1pu_3$  and so,  $\text{ssdr}_{p,q}(u) \equiv \text{ssdr}_{q,s}(u)$ .  $\square$

We can now extend the results above for strategies using all three operations.

**Theorem 2** *Let  $u$  be a legal string over  $\Sigma_n$  and  $\phi, \psi \in \text{Snr} \cup \text{Sspr} \cup \text{Ssdr}$  be two operations applicable to  $u$ . Then either  $\phi \circ \psi(u) = \psi \circ \phi(u)$ , or  $\psi(u) \equiv \phi(u)$ .*

*Proof.* Based on the previous three lemmata, we only need to prove the claim in the case  $\phi \in \text{Sspr}$ ,  $\psi \in \text{Ssdr}$ . Now, if  $\phi = \text{sspr}_p$ ,  $\psi = \text{ssdr}_{q,r}$ , and  $p \neq q, r$ , then  $\phi \circ \psi(u) = \psi \circ \phi(u)$ .

Assume then that  $p = q$ . Then in  $\text{sspr}_p$ , the occurrences of  $p$  must have different signatures, and in  $\text{ssdr}_{p,r}$ , all occurrences of  $p$  and  $r$  must have the same signature, a contradiction.  $\square$

**Example 10** *Let  $u = b2345623456e$ .*

- i. Both  $\text{ssdr}_{2,3}$  and  $\text{ssdr}_{5,6}$  are applicable to  $u$ . Moreover,  $\text{ssdr}_{2,3} \circ \text{ssdr}_{5,6}$  and  $\text{ssdr}_{5,6} \circ \text{ssdr}_{2,3}$  are also applicable to  $u$  and*

$$\text{ssdr}_{2,3} \circ \text{ssdr}_{5,6}(u) = \text{ssdr}_{5,6} \circ \text{ssdr}_{2,3}(u) = b44e.$$

- ii. Both  $\text{ssdr}_{2,3}$  and  $\text{ssdr}_{3,4}$  are applicable to  $u$ . Moreover, applying either operation gives an equivalent result:*

$$\sigma(\text{ssdr}_{2,3}(u)) = b456456e \equiv b256256e = \sigma(\text{ssdr}_{3,4}(u)).$$

We now have the necessary information to show that the result for signed permutations described in [9] also applies on legal strings, namely that simple operations on legal strings are confluent.

**Theorem 3** *Let  $u$  be a legal string over  $\Sigma_n$  and  $\phi, \psi$  be two strategies for  $u$ . Then either  $\phi$  and  $\psi$  are both sorting strategies for  $u$ , or they are both unsuccessful strategies. Moreover,  $\phi(u) \equiv \psi(u)$ .*

*Proof.* Assume that the claim of the theorem is not true and consider a legal string  $u$  of minimal length such that  $\phi = \phi_1 \dots \phi_k$  is a successful strategy for  $u$ , while  $\psi = \psi_1 \dots \psi_l$  is an unsuccessful one,  $\phi_i, \psi_j \in \text{Snr} \cup \text{Sspr} \cup \text{Ssdr}$ ,  $1 \leq i \leq k$ ,  $1 \leq j \leq l$ .

It follows from Theorem 2 that either  $\phi_k(u) \equiv \psi_l(u)$ , or  $\phi_k \circ \psi_l(u) = \psi_l \circ \phi_k(u)$ . If they are equivalent, then  $\phi_k(u)$  or  $\psi_l(u)$  would be a smaller counterexample than  $u$  contradicting the minimality of  $u$ . In the latter case note that due to the minimality of  $u$ , it follows that  $\phi_k(u)$  has only successful strategies and  $\psi_l(u)$  has only unsuccessful strategies. Consequently,  $\psi_l \circ \phi_k(u)$  has both successful and unsuccessful strategies, contradicting the minimality of  $u$ .  $\square$

**Example 11** *The legal string  $u = 45\overline{76}7eb2\overline{3}2\overline{3}456$  has several sorting strategies. some of them are shown below.*

$$\begin{aligned} \phi_1(u) &= \text{snr}_6 \circ \text{snr}_3 \circ \text{sspr}_2 \circ \text{sspr}_7 \circ \text{ssdr}_{4,5}(u) = eb, \\ \phi_2(u) &= \text{snr}_4 \circ \text{ssdr}_{5,6} \circ \text{sspr}_7 \circ \text{snr}_2 \circ \text{sspr}_3(u) = eb, \\ \phi_3(u) &= \text{snr}_6 \circ \text{ssdr}_{4,5} \circ \text{sspr}_7 \circ \text{snr}_3 \circ \text{sspr}_2(u) = eb, \\ \phi_4(u) &= \text{snr}_4 \circ \text{ssdr}_{5,6} \circ \text{sspr}_7 \circ \text{snr}_3 \circ \text{sspr}_2(u) = eb. \end{aligned}$$

**Example 12** *The legal string  $v = b234678e56782345$  has several unsuccessful strategies. Some of them are shown below.*

$$\begin{aligned}\psi_1(v) &= \text{ssdr}_{6,7} \circ \text{ssdr}_{2,3}(v) = b48e5845, \\ \psi_2(v) &= \text{ssdr}_{7,8} \circ \text{ssdr}_{3,4}(v) = b26e5625, \\ \psi_3(v) &= \text{ssdr}_{6,7} \circ \text{ssdr}_{3,4}(v) = b28e5825, \\ \psi_4(v) &= \text{ssdr}_{2,3} \circ \text{ssdr}_{7,8}(v) = b46e5645.\end{aligned}$$

*Note that  $\psi_1(v) \equiv \psi_2(v) \equiv \psi_3(v) \equiv \psi_4(v)$ .*

## References

1. Cavalcanti, A., Clarke, T.H., Landweber, L., MDS\_IES\_DB: a database of macronuclear and micronuclear genes in spirotrichous ciliates. *Nucleic Acids Research* **33** (2005) 396–398.
2. Ehrenfeucht, A., Harju, T., Petre, I., Prescott, D. M., and Rozenberg, G., *Computation in Living Cells: Gene Assembly in Ciliates*, Springer (2003).
3. Ehrenfeucht, A., Prescott, D. M., and Rozenberg, G., Computational aspects of gene (un)scrambling in ciliates. In: L. F. Landweber, E. Winfree (eds.) *Evolution as Computation*, Springer, Berlin, Heidelberg, New York (2001) pp. 216–256.
4. Harju, T., Li, C, Petre, I., and Rozenberg, G., Parallelism in gene assembly. *Natural Computing*, (2006).
5. Harju, T., Li, C, Petre, I., and Rozenberg, G., Complexity Measures for Gene Assembly In: K. Tuyls (Eds.), *Proceedings of the Knowledge Discovery and Emergent Complexity in Bioinformatics workshop*, Springer, Lecture Notes in Bioinformatics, 4366, 2007.
6. Harju, T., Li, C, Petre, I., and Rozenberg, G., Modelling simple operations for gene assembly. In: Junghuei Chen, Natasha Jonoska, Grzegorz Rozenberg (Eds), *Nanotechnology: Science and Computation*, 361-376, Springer, 2006.
7. Landweber, L. F., and Kari, L., The evolution of cellular computing: Nature’s solution to a computational problem. In: *Proceedings of the 4th DIMACS Meeting on DNA-Based Computers*, Philadelphia, PA (1998) pp. 3–15.
8. Landweber, L. F., and Kari, L., Universal molecular computation in ciliates. In: L. F. Landweber and E. Winfree (eds.) *Evolution as Computation*, Springer, Berlin Heidelberg New York (2002).
9. Langille, M., Petre, I., Simple gene assembly is deterministic. *Fundamenta Informaticae* **72** (2006) 1–12, IOS Press.
10. Prescott, D. M., Ehrenfeucht, A., and Rozenberg, G., Molecular operations for DNA processing in hypotrichous ciliates. *Europ. J. Protistology* **37** (2001) 241–260.