

The evolution of cellular computing: nature's solution to a computational problem

Laura F. Landweber^{a,*}, Lila Kari^b

^a Ecology and Evolutionary Biology, Princeton University, Princeton, NJ 08544-1003, USA

^b Department of Computer Science, University of Western Ontario, London, Ont., N6A 5B7, Canada

Abstract

How do cells and nature 'compute'? They read and 'rewrite' DNA all the time, by processes that modify sequences at the DNA or RNA level. In 1994, Adleman's elegant solution to a seven-city directed Hamiltonian path problem using DNA launched the new field of DNA computing, which in a few years has grown to international scope. However, unknown to this field, two ciliated protozoans of the genus *Oxytricha* had solved a potentially harder problem using DNA several million years earlier. The solution to this problem, which occurs during the process of gene unscrambling, represents one of nature's ingenious solutions to the problem of the creation of genes. RNA editing, which can also be viewed as a computational process, offers a second algorithm for the construction of functional genes from encrypted pieces of the genome. © 1999 Elsevier Science Ireland Ltd. All rights reserved.

Keywords: DNA computing; Scrambled gene; Molecular evolution; Ciliate; Hypotrich

1. Gene unscrambling as a computational problem

1.1. Introduction

Ciliates are a diverse group of 8000 or more unicellular eukaryotes (nucleated cells) named for their wisp-like covering of cilia. They possess two types of nuclei: an active *macronucleus* (soma) and a functionally inert *micronucleus* (germline), which contribute only to sexual reproduction. The somatically active macronucleus forms from the

germline micronucleus after sexual reproduction, during the course of development. The genomic copies of some protein-coding genes in the micronucleus of hypotrichous ciliates are obscured by the presence of intervening non-protein-coding DNA sequence elements (internally eliminated sequences or IESs). These must be removed before the assembly of a functional copy of the gene in the somatic macronucleus. Furthermore, the protein-coding DNA segments (macronuclear destined sequences or MDSs) in *Oxytricha* species are sometimes present in a permuted order relative to their final position in the macronuclear copy. For example, in *O. nova*, the micronuclear copy of three genes (Actin I, α -telomere binding protein, and DNA polymerase α) must be re-

* Corresponding author. Tel.: +1-609-258-1947; fax: +1-609-258-1682.

E-mail addresses: lfl@princeton.edu (L.F. Landweber), lila@csd.uwo.ca (L. Kari)

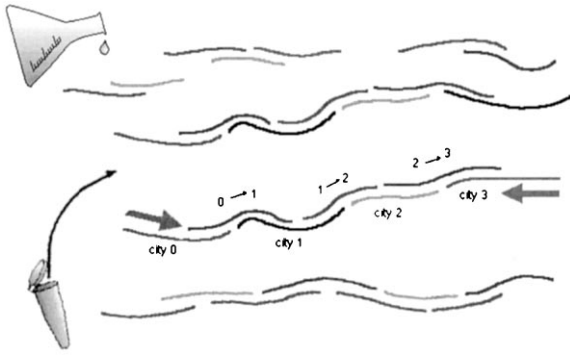


Fig. 1. DNA hybridization in a molecular computer. PCR primers are indicated by arrows.

ordered and intervening DNA sequences removed in order to construct functional macronuclear genes. Most impressively, the gene encoding DNA polymerase α (DNA pol α) in *O. trifallax* is apparently scrambled in 50 or more pieces in its germline nucleus (Hoffman and Prescott, 1997). Destined to unscramble its micronuclear genes by putting the pieces together again, *O. trifallax* routinely solves a potentially complicated computational problem when rewriting its genomic sequences to form the macronuclear copies.

This process of unscrambling bears a remarkable resemblance to the DNA algorithm used by Adleman (1994) to solve a seven-city instance of the Directed Hamiltonian Path problem. Section

1.4 introduces a formal model of gene unscrambling. (Adleman's algorithm involves the use of edge-encoding sequences as splints to connect city-encoding sequences, allowing the formation of all possible paths through the graph (Fig. 1). Afterwards, a screening process eliminates the paths that are not Hamiltonian, i.e. ones which either skip a city, enter a city twice, or do not start and end in the correct origin and final destinations.)

The developing ciliate macronuclear 'computer' (Figs. 2 and 3) apparently relies on the information contained in short repeat sequences to act as guides in a series of homologous recombination events (Table 1). These guide sequences provide the splints analogous to the edges in Adleman's graph, and the process of recombination results in linking the protein-encoding segments (MDSs, or 'cities') that belong next to each other in the final protein coding sequence ('Hamiltonian path'). As such, the unscrambling of sequences that encode DNA polymerase α accomplishes an astounding feat of cellular computation, especially as 50-city Hamiltonian path problems are often considered hard problems in computer science and present a formidable challenge to a biological computer. Other structural components of the ciliate chromatin presumably play a significant role, but the exact details of the mechanism are still unknown.

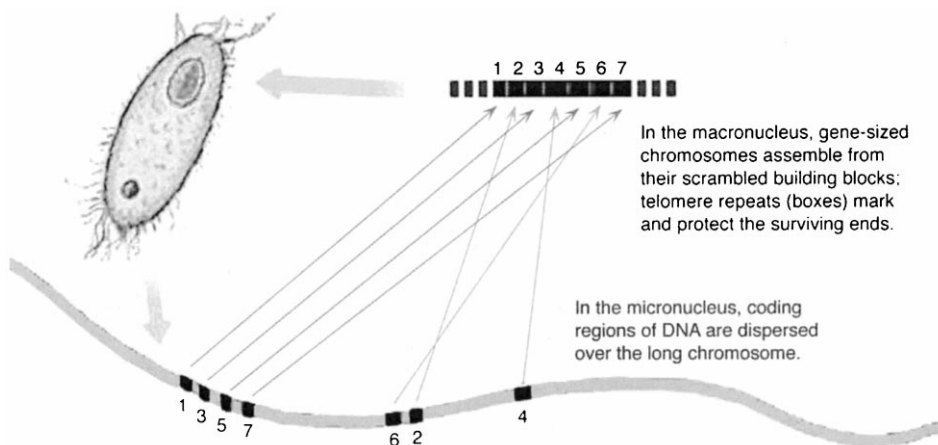


Fig. 2. Overview of gene unscrambling. Dispersed coding MDSs 1–7 reassemble during macronuclear development to form the functional gene copy (top), complete with telomere addition to mark and protect both ends of the gene.

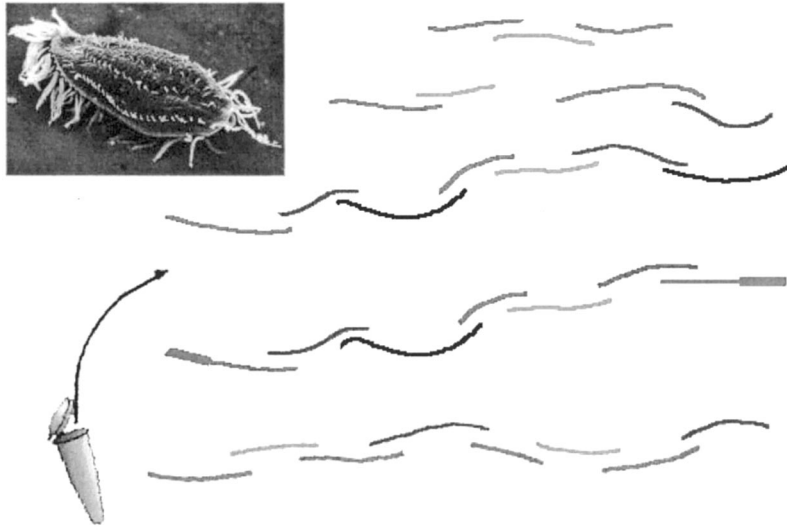


Fig. 3. A ciliate molecular computer? Correct gene assembly in *Styloynchia* (inset) (Lynn and Corliss, 1991) requires the joining of many segments of DNA guided by short sequence repeats, only at the ends. Telomeres, indicated by thicker lines, mark the termini of correctly assembled gene-sized chromosomes. Note the striking similarities to DNA computations that specifically rely on pairing of short repeats at the ends of DNA fragments, as in the experiment of Adleman (1994).

1.2. The path towards unscrambling

Typical IES excision in ciliates involves the removal of short (14 ~ 600 bp) A–T rich sequences, often released as circular DNA molecules (Tausta and Klobutcher, 1989). The choice of which sequences to remove appears to be minimally ‘guided’ by recombination between direct repeats of only 2–14 bp (Table 1).

Unscrambling is a particular type of IES removal in which the order of the MDSs in the MIC is often radically different from that in the MAC. For example, in the micronuclear genome of *Oxytricha nova*, the MDSs of α -telomere binding protein (α -TP) are arranged in the cryptic order 1–3–5–7–9–11–2–4–6–8–10–12–13–14 relative to their position in the ‘clear’ macronuclear sequence 1–2–3–4–5–6–7–8–9–10–11–12–13–14. This particular arrangement predicts a spiral mechanism in the path of unscrambling, which links odd and even segments in order (Fig. 4; Mitcham et al., 1992).

Homologous recombination between identical short sequences at appropriate MDS–IES junctions is implicated in the mechanism of gene unscrambling, as it could simultaneously remove

the IESs and reorder the MDSs. For example, the DNA sequence present at the junction between MDS n and the downstream IES is generally the same as the sequence between MDS $n + 1$ and its upstream IES, leading to correct ligation of MDS n to MDS $n + 1$ over a distance (Table 1). However, the presence of such short repeats (average length four bp between non-scrambled MDSs, nine bp between scrambled MDSs (Prescott and Dubois, 1996)) suggests that although these guides are necessary, they are certainly not sufficient to guide accurate splicing. Hence, it is likely that the repeats satisfy more of a structural requirement for MDS splicing, and less of a role in substrate recognition. Otherwise, incorrectly spliced sequences (the results of promiscuous recombination) would dominate, especially in the case of very small (2–4 bp) repeats that would be present thousands of times throughout the genome. This incorrect hybridization could be a driving force in the production of newly scrambled patterns in evolution. However during macronuclear development only unscrambled molecules that contain 5' and 3' telomere addition sequences would be selectively retained in the macronucleus, ensuring that most promiscuously ordered genes would be lost.

Table 1

O. trifallax DNA polymerase α (data modified from Hoffman and Prescott, 1997)

5' MDS/IES junction sequence	MDS	3' MDS/IES junction sequence	Number repeats in Mac sequence*
5' Telomere addition site	1	AGATA	8
AGATA	2	ATT	*
ATT	3	ATA	*
ATA	4	ATGATGAGTGAAT	1
ATGATGAGTGAAT	5	AACAGAAC	1
AACAGAAC	6	AGAAATATG	1
AGAAATATG	7	n.d.	
n.d.	9	TTATCATT	2
TTATCATT	10	AAAATAAT	1
AAAATAAT	11	GTTTCTTG	1
GTTTCTTG	12	ATGCAA	1
ATGCAA	13	TAAAATGA	1
TAAAATGA	14	AGAGGAG	1
AGAGGAG	15	TAATGATGG	1
TAATGATGG	16	ATGGTGAG	1
ATGGTGAG	17	AAAATCAA	3
AAAATCAA	18	AAAGCATGCTTG	1
AAAGCATGCTTG	19	GATTTCAAGAAAA	1
GATTTTAAAGAAAA	20	GTTACTCTTG	1
GTTACTCTTG	21	GCTCAATAAAAA	1
GCTCAATAAAAA	22	ATCTTG	2
ATCATG	23	AAAACTT	1
AAAACTT	24	GAGAGATAGA	1
GAGAGATAGA	25	TAGTTGCTC	1
TAGTTGCTC	26	AAGCTAGATTTT	1
AAGCTAGATTTT	27	GGAGGATC	1
GGAGGATC	28	CAAGATAA	1
CAAGATAA	29	GTTCAACT	1
GTTCAACT	30	ATAAGACTTTGATGA	1
ATAAGACTTTGATGA	31	CTAATGAA	1
CTAATGAA	32	n.d.	
n.d.	36	CTTGAGAT	1
CTTGAGAT	37	AAAGTAGTTTAG	1
AAAGTAGTTTAG	38	CACTTTCAA	1
CACTTTCAA	39	ATGAAAAATAA	1
ATGAAAAATAA	40	CCTTGGATCA	1
CCTTGGATCA	41	AAGAGTGAAT	1
AAGAGTGAAT	42	TGAACAACCTT	1
TGAACAACCTT	43	GTGCTTAG	1
GTGCTTAG	44	n.d.	
n.d.	49	ATAAAA	4
ATAAAA	50	AT	*
AT	51	3' Telomere addition site	

* The number of occurrences of di- and trinucleotides was not determined since they would be extremely represented in any gene sequence. Note the values shown in this column only represent the number of occurrences of these sequence motifs in the macronuclear copy of the gene. There are also several occurrences of these repeats throughout the non-coding portion of the micronuclear copy of this gene as well as throughout the entire genome; each such occurrence offers the opportunity for incorrect pairing, which would lead to the production of 'dead-end' copies of the gene. These would, however, be unlikely to contain telomere addition sequences at both ends. Junction sequences for MDSs 32–36 and 44–49 are unknown because of missing micronuclear sequence data.

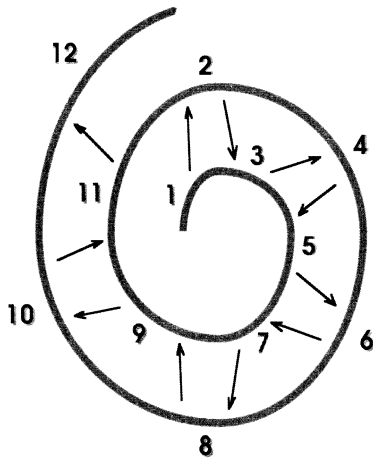


Fig. 4. Model for unscrambling in α -TP. (adapted from Mitcham et al., 1992).

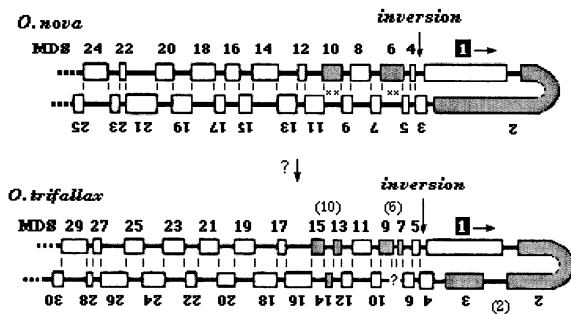


Fig. 5. Model for scrambling of DNA pol α . Vertical lines indicate recombination junctions between scrambled MDSs, guided by direct repeats. MDS 1 contains the start of the gene. MDS 10 in *O. nova* can also give rise to three new MDSs (13–15) in *O. trifallax*, one scrambled on the inverted strand, by two spontaneous intramolecular recombination events (\times) in the folded orientation shown. *O. nova* MDS 6 can give rise to *O. trifallax* MDSs 7–9 (MDS 8, shaded, is only 6 bp and was not identified in (Hoffman and Prescott, 1997)). *O. trifallax* non-scrambled MDSs 2 and 3 could be generated by the insertion of an IES in *O. nova* MDS 2 (similar to a model suggested by M. DuBois in Hoffman and Prescott, 1997).

1.3. Inversions as catalysts of DNA rearrangements

The micronuclear actin I gene has a scrambled MDS order 3–4–6–5–7–9–2–1–8 in *O. nova*, with MDS 2 inverted (present on the opposite strand and in the opposite direction) relative to the others (Dubois and Prescott, 1995). DNA

polymerase α has at least 44 MDSs in *O. nova* and 51 in *O. trifallax* (Table 1), scrambled in a non-random order with an inversion in the middle, and some MDSs located at least several kb away from the main gene (in an unmapped PCR fragment). The resulting hairpin structural constraint predicted in Fig. 5 equips the ciliate with a dramatic shortcut to finding the correct solution to its DNA polymerase α unscrambling problem.

Figs. 5 and 6 outline a model for the origin and accumulation of scrambled MDSs. The appearance of an inversion is likely to encourage the formation of new MDSs in a nonrandomly scrambled pattern. By Muller's Ratchet, an inversion makes the addition of new MDSs much more likely, given that the hairpin structure, which juxtaposes coding and non-coding DNA sequences, would promote recombination, possibly between short arbitrary repeats. For example, the arrangement of MDSs 2, 6 and 10 in *O. nova* could have given rise to the arrangement of eight new MDSs in *O. trifallax* (Fig. 5).

We have recently discovered scrambling in the gene encoding DNA polymerase α in the micronucleus of a different ciliate, *Stylonychia lemnae*, which enjoys the benefit of a working transfection system (Wen et al., 1996). The scrambled gene in *S. lemnae* appears to share the presence of an inversion with the two *Oxytricha* species. These scrambled genes in ciliates thus offer a unique system in which to study the origin of a complex genetic mechanism and the role of inversions as catalysts of acrobatic DNA rearrangements during evolution (Fig. 6). DNA polymerase α 's complex scrambling pattern is possibly the best analog equivalent of a hard path finding problem in nature. Alternate splicing at the RNA level, as well as other forms of programmed DNA rearrangements, could also be viewed as solutions to path finding problems in nature. Dynamic processes, such as maturation of the immune response, provide examples of genuine evolutionary computation in cells, whereas the path finding problems here may follow a more deterministic algorithm. Current effort is directed toward (1) recoding DNA in the laboratory and (2) understanding how cells unscramble DNA, how this process has arisen, and how the 'programs' are

written and executed. Do they decode the message by following the shortest unscrambling path or by following a more circuitous but equally effective route, as in the case of RNA editing (below)? Also, how error prone is the unscrambling process? Does it actually search through several plausible unscrambled intermediates or follow a strictly deterministic pathway? The isolation of functional nucleic acid molecules, such as RNA catalysts (ribozymes), from large pools of random sequence offers yet a different paradigm for molecular computation (Bartel and Szostak, 1993; Landweber, 1997).

1.4. The formal model

Before introducing the formal model, we summarize our notation. An alphabet Σ is a finite, non-empty set. In our case $\Sigma = \{A, C, G, T\}$. A sequence of letters from Σ is called a string (word) over Σ and in our interpretation corresponds to a linear strand. The words are denoted by lowercase letters such as u, v, α_i, x_{ij} . A word that has no letters in it is called an empty word. The set of all possible words consisting of letters from Σ is denoted by Σ^* . We also define circular words over Σ by declaring two words to be equivalent if and only if one is a cyclic permutation of the other. In other words, w is equivalent to w' if and only if

they can be factored as $w = uv$ and $w' = vu$, respectively. We denote a representative of the equivalence class of w by $\bullet w$. Such a circular word $\bullet w$ refers to any of the circular permutations of the letters in w .

With this notation, we introduce two operations that model the processes that occur during the homologous recombination.

Operation (1), intramolecular recombination, is unary:

$$uxwxv \Rightarrow uxv + \bullet wx$$

where u, w, x and v are words in Σ^* , and x is non-empty. Here ‘+’ is interpreted as the union of the two resulting strands.

Operation (1) is reversible. Note that op1 in the forward direction is formally intramolecular recombination, whereas op1 in the reverse direction is intermolecular recombination.

Thus, operation (1) models the process of intramolecular recombination that occurs during unscrambling of the gene. x is a repeated sequence that guides the homologous recombination. After x finds its second occurrence in $uxwxv$, the molecule undergoes a strand exchange in x that leads to the formation of two new molecules: uxv and a circular DNA molecule $\bullet wx$.

Operation (1) also accomplishes the deletion of either sequence wx or xw from the original

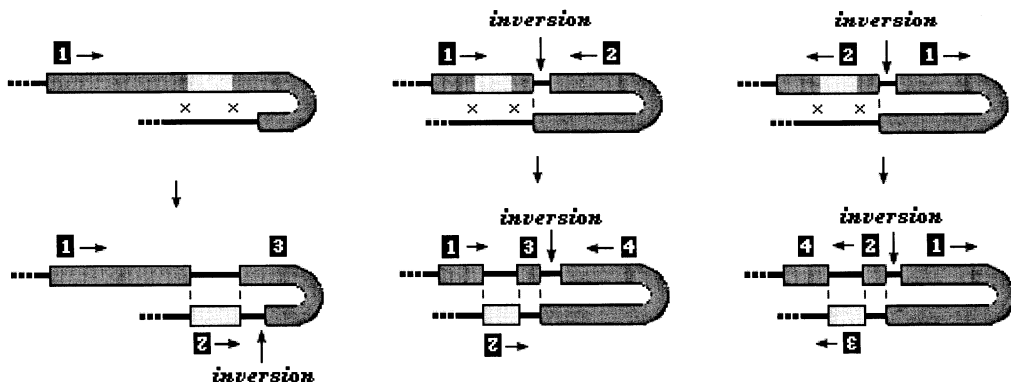


Fig. 6. Proposed model for the origin of a scrambled gene. Left: birth of a scrambled gene from a non-scrambled gene by a double recombination with an IES or any non-coding DNA (new MDS order 1–3–2 with an inversion between MDSs 3 and 2). Middle: generation of a scrambled gene with a non-random MDS order, from a non-scrambled gene with an inversion between two MDSs. Right: creation of new scrambled MDSs in a scrambled gene containing an inversion. Inversions may dramatically increase the production of scrambled MDSs, by stabilizing the folded conformation that allows reciprocal recombination across the inversion.

molecule $uxwxv$. The fact that $\bullet wx$ is circular implies that we can use any circular permutation of its sequence as an input for a subsequent operation.

Operation (1 \bullet) is also unary:

$$\bullet uxwxv \Rightarrow \bullet uxv + \bullet wx$$

Operation (1 \bullet) is similar to op1, the only difference in that the input is circular, which results in the output of two circular strands. Like op1, op1 \bullet is reversible.

Operation (2), intermolecular recombination, is binary:

$$uxv + u'xv' \Rightarrow uxv' + u'xv$$

where u, x, v, u', v' are words in Σ^* , and x is non-empty. Operation 2 is also reversible.

Operation (2) models most processes of intermolecular recombination. Given two molecules uxv and $u'xv'$, both of which contain a homologous (identical) subsequence x , the molecules undergo a strand exchange (homologous recombination) in x that leads to the formation of molecules uxv' and $u'xv$. This operation effectively rewrites the input sequences, for example by replacing the suffix v from uxv with v' , a process analogous to *trans*-splicing (Sullenger and Cech, 1994).

Note that each operation and its reverse conserves the number of ‘ends’ (a linear strand has two ends while a circular strand has none). Having defined the operations modeling intra- and intermolecular recombinations, we now remark that the recombination events predicted to occur during gene unscrambling are capable of generating a variety of products. These include the production of one circular output strand from two circular inputs, and vice versa (by op1 \bullet and its reverse), the generation of both a linear and a circular output strand from one linear input strand (by op1), the generation of a linear strand by the combination of a linear and circular input (by reversed op1), and finally the generation of two recombined linear output strands from two linear input strands (by op2 and its reverse).

These operations resemble the splicing operation introduced by Head (1987) as a model of DNA recombination and the splicing on circular

strands studied by Siromoney et al. (1992) and Pixton (1995). Paun (1995) and Cshuhaj-Varju et al. (1996) subsequently showed that this model has the computational power of a universal Turing Machine.

The process of gene unscrambling entails a series of successive or possibly simultaneous intra- and intermolecular homologous recombination events. This is followed by the excision of all sequences $\tau_s y \tau_e$, where the sequence y is marked by the presence of τ_s , a telomere ‘start’ (at its 5’ end), and τ_e , a telomere ‘end’ (at its 3’ end). Thus, from a long sequence $u \tau_s y \tau_e v$, this step retains only $\tau_s y \tau_e$ in the macronucleus. Lastly, the enzyme telomerase extends the length of the telomeric sequences (usually double-stranded $\{\text{TTTTGGGG}\}_n$ repeats in these organisms) from the ‘telomere addition sequences’, τ_s and τ_e , to protect the ends of the DNA molecule; however the telomere addition step is not present in our formal model.

We now make the assumption that, by a clever structural alignment, such as the one depicted in Fig. 4, and numerous other biological factors, the cell decides which sequences are non-protein-coding (IESs) and which are ultimately protein-coding (MDSs), as well as which sequences x guide homologous recombination. Moreover, such biological shortcuts are presumably essential to bring into proximity the guiding sequences x .

Each MDS, denoted primarily by α_i , $1 < i < n$ (where n is the number of pieces sparsely present in the micronucleus that assemble to form the functional gene in the macronucleus) is flanked by the guiding sequences $x_{i-1,i}$ and $x_{i,i+1}$. Each guiding sequence points to the MDS that should precede or follow α_i in the final sequence. The only exceptions are α_1 , which is preceded by τ_s , and α_n which is followed by τ_e (Table 1). Note that although present generally once in the final macronuclear copy, each $x_{i,i+1}$ occurs at least twice in the micronuclear copy, once after α_i and once before α_{i+1} .

We denote by ε_k an internal sequence that is deleted; ε_k does not occur in the final sequence. Thus, since unscrambling leaves one copy of each $x_{i,i+1}$ between α_i and α_{i+1} , an IES is non-deterministically either $\varepsilon_k x_{i,i+1}$ or $x_{i,i+1} \varepsilon_k$, depending

on which guiding sequence x_{i+1} is eliminated. Similarly an MDS is technically either $\alpha_i x_{i,i+1}$ or $x_{i-1,i} \alpha_i$. For the purposes of this model, either choice is equivalent.

On a technical note, removal of simple (non-scrambled) IES's in *Euplotes* leaves extra sequences (including a duplication of x_{ij}) at the junctions between ϵ_k 's in the resulting non-protein-coding products. This may result when the x_{ij} values are as short as two nucleotides (Klobutcher et al., 1993). It is unknown whether unscrambling also introduces extra sequences, since it uses considerably longer x_{ij} values on average. However, since the extra sequences have always been found at junctions between ϵ_k values, this would not affect our unscrambling model.

The following example models unscrambling of a micronuclear gene that contains MDSs in the scrambled order 2–4–1–3:

$$\begin{aligned} & UX_{12}\alpha_2 X_{23}\epsilon_1 X_{34}\alpha_4 \tau_e \epsilon_2 \tau_s \alpha_1 X_{12}\epsilon_3 X_{23}\alpha_3 X_{34}V \Rightarrow \\ & UX_{12}\epsilon_3 X_{23}\alpha_3 X_{34}V + \bullet \alpha_2 X_{23}\epsilon_1 X_{34}\alpha_4 \tau_e \epsilon_2 \tau_s \alpha_1 X_{12} = \\ & UX_{12}\epsilon_3 X_{23}\alpha_3 X_{34}V + \bullet \epsilon_1 X_{34}\alpha_4 \tau_e \epsilon_2 \tau_s \alpha_1 X_{12}\alpha_2 X_{23} \Rightarrow \\ & UX_{12}\epsilon_3 X_{23}\epsilon_1 X_{34}\alpha_4 \tau_e \epsilon_2 \tau_s \alpha_1 X_{12}\alpha_2 X_{23}\alpha_3 X_{34}V \Rightarrow \\ & UX_{12}\epsilon_3 X_{23}\epsilon_1 X_{34}V + \bullet \alpha_4 \tau_e \epsilon_2 \tau_s \alpha_1 X_{12}\alpha_2 X_{23}\alpha_3 X_{34} = \\ & UX_{12}\epsilon_3 X_{23}\epsilon_1 X_{34}V + \bullet \tau_s \alpha_1 X_{12}\alpha_2 X_{23}\alpha_3 X_{34}\alpha_4 \tau_e \epsilon_2 \Rightarrow \\ & \tau_s \alpha_1 X_{12}\alpha_2 X_{23}\alpha_3 X_{34}\alpha_4 \tau_e + \epsilon_2 + UX_{12}\epsilon_3 X_{23}\epsilon_1 X_{34}V \end{aligned}$$

Note that the process is non-deterministic in that, for example, one could start by replacing the first step, which was the recombination between homologous sequences x_{12} , by recombination between the homologous sequences x_{34} instead, obtaining thus

$$\begin{aligned} & UX_{12}\alpha_2 X_{23}\epsilon_1 X_{34}\alpha_4 \tau_e \epsilon_2 \tau_s \alpha_1 X_{12}\epsilon_3 X_{23}\alpha_3 X_{34}V \Rightarrow \\ & UX_{12}\alpha_2 X_{23}\epsilon_1 X_{34}V + \bullet \alpha_4 \tau_e \epsilon_2 \tau_s \alpha_1 X_{12}\epsilon_3 X_{23}\alpha_3 X_{34} = \\ & UX_{12}\alpha_2 X_{23}\epsilon_1 X_{34}V + \bullet \epsilon_3 X_{23}\alpha_3 X_{34}\alpha_4 \tau_e \epsilon_2 \tau_s \alpha_1 X_{12} \Rightarrow \\ & UX_{12}\epsilon_3 X_{23}\alpha_3 X_{34}\alpha_4 \tau_e \epsilon_2 \tau_s \alpha_1 X_{12}\alpha_2 X_{23}\epsilon_1 X_{34}V \Rightarrow \\ & \bullet \alpha_3 X_{34}\alpha_4 \tau_e \epsilon_2 \tau_s \alpha_1 X_{12}\alpha_2 X_{23} + UX_{12}\epsilon_3 X_{23}\epsilon_1 X_{34}V = \\ & \bullet \tau_s \alpha_1 X_{12}\alpha_2 X_{23}\alpha_3 X_{34}\alpha_4 \tau_e \epsilon_2 + UX_{12}\epsilon_3 X_{23}\epsilon_1 X_{34}V \Rightarrow \\ & \tau_s \alpha_1 X_{12}\alpha_2 X_{23}\alpha_3 X_{34}\alpha_4 \tau_e + \epsilon_2 + UX_{12}\epsilon_3 X_{23}\epsilon_1 X_{34}V \end{aligned}$$

in the same number of steps.

In effect, the above examples show that, as the input strand is always linear, albeit scrambled, the intermediate steps will generally produce at most one linear strand as output (apart from telomere addition or other mechanisms that may lead to chromosome breakage or fragmentation). Indeed, in the most basic case, the output of an operation that has only one linear strand as input can never be two linear strands. Consequently, the process involves only iterative application of op1 and op1*. Formally, op2 can only occur after the telomere addition phase (the last step) which we do not define as a separate operation in our model. This does not reduce the generality of the model, as telomere addition happens only once at the end of the process.

Note that, once we assume that the cell has 'decided' which are the α_i , $x_{i,i+1}$ and ϵ_i values, the process that follows is simply sorting, requiring $O(n)$ steps (possibly fewer than n if some of the recombination events take place simultaneously).

However, this 'decision' process, the details of which are still unknown, amounts to finding the correct 'path' linking the pieces of protein-coding regions in the correct order. Indeed, the occurrence of $\alpha_i x_{i,i+1}$ and $x_{i,i+1} \alpha_{i+1}$ in the micronuclear sequence provides the link between α_i and α_{i+1} , to indicate that they belong next to each other in the macronuclear sequence. The junction sequences $x_{i,i+1}$ thus serve the role of the 'edge' sequences in Adleman's graph.

A computational difficulty is the presence of multiple copies of the sequences $x_{i,i+1}$ (Table 1) which may direct the formation of incorrect 'paths'. Indeed, throughout the genome, such simple sequences may be present in extremely high redundancy. Some of the $x_{i,i+1}$ even overlap with each other. For example, in Table 1, $x_{24,25} = \text{GAGAGATAGA}$ contains $x_{1,2} = \text{AGATA}$ as a subsequence.

The search for the proper junction sequences thus amounts to finding the correct 'path' and is potentially the most costly part of the computation. Production of incorrect paths will not necessarily lead to the production of incorrect proteins unless the path sequences start and end with the correct telomere addition sites (τ_s and τ_e), since these ensure survival of the genes in the macronucleus.

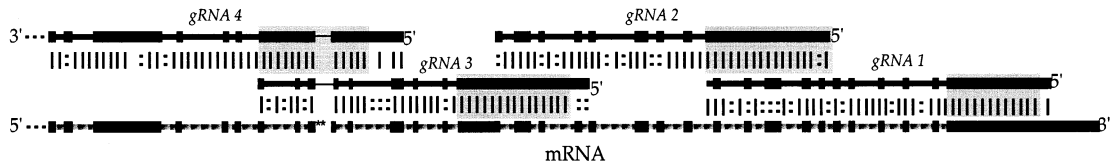


Fig. 10. Editing by four overlapping gRNAs. Thick lines in the mRNA are encoded in the mitochondrial DNA. Thin shaded lines are inserted U; the two asterisks are deleted U (Maslov and Simpson, 1992). Thin lines in the gRNAs are guide nucleotides (A or G) that pair with inserted U. Vertical lines indicate Watson–Crick base pairs; colons indicate G:U wobble base pairs, illustrating formation of well-paired ‘anchors’ between the 5' ends of gRNAs and the corresponding region of the mRNA.

introduced by fixed changes in the number and position of U's inserted or deleted by editing and then compensated by editing at another site that restores the reading frame (Fig. 8). Editing therefore allows the production of combinatorially diverse protein products from a single gene, either within an individual cell (Sommer et al., 1991) or over evolutionary time (Landweber and Gilbert, 1993). The surprising result that RNA editing provides an additional level of sequence variation, rather than a faithful ‘editing’ or correcting mechanism, underscores the importance of the question of why is it still used by some organisms to generate a sequence that encodes a single conserved protein (Landweber and Gilbert, 1994).

The genetic information for editing is stored in the form of ‘guide’ RNA molecules (gRNAs), very small (50–70 nt) transcripts that mediate editing by base-pairing with specific regions of the edited transcript, exploiting G:U base-pairs in RNA (Blum et al., 1990). Each gRNA contains the sequence information to edit approximately 30 nt of edited RNA (Landweber et al., 1993) and pairs more efficiently with the final product than with the pre-edited substrate. For every inserted U in the messenger RNA sequence, there is a corresponding A or G in the gRNA which pairs with the fully edited product (Fig. 9). Complete editing proceeds 3' to 5' on the mRNA and requires a full set of overlapping guide RNAs. Editing by each guide RNA creates an anchor sequence for binding the next guide RNA (Fig. 10, Blum et al., 1990; Maslov and Simpson, 1992) leading to an ordered cascade of insertion and deletion editing events. RNA editing is thus a cellular process which uses RNA sequences as

guides to convert seemingly disordered RNA sequences into a final messenger RNA molecule: a truly RNA-based computer. Together, the stunning acrobatics of DNA, such as scrambling or editing, give proof to the versatility of nucleic acids and their potential use in solving computational problems that occur in biological systems.

Acknowledgements

The authors thank Richard Lipton, David Prescott and Ed Curtis for discussion and Hans Lipps for DNA. L.F.L. is a Burroughs-Wellcome Fund New Investigator in Molecular Parasitology.

References

- Adleman, L.M., 1994. Molecular computation of solutions to combinatorial problems. *Science* 266, 1021–1024.
- Bartel, D.P., Szostak, J.W., 1993. Isolation of new ribozymes from a large pool of random sequences. *Science* 261, 1411–1418.
- Beaver, D., 1996. A universal molecular computer. In: Lipton, R.J., Baum, E.B. (Eds.), *DNA Based Computers*. DIMACS Series in Discrete Mathematics and Theoretical Computer Science, vol. 27. AMS, pp. 29–36.
- Blum, B., Bakalara, N., Simpson, L., 1990. A model for RNA editing in kinetoplast mitochondria: ‘Guide’ RNA molecules transcribed from maxicircle DNA provide the edited information. *Cell* 60, 189–198.
- Csuhaj-Varju, E., Freund, R., Kari, L., Paun, G., 1996. DNA computing based on splicing: universality results. In: Hunter, L., Klein, T. (Eds.), *Proceedings of 1st Pacific Symposium on Biocomputing*. World Scientific, Singapore, pp. 179–190.
- Dubois, M., Prescott, D.M., 1995. Scrambling of the actin I gene in two *Oxytricha* species. *Proc. Natl. Acad. Sci. USA* 92, 3888–3892.

- Feagin, J.E., Abraham, J.M., Stuart, K., 1988. Extensive editing of the cytochrome c oxidase III transcript in *Trypanosoma brucei*. *Cell* 53, 413–422.
- Head, T., 1987. Formal language theory and DNA: an analysis of the generative capacity of specific recombinant behaviors. *Bull. Math. Biol.* 49, 737–759.
- Hoffman, D.C., Prescott, D.M., 1997. Evolution of internal eliminated segments and scrambling in the micronuclear gene encoding DNA polymerase α in two *Oxytricha* species. *Nucl. Acids Res.* 25, 1883–1889.
- Kari, L., Thierrin, G., 1996. Contextual insertion/deletions and computability. *Inform. Comput.* 131, 47–61.
- Klobutcher, L.A., Turner, L.R., LaPlante, J., 1993. Circular forms of developmentally excised DNA in *Euplotes crassus* have a heteroduplex junction. *Genes Dev.* 7, 84–94.
- Landweber, L.F., 1997. RNA based computing. DIMACS Technical Report 97-83.
- Landweber, L.F., Gilbert, W., 1993. RNA editing as a novel source of genetic variation. *Nature* 363, 179–182.
- Landweber, L.F., Gilbert, W., 1994. Phylogenetic analysis of RNA editing: a primitive genetic phenomenon. *Proc. Natl. Acad. Sci. USA* 91, 918–921.
- Landweber, L.F., Fiks, A.G., Gilbert, W., 1993. The boundaries of partially edited cytochrome c oxidase III transcripts are not conserved in kinetoplasts: implications for the guide RNA model of editing. *Proc. Natl. Acad. Sci. USA* 90, 9242–9246.
- Lynn, D.H., Corliss, J.O., 1991. Ciliophora. In: *Microscopic Anatomy of Invertebrates*. Volume 1. Protozoa. Wiley-Liss, New York, pp. 333–467.
- Maslov, D.A., Simpson, L., 1992. The polarity of editing within a multiple gRNA-mediated domain is due to formation of anchors for upstream gRNAs by downstream editing. *Cell* 70, 459–467.
- Mitcham, J.L., Lynn, A.J., Prescott, D.M., 1992. Analysis of a scrambled gene: the gene encoding α -telomere-binding protein in *Oxytricha nova*. *Genes Dev.* 6, 788–800.
- Paun, G., 1995. On the power of the splicing operation. *Int. J. Comp. Math.* 59, 27–35.
- Pixton, D., 1995. Linear and circular splicing systems. In: *Proceedings of the First International Symposium on Intelligence in Neural and Biological Systems*. IEEE Computer Society Press, Los Alamos, pp. 181–188.
- Prescott, D.M., Dubois, M.L., 1996. Internal eliminated segments (IESs) of oxytrichidae. *J. Euk. Microbiol.* 43, 432–441.
- Siromoney, R., Subramanian, K.G., Rajkumar Dare, V., 1992. Circular DNA and splicing systems. In: *Parallel Image Analysis*. Lecture Notes in Computer Science 654. Springer, Berlin, pp. 260–273.
- Smith, W.D., 1996. DNA computers in vitro and vivo. DIMACS Series in Discrete Mathematics and Theoretical Computer Science 27, 121–185.
- Sommer, B., Kohler, M., Sprengel, R., Seeburg, P.H., 1991. RNA editing in brain controls a determinant of ion flow in glutamate-gated channels. *Cell* 67, 11–19.
- Sullenger, B.A., Cech, T.R., 1994. Ribozyme-mediated repair of defective mRNA by targeted, trans-splicing. *Nature* 371, 619–622.
- Tausta, S.L., Klobutcher, L.A., 1989. Detection of circular forms of eliminated DNA during macronuclear development in *E. crassus*. *Cell* 59, 1019–1026.
- Wen, J., Maercker, C., Lipps, H.J., 1996. Sequential excision of internal eliminated DNA sequences in the differentiating macronucleus of the hypotrichous ciliate *Stylonychia lemnae*. *Nucl. Acids Res.* 24, 4415–4419.