

## Genetic Fusion

Takashi Ikegami<sup>(a)</sup>

*The Graduate School of Science and Technology, Kobe University, Rokkodai, Nada-ku, Kobe 657, Japan*

Kunihiko Kaneko<sup>(b)</sup>

*Institute of Physics, College of Arts and Sciences, University of Tokyo, Komaba, Meguro-ku, Tokyo 153, Japan*

(Received 25 June 1990)

Genetic fusion is introduced as a model for evolution. In fusion two genomes are combined to generate a longer genome. Representing each species by a binary genetic sequence, we introduce a fitness function on the bit sequence. As the evolutionary dynamics, we incorporate mutation, genetic fusion, and reproduction in proportion to fitness. It is found that genetic fusion leads to the appearance of module-type sequences and duplicated genes. The time necessary to find a sequence with larger fitness is largely reduced by the inclusion of genetic fusion, which suggests the application of our algorithm to optimization problems.

PACS numbers: 87.10.+e

Theoretical studies on mechanisms of evolutionary processes are of importance not only for theoretical biology but also from the viewpoint of information processing. Fitness in evolution is believed to depend on the genetic sequence. The dependence of this fitness on sequences can be complicated, and has been referred to as a rugged landscape.<sup>1</sup> A spin-glass model<sup>2</sup> provides a simple example which generates this class of landscape. A typical model is the Sherrington-Kirkpatrick (SK) model.<sup>2</sup> In the model, a spin  $S_i$  ( $=1$  or  $-1$ ) is assigned to a site  $i$ . Energy is defined on a binary sequence  $[S_i]$ , given by  $E = \sum_{i,j} J_{i,j} S_i S_j$ , where a spin  $S_i$  interacts with all the other spins  $S_j$ . The coupling  $J_{i,j}$  is set at a fixed random value distributed over positive and negative values. It is known that the model provides many metastable states and a rugged landscape.<sup>2</sup> The application of this class of model to evolution has already been discussed by Anderson,<sup>3</sup> where the SK model is applied to prebiotic evolution.

Another important mechanism in evolution is sexual recombination of sequences by mating. This mechanism has frequently been used in the so-called genetic algorithm.<sup>4</sup> Here two genetic sequences are spliced and joined; for example, parents  $[\sigma_1, \sigma_2, \sigma_3, \sigma_4, \dots]$  and  $[\tau_1, \tau_2, \tau_3, \tau_4, \dots]$  lead to the offsprings  $[\sigma_1, \sigma_2, \tau_3, \tau_4, \dots]$  and  $[\tau_1, \tau_2, \sigma_3, \sigma_4, \dots]$ . This recombination is especially useful to escape from a local minimum in configuration space and has been applied to optimization problems.

In the present Letter, we focus on another mechanism of evolution; genetic incorporation involving different species, which also changes the length of the genetic sequence itself. This mechanism may have been especially important in the early stages of evolution and it is probably still being used in present-day organisms.<sup>5</sup>

There is much evidence to confirm this genetic incorporation. Ohno has put forward a theory of gene duplication.<sup>6</sup> In his theory, some parts of a genetic sequence may incur duplication error during the process of repro-

duction, as

$$[\sigma_1, \sigma_2, \sigma_3, \sigma_4, \sigma_5, \sigma_6, \dots] \\ \rightarrow [\sigma_1, \sigma_2, \sigma_3, \sigma_1, \sigma_2, \sigma_3, \sigma_4, \sigma_5, \sigma_6, \dots],$$

where  $\sigma_i$  represents a gene at a site  $i$ . The duplication theory is confirmed in real DNA. Multiple repetition of short DNA sequences is commonly observed.<sup>5</sup>

Besides the above gene duplication, a biological system has transposable genetic elements. They are self-reproductive elements such as phages, plasmids, and transposons. For example, bacteria can evolve to be antibiotic resistant with the use of  $R$  plasmids. Plasmids can insert a particular sequence into other plasmids. Even in eukaryotes, the existence of a cell symbiotic partner is known to have escalated the tempo of evolution, as has been discussed in the cell symbiosis of mitochondria.<sup>7</sup>

Then general questions arise: Why are the above genetic incorporation algorithms adopted in real evolution? Can we construct a simple model for these processes? If so, does the model give rise to practical merits over simple genetic algorithms? In the present Letter, we give a partial answer to these questions, by introducing an abstract model of the above cell symbiosis and gene duplication. The process in our model is called "genetic fusion," here borrowing from molecular biology.

First, we represent each species by a bit string of genes  $\sigma_i$ , which take the value 0 or 1. The length of this sequence is not fixed. Species are assumed to evolve in a fixed environment. For the sake of simplicity, we adopt a spin-glass-type energy for our fitness function. Instead of the SK model, we use the following translational-invariant version of it:

$$E = \sum_{i,j} J_{|i-j|} S_i S_j, \quad (1)$$

where  $S_i = 2\sigma_i - 1$ , and  $J_m$  is set at a fixed random value distributed over  $[-1, 1]$ .<sup>8</sup>

Introducing an ensemble of species, we take the fol-



TABLE I. Bit sequence of module and remaining species. Examples of module species are depicted. Each of these species is expected to be in a local minimum. Each species ( $l-m$ ) is specified by a decimal code  $l$  and its length  $m$ . The decimal code is obtained by the conversion of its binary sequence. Examples of remaining species after 30 time steps are also exhibited. Note that they are composed of the modules in tandem. The symbol  $[A']$  denotes a sequence which is different from  $[A]$  by one bit.

	Species	Sequence	Energy
Modules			
$A$	4923-13	1001100111011	-13.796
$B$	4095-12	111111111111	-13.066
$C$	4064-13	0111111100000	-7.475
Remaining species			
	645100347-30	$[A']$ $[1011]$ $[A]$	-55.714
	577139515-30	$[1000]$ $[A']$ $[A]$	-46.592
	16777215-29	$[00000]$ $[B]$ $[B]$	-43.532
	40334134-26	$[A]$ $[A]$	-36.634
	16777215-24	$[B]$ $[B]$	-31.838
	33550304-25	$[B]$ $[C]$	-26.564

species easily leads to species of better fitness. For this condition, the energy of a module species should not be too small. In our simulation, module species have small fitness and thus have small population size.<sup>11</sup>

There can simultaneously be several different module species. If this is the case, species are categorized into phylogenetic groups according to whether they share a common module or not. Examples of such module sequences in our simulation are listed in Table I. All of the persistent species originate in one of these module species.

It often happens that a module sequence  $M_1 = [\sigma_1, \dots, \sigma_k]$  combines with itself. This leads to gene duplication, as  $M_1 \rightarrow [M_1, M_1]$ . The chain-duplication process also often occurs as  $[M_1] \rightarrow [M_1, M_1] \rightarrow [M_1, M_1, M_1]$  or as  $[M_1] \rightarrow [M_1, M_1] \rightarrow [M_1, M_1, M_1, M_1]$ . This class of gene duplication is frequently seen in biological evolution, as is stressed by Ohno.<sup>6</sup>

In our simulation, the evolutionary process includes duplication, fusion, and mutations. Thus a genetic sequence does not have a complete repetitive structure, but consists of fragments of repetitive parts. This is true in our simulation (see Table I) and in real data.<sup>6</sup>

If the mutation rate is small, few (often a single) module species are allowed in our simulation. Possible fusions consist of those between the same module species. The identical module sequence is repetitively made use of in the fusion process. In other words, the genetic duplication is frequently seen in a low-mutation regime. If the mutation rate is higher, fusion with many other species can happen. The obtained species consist of a combination of a variety of sequences.

Our genetic fusion can provide an algorithm for the search of a lower energy (or larger fitness) in a rugged landscape, faster than the conventional genetic algorithm or simulated annealing.

Let us compare our genetic fusion and a mere mutation process. The simple mutation process corresponds to a traditional Monte Carlo method.

In Fig. 3, we have plotted the minimum energy of the species versus CPU time for simulations with and without genetic fusion (both have the same mutation rates). In the simulation of mutation only, we have started from ten randomly chosen sample species with bit length of 30. The advantage of our genetic fusion is clearly seen.<sup>12</sup>

The advantage comes from the following two features.

(i) Our fusion process tries to find a minimum by combining small fragments. Thus it is effective when an optimal solution is well approximated by a combination of smaller parts and modules (e.g., a partial solution).

(ii) With mere mutation, our system is easily trapped by a local minimum. It may require a long time to hit a global minimum. The present genetic fusion includes a global jump in the configuration space. Our system can skip out of the local minimum.

The second merit is also present in the conventional genetic algorithm, where sexual recombination leads to a global jump. Owing to the first merit, our fusion clearly shows predominance over the conventional algorithm.<sup>4</sup>

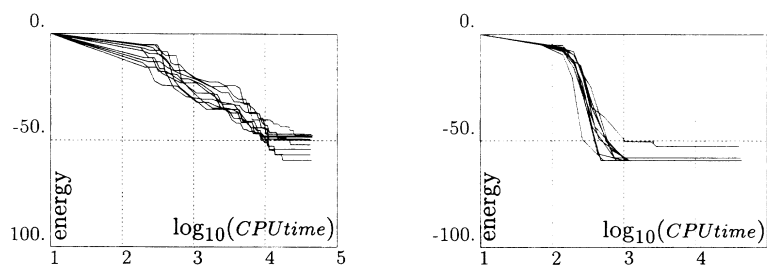


FIG. 3. The lowest-energy values at each CPU time are overlaid for ten different samples. With fusion (right-hand panel) lower energy is attained faster than without fusion (left-hand panel). Both processes are simulated with a mutation rate of 0.1 with  $\beta=0.1$  and  $g=0.1$ .

The fusion process is powerful, if the optimal solution can be constructed as a combination of partial solutions. It is promising to apply our approach to other optimization problems. We believe that it is better than the mere mutational process in general.

The evolutionary process is often believed to be organized in a genealogical tree. This is true, if the process consists only of mutations and sexual recombination within species. If genetic fusion is a principal mechanism in evolution, the genealogy is organized in a rhizomelike structure.

The fusion process introduces deviation from a tree-like structure, as is shown in Fig. 2. A descendant can have more than one ancestor. The mixture of ancestors gives rhizome complexity. For example, the species 645337915-30 consists of the ancestors 10 and 00 and two occurrences of 1001100111011 (the module species). The rhizome complexity can be defined through the frequency of ancestor sequences, which will be useful to discuss the deviation from ultrametricity.<sup>2</sup>

In the present Letter, we have introduced a novel dynamical model for evolution. We have discussed the appearance of the module sequence, gene duplication, and a possible practical merit of our dynamics in optimization problems.

For future studies, it may be important to include splitting of a sequence into two subsequences.<sup>13</sup> If we include this process, some "good" parts in a long sequence can be added to other creatures, as is often used in genetic technology.

In our present model fusion always occurs if the energy condition is satisfied. Since fusion itself is an error in the maintenance of genes, it must be more realistic to introduce a small rate of fusion error. With this extension, the evolutionary process exhibits much clearer stepwise changes than in Fig. 1: rapid-change eras and quasistationary phases, as has been discussed in Ref. 14 as punctuated equilibrium theory.<sup>15</sup>

Evolution provides a novel viewpoint for biological information processing such as autocatalytic,<sup>16</sup> immune, and neural networks.<sup>17</sup> The idea of genetic fusion will be useful to these problems.<sup>18</sup>

The authors would like to thank Y. Oono for critical comments, P. Davis for critical reading of the manuscript, and D. Farmer, C. Langton, and S. Rasmussen for useful discussions. This work was partially supported by Grant-in-Aids for Scientific Research from the Ministry of Education, Science and Culture of Japan. One of

the authors (T.I.) is indebted to the Japan Society for the Promotion of Science for financial support.

<sup>(a)</sup>Electronic address: ikegami@gradis.scitec.kobe-u.ac.jp.

<sup>(b)</sup>Electronic address: d34205@tansei.cc.u-tokyo.ac.jp.

<sup>1</sup>See, for the evolution in a rugged landscape, S. A. Kaufmann and S. Levin, *J. Theor. Biol.* **128**, 11 (1987).

<sup>2</sup>*Spin Glass Theory and Beyond*, edited by M. Mezard, G. Parisi, and M. A. Virasoro (World Scientific, Singapore, 1988).

<sup>3</sup>P. W. Anderson, *Proc. Natl. Acad. Sci. U.S.A.* **80**, 3386 (1983).

<sup>4</sup>J. Holland, in *Escaping Brittleness in Machine Learning II*, edited by R. S. Mishalski, J. Carbonell, and T. M. Mitchell (Kaufman, Ohio, 1986).

<sup>5</sup>See, e.g., J. M. Smith, *Evolutionary Genetics* (Oxford, Univ. Press, Oxford, 1989), Chap. 11.

<sup>6</sup>S. Ohno, *Evolution by Gene Duplication* (Springer-Verlag, Berlin, 1970).

<sup>7</sup>L. Margulis, *Symbiosis in Cell Evolution* (Freeman, San Francisco, 1981).

<sup>8</sup>We choose this model, not because of direct biological significance, but because it seems to have many local minima in a rugged landscape with translational invariance of gene sites. Detailed studies of the landscape of our model are left for future work.

<sup>9</sup>In the model, nothing stops the increase of bit length. We put a limit here externally at 30.

<sup>10</sup>Here we use the first algorithm for the reproduction, except for Fig. 3. Results, however, are essentially the same for either algorithm.

<sup>11</sup>Although the second condition itself can depend on the distribution of other species, it is rather robust with respect to the change of initial distribution in our simulation.

<sup>12</sup>After the fusion process hits the maximal length (30), mere mutation can occur to decrease the energy. The large plateau in CPU time comes from this saturation.

<sup>13</sup>This introduction of rearrangement is important to discuss the evolution of chromosomes (partitioned units of genetic codes in eukaryotes), since they split or join through evolution.

<sup>14</sup>N. Eldredge and S. J. Gould, in *Models in Paleobiology*, edited by T. J. M. Schopf (Freeman, San Francisco, 1972).

<sup>15</sup>See also K. Lindgren (to be published), for the application of gene duplication to evolution in the iterated prisoner's dilemma.

<sup>16</sup>J. D. Farmer, S. A. Kauffman, and N. H. Packard, *Physica (Amsterdam)* **22D**, 187 (1986).

<sup>17</sup>G. M. Edelman, *Neural Darwinism* (Basic Books, New York, 1987).

<sup>18</sup>T. Ikegami, Ph.D. thesis, University of Tokyo, 1989 (unpublished).