



# A trade-off relationship between energetic cost and entropic cost for *in vitro* evolution

Takuyo Aita

Graduate School of Science and Engineering, Saitama University, Saitama 338-8570, Japan

## ARTICLE INFO

### Article history:

Received 3 June 2010

Received in revised form 8 July 2010

Accepted 9 July 2010

### Keywords:

Artificial evolution

Fitness landscape

Optimization

Information

Search algorithm

Sequence space

## ABSTRACT

In this paper, we consider two complementary cost functions for the landscape exploring processes to obtain the global optimum sequence through *in vitro* evolution protocol: one is the entropic cost  $C_{\text{etp}}$ , which is based on the deviation from the uniformity of a mutant distribution in sequence space, and the other is the energetic cost  $C_{\text{eng}}$ , which is based on the total number of sequences to be generated and evaluated. Based on a prior knowledge about the structure of a given fitness landscapes, the conductor of the experiment can think up the efficient search algorithm (ESA), which requires the minimum number of points (=sequences) to be searched up to the global optimum. For five typical fitness landscapes, we considered their respective (putative) ESA,  $C_{\text{etp}}^*$  and  $C_{\text{eng}}^*$  based on the ESA. As a result, we found a trade-off relationship between  $C_{\text{etp}}^*$  and  $C_{\text{eng}}^*$  for every case, that is,  $C_{\text{eng}}^* + C_{\text{etp}}^*$  is approximately equal to the logarithm of the volume of the sequence space.  $C_{\text{etp}}^*$  and  $C_{\text{eng}}^*$  are interpreted in terms of the information-theoretic concepts.

© 2010 Elsevier Ireland Ltd. All rights reserved.

## 1. Introduction

*In vitro* evolution is comprehended as an optimization process on a given fitness landscape, where “fitness landscape” is defined by the mapping from sequences (=genotype) to the corresponding fitness values (Wright, 1932; Maynard-Smith, 1970). Recent development of this field has been reviewed in some books (e.g. Arnold, 2000) and review papers (e.g. Romero and Arnold, 2009). Ideally, it is desirable to find the global optimum sequence at the summit. The efficient search algorithm on fitness landscapes is deeply dependent on our knowledge about statistical orders of the structures of fitness landscapes (Wolpert and Macready, 1997), where “the efficient search algorithm (which is abbreviated to ESA)” is defined as a particular strategy that requires the minimum number of points (=sequences) to be searched up to the global optimum. For example, we can say that the fitness landscapes with high order, such as Fujiyama landscape, have a lot of information about their structures. That is, we know the mutational effect is additive for Fujiyama landscape. Therefore, we should think up the efficient search algorithm based on the available information.

Generally, the search process is conducted by two processes: mutagenesis and selection. In the mutagenesis process, controlling the distribution of mutant sequences in sequence space requires a “cost” that originates from the negative entropy. Therefore, we

introduced the entropic cost:  $C_{\text{etp}}$ . If a site-directed mutagenesis is applied and the mutant sequences are distributed in a local area in sequence space,  $C_{\text{etp}}$  is large due to strong control of mutations. This control is done by our intelligence or replication fidelity of replication enzyme (e.g. DNA polymerase) (Ishii et al., 1989; Cady and Qian, 2009). If mutant sequences are distributed randomly in sequence space,  $C_{\text{etp}}$  is small due to no control of mutations. On the other hand, the energy consumed through all the processes is proportional to the total number  $N$  of sequences to be generated and evaluated. Therefore, we introduced the energetic cost defined by  $C_{\text{eng}} = \ln N$ .

As a whole, there seems a trade-off relationship between energetic cost  $C_{\text{eng}}$  and entropic cost  $C_{\text{etp}}$  based on the corresponding ESA. For example, it is easy to find the global optimum sequence of Fujiyama landscape, because one has only to identify the fittest letter for each site by positional scanning method. In this case,  $C_{\text{eng}}$  is small due to small number of  $N$ , while  $C_{\text{etp}}$  is large due to strong control of mutations. On the other hand, as for random rugged landscape, an exhaustive search over the sequence space is necessary by performing random synthesis of all possible sequences. In this case,  $C_{\text{eng}}$  is large due to vast numbers of  $N$ , while  $C_{\text{etp}}$  is small due to no control of mutations.

In this paper, focusing on five typical fitness landscapes, we considered their respective (putative) ESA,  $C_{\text{etp}}$  and  $C_{\text{eng}}$  based on the ESA. Our aim is to describe quantitatively a trade-off relationship between  $C_{\text{eng}}$  and  $C_{\text{etp}}$ , and to interpret them in terms of the information-theoretic concepts.

E-mail address: [taita@mail.saitama-u.ac.jp](mailto:taita@mail.saitama-u.ac.jp).

## 2. Definition

### 2.1. Deviation from the uniformity of a mutant distribution in sequence space

We consider all conceivable sequences of  $\lambda^v$ , where  $v$  is the number of all sites and  $\lambda$  is the number of available letters for each site. The  $\lambda$  is much larger than one ( $\lambda \gg 1$ ). Then, an arbitrary sequence  $s$  ( $s = 1, 2, \dots, \lambda^v$ ) is mapped into the corresponding point in the  $\lambda$ -valued  $v$ -dimensional sequence space.

Consider that sequences with an infinite population are distributed in sequence space according to the underlying probability distribution  $Q(s)$  ( $\sum_{s=1}^{\lambda^v} Q(s) = 1$ ), where  $Q(s)$  is a probability of being occupied by a sequence  $s$  in sequence space. Let  $d$  be the Hamming distance between two arbitrary sequences among the population, and let  $P(d)$  ( $\sum_{d=0}^v P(d) = 1$ ) be the probability distribution of  $d$  over all pairs of them. That is, the probability distribution  $Q(s)$  is converted to the probability distribution  $P(d)$  by

$$P(d) = \sum_{s=1}^{\lambda^v} \sum_{s'=1}^{\lambda^v} \delta(d(s, s'), d) Q(s) Q(s'), \quad (1)$$

where  $d(s, s')$  is the Hamming distance between sequences  $s$  and  $s'$ , and  $\delta(x, x_0)$  is the Kronecker's delta defined by

$$\delta(x, x_0) \equiv \begin{cases} 1, & \text{if } x = x_0 \\ 0, & \text{if } x \neq x_0. \end{cases}$$

Let  $B(d)$  be a back ground distribution where sequences are distributed randomly over the sequence space. As the deviation from the uniformity of a mutant distribution in the sequence space, we introduce the following relative entropy  $D$ :

$$D \equiv \sum_{d=0}^v P(d) \ln \frac{P(d)}{B(d)}, \quad (2)$$

$$B(d) = \binom{v}{d} \left(1 - \frac{1}{\lambda}\right)^d \left(\frac{1}{\lambda}\right)^{(v-d)}. \quad (3)$$

Eq. (3) is derived in Appendix A, based on the “profile” of given sequences (Gribskov et al., 1987).

### 2.2. Energetic cost and entropic cost for in vitro evolution

In this paper, we consider landscape exploring processes to obtain the global optimum sequence through *in vitro* evolution protocol. *In vitro* evolution is conducted by iterating the evolution cycle, which consists of generating mutant sequences of  $N_t$ , measuring fitness values of them and selecting the fittest one from among them. The  $N_t$  sequences consist of heterogeneous sequences. The step number of the iteration process is denoted by  $t$  ( $t = 1, 2, 3, \dots, t_e$ ). Suppose that the following two different costs are generated in the *in vitro* evolution: one is the “energetic cost”  $C_{\text{eng}}$ , which is based on the total number of sequences to be generated and evaluated, the other is the “entropic cost”  $C_{\text{etp}}$ , which is based on the deviation from the uniformity of a mutant distribution in sequence space.

Consider that, in the  $t$  th step, mutant sequences of  $N_t$  are generated according to the underlying probability distribution  $Q_t(s)$  ( $\sum_{s=1}^{\lambda^v} Q_t(s) = 1$ ), where  $Q_t(s)$  is a probability of being occupied by a sequence  $s$  in sequence space. Note that the shape of  $Q_t(s)$  is set up by the “conductor”, who sets up the experimental algorithm and implements it. The probability distribution  $Q_t(s)$  is converted to the probability distribution of the Hamming distance  $d$  between arbitrary two sequences,  $P_t(d)$  ( $\sum_{d=0}^v P_t(d) = 1$ ), by using

Eq. (1). Therefore, we apply the measure  $D$  defined in Eq. (2) to the probability distribution for the  $t$  th step,  $P_t(d)$ :

$$D_t \equiv \sum_{d=0}^v P_t(d) \ln \frac{P_t(d)}{B(d)}. \quad (4)$$

We note the following two similar but different cases: one is the case where a single particular sequence  $s^*$  is generated ( $N_t = 1$ ), and the other is the case where a single arbitrary sequence  $s$  is generated ( $N_t = 1$ ). For the former case,  $D_t = v \ln \lambda$ , because  $Q_t(s)$  is given by the Kronecker's delta,  $Q_t(s) = \delta(s, s^*)$ , and then  $P_t(d) = \delta(d, 0)$ . For the latter case,  $D_t = 0$ , because  $Q_t(s)$  is the uniform distribution,  $Q_t(s) = 1/\lambda^v$ , and then  $P_t(d) = B(d)$ . Then, we define the “entropic cost” through the whole evolution process by

$$C_{\text{etp}} \equiv \langle D \rangle \equiv \frac{1}{t_e} \sum_{t=1}^{t_e} D_t. \quad (5)$$

In this paper, we use  $\langle X \rangle$  as the arithmetic mean of a  $t$ -dependent quantity  $X_t$  over the whole process from  $t = 1$  to  $t = t_e$ .

The energy consumed through all the processes is proportional to the total number of sequences to be generated and evaluated. Let  $N_{\text{tot}}$  be the number of points (=sequences) to be searched up to reaching the global optimum sequence:  $N_{\text{tot}} = \sum_{t=1}^{t_e} N_t$ . We define the “energetic cost” by

$$C_{\text{eng}} \equiv \ln N_{\text{tot}}. \quad (6)$$

The total cost is defined by

$$C_{\text{tot}} \equiv C_{\text{eng}} + C_{\text{etp}}. \quad (7)$$

Both  $C_{\text{etp}}$  and  $C_{\text{eng}}$  range from 0 to  $v \ln \lambda$ .

### 2.3. Typical fitness landscape and efficient search algorithm

Here, we define a “typical fitness landscape” as a set of all possible landscapes described by a given mathematical model with several parameters (e.g. Fujiyama landscape, NK landscape). A typical fitness landscape has its characteristic order and then the conductor knows it as a prior knowledge about the landscape. The conductor's job is to obtain the global optimum sequence of the given typical landscape by implementing an “efficient search algorithm”, where the efficient search algorithm (which is abbreviated to ESA) is, in this paper, defined as a particular algorithm that requires the minimum number of points (=sequences) to be searched up to reaching the global optimum. This minimum number is denoted by  $N_{\text{tot}}^*$ . Based on the prior knowledge about the given landscape, the conductor can think up the corresponding ESA.  $C_{\text{etp}}$ ,  $C_{\text{eng}}$ , and  $C_{\text{tot}}$  based on the corresponding ESA are denoted by  $C_{\text{etp}}^*$ ,  $C_{\text{eng}}^*$  ( $= \ln N_{\text{tot}}^*$ ), and  $C_{\text{tot}}^*$ , respectively.

The ESA gives the minimum value of  $C_{\text{eng}}$ . We can define another search algorithm that gives the minimum value of  $C_{\text{tot}}$ . However, determining this algorithm for a given landscape is so difficult. Therefore, we leave it and focus on the ESA in this paper.

## 3. Application to typical fitness landscapes

In this section, as for five typical fitness landscapes, we consider their respective ESA,  $C_{\text{etp}}^*$ ,  $C_{\text{eng}}^*$  and  $C_{\text{tot}}^*$  based on the ESA. In some cases, although the conductor can think up a putative ESA for a given landscape, the putative ESA is however difficult to be proven true rigorously. Therefore, the ESAs we described below are putative ones. The results are compiled in Table 1.

### 3.1. Additive fitness landscape (Fujiyama landscape)

In this case, each letter at every site in a given sequence contributes to the fitness independently and additively, that is,

Download English Version:

<https://daneshyari.com/en/article/2076301>

Download Persian Version:

<https://daneshyari.com/article/2076301>

[Daneshyari.com](https://daneshyari.com)