

# On the Rate of Molecular Evolution\*

MOTOO KIMURA and TOMOKO OHTA

National Institute of Genetics, Mishima, 411, Japan

Received September 15, 1970

*Summary.* There are at least two outstanding features that characterize the rate of evolution at the molecular level as compared with that at the phenotypic level. They are; (1) remarkable uniformity for each molecule, and (2) very high overall rate when extrapolated to the whole DNA content.

The population dynamics for the rate of mutant substitution was developed, and it was shown that if mutant substitutions in the population are carried out mainly by natural selection, the rate of substitution is given by  $k = 4N_e s_1 v$ , where  $N_e$  is the effective population number,  $s_1$  is the selective advantage of the mutants, and  $v$  is the mutation rate per gamete for such advantageous mutants (assuming that  $4N_e s_1 \gg 1$ ). On the other hand, if the substitutions are mainly carried out by random fixation of selectively neutral or nearly neutral mutants, we have  $k = v$ , where  $v$  is the mutation rate per gamete for such mutants.

Reasons were presented for the view that evolutionary change of amino acids in proteins has been mainly caused by random fixation of neutral mutants rather than by natural selection.

It was concluded that if this view is correct, we should expect that genes of "living fossils" have undergone almost as many DNA base replacements as the corresponding genes of more rapidly evolving species.

*Key-Words:* Molecular Evolutionary Rate — Population Genetics Theory.

## Introduction

The rate at which new mutant genes are incorporated into the species in the course of evolution should characterize the speed of evolution more unambiguously than any other criterion.

Yet, there had been no ways of measuring this until comparative studies of amino acid sequence among homologous proteins became available, as was so clearly pointed out by Crow (1969).

The pioneering work by Zuckerkandl and Pauling (1965) on the evolution of informational macromolecules, especially their systematic and quantitative treatment of evolutionary divergence of proteins has made it clear that such measurement is indeed possible.

In the present paper, we intend to discuss some problems relating to the rate of molecular evolution from the standpoint of population genetics.

---

\* Contribution No. 789 from the National Institute of Genetics, Mishima, Shizuoka-ken 411 Japan. Aided in part by a grant-in-aid from the Ministry of Education, Japan.

## Mutation and Mutant Substitution

First of all we must emphasize the basic difference between "gene mutation" and "gene substitution", or more generally, between "mutation" and "mutant substitution".

The former refers to change of genetic material at the *individual level*, while the latter refers to that at the *population level*. Without keeping these two events conceptually distinct, there can be no meaningful discussions in population genetics. For example, to ascribe the difference of amino acids between homologous proteins simply to "randomly distributed point mutations" while neglecting the population dynamics underlying such a difference only obscures this basic distinction.

Although a large number of mutants arise in each generation within any reasonably large population, the majority of them are lost by chance within a few generations (cf. Fisher, 1930; Kimura and Ohta, 1969b). It is not often realized that this is true not only for deleterious and selectively neutral mutants but also for advantageous mutants unless the advantage is very large. For example, if a mutant has one per cent selective advantage, the chance is only about two per cent that it will eventually spread into the whole population (Haldane, 1927; Fisher, 1930). In the remaining 98% of the cases, it will be lost by chance from the population without being used in evolution. Thus, a vast difference exists between the total number of advantageous mutants that have ever occurred in any species in the course of evolution and the number that have actually been incorporated into the species.

Actually, it is a lucky minority that manage to increase their frequencies and spread in the species reaching the state of fixation. Each such event usually takes a large number of generations.

Throughout this paper, we will use the term "mutant incorporation" or "mutant substitution" to represent an event in which an originally rare mutant increases its frequency and spreads through the population. A word of caution may be appropriate here if we want to avoid a confusion that sometimes occurs when we talk about "the rate of gene substitution". Not infrequently, this is misinterpreted to mean the rate at which an individual mutant or mutant combination increases its frequency within the population, rather than the rate at which mutants are incorporated one by one into the population in the course of time. The former (rate of increase of mutant frequency) can be increased proportionally by increasing the selective advantage. The effect of selection on the latter is more complicated; the selective advantage influences the rate of mutant incorporation through modifying the rate of eventual fixation of the mutants. Moreover, this rate is determined to a large extent by the very rate at which new mutants appear in the population in each generation. If, by increasing the

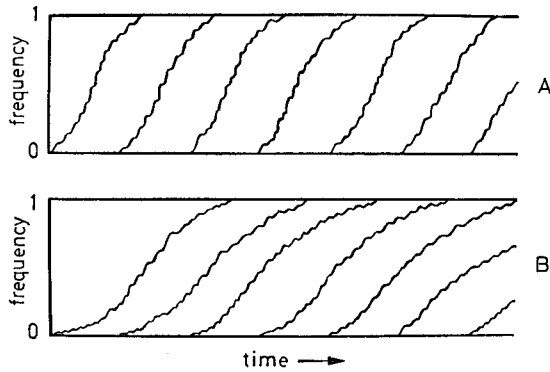


Fig. 1. Diagrams showing the processes of mutant substitution in the population. Of the two cases illustrated, individual mutants increase much more rapidly within a population in A than in B, yet the rate of mutant substitution is the same

selective advantage of mutants, the rate of appearance of such mutants were reduced, there might be no increase in the rate of mutant incorporation. By the term rate of mutant substitution (or incorporation) we therefore mean the long term average of the number of mutants that become fixed in the population (species). In other words, the rate of substitution ( $k$ ) per unit time (year, generation etc.) may be defined by

$$k = \lim_{T \rightarrow \infty} n(T)/T, \quad (1)$$

where  $n(T)$  is the cumulative number of mutants that have been fixed in the population during the time of length  $T$ , assumed to be very long. Note that the length of time involved for each mutant substitution ( $\bar{t}_1$ ) does not come into this rate. We wait long enough that  $T$  is much larger than  $\bar{t}_1$ . Fig. 1 illustrates two cases where the average length of time for each substitution is very different, yet the rate of gene substitution is the same.

### Measurement of the Rate of Amino Acid Substitution

The simplest and surest way of measuring the rate of mutant substitution is through comparison of the amino acid sequence in homologous proteins of related organisms, coupled with paleontological information about the time of their divergence. The extensive compilation of protein sequences recently made by Dayhoff (1969) helps greatly in such estimation.

If we denote by  $n_{aa}$  the total number of amino acid sites in each of the two polypeptide chains compared and if they differ at  $d_{aa}$  sites, then the mean number of substitutions per amino acid site over the whole evolutionary period that separates these two chains may be estimated from

$$K_{aa} = -\log_e(1 - p_d) \approx -2.3 \log_{10}(1 - p_d), \quad (2)$$

where  $p_d = d_{aa}/n_{aa}$  is the fraction of differing sites. In counting the number of differing sites, we exclude deletions and insertions if any exist, since we are here concerned with the rate of evolution due to amino acid substitutions. Actually, mutant substitutions involving such structural changes are much less frequent than those involving amino acid replacement and they will not influence our estimates very much.

On the other hand, in reconstructing phylogeny from amino acid sequences, structural changes, because of their rarity, are particularly informative as to the actual evolutionary history, and should therefore not be excluded.

Formula (2) is originally due to Zuckerkandl and Pauling (1965), and it involves a correction for the undetected occurrence of two or more substitutions at one site, assuming that amino acid substitutions are random. We note that the actual number of nucleotide substitutions per site can often be determined directly by using knowledge of the DNA code, if a number of homologous proteins among related organisms are analysed simultaneously (cf. Jukes, 1966; King and Jukes, 1969). This is one of the places where a computer is of real help (see for example, Fitch and Margoliash, 1970). It appears that the assumption of randomness in the pattern of substitution is valid as a first approximation, if we exclude certain "invariant" sites.

The standard error of the estimate of  $K_{aa}$  is given by

$$\sigma_K = \sqrt{\frac{p_d}{(1-p_d)n_{aa}}} \quad (3)$$

(Kimura, 1969b).

The quantity of particular interest to us is the rate of substitution per amino acid site *per year* and this can be computed by

$$k_{aa} = K_{aa}/(2T), \quad (4)$$

where  $T$  is the number of years that have elapsed since the evolutionary divergence of the two polypeptides from their common ancestor. The factor 2 in the denominator arises from the fact that each polypeptide has evolved independently for  $T$  years from their common ancestor.

As an example, let us compare the  $\alpha$  hemoglobin chains of man and carp. From the sequence of these chains listed in Dayhoff, we find that they have the same amino acids at 72 of the sites and are differentiated by 68 of the sites, if we exclude insertions or deletions that amount to 3 amino acids. Furthermore the number of invariant sites in hemoglobin is small enough to be neglected. Thus, we have  $d_{aa} = 68$  and  $n_{aa} = 140$ . Substituting these in formulae (2) and (3) we obtain  $K_{aa} = 0.665$  and  $\sigma_K = 0.082$ . In order to apply formula (4) to obtain the rate of substitution per year, we must know the length of time,  $T$ , since divergence. Fortunately,

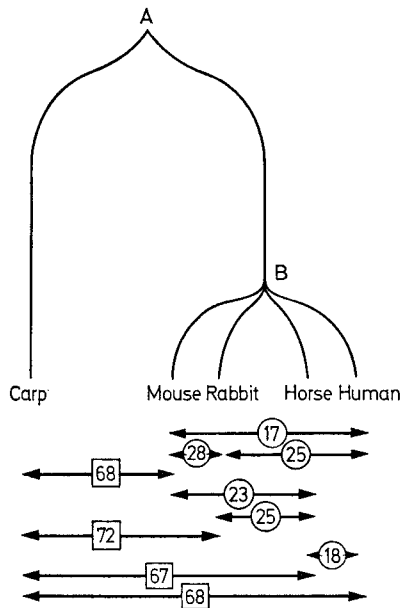


Fig. 2. A phylogenetic tree involving carp and some mammals. Numerals with arrows indicate the numbers of amino acid differences among them with respect to the hemoglobin  $\alpha$  chain

the paleontology of vertebrate evolution, especially the evolution of fishes is well documented (cf. Simpson, Pittendrigh and Tiffany, 1958; Romer, 1968), and coupled with modern methods of isotope dating, we are reasonably sure that the common ancestor of carp and man lived in the Devonian period or little earlier, dating back some 400 million years ago. So we may take  $T = 400 \times 10^6$  (years) in formula (4), giving the rate of substitution,

$$k_{aa} = 8.3 \times 10^{-10}$$

per year per amino acid site, or roughly  $k_{aa} = 10^{-9}$ .

Applying this method to compare the  $\alpha$  hemoglobin chain of carp with the  $\alpha$  chains of man, mouse, rabbit, horse and cattle (see Fig. 2), we obtain  $K_{aa}$  values that agree with each other within the limit of statistical variation despite the fact that the members of the latter group (mammals) differ from each other on the average at some 20 amino acid sites. Furthermore, by comparing  $\alpha$  chains among the mammals, noting that they diverged from their common ancestor some 80 million years ago, we again obtain a rate of substitution per year per site about  $k_{aa} = 10^{-9}$ . Extensive calculations of this sort including  $\beta$  hemoglobin and lamprey globin as well as  $\alpha$  hemoglobin reveal amazing uniformity of the rate of amino acid

substitution among diverse lines of vertebrate evolution, always giving a value of approximately  $10^{-9}$  per year per amino acid site (Kimura, 1969b).

One of the best examples showing uniformity of the rate of substitution comes from the comparison of hemoglobin  $\alpha$  and  $\beta$  chains. By comparing the  $\beta$  chain of man with the  $\alpha$  chains of man and carp, we find that  $K_{aa} = 0.776 \pm 0.092$  for human  $\beta$ -human  $\alpha$  comparison, and  $K_{aa} = 0.807 \pm 0.094$  for human  $\beta$ -carp  $\alpha$  comparison. These two estimates are well within the limit of statistical variation, despite the fact that human  $\alpha$  and carp  $\alpha$  are so greatly differentiated as to give  $K_{aa} = 0.665$  as mentioned already. This shows that the two structural genes corresponding to  $\alpha$  and  $\beta$  chains have diverged from each other independently and to the same extent since their origin by duplication back possibly in Ordovician period (about 450 million years ago). It may be hard to imagine, from the traditional evolutionary point of view, that mutant substitutions at gene loci coding for  $\alpha$  and  $\beta$  hemoglobins have proceeded at the same speed in two separate lines that have evolved independently over almost a half billion years, one leading to present day carp and one leading to man, while the speed of evolution in terms of morphology is so utterly different.

Fig. 3 illustrates similar comparisons for cytochrome c. Note that although we know very little about the times of divergence among various organisms in the figure, we are reasonably sure that divergence (A) between plants and animals is the oldest, followed by divergence (B) between insects and the rest of the animals. If the rate of mutant substitution is the same, we should expect the various animals to differ from wheat to the same extent (except for statistical fluctuations) and this is indeed the case.

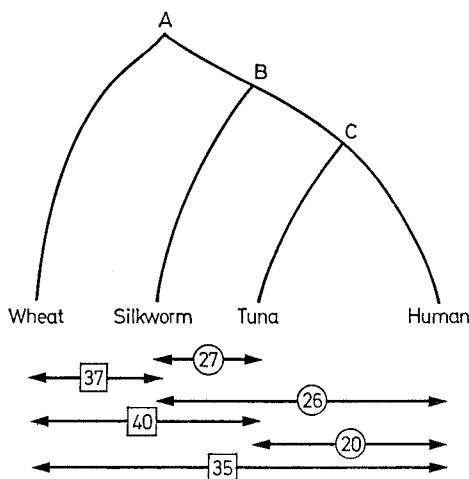


Fig. 3. A phylogenetic tree involving wheat, silkworm, tuna and man. Numerals with arrows indicate the numbers of amino acid differences among their cytochrome c molecules

## Units of Evolutionary Rates: Darwin and Pauling

The rate of evolution may be considered at various levels. Traditionally, at least two measures have been used. One is the rate at the "organism" level, as extensively studied by Simpson (1944). For example, the line leading to horse (*Equus*) from *Eohippus* went through eight successive genera during some 45 million years, giving about 0.18 genera per million years. This is within the range of "standard" rates which Simpson called "horotelic". It is well known among students of evolution that there are enormous differences among evolution rates. Some forms have evolved very rapidly while others have changed so slowly as to be hardly evolving at all.

Less subjective is the rate of change in quantitative measurements such as body size and length of tooth as studied by Haldane (1949). For example, in the evolution of the horse, he found that tooth length changed on the average at the rate of about  $4 \times 10^{-8}$  per year, or 4% per million years. He proposed the term *darwin* as a unit of evolutionary rate, representing a change in measurement at the rate of  $10^{-6}$  per year. In these terms the horse rate is about 40 millidarwins. Also according to Haldane, the rate of evolutionary increase in body length of Dinosaurs during the mesozoic era is roughly half as large. On the other hand, it is known that in hominid evolution, cranial capacity increased by the factor of 2 or 3 in million years (since *Australopithecus*), so the rate is near one darwin, an example of rapid evolution.

It is evident that an entirely different measure is needed to represent the rate of evolution at the molecular level, since the rate must be measured in terms of the number of mutant substitutions. Kimura (1969b) proposed the term *pauling* to represent the rate of substitution of  $10^{-9}$  per amino acid site per year. In terms of this, the hemoglobin rate is very near to one *pauling*, while the rate of cytochrome c is about 25 centipaulings. The most rapid rate known is 4 paulings for Fibrinopeptide A. The rate for 7 proteins were computed by King and Jukes (1969). The average rate they obtained is 1.6 paulings. Also Dayhoff (1969) listed the rates for 16 proteins. On the whole, difference of rates among molecules are much less pronounced than "organism" rate. The only exception is histones which had only 2 changes in approximately 1.5 billion years per 100 amino acids, or roughly one centipauling.

## Population Dynamics of Mutant Substitution

We will now consider the sequence of events within a population whereby a rare molecular mutant increases in frequency and spreads through the entire population, eventually reaching the state of fixation (see Fig. 1). Each event represents one mutant substitution or incorporation. By mole-

cular mutant we mean a mutant produced by base replacement at one or more nucleotide site, the great majority of which involve a single base replacement. Each amino acid site corresponds to a codon consisting of 3 nucleotide sites. By correcting for synonymy of codons, which amounts to some 25 % (cf. Kimura, 1968b), the rate of nucleotide substitution can be estimated from the rate of amino acid substitution. Roughly speaking, however, the former is  $1/3$  of the latter.

In the following treatment, "site" refers to a nucleotide site, but it applies equally to a codon or amino acid site corresponding to 3 nucleotide sites.

Although most biologists will have no difficulty in understanding the nature of molecular mutants, some applied mathematicians working today on population genetic theory seem to be still preoccupied by a classical gene concept, with reversible mutation between a pair of alleles, say  $A$  and  $a$ , at a comparable rate at each genetic locus. Therefore, it may be pertinent here to mention some salient features of molecular mutants which are essential to the understanding of their population consequences. First, the total number of nucleotide sites making up a haploid genome of higher organisms is very large. For mammals, this amounts to some 4 billion ( $4 \times 10^9$ ). This is several orders of magnitude higher than the conventional "gene number" which is estimated to be some 5 thousands for *Drosophila*, and, at most, a few times as numerous for men. The mutation rate per site per generation is of the order of  $10^{-8}$  or  $10^{-9}$  rather than  $10^{-5}$  or  $10^{-6}$  for classical genes (cf. Kimura, 1968b). Secondly, each cistron or a structural gene consists of a fairly large number of nucleotide sites. For example, the cistron coding for the hemoglobin  $\alpha$  chain in mammals consists of 141 codons or 423 nucleotide sites, any one of which may change to produce a molecular mutant. With 4 kinds of base pairs that can occur in each site, the total number of possible alleles in this cistron is  $4^{423}$  or roughly  $10^{254}$ , truly an astronomical number. For each allele, there are  $3 \times 423$  or 1269 other alleles that can be reached by single base substitution. They will in turn change into other alleles by additional single base replacement, but the chance is now 1 in 1269 that it comes back to the original allele. This means that the "back mutation" in the strict sense is so rare to be negligible at the molecular level.

To formulate the process of mutant incorporation mathematically, we use the following model. We assume that the total number of sites making up the genome is so large while the mutation rate per site is so low that whenever a mutant appears it represents a new, previously homallelic site, in which mutant forms are not segregating. The same model was used by Kimura (1969a) to calculate the number of heterozygous nucleotide sites per individual under a steady flux of mutations in a finite population, and



by Ohta and Kimura (1971) to study linkage disequilibrium between two segregating sites.

Let  $\nu_m$  be the number of sites at which new mutants appear each generation in the entire population. We assume that at each site, a new mutant is represented only once at the moment of appearance. Since in any reasonably large population, deleterious mutants will eventually be lost from the population, we will restrict our consideration to mutants that are either advantageous or selectively neutral.

Let  $u$  be the probability of a single mutant ultimately reaching fixation in a population of actual size  $N$  and "effective" size  $N_e$  (for the meaning of effective population number, readers may refer to Kimura and Crow (1963), and Crow and Kimura (1970)).

Then the rate of mutant substitution or incorporation defined by formula (1) is given by

$$k = \nu_m u = 2Nv u, \quad (5)$$

where  $v = \nu_m / (2N)$  is the mutation rate per gamete per generation. Here we assume that mutants at different sites behave independently with free recombination and without epistasis.

The probability,  $u$ , of ultimate fixation of an individual mutant is in general a function of its selective advantage, and also of  $N$  and  $N_e$  of the population where it occurs (cf. Kimura, 1957, 1962, 1964).

So, we will consider two important cases. First, if the mutants have a small but definite selective advantage  $s$  over their preexisting form, then we have approximately

$$u = 2s_1(N_e/N) \quad (6)$$

(Kimura, 1964). This is valid as long as  $s_1$  is small but  $4N_e s_1$  is large, or more precisely,  $0 < s_1 \ll 1$  and  $\exp(-4N_e s_1) \ll 1$ . If the actual and effective sizes of the population are the same ( $N_e \equiv N$ ), the formula reduces to  $u = 2s$ , a result first obtained by Haldane (1927); namely, the probability of ultimate fixation of a single mutant is roughly twice its selective advantage. For example, if the mutant has one per cent selective advantage the chance of its ultimate fixation is about 2 per cent. However, in natural populations, several factors are at work which make the effective number ( $N_e$ ) smaller than the actual number ( $N$ ). The most important of these is probably occasional reduction in the number of individuals. Also, inequality in the numbers of breeding males and females, and deviation of progeny distribution from the Poisson law in the direction of larger variance contribute to make  $N_e$  small. It is quite likely therefore that in the evolution of species the effective number is usually very much less than the total number of individuals comprising the species (summed over all local subpopulations), and the ratio  $N_e/N$  is much smaller than unity. It is even probable that  $N_e$  is often an order of magnitude less than  $N$ .

Substituting formula (6) in (5), we obtain

$$k = 4N_e s_1 v, \quad (7)$$

showing that in this case the rate of mutant substitution in the course of evolution depends on the effective population number, and on the selective advantage, as well as on the rate at which mutants having such selective advantage are produced each generation. One should expect then that  $k$  depends strongly on the environment in which the species is placed, being high for a species offered a new "ecologic opportunity" (cf. Wright, 1950) but low for those kept in a stable environment.

Secondly, if the mutant is selectively neutral, the probability of ultimate fixation of a single mutant is equal to its initial frequency so that

$$u = 1/(2N), \quad (8)$$

where  $N$  is the actual population number. Substituting this in (5), we obtain

$$k = v \quad (9)$$

showing that for a neutral mutant, the rate of evolution in terms of mutant incorporation is equal to the mutation rate per gamete—a remarkably simple principle (Kimura, 1968a; King and Jukes, 1969; Crow, 1969). This formula is valid as a good approximation as long as

$$|4N_e s_1| \ll 1, \quad (10)$$

where  $s_1$  is the selection coefficient of the mutant (positive or negative) relative to the preexisting form. Mutants having such a small selective effect have been called "almost neutral" by Kimura (1968b). They have the definite characteristic that their rate of substitution is independent of the effective population number, but depends only on the mutation rate. This should not be confused with the rate at which an individual mutant increases within the population, which does depend on the effective population number. Actually, it was shown by Kimura and Ohta (1969a) that it takes on the average  $4N_e$  generations for a selectively neutral mutant to reach fixation in a population of effective size  $N_e$  if we exclude the cases in which it is lost from the population by chance. They also showed that the average number of generations until fixation is less for definitely advantageous mutants (for the general theory on the subject, readers may refer to Kimura and Ohta, 1969a).

An additional point regarding formula (9) for neutral mutants is that it is independent of the mode of reproduction; it is equally applicable to sexual and asexual organisms. On the other hand, formula (7) for selected mutants is valid only for sexually reproducing species in which enough recombination can take place between mutants in the course of substitution.

Under asexual reproduction, the rate of evolution can be lower (for details see Crow and Kimura, 1965).

There is one more point which needs consideration when we discuss the dynamics of gene substitution, especially under natural selection. This is the amount of selective elimination that accompanies gene substitution by natural selection. In his pioneering work Haldane (1957) showed that the number of selective deaths required to carry out one gene substitution by natural selection is independent of the selection coefficient ( $s_1$ ) but depends on the initial frequency of the mutant gene used for the substitution. He used the term "cost" for the amount of selective elimination and obtained an elegant result;  $D = -2 \log_e p$  for a semidominant mutation, where  $D$  is the fraction of cumulative selective deaths and  $p$  is the initial mutant frequency. For example, if the initial frequency of the advantageous mutant is one in million ( $p = 10^{-6}$ ),  $D = 27.6$ . If the mutant is completely recessive,  $D$  becomes much larger. Haldane took  $D = 30$  as a representative value in the actual course of evolution. If selection is taking place slowly at a number of loci with an average rate of one mutant substitution every  $n$  generations, the fitness of the species will fall below the optimum by a factor of  $30/n$ , namely, selection intensity  $I = 30/n$ . He conjectured that  $I = 0.1$  is a reasonable figure for horotelic evolution and suggested that mutants are substituted at the rate of about  $1/300$ . He also considered that this accords with the observed slowness of evolution at the phenotypic level.

In his calculation on the cost of natural selection, or "substitutional load" as later called by Kimura (1960), Haldane assumed an infinitely large population. However, actual populations are all finite and it may be desirable to take stochastic elements into account. This was done by Kimura and Maruyama (1969) using diffusion models. They showed that the load for one mutant substitution is approximately

$$L(p) = -2 \log_e p + 2, \quad (11)$$

where  $s_1$  is the selective advantage of the mutant assumed to be semidominant in fitness, provided if  $4N_e s_1 p$  is small but  $4N_e s_1$  is large. Thus, this formula is valid under the same conditions for which formula (7) is valid. The formula shows that the load for one mutant substitution in a finite population is larger by 2 than the corresponding Haldane formula. In a finite population, a large fraction of slightly advantageous mutants are lost by chance, never contributing to the actual mutant substitution, and this inflates the load. However, the correction to be added is relatively small.

On the other hand, if the mutant is almost neutral ( $4N_e s_1 \ll 1$ ), we have  $L(p) = -8N_e s_1 \log_e p$ , and therefore the substitutional load can be very small. For such a class of mutants, the Haldane principle of cost does not set an upper limit to the rate of substitution in evolution.

Haldane had insights that were deep, both mathematically and biologically. In particular he was fully aware of the various means of population regulation. In our opinion a great deal has been written that adds little or nothing to Haldane's original papers (1957, 1960) and in some cases have led to semantic confusion. It is encouraging, however, that further developments along the line originally formulated by Haldane are being published by Crow (1970), Felsenstein (1971) and Nei (1971) (see also Kimura and Crow, 1969).

In view of the importance of formula (9) giving the rate of neutral mutant substitution for our consideration of the molecular evolution, we have performed Monte Carlo experiments to check the validity of this formula.

The results are summarized in the Table. The experiments were performed by using the TOSBAC 3400 computer. The population consists of  $2N$  chromosomes, each represented by a binary integer number. Number "0" at any digit represents the normal allele and "1", the mutant allele. Each experiment starts from a population of individuals having all 0's in their chromosomes. In each generation, if the uniform pseudo-random number (RAND 20 in Fortran IV) is less than  $2Nv$  (the total number of mutants in the population per generation) one mutant is introduced into a randomly chosen chromosome and a randomly chosen locus at which no mutants are currently segregating within the population (we assumed  $2Nv \leq 1$ ). Recombination and sampling were also performed by generating uniform random numbers, and the experiments were repeated over many generations. Whenever, the mutant ("1") becomes fixed in one locus, all the numbers at that locus are set back to zero, ready to be used for next mutation. The agreement between the expected and the observed numbers of substitutions is satisfactory.

Table. Results of Monte Carlo experiments on the rate of substitution of neutral mutants. A total of 20 loci were assumed in each chromosome, and 4 levels of recombination fraction between adjacent loci ( $C_1$ ) were tried ( $C_1 = 0.0, 0.01, 0.02$  and  $0.03$ ). Population size ( $N$ ) was 10 (diploid population) and 3 levels of mutation rates i.e.,  $2Nv = 0.1, 0.5$  and  $1.0$  per generation were used. The table lists the number of substitutions actually occurred during 2000 generations, together with expected numbers in parenthesis

$C_1$	$2Nv$		
	0.1	0.5	1.0
0.0	10	42	107
0.01	6	54	84
0.02	7	51	90
0.03	9	45	104
(Expected)	(10)	(50)	(100)

## Molecular Evolution by Neutral Mutation and Random Frequency Drift

There are at least two outstanding features that characterize the rate of evolution at the molecular level as compared with that at the phenotypic level. They are (1) a very high overall rate, and (2) remarkable uniformity for each molecule.

The first becomes evident if we extrapolate the observed rate per amino acid site to the whole genome of higher organisms. In man, the total number of nucleotide sites making up the genome is about  $4 \times 10^9$  (Muller, 1958; Vogel, 1964) and this figure is roughly the same among different species of mammals. If the rate of mutant substitution is on the average about 1.6 paulings per nucleotide triplet as estimated by King and Jukes (1969), the rate of mutant substitution per genome is roughly 2 per year. For a mammalian species which takes 3 years for one generation, the mutant substitution proceeds in the species at the rate of half dozen per generation. This is a surprisingly high rate, especially in view of Haldane's (1957) conjecture that new mutants may be substituted by natural selection at the rate of  $1/300$  in standard rate evolution. Actually, if we calculate the substitutional load using formula (11), assuming that the majority of mutant substitutions at the molecular level are carried out by natural selection, the substitutional load in each generation is so large that no mammalian species could tolerate it. For example, for a species consisting of a half million individuals,  $L(p) = 29.6$  with  $p = 1/(2N) = 10^{-6}$ , and the load per generation is 6 times this figure or roughly 180. This means that to maintain the same population number and still to carry out mutant substitution at the above rate, each parent must leave  $e^{180} \approx 10^{78}$  offspring for only one of the offspring to survive.

This was the main reason why random fixation of selectively neutral mutants was first proposed by one of us (Kimura, 1968a) as the main factor in molecular evolution. An additional reason not mentioned but realized at that time is that an unusually high rate of production of "advantageous" mutants is required to explain the high rate of molecular evolution by natural selection, especially assuming mutants with slight selective advantage such as  $s_1 = 0.001$  or less. To see this, consider a mammalian species having a large body size and a generation time of three years. For such a species, the effective population number (but not the actual number) may be  $10^5$  or less. From formula (7), we have

$$v = k/(4N_e s_1) \quad (12)$$

and if we put  $k = 6$ ,  $2N_e = 10^5$  and  $s_1 = 10^{-3}$ , we get  $v = 3/100$  or 3 per cent per gamete per generation. This is a very high rate of production for advantageous mutations, since it is the same order of magnitude as the total rate for lethal and semilethal mutations per gamete. If  $k$  is larger, or both  $N_e$  and  $s_1$  are smaller,  $v$  becomes higher. The assumption that

“advantageous” mutations are being produced at such a high and constant rate throughout the long history of evolution seems to contradict the principle of adaptive evolution by natural selection, for one should expect that substitution of each advantageous mutation must generally decrease the probability of subsequent mutations being advantageous unless the environment changes drastically. Also, advantageous mutations should be much less frequent than deleterious mutations.

The second characteristic of the rate of molecular evolution, namely its remarkable uniformity, is also consistent with the theory of neutral mutant substitution rather than that of selection, as pointed out by King and Jukes (1969), Kimura (1969b) and Crow (1969). In the previous section, we have demonstrated this point in the form of formulas (9) versus (7). According to the neutral mutation-random fixation theory, the remarkable uniformity of the rate of evolution of hemoglobins can be explained simply by assuming that the neutral molecular mutations due to base replacement in the hemoglobin genes occur at a constant rate per year throughout the diverse lines of vertebrate evolution. On the other hand, if we adopt the theory of mutant substitution by natural selection, we must assume that values of  $N_e$  (effective population number),  $s_1$  (selective advantage of the mutants), and  $v$  (rate of production of such mutants) are adjusted in such a way that their product always stays constant and is equal (per year, not per generation) throughout diverse lines of vertebrate evolution, irrespective of whether the evolution at the phenotypic level is very rapid (as in the line leading to man) or practically stopped (as in the line leading to carp). Uniformity of the rate of molecular evolution has also been found by Sarich and Wilson (1967), using an immunological approach, for albumins in primates. More recently, McLaughlin and Dayhoff (1970) obtained a result showing that the rate of molecular evolution is also quite uniform among prokaryote and eukaryote lines with respect to differentiation of several transfer RNA's.

Although the rate of evolution usually is remarkably uniform for each protein, there is a clear evidence for small but definite variation of evolutionary rate among different proteins. This can most readily be explained by assuming that the fraction of neutral change is different among different proteins, depending on their biological function, as pointed out by King and Jukes (1969), in addition to the possibility of a difference in mutation rate. Recently, Fitch and Markovitz (1970) obtained an important result suggesting that in cytochrome c only about 10% of amino acid replacements are neutral at any moment in the course of evolution, while in fibrinopeptide A most of the change is neutral. This is consistent with the observed fact that fibrinopeptide A evolved roughly 10 times as fast as cytochrome c.

There are several additional pieces of evidence that support our theory of neutral mutant substitution by random drift.

One is the observation that the average amino acid composition of vertebrate proteins can be predicted by the DNA base frequencies and the knowledge of the genetic code, assuming complete randomness in the arrangement of DNA bases within cistrons (Kimura, 1968b; King and Jukes, 1969; Ohta and Kimura, 1970).

The agreement between observed and expected amino acid frequencies is remarkably good except for arginine, suggesting that amino acid replacements and underlying DNA base replacements have occurred mostly at random in the course of evolution. If a majority of mutant substitutions are due to natural selection acting on advantageous mutations, it is unlikely that an expectation based on randomness would give such a good fit. The alternative argument is sometimes put forward that the code evolved to fit the amino acid composition that are essential to biological activity of the proteins. However, this is against the fact that the origin of present code is very old and it was established long before the amino acid composition of the vertebrate proteins have evolved (King, 1971).

Probably the strongest evidence for the prevalence of neutral mutation at the molecular level comes from the experiment by Cox and Yanofsky (1967) on Treffers' mutator gene in *E. coli*.

This mutator gene has an unusual property of causing preferentially the transversion from an *AT* pair to a *CG* pair at the enormous rate of  $3.5 \times 10^{-6}$  per *AT* pair per generation, with a negligible rate of change in the reverse direction. Since the total number of nucleotide pairs making up the genome of *E. coli* is estimated to be about  $3.7 \times 10^6$ , of which about half are *A-T* pairs, it is expected that roughly 7 molecular mutations occur on the average per bacterium each generation. In a strain containing this mutator, they observed an 0.2–0.5 per cent increase in *G-C* content of the bacterial DNA after 80 subcultures which corresponds to 1200–1600 cell generations. This is in good agreement with the theoretical expectation that *G-C* content should increase by 0.21–0.28 per cent under mutational pressure if unopposed by natural selection. This agreement suggests that majority of base substitutions caused by the mutator are neutral and not eliminated by selection. Note here that bacteria can not tolerate more than 50% selective death per division (generation) as pointed out by Crow (1969). Cox and Yanofsky's observation that the strain is fully viable after accumulation of more than 7000 base substitutions is also consistent with selective neutrality of the base substitution. In addition, competition experiments between the mutator strain and a coisogenic normal strain have since been made by Gibson, Scheppe and Cox (1970) using chemostats, and they obtained the result that the former outcompetes the latter. They consider that their result is consistent with the view that a majority of base substitutions have been almost neutral.

The hypothesis that molecular evolution has mainly been carried out through random fixation of selectively neutral mutants is in sharp contrast to the prevailing view in the evolution theory today maintaining that neutral mutant genes must be very rare if they ever occur, and random gene frequency drift is negligible in determining the genetic structure of biological populations (cf. Mayr, 1965). It is not surprising, therefore, that several criticisms have been published against this hypothesis, some of which are based on misunderstandings, as exemplified by the critique of Richmond (1970) which in turn has been rebutted by King, Jukes and Arnheim (1970).

Whether the hypothesis of random fixation of neutral mutants will survive or not will be judged in future through new data on molecular evolution together with their meticulous analyses based on population genetics theory. Verbal discussions based on evolutionary data at the phenotypic level will find less and less place in this field.

Finally, if this hypothesis turns out to be correct, we should expect that every species, including those that remain unchanged at the phenotypic level, are undergoing constant change at the molecular (DNA) level. Thus, genes in "living fossils" such as coelacanths, horseshoe crabs, and *Lingula* may be expected to have undergone as many base substitutions as corresponding genes in more rapidly evolving species (Kimura, 1969b). In other words, underneath their unchanged morphology (and probably physiology as well) that have been kept remarkably constant by incessant action of *natural selection* for hundred million years, a steady stream of almost neutral genetic variations has flowed through *random gene frequency drift*, transforming their informational macromolecules tremendously.

This will make comparative studies of amino acid sequences much more useful than if their evolutionary changes were mainly caused by natural selection.

We may look forward to the days when computer-assisted analyses of the phylogeny of information macromolecules supply much more useful and reliable information on the genealogy of existing organisms than can be attained by any other means based on phenotypes.

We would like to thank Dr. J. F. Crow for reading the manuscript and making valuable suggestions.

### References

- Cox, E. C., Yanofsky, Ch.: Proc. nat. Acad. Sci. (Wash.) **58**, 1895-1902 (1967).
- Crow, J. F.: Proc. XII Intern. Congr. Genetics **3**, 105-113 (1969).
- Genetic loads and the cost of natural selection. In: Biomathematics 1, mathematical topics in population genetics. K. Kojima, ed., p. 127-177. Berlin-Heidelberg-New York: Springer 1970.
- Kimura, M.: Amer. Natur. **99**, 439-450 (1965).
- — An introduction to population genetics theory. New York: Harper & Row 1970.



- Dayhoff, M. O. (ed.): Atlas of protein sequence and structure. Silver Spring, Maryland: National Biomedical Research Foundation 1969.
- Felsenstein, J.: Amer. Natur. **105**, 1-11 (1971).
- Fisher, R. A.: Proc. roy. Soc. Edinb. **50**, 205-220 (1930).
- Fitch, W. M., Margoliash, E.: The usefulness of amino acid and nucleotide sequences in evolutionary studies. In: Evolutionary biology, Steere, Dobzhansky & Hecht, eds. (in press).
- Markowitz, E.: Biochemical Genetics **4**, 579-593 (1970).
- Gibson, T. C., Scheppe, M. L., Cox, E. C.: Science **169**, 686-688 (1970).
- Haldane, J. B. S.: Proc. Camb. Phil. Soc. **23**, 838-844 (1927).
- Evolution **3**, 51-56 (1949).
- J. Genet. **55**, 511-524 (1957).
- J. Genet. **57**, 351-360 (1960).
- Jukes, T. H.: Molecules and evolution. New York: Columbia Univ. Press 1966.
- Kimura, M.: Ann. Math. Stat. **28**, 882-901 (1957).
- J. Genet. **57**, 21-34 (1960).
- Genetics **47**, 713-719 (1962).
- J. appl. Probability **1**, 177-232 (1964).
- Nature (Lond.) **217**, 624-626 (1968a).
- Genet. Res. Camb. **11**, 247-269 (1968b).
- Genetics **61**, 893-903 (1969a).
- Proc. nat. Acad. Sci. (Wash.) **63**, 1181-1188 (1969b).
- Crow, J. F.: Evolution **17**, 279-288 (1963).
- — Genet. Res. **13**, 127-141 (1969).
- Maruyama, T.: Heredity **24**, 101-114 (1969).
- Ohta, T.: Genetics **61**, 763-771 (1969a).
- — Genetics **63**, 701-709 (1969b).
- King, J. L.: The influence of the genetic code on protein evolution. In: Biochemical evolution and the origin of life. E. Schoffeniels, ed. North Holland (in press).
- Jukes, T. H.: Science **164**, 788-798 (1969).
- — Arnheim, N.: Nature (submitted).
- Mayr, E.: Animal species and evolution. Cambridge: Harvard University Press 1965.
- McLaughlin, P. J., Dayhoff, M. O.: Science **168**, 1469-1471 (1970).
- Muller, H. J.: Bull. Amer. Math. Soc. **64**, 137-160 (1958).
- Nei, M.: Genetics (in press).
- Ohta, T., Kimura, M.: Genetics **64**, 387-395 (1970).
- — Genetics (in press).
- Richmond, R. C.: Nature (Lond.) **225**, 1025-1028 (1970).
- Romer, A. S.: The procession of life. London: Weidenfeld and Nicolson 1968.
- Sarich, V. M., Wilson, A. C.: Proc. nat. Acad. Sci. (Wash.) **58**, 142-148 (1967).
- Simpson, G. G.: Tempo and mode in evolution. New York: Columbia Univ. Press 1944.
- Pittendrigh, C. S., Tiffany, L. H.: Life: An introduction to biology. London: Routledge & Kegan Paul 1958.
- Vogel, F.: Nature (Lond.) **201**, 847 (1964).
- Wright, S.: Population structure as a factor in evolution. In: Moderne Biologie. F. W. Peter, Editor, p. 274-287. Festschrift für Hans Nachtsheim, Berlin (1950).
- Zuckerkindl, E., Pauling, L.: Evolutionary divergence and convergence in proteins. In: Evolving genes and proteins. V. Bryson, & H. J. Vogel eds., p. 97-166. New York: Academic Press 1965.

Motoo Kimura  
 Tomoko Ohta  
 National Institute of Genetics  
 Yata 1, 111, Mishima  
 Shizuoka-ken, 411 Japan