# FEATURE ARTICLE

## Molecular Quasi-Species[†]

### Manfred Eigen,* John McCaskill,

*Max Planck Institut für biophysikalische Chemie, Am Fassberg, D 3400 Göttingen-Nikolausberg, BRD*

### and Peter Schuster*

*Institut für theoretische Chemie und Strahlenchemie, der Universität Wien, Währinger Strasse 17, A-1090 Wien, Austria (Received: June 9, 1988)*

The molecular quasi-species model describes the physicochemical organization of monomers into an ensemble of heteropolymers with combinatorial complexity by ongoing template polymerization. Polynucleotides belong to the simplest class of such molecules. The quasi-species itself represents the stationary distribution of macromolecular sequences maintained by chemical reactions effecting error-prone replication and by transport processes. It is obtained deterministically, by mass-action kinetics, as the dominant eigenvalue of a *value* matrix, **W**, which is derived directly from chemical rate coefficients, but it also exhibits stochastic features, being composed to a significant fraction of unique individual macromolecular sequences. The quasi-species model demonstrates how macromolecular information originates through specific nonequilibrium autocatalytic reactions and thus forms a bridge between reaction kinetics and molecular evolution. Selection and evolutionary optimization appear as new features in physical chemistry. Concentration bias in the production of mutants is a new concept in population genetics, relevant to frequently mutating populations, which is shown to greatly enhance the optimization properties. The present theory relates to asexually replicating ensembles, but this restriction is not essential. A sharp transition is exhibited between a drifting population of essentially random macromolecular sequences and a localized population of close relatives. This transition at a threshold error rate was found to depend on sequence lengths, distributions of selective values, and population sizes. It has been determined generically for complex landscapes and for special cases, and, it was shown to persist generically in the presence of nearly neutral mutants. Replication dynamics has much in common with the equilibrium statistics of complex spin systems: the error threshold is equivalent to a magnetic order–disorder transition. A rational function of the replication accuracy plays the role of temperature. Experimental data obtained from test-tube evolution of polynucleotides and from studies of natural virus populations support the quasi-species model. The error threshold seems to set a limit to the genome lengths of several classes of RNA viruses. In addition, the results are relevant even in eucaryotes where they contribute to the exon–intron debate.

## 1. Molecular Selection

Our knowledge of physical and chemical systems is, in a final analysis, based on models derived from repeatable experiments. While none of the classic and rather besieged list of properties rounded up to support the intuition of a distinction between the living and nonliving—metabolism, self-reproduction, irritability, and adaptability, for example—intrinsically limit the application of the scientific method, a determining role by unique or *individual* entities comes into conflict with the requirement of repeatability. Combinatorial variety, such as that in heteropolymers based on even very small numbers of different bases, even just two, readily provides numbers of different entities so enormous that neither consecutive nor parallel physical realization is possible. The physical chemistry of finite systems of such macromolecules must deal with both known regularities and the advent of unique copolymeric sequences. Normally this would present no difficulty in a statistical mechanical analysis of typical behavior, where rare events play no significant role, but with autocatalytic polymerization processes even unique single molecules may be amplified to determine the fate of the entire system. Potentially creative, self-organizing around unique events, the dynamics of this simplest living chemical system is invested with regularities that both allow and limit efficient adaptation. The quasi-species model is a study of these regularities.

The fundamental regularity in living organisms that has invited explanation is adaptation. Why are organisms so well fitted to their environments? At a more chemical level, why are enzymes optimal catalysts? Darwin's theory of natural selection has provided biologists with a framework for the answer to this question. The present model is constructed along Darwinian lines but in terms of specific macromolecules, chemical reactions, and physical processes that make the notion of survival of the fittest precise. Not only does the model give an understanding of the physical limitations of adaptation, but also it provides new insight into the role of chance in the process. For an understanding of the structure of this minimal chemical model it is first necessary to recall the conceptual basis of Darwin's theory.

Darwin recognized that new inheritable adaptive properties were not induced by the environment but arose independently in the production of offspring. Lasting adaptive changes in a population could only come about by natural selection of the heritable traits or *genotype* based on the full characteristics or *phenotype* relevant for producing offspring. A process of chance, i.e., uncorrelated with the developed phenotype, controls changes in the genotype from one generation to the next and generates the diversity necessary for selection. Three factors have probably prevented chemists from gaining a clear insight into these phenomena in the past, despite the discovery of the polymeric nature of the genotype (DNA): the complexity of a minimum replication phenotype, the problem of dealing with a huge number of variants, and the nonequilibrium nature of these ongoing processes.

The formulation of a tractable chemical model based on Darwin's principle may be understood in several steps:

1. The major constituents of the system have to be inherently self-reproductive. Only two classes of molecules are presently

---

(1) Eigen, M.; McCaskill, J. S.; Schuster, P. *Adv. Chem. Phys.*, in press.

known to have this property: RNA and DNA. In the kinetic equations the capability of self-replication is expressed through autocatalytic terms in the reaction scheme. These each have an overall rate that in detail may depend on a great many elementary reaction steps involving intermediates and small molecules (see below). The essential rate dependence is the linear one on template concentration so that the rate coefficients are phenotypic properties determined by the polymeric sequence, the genotype, and the chemical environment.

2. Replication is never entirely precise but rather prone to mutation. This is expressed by heterocatalysis of template in the polymerization of nearby sequences. The concept of neighborhood in the space of polymeric sequences plays an important role in what follows. The simplest type of mutation is the substitution of one monomer by another. Others include insertions and deletions. In this simplest model we restrict mutational variation to substitution, noting that any sequence of a given length may be obtained via subsequent mutations from a starting polymer. The Hamming distance $d(i,k)$, which counts the number of positions in which two sequences, $I_i$ and $I_k$, differ,[2] appropriately arranges the multitude of different copolymers into a (high-dimensional) *sequence space*.

3. The system has to stay far from equilibrium where neither the formation nor the decomposition process is microscopically reversible and detailed balance does not hold. This is necessary to allow concentrations to depend more strongly on the rate coefficients than expressed by the simple equilibrium ratios and hence to make selection possible. This is achieved in the kinetic equations by assigning irreversibility to all processes and including a process of dilution to constrain the total concentration. In reality, *far from equilibrium conditions* are sustained by transport processes that in the simplest case maintain a steady supply of energy-rich monomeric building blocks of the replicating biopolymers and a steady removal of polymerized products.

Apart from these three prerequisites, no specific interactions among the constituents—such as mutual enhancement or suppression, recombinative exchange, or other regulative couplings—have been assumed. Their presence would certainly alter the dynamics and add new features to the Darwinian scenario, but the above introduces a minimal model that is made precise in section 2 and solved in sections 3–5.

What is gained by this formulation? Perhaps most importantly, the physical chemistry makes explicit what is meant by the phrase *survival of the fittest*. Without defining *survival* and *fittest* independently of Darwin's principle we get inevitably caught in the tautological loop of *survival of the survivor*. Survival is appropriately measured in terms of population numbers or, in the language of the physical chemist when these are large, concentrations. A type survives when it is present at nonzero concentration in the long time limit. Fitness is a property of the individual type in a given environment: it is a value parameter measuring the efficiency in producing offspring. If this can be quantified, independently of concentrations, as a chemical property $W_i$ of a given macromolecule $I_i$, then the selection problem boils down to the establishment of quantitative correlations between selective values $W_i$ and population numbers $N_i$. *Survival of the fittest* would then imply that this correlation is of a rather simple all or none form: all types except that with the highest selective value are bound to die out. However, if significant cross-catalysis of neighboring sequence polymerization occurs, i.e., if mutations happen much more frequently than it is commonly assumed in population genetics, then not a single fittest type but an optimal ensemble—a *master sequence* together with its frequent mutants—will survive. Such stationary mutant distributions, usually distinguished by a unique consensus sequence, were named *quasi-species* to point out the analogy with the species concept in biology.

The most intriguing feature introduced by the quasi-species concept is the ability to describe a transition between a *living*

chemical system in which, qualitatively, heredity is strong enough to maintain adaptive information in a well-defined stationary distribution of polymers and a *nonliving* system in which the distribution is neither stationary nor well localized. To the physical chemist such a sharp transition, although a property of a nonequilibrium system, suggests an analogy with the equilibrium phase transitions that connect ordered and disordered phases of matter. Indeed, techniques developed in the study of localization in spin glasses allowed a calculation of the generic properties of this transition as discussed in section 5. The transition persists even when *neutral* mutants, with nearly equal rate coefficients, occur (cf. the neutral theory[3]). A formal analogy with an equilibrium phase transition in a heterogeneous Ising model was established more recently and forms the subject of section 6.

In addition to this change to the Darwinian picture of the outcome of natural selection, the dynamics of the adaptive process in this chemical model is also made very different by the replacement of single species by quasi-species. The quasi-species is shown to consist of a nonsymmetric cloud of mutants, replicating at different rates, that provide a mass-action bias in the production of further mutants. The higher concentrations of well-adapted macromolecular sequences result in more of their mutants, with correlated high values, being produced than for other sequences. This runs contrary to the classical belief that mutants of equal distance appear stochastically at a rate independent of their selective value and is shown in section 4 to radically alter the possibilities of evolutionary adaptation which are otherwise limited by the occurrence of local minima in the fitness function over the (Hamming) space of mutant copolymers.

What is the basis for our confidence that we have captured the essential features of a living system in the quasi-species model? Complex living systems are characterized by a great variety of mutual interferences, and then it is exceedingly difficult to evaluate fitness in qualitative terms. Indeed the environment, as introduced above, comprises not only all external parameters but also the influences exerted by competing types. This basic complexity has been confirmed for DNA replication: more than 10 different enzyme molecules and many intermediates are involved. Moreover, although most of the basic features of this process are presently known, it is impossible as yet to conceive a DNA replication process that can be studied in detail with the techniques of chemical reaction kinetics. Template-induced synthesis of RNA molecules, on the other hand, follows a much simpler mechanism, being normally accomplished by a single enzyme, and indeed it is experiments with this system in isolation from other cellular processes that laid the basis for the quasi-species model. Owing to the pioneering works of Spiegelman,[4] Charles Weissmann,[5] and others, in vitro replication assays were developed that use the single enzyme of the *Escherichia coli* bacteriophage $Q\beta$—the RNA polymerase $Q\beta$-replicase. More recently this system was studied in great detail.[6-8] It turned out to be well suited for systematic kinetic studies, and within the past 5 years a mechanism of RNA replication in the $Q\beta$ system was established and examined experimentally.[9-11] A concentration regime indeed occurs where catalysis is linear in the template concentration as in the simplest model. Error rates of RNA replication were also determined for several viruses,[12-14] and, as in the $Q\beta$ system, the frequency of

(2) Hamming, R. W. *Coding and Information Theory*; Prentice-Hall: Englewood Cliffs, NJ, 1980.

(3) Kimura, M. *The Neutral Theory of Molecular Evolution*; Cambridge University Press: Cambridge, UK, 1983.

(4) Spiegelman, S. *Q. Rev. Biol.* **1971**, *4*, 213.

(5) Weissmann, C.; Billeter, M. A.; Goodman, H. M.; Hindley, J.; Weber, H. *Ann. Rev. Biochem.* **1973**, *42*, 303.

(6) Biebricher, C. K.; Eigen, M.; Luce, R. *J. Mol. Biol.* **1981**, *148*, 369.

(7) Biebricher, C. K.; Eigen, M.; Luce, R. *J. Mol. Biol.* **1981**, *148*, 391.

(8) Biebricher, C. K. In *Evolutionary Biology*, Hecht, M. K., Wallace, B., Prance, G. T. Eds.; Plenum: New York, 1983; Vol. 16, pp1–52.
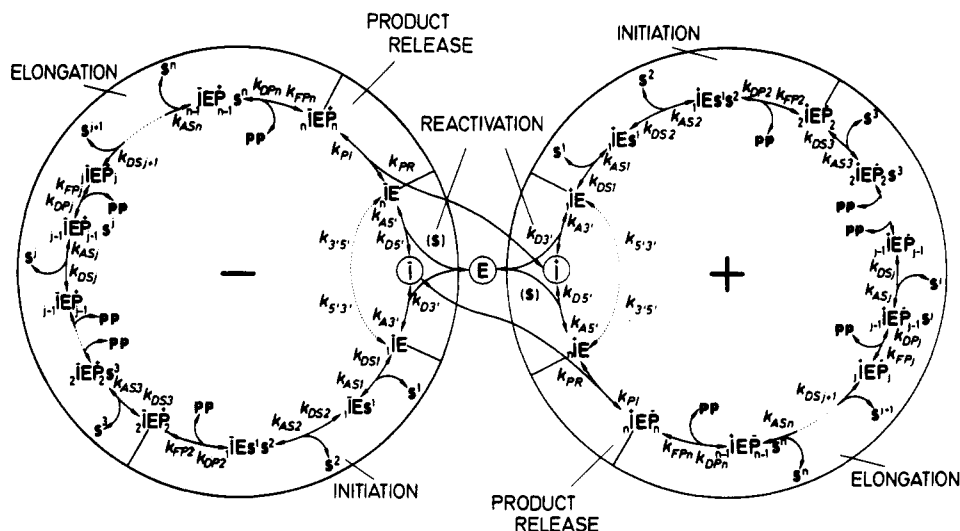
(9) Biebricher, C. K.; Eigen, M.; Gardiner, W. C., Jr. *Biochemistry* **1983**, *22*, 2544.

(10) Biebricher, C. K.; Eigen, M.; Gardiner, W. C., Jr. *Biochemistry* **1984**, *23*, 3168.

(11) Biebricher, C. K.; Eigen, M.; Gardiner, W. C., Jr. *Biochemistry* **1985**, *24*, 6550.

(12) (a) Domingo, E.; Flavell, A.; Weissmann, C. *Gene* **1976**, *1*, 3. (b) Batschelet, E.; Domingo, E.; Weissmann, C. *Gene* **1976**, *1*, 27.
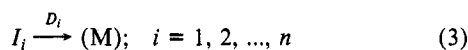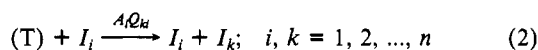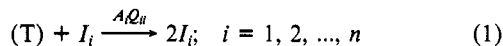
**Figure 1.** Reaction scheme of complementary replication of single-stranded RNA. The reaction consists of four phases: initiation, elongation, product release, and template activation. The reaction product—the replica—is complementary to the template. The substrates are the four nucleoside triphosphates: ATP, GTP, UTP, and CTP. The waste product at each step of incorporation is pyrophosphate: pp. The symbol I refers to the RNA template chain, E to the replicase enzyme, and P to the growing RNA replica chain. The indexes are as follows: A = association, D = dissociation, S = substrate, F = phosphate diester bond formation, and RP = product release. The numbers 3' or 5' refer to the particular end of the RNA chain to which the enzyme binds or from which it dissociates.[8]

mutations is high enough to give significant concentrations to a large spectrum of mutants.

A dynamical theory of molecular evolution was originally derived in 1971 from the kinetics of replication and mutation.[15] Later developments are found in ref 16–21. Here we shall briefly review the basic concept and then present the results of recent work which are described in detail in ref 1.

## 2. Kinetics of Replication and Mutation

Template-induced, enzymatic synthesis of RNA follows a multistep polymerization mechanism which is sketched in Figure 1.[8-10] If the activated monomers (the ribonucleoside triphosphates GTP, ATP, CTP, and UTP) and the polymerizing enzyme are applied in excess, the intermediates do not play any specific role in the production of polymer and hence can be subsumed in overall rate constants. Replication, mutation, and degradation can then be modeled directly by single-step processes:

$$(T) + I_i \xrightarrow{A_iQ_{ii}} 2I_i; \quad i = 1, 2, ..., n \tag{1}$$

$$(T) + I_i \xrightarrow{A_iQ_{ki}} I_i + I_k; \quad i, k = 1, 2, ..., n \tag{2}$$

$$I_i \xrightarrow{D_i} (M); \quad i = 1, 2, ..., n \tag{3}$$

Nucleoside triphosphates and nucleoside monophosphates are denoted here by T and M respectively—they are put in parentheses since their concentrations are not considered as variables here.

It is essential to be able to describe the variety and relatedness of the many different macromolecules in this model. They are visualized as strings of digits from a finite alphabet of $\kappa$ symbols:

$$I_i = (s_1^{(i)}s_2^{(i)}s_3^{(i)}...s_\nu^{(i)}) \tag{4}$$

(13) Holland, J.; Spindler, K.; Horodyski, F.; Grabau, E.; Nichol, S.; VandePol, S. *Science (Washington, D.C.)* **1982**, *215*, 1577.

(14) Parvin, J. D.; Moscona, A.; Pan, W. T.; Lieder, J.; Palese, P. *J. Virol.* **1986**, *59*, 377.

(15) Eigen, M. *Naturwissenschaften* **1971**, *58*, 465.

(16) Eigen, M.; Schuster, P. *The Hypercycle—A Principle of Natural Self-Organization*; Springer Verlag: Berlin 1979. (a) Eigen, M.; Schuster, P. *Naturwissenschaften* **1977**, *64*, 541. (b) Eigen, M.; Schuster, P. *Naturwissenschaften* **1978**, *65*, 7. (c) Eigen, M.; Schuster, P. *Naturwissenschaften* **1978**, *65*, 341.
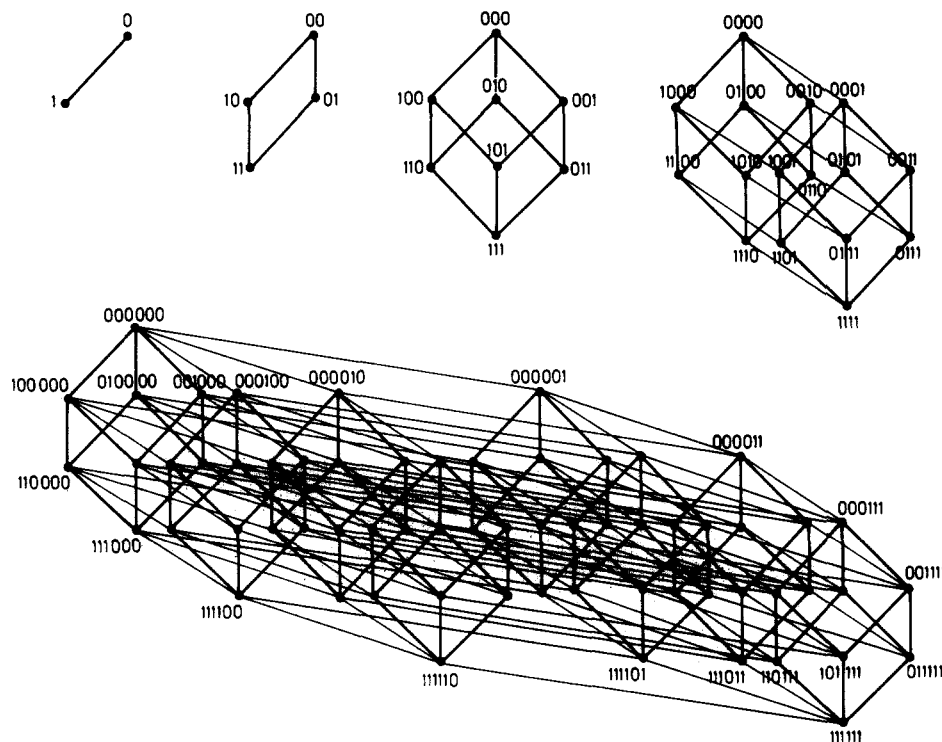
(17) McCaskill, J. S. *J. Chem. Phys.* **1984**, *80*, 5194.

(18) Eigen, M. *Ber. Bunsen-Ges. Phys. Chem.* **1985**, *89*, 658.

(19) Schuster, P.; Sigmund, K. *Ber. Bunsen-Ges. Phys. Chem.* **1985**, *89*, 668.

(20) Eigen, M. *Chem. Scr.* **1986**, *26B*, 13.

(21) Schuster, P. *Phys. Scr.* **1987**, *35*, 402.

with $s_j^{(i)} \in \{G,A,C,U\}$ for polynucleotides ($\kappa = 4$) or $s_j^{(i)} \in \{0,1\}$ for binary sequences ($\kappa = 2$). Clearly, the number of possible different sequences, $\kappa^\nu$, is hyperastronomically large—$4^\nu$ or $2^\nu$ for polynucleotides or binary sequences, respectively—even for moderate chain lengths $\nu$. In most studies carried out so far $\nu$ was kept constant. In this case the sequences are naturally arranged into an abstract point space, the sequence space originally introduced by Hamming,[2] which places mutational neighbors adjacently. It proved particularly useful for the visualization of evolutionary optimization processes. An example of the sequence space of binary sequences is shown in Figure 2.

Ordinary chemical reactions involve but a few molecular species, each of which is present in a very large number of copies. The converse situation also occurs in molecular evolution: the numbers of different polynucleotide sequences that may be mutually interconverted through replication and mutation exceed by far the number of molecules actually present in any experiment or even the total numbers of molecules available on earth or in the entire universe. Hence, the applicability of conventional kinetics to problems of evolution is a subtle question that has to be considered carefully whenever the deterministic approach is used.

The quality factors ($Q_{ii}$) measure the fraction of correct replicas synthesized on the template $I_i$. They are readily derived from single-digit accuracies $q_j^{(i)}$ by multiplication:

$$Q_{ii} = q_1^{(i)}q_2^{(i)}...q_n^{(i)} \tag{5}$$

The single-digit accuracy $q_j^{(i)}$ represents the fraction of correct digit incorporations at the position $j$ of the template $I_i$. In reality the $\nu$ factors in eq 5 will be all different depending on the nature of the digit to be incorporated as well as on its neighborhood in the sequence. For many purposes it is sufficient to assume uniform replication fidelities for all sequences:

$$q_1^{(i)} = q_2^{(i)} = ...q_\nu^{(i)} = q^{(i)} = q \tag{6}$$

Then, the quality factors depend only on chain lengths and (mean) single-digit accuracies: $Q_{ii} = q^\nu$.

The computation of mutation frequencies (which form the off-diagonal elements of the mutation matrix Q: $Q_{ki}$, $k \neq i$) depends, in addition, on the Hamming distance of the two sequences $I_i$ and $I_k$, which we denote by $d(i,k)$:

$$Q_{ki} = \epsilon^{d(i,k)}Q_{ii} \tag{7}$$

with $\epsilon = (q^{-1} - 1)/(\kappa - 1)$. We make here the implicit assumption that the probabilities of all ($\kappa - 1$) mutations to the different *error digits* are equal. This is consistent with the assumption of uniform replication fidelities.

**Figure 2.** *Lexicographic* ordering of sequences through successive duplications of sequence space. As shown in the figure for binary sequences the sequence space of dimension $\nu$ can be constructed by duplication of the sequence space of dimension $(\nu - 1)$. Each of the $2^\nu$ points specifies a binary $(0,1)$ sequence. If, in addition, the two alternative base classes $(0 = G$ or $A$ and $1 = C$ or $U)$ are specified, then to each of the $2^\nu$ points in binary sequence space another subspace of binary specification is added, yielding a total of $4^\nu$ points or a dimension of the hypercube of $2\nu$.

The replication rate constant $A_i$ is a measure of all copies—correct and erroneous—synthesized on the template $I_i$. Hence, $\mathbf{Q} = \{Q_{ik}\}$ is a $(b_i)$ stochastic matrix:

$$\sum_{k=1}^{n} Q_{ki} = 1 \qquad (8)$$

It is useful to define an excess production as the difference between the replication and the degradation rate constants of the individual templates:

$$E_i = A_i - D_i; \quad i = 1, 2, ..., n \qquad (9)$$

The combinations of rate constants that will turn out to be relevant for selection dynamics are summarized in the value matrix $\mathbf{W} = \{W_{ki}\}$. Its diagonal elements, the selective values

$$W_{ii} = A_i Q_{ii} - D_i \qquad (10)$$

correspond to the fitness factors of conventional population dynamics. The off-diagonal elements

$$W_{ki} = A_i Q_{ki} \qquad (11)$$

represent mutation rates $I_i \rightarrow I_k$—these are rate constants for the synthesis of $I_k$ on the template $I_i$ as error copy.

The variables of the dynamical system are the concentrations of individual polynucleotide sequences: $[I_i] = c_i(t)$. We are interested, essentially, in the relative concentrations of the different species

$$x_i(t) = c_i(t) / \sum_{i=1}^{n} c_i(t); \quad i = 1, 2, ..., n \qquad (12)$$

The resulting kinetic equations, around which quasi-species theory centers, are then

$$dx_i(t)/dt \equiv \dot{x}_i(t) = (W_{ii} - \bar{E}(t))x_i(t) + \sum_{k \neq i} W_{ik}x_k(t);$$
$$i, k = 1, 2, ..., n \qquad (13)$$

The mean excess production

$$\bar{E}(t) = \sum_{i=1}^{n} x_i(t)E_i \qquad (14)$$

of the population may be physically compensated by a dilution

flux $\Phi(t)$—corresponding to physical transport—which keeps the total concentration constant:

$$\sum_{i=1}^{n} c_i(t) = c_0 = \text{constant}$$

without change to eq 13. In a finite volume, the consequent finite population size is important in evolution, since low concentration sequences may fall below the single-molecule level where stochastic effects predominate.

To illustrate the meaning of this differential equation, we neglect first the sum of off-diagonal coupling terms which describe the *mutational* backflow. The first term on the right-hand side weighs the net production of correct copies on a given template—here $I_i$—against the mean excess production of the population. Sequences with selective values above the mean excess production—$W_{ii} > \bar{E}(t)$—will increase in frequency. The relative concentrations of those sequences that produce less than the average—$W_{ii} < \bar{E}(t)$—will decrease. The resulting shift in the distribution of relative concentrations causes an increase in $\bar{E}(t)$, which implies that more and more sequences fall below the critical threshold $W_{ii} = \bar{E}(t)$. If there were no mutations, this process would continue until the sequence with the largest $W_{ii}$ value had been selected. Backflow mutations become important when the relative concentration of a sequence becomes so small such that the first term on the right-hand side of eq 13 is comparable to the sum of the coupling terms. Then the sequence is protected from dying out by its production as error copy on the other templates. To be more precise, we have to cast these qualitative arguments now into quantitative predictions derived from the kinetic equations.

### 3. Solutions of Kinetic Equations and Stationary Mutant Distributions

Solutions of the kinetic equations (13) can be obtained in terms of eigenvalues and eigenvectors of the value matrix $\mathbf{W}$. After transformation[22,23] the off-diagonal coupling terms vanish:

$$\dot{y}_i(t) = (\lambda_i - \bar{\lambda}(t))y_i(t); \quad i = 1, 2, ..., n \qquad (15)$$

(22) Thompson, C. J.; McBride, J. L. *Math. Biosci.* **1974**, *21*, 127.

Here we denote the eigenvalues of $\mathbf{W}$ by $\lambda_i$. The mean eigenvalue

$$\bar{\lambda}(t) = \sum_{i=1}^{n} y_i(t)\lambda_i / \sum_{i=1}^{n} y_i(t)$$

is identical with the mean excess production: $\bar{E}(t) = \bar{\lambda}(t)$. The variables $y_i(t)$ are the components in the direction of the corresponding *right-hand* eigenvectors $\bar{l}_i$:

$$\vec{x}(t) = \sum_{k=1}^{n} x_k(t)\vec{e}_k = \sum_{k=1}^{n} y_k(t)\vec{l}_k \qquad (16)$$

with $\vec{x} = (x_1, x_2, ..., x_n)$ and $\vec{e}_k$ ($k = 1, ..., n$) being the conventional Cartesian unit vectors in concentration space. The equation

$$\mathbf{W}\vec{l}_k = \lambda_k\vec{l}_k; \quad k = 1, 2, ..., n \qquad (17)$$

defines the new basis $\vec{l}_k$ ($k = 1, ..., n$). The components of all *right-hand* (column) eigenvectors can be subsumed in the matrix

$$\mathbf{L} = (\vec{l}_1...\vec{l}_k...\vec{l}_n) = \{l_{ik}\} \qquad (18)$$

The *left-hand* (row) eigenvectors of the value matrix are given by

$$\mathbf{H} = \mathbf{L}^{-1} = \begin{pmatrix} \vec{h}_1 \\ \vdots \\ \vec{h}_k \\ \vdots \\ \vec{h}_n \end{pmatrix} = \{h_{ki}\} \qquad (19)$$

The general solution of the kinetic equations can be expressed now in terms of elements of the two matrices $\mathbf{L}$ and $\mathbf{H}$:

$$x_i(t) = \frac{\sum_{k=1}^{n} l_{ik} \exp(\lambda_k t)\sum_{j=1}^{n} h_{kj}x_j(0)}{\sum_{m=1}^{n}\sum_{k=1}^{n} l_{mk} \exp(\lambda_k t)\sum_{j=1}^{n} h_{kj}x_j(0)}; \quad i = 1, 2, ..., n \qquad (20)$$

Two conclusions can be drawn directly:

(1) The mean excess production—which is identical with the mean eigenvalue of $\mathbf{W}$—converges asymptotically to the largest eigenvalue of the value matrix, $\lambda_0 > \lambda_1 \geq \lambda_2 \geq ... \geq \lambda_n$:

$$\lim_{t\to\infty} \bar{E}(t) = \lim_{t\to\infty} \bar{\lambda}(t) = \lambda_0 \qquad (21)$$

This optimization principle was derived independently by application of techniques from statistical mechanics of spin lattices.[24] We shall come back to analogies between replication dynamics and spin statistics in section 6.

(2) The distribution of sequences converges asymptotically to the dominant *right-hand* eigenvector, $\vec{l}_0$, of the value matrix $\mathbf{W}$:

$$\lim_{t\to\infty} x_i(t) = \bar{x}_i = l_{i0}/\sum_{m=1}^{n} l_{m0} \qquad (22)$$

The stationary sequence distribution, thus determined by the dominant eigenvector $l_0$, is called the *quasi-species*. It consists of a *master sequence* $I_0$, which is the most frequent sequence and commonly has the maximum selective value, and a mutant distribution centered around the master. The value matrix $\mathbf{W}$, according to eq 10 and 11, has exclusively positive entries, and the Perron–Frobenius theorem[25] applies: the largest eigenvalue $\lambda_0$ is nondegenerate, and all components of the vector $\vec{l}_0$ are positive. Hence, it fulfills all requirements to describe a mixture of chemical compounds—here the distribution of the master sequence and its mutants.

Equation 20, in principle, provides a solution to replication–mutation dynamics, but it is of little use for most practical purposes. The value matrix $\mathbf{W}$ is huge—of dimension $\kappa^\nu \times \kappa^\nu$—and numerical methods to determine eigenvalues and eigenvectors fail. If all sequences in a given error class—master sequence $I_0$, all one-error mutants, all two-error mutants, etc.—have identical

replication rate constants, the dimension of the eigenvalue problem can be reduced to $\nu \times \nu$.[26] Analytical solutions are available only for the case of selective neutrality.[27] Therefore, we have to search for approximations. Perturbation theory turned out to be particularly useful for this goal.

## 4. Perturbational Approach to Replication–Mutation Dynamics

Conventional Rayleigh–Schrödinger perturbation theory can be applied directly to compute approximate expressions for the largest eigenvalue and the dominant eigenvector of the value matrix $\mathbf{W}$. To second order we obtain for the eigenvalue—$I_0$ being the master sequence ($W_{00} > W_{kk}, k \neq 0$)

$$\lambda_0 = W_{00} + \sum_{k\neq0} \frac{W_{0k}W_{k0}}{W_{00} - W_{kk}} \qquad (23)$$

The stationary fraction of the master sequence is given by

$$\bar{x}_0 = \frac{W_{00} - \bar{E}_{k\neq0}}{E_0 - \bar{E}_{k\neq0}} \qquad (24)$$

where we have introduced the average over all sequences except the master, in analogy to eq 9:

$$\bar{E}_{k\neq0} = \sum_{k\neq0} x_k(t)E_k / \sum_{k\neq0} x_k(t) \qquad (25)$$

The stationary concentrations of all other sequences are derived from the value $\bar{x}_0$ and the perturbation expression

$$\frac{\bar{x}_k}{\bar{x}_0} = \frac{W_{k0}}{W_{00} - W_{kk}} + \sum_{j\neq0,k} \frac{W_{kj}W_{j0}}{(W_{00} - W_{kk})(W_{00} - W_{jj})} \qquad (26)$$

To be more specific, we shall adopt the uniform error assumption and consider the stationary frequency of mutants $I_k$ as a function of the Hamming distance to the master sequence, $d(0,k)$. The probability of producing a mutant with Hamming distance $d$ is given by

$$Q_{k0} = q^\nu \epsilon^{d(0,k)} = q^\nu \left(\frac{q^{-1} - 1}{\kappa - 1}\right)^{d(0,k)} \qquad (27)$$

with $\epsilon$ as in eq 7. This probability decreases exponentially with increasing Hamming distance and, since $\epsilon$ often is very small ($\epsilon \approx 10^{-4}$ in viral RNA replication), becomes extremely small for the higher error classes ($d > 5$).

To give an estimate of the stationary concentrations of mutants, we must take into account also the selective values of the master sequence, the mutant under consideration, and all possible intermediates. This can be done directly within the perturbational approach when we neglect mutational backflow from the mutants to the master sequence. Then ratios of relative concentrations may be computed through recursion

$$\frac{\bar{x}_{dk}}{\bar{x}_0} = \epsilon^d \frac{W_{00}}{W_{kk}} f_{dk} \qquad (28)$$

with

$$f_{1i} = \frac{W_{ii}}{W_{00} - W_{ii}}$$

$$f_{2j} = \frac{W_{jj}}{W_{00} - W_{jj}}\left(1 + \sum_{i=1}^{\binom{2}{1}} f_{1i}\right)$$

$$f_{3m} = \frac{W_{mm}}{W_{00} - W_{mm}}\left(1 + \sum_{i=1}^{\binom{3}{1}} f_{1i} + \sum_{j=1}^{\binom{3}{2}} f_{2j}\right)$$

$$\vdots$$

$$f_{dk} = \frac{W_{kk}}{W_{00} - W_{kk}}\left(1 + \sum_{i=1}^{\binom{k}{1}} f_{1i} + ... + \sum_{l=1}^{\binom{k}{k-1}} f_{k-1,l}\right)$$

The double index $dk$ used here refers to the mutant $I_k^{(d)}$ in the

(23) Jones, B. L.; Enns, R. H.; Rangnekar, S. S. *Bull. Math. Biol.* 1976, 38, 12.
(24) Demetrius, L. *Phys. Scr.* 1987, 36, 693.
(25) Perron, O. *Math. Ann.* 1907, 64, 248.

(26) Swetina, J.; Schuster, P. *Biophys. Chem.* 1982, 16, 329.
(27) Rumschitzki, D. S. *J. Math. Biol.* 1987, 24, 667.

$d$-error class of the master sequence $I_0$: $d(0,k) = d$. The first term in the recursion for the factor $f_{dk}$ corresponds to the first term on the right-hand side of eq 26 expressed in the uniform error model (6). It describes single-step mutant formation:
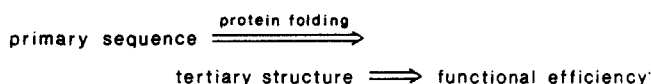
$$I_0 \rightarrow I_k^{(d)}$$

All other terms refer to mutant formation via intermediates

$$I_0 \rightarrow I_i^{(d_1)} \rightarrow I_k^{(d)}; \quad d_1 < d$$

$$I_0 \rightarrow I_i^{(d_1)} \rightarrow I_j^{(d_2)} \rightarrow I_k^{(d)}; \quad d_1 < d_2 < d$$

$$\vdots$$

$$I_0 \rightarrow I_i^{(1)} \rightarrow I_j^{(2)} \rightarrow ... \rightarrow I_l^{(d-1)} \rightarrow I_k^{(d)}$$

The cases in the last row pass through the maximum number of intermediates: $d - 1$. There are $d!$ such routes which contain $d$-fold products of $W_{ii}/(W_{00} - W_{ii})$ terms. Most mutants will have very small selective values, $W_{ii} \ll W_{00}$, or may even be lethal, $W_{ii} = 0$, and hence cannot contribute as intermediates. If all mutants were of this kind, mutant appearance would resemble just the Poissonian distribution, yielding a probability of $\epsilon^d$ for every individual in the error class $d$.

If, however, selective values are not distributed randomly in sequence space—if fitness landscapes show cohesive domains like landscapes do on earth where lowlands and mountainous areas are clustered—mutant frequencies may be larger than the Poissonian values by many orders of magnitude. In particular, mutant frequencies are large along ridges of high selective values—these are the paths along which all factors $W_{ii}/(W_{00} - W_{ii})$ are fairly large. Stationary mutant distributions thus will protrude far into sequence space along such ridges. As an example, consider the difference between a 12-error mutant that is connected to the wild-type by a ridge of higher valued intermediate mutants and one that is not. If the sequence length is 100 and the directly intervening mutants have an inferiority dipping smoothly to 0.5 at $d = 6$ and 7, whereas other mutants have value 0.2, then the first particular 12-error mutant appears by 15–20 orders of magnitude more frequently than other 12-error mutants that are not connected to the wild-type in this way!

There is indeed evidence for at least peak connectivity in value landscapes. Functional efficiency in proteins is clustered around certain sequences. The catalytic activity of an enzyme depends on the correct spatial arrangement of those amino acid residues that constitute the active center. This is achieved by three-dimensional folding of the polypeptide chain:

$$\text{primary sequence} \xrightarrow{\text{protein folding}}$$

$$\text{tertiary structure} \Longrightarrow \text{functional efficiency}$$

If similarities in primary sequences map into similarities of folding which in turn lead to similarities in tertiary structures and functional efficiency, we are dealing with a highly correlated value landscape. This indeed holds for many changes in primary sequences of enzyme molecules that were produced by site-directed mutagenesis—provided these changes do not occur at certain strategic positions that are vital for the structure and function of the active site. The same conclusion is reached from phylogenetic sequence relations: in a particularly well-studied example, the enzyme cytochrome c, sequences from different species may be nonhomologous in more than 70% of the positions and yet the same reaction is catalyzed (almost) equally well. There is no doubt that peaks in the fitness landscape do cluster, but the detailed relation between changes in sequence and their consequences in functional performance is anything but simple. Recent work has exposed essential features of this relationship based on the sequence to secondary structure mapping in RNA. One of us has recently calculated the full sequence-dependent equilibrium distribution of structures, which places the optimal structure in context as a member of an ensemble of related structures.[28] Secondary-

structure-based sequence-to-value mappings confirm the existence of correlations.[29] For the quasi-species model the existence of clustered value distributions is of paramount importance. The theory shows how value distributions map into population structures, and therefore quantitative analysis of mutant distributions in selection experiments could provide direct evidence for connectivity in landscapes and guidance of mutations in evolution.

The comparison of value landscapes with landscapes on earth may be misleading, since the sequence of realistic cases is always high-dimensional. As shown in Figure 2 the states are highly interconnected and, hence, distances remain relatively small: the maximum Hamming distance of sequences with chain length $\nu$, i.e., the maximum distance in the $\nu$-dimensional sequence space, is $\nu$ in an ensemble of $\kappa^\nu$ different states.

Perturbation theory, in the form applied so far, faces a problem when the selective values of a mutant are very close or even equal to that of the master sequence. In this case of selective neutrality, denominators in the perturbation expansion become zero and the series diverges. In contrast to the previous examples, mutational backflow becomes important too. Then perturbation theory of degenerate states has to be applied. The case of two sequences with identical selective values was studied in detail.[30] The structure of the quasi-species depends characteristically on the Hamming distance $d(1,2)$ of the two *best* sequences, $I_1$ and $I_2$. Three cases are distinguished:

(1) $d(1,2) = 1$. The mutant distribution formed at small error rates centers around the two sequences $I_1$ and $I_2$ which are present in equal amounts, $\bar{x}_1/\bar{x}_2 = 1$.

(2) $d(1,2) = 2$. The mutant distribution formed at small error rates centers around the two sequences $I_1$ and $I_2$ whose concentrations are different, $\bar{x}_1/\bar{x}_2 = \alpha$ with $0 < \alpha < 1$. The ratio of stationary concentrations is determined by the neighborhood of the two *best* sequences. The one that is situated in the *fitter* neighborhood is present in higher concentrations.

(3) $d(1,2) \geq 3$. Selection occurs in the limit of vanishing error rates, $\epsilon \rightarrow 0$. The sequence with the fitter neighborhood is selected. These results demonstrate that the neutral theory has to be modified in view of the results of the quasi-species model. It is valid only in cases where mutations are so rare that mutational backflow can be completely ruled out. This limiting case may well be fulfilled in higher organisms with very large genomes and relatively small populations. There is no mutational correlation between different selectively neutral sequences, and their relative concentrations drift randomly. In test-tube evolution experiments in natural populations of viruses or bacteria, however, mutations occur frequently and populations are so large that many selectively neutral polymer sequences are mutationally interconnected. Then the ratios of their concentrations converge to values that are determined according to the kinetic equations (13) and are not subject to random drift.
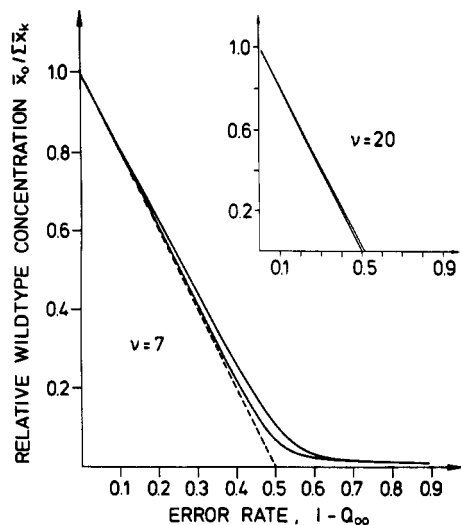
In addition, cases were investigated where two sequences have slightly different selective values, but the less efficient sequence *lives* in a more efficient neighborhood.[30] We consider the structure of the quasi-species as a function of increasing error rate: at small error rates the mutant distribution centers around the most efficient sequence. However, at some critical replication accuracy $q_{cr}$ it may happen that the less efficient sequence becomes the master sequence, because the difference in mutational backflow overcompensates the small difference in selective values. Such a case is discussed below in relationship to the error threshold and is displayed in Figure 4, bottom.

These examples demonstrate an important feature of the replication–mutation system: the fitness of a given sequence is not only determined by its own selective value. Rather, neighboring sequences contribute too, through mutations, and their influence becomes more important as the error rate increases. In conventional natural selection theory, advantageous mutations drove the evolutionary process. The neutral theory introduced selectively neutral mutants, in addition to the advantageous ones, which

(28) McCaskill, J. S. *The Equilibrium Partition Function and Base Pair Binding Probabilities for RNA Structure, Biopolymers*, in press.

(29) Fontana, W.; Schuster, P. *Biophys. Chem.* **1987**, *26*, 123.

(30) Schuster, P.; Swetina, J. *Bull. Math. Biol.*, in press.

**Figure 3.** Relative wild-type concentrations as a function of the error rate $1 - Q_{00}$. We show the exact solution curve (upper full line) and compare it with the result of perturbation theory (according to (29), broken line) and the exact phenomenological result (lower full line). The following parameters and rate constants were applied: $A_0 = 10$, $\bar{A}_{k\neq0} = 5$; $D_0 = D_1 = ... = D_n$ with $n = 2^\nu - 1$ and $\nu = 7, 20$.

contribute to evolution through random drift. The concept of quasi-species shows that much weight is attributed to those slightly deleterious mutants that are situated along high ridges in the value landscape. They guide populations toward the peaks of high selective values.

## 5. Error Thresholds and Localization of Quasi-Species

The kinetic model of replication and mutation gives rise to a sharp error threshold which depends on the value landscape ($W_{ii}$; $i = 1, ..., n$), on the mean fidelity of single-digit replication, $q$, and on the sequence length, $\nu$. The error threshold confines the range of parameters within which evolution proceeds to selection of a stable quasi-species which is localized—i.e., concentrated around one, two, or eventually more *best* sequences—in some region of the sequence space. The error threshold defines a critical minimum accuracy of replication, $(Q_{00})_{min}$. If the error rate exceeds the threshold value, the class of sequences that are populated in amounts comparable to the fittest types—in the stationary solution of the deterministic eq 13—becomes so large that it cannot be sampled by any realistic population. A stochastic interpretation of the dynamics as random drift over this class in sequence space is then called for. However, in describing the approach to the error threshold in large populations the deterministic approach suffices.[31]

The limiting nature of error propagation is best understood by first adopting a phenomenological formulation. The master sequence $I_0$ competes successfully with its mutants if the net production of correct copies—the selective value—exceeds the mean production of all others:
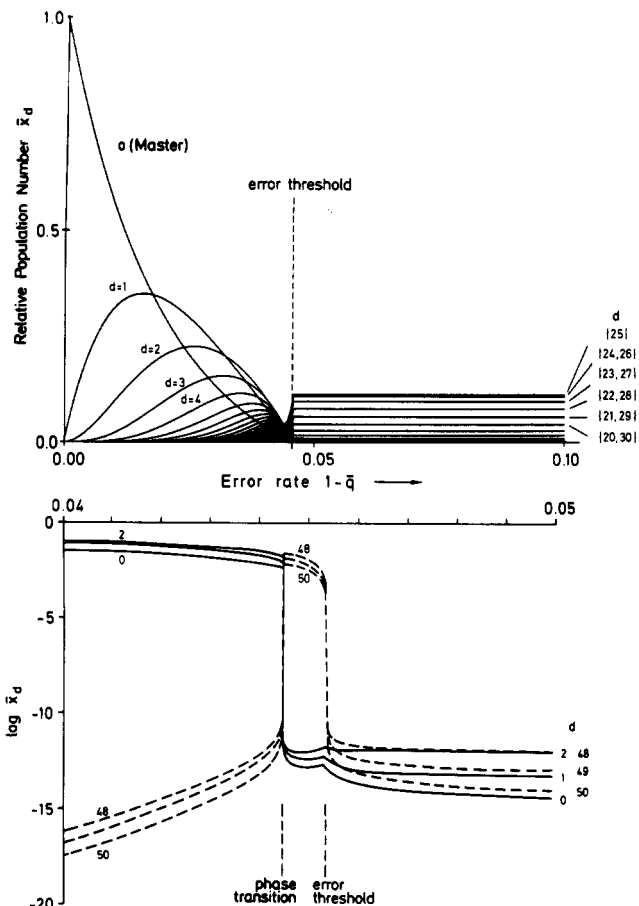
$$W_{00} > \bar{E}_{k\neq0}$$

Neglecting backflow of mutations—identifying $W_{00}$ with the eigenvalue $\lambda_0$—we find for the relative concentration of the master sequence

$$\bar{x}_0 = \frac{W_{00} - \bar{E}_{k\neq0}}{E_0 - \bar{E}_{k\neq0}} = \frac{Q_{00} - \sigma_0^{-1}}{1 - \sigma_0^{-1}} \qquad (29)$$

which is equivalent to the second-order perturbation result. To isolate the dependence on copying fidelity, we introduced a *superiority* parameter of the master sequence

$$\sigma_0 = \frac{A_0}{D_0 + \bar{E}_{k\neq0}} \qquad (30)$$

(31) McCaskill, J. S. *Biological Cybernetics* **1984**, *50*, 63.



**Figure 4.** Quasi-species as a function of the single-digit accuracy ($q$) for chain length $\nu = 50$. In the upper part we plot the relative concentration of the master sequence ($\bar{x}_0$), the sum of the relative stationary concentrations of all one-error mutants ($\bar{x}_1$), of all two-error mutants ($\bar{x}_2$), etc. We observe selection of the master sequence at $q = 1$. Then the relative concentration of the master sequence decreases with decreasing $q$. At the same time the relative concentrations of error copies increase. Above the critical error rate, $1 - q > 0.046$, the organized quasi-species is replaced by the uniform distribution, and hence the sum of the concentrations of the statistically (most probable) 25-error mutants is largest, followed by the 24- and 26-error mutants, etc. At this chain length ($\nu = 50$) the transition is very sharp already. In the lower part of the plot, we show a characteristic example where a sharp change in the quasi-species distribution occurs as $q$ is decreased prior to the error threshold. At intermediate $q$, the wild-type is no longer the sequence with the maximum replication rate. The transition is still sharp on the logarithmic scale applied. We show the region $0.96 > q > 0.95$. At the critical value $q_{tr} = 0.9555$ the new master sequence ($I_{50}$) and its neighbors become dominant. At $q_{min} = 0.9546$ the usual error threshold is observed. The following selective values were used in the two examples: in the upper part $A_0 = 10$ and $A_k = 1$ for all $k \neq 0$ and in the lower part $A_0 = 10$, $A_{50} = 9$, $A_{49} = 5$, and $A_k = 1$ for all $k = 1, ..., 48$. The index is given in parentheses when it refers to the whole error class.

Accordingly, the relative stationary concentration of the master sequence vanishes for some critical error rate that fulfils

$$(Q_{00})_{min} = \sigma_0^{-1} \qquad (31)$$

This, of course, relies on the neglect of mutational backflow. The phenomenological calculation of $\bar{x}_0$—based on the use of $\bar{D}_{k\neq0}$—can also be carried out exactly, leading to the results shown in Figure 3.

To make eq 31 more easily applicable to realistic cases, we adopt the uniform error assumption (6) involving the single-digit accuracy $q$:

$$q_{min} = 1/\sigma_0^{1/\nu} \qquad (32)$$

The nature of the error threshold is illustrated in Figure 4. There we use an extremely simple *single-peak* value landscape that is useful for mathematical analysis: a replication rate constant $A_0$

is assigned to the master sequence; all other types have the same selective values ($A_1 = A_2 = ... = A_{2^\nu-1} = A$) and all degradation rate constants are equal ($D_0 = D_1 = ... = D_{2^\nu-1} = D$). A mutant distribution that is localized around the master sequence $I_0$ changes at the critical minimum single-digit accuracy $q_{min}$ into the delocalized, almost uniform distribution. This transition, already abrupt at the chain length $\nu = 50$ shown in Figure 4, top, sharpens with increasing chain length. The error threshold thus exhibits apparent similarity to cooperative transitions in biopolymers or to equilibrium phase transitions. We shall take up this point again in section 6. Additional computations based on other, more complicated value landscapes[29] show the same qualitative behavior. Figure 4, bottom, shows how changing the error rate can cause a shift in focus of the quasi-species before a final error catastrophe. The former is marked as a phase transition in the figure to underline the criticality of the reorganization of concentrations, but it should be distinguished from the final error threshold transition. The analogy of the latter with equilibrium phase transitions is discussed in section 6. Clearly, the sharpness of the transition (especially apparent in the logarithmic scale of Figure 4, bottom) and its dependence on chain lengths are not consequences of the simplifications applied.

With the uniform error assumption, eq 31 may also be rewritten to work out a maximum chain length for constant single-digit accuracy $q$:

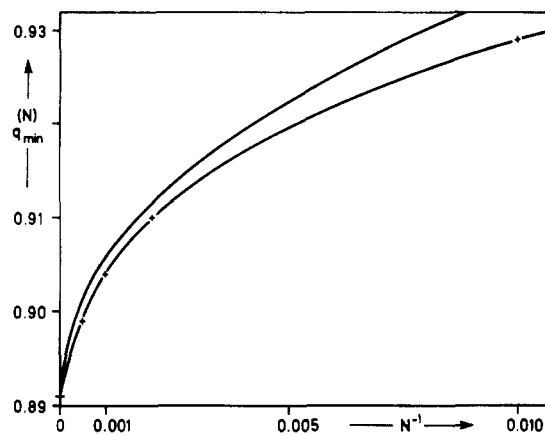$$\nu_{max} = -\frac{\ln \sigma_0}{\ln q} \approx \frac{\ln \sigma_0}{1 - q} \qquad (33)$$

Below the critical length the quasi-species is localized. The limitation in chain length, as predicted by eq 33, is particularly important for organisms that operate with a *replication machinery* of low accuracy. This is true for most RNA viruses. Indeed, we find that their genomes rarely exceed 10000 bases (see also section 7).

Fitness landscapes are extremely complicated objects, and their complete computation—for all sequences—is not possible. So it is highly desirable to be able to handle also cases with limited information, for example, systems in which only the statistics of the distribution of selective values is known. Such a distribution may be characterized by a probability density $\rho(W)$ which covers selective values up to the master sequence ($\rho(W) = 0$ if $W > W_{00}$).[17] Brillouin–Wigner perturbation theory up to infinite order was applied to obtain a formal solution for the wild-type fraction in the population for any particular set of replication rates. Statistical analysis of the perturbation expansion, over an ensemble of choices for the replication rates, was aided by first applying the Watson renormalization procedure.[32] The result showed that eq 33 can still be used with probability 1, provided there are no long-range correlations between the values of distant mutants, if the superiority parameter $\sigma_0$ is replaced by an effective superiority:

$$\sigma_{eff} = 1 + \left\langle \frac{W_{00} - W}{W} \right\rangle_{ln} \qquad (34)$$

where $\langle \ \rangle_{ln}$ is the geometric mean—expressed here as a *logarithmic average*: $\langle \alpha \rangle_{ln} = \exp(\langle \ln \alpha \rangle)$.

Most interestingly, the result demonstrates that the error threshold persists even in the near neutral limit where the distribution includes mutants with replication rates arbitrarily close to the maximum. Actually, a great deal more can be said if one notes that the maximum value $W_0$ is actually an extreme value of many samples from the underlying distribution of replication rates. For a large number of samples, as occurs in even a short period of quasi-species evolution, this extreme value is relatively independent of the form of the underlying distribution and depends only very weakly on the number of samples, $n$ say, that may be estimated to sufficient accuracy as the product of population size and evolution time in generations. The effective superiority may then be shown to have the form (for an underlying normal dis-

(32) Watson, K. M. *Phys. Rev.* **1957**, *107*, 1388.



**Figure 5.** Error threshold as a function of population size. Stochastic replication–mutation dynamics in an ensemble of binary sequences with chain length $\nu = 20$ was simulated by Gillespie's algorithm.[32] The critical single-digit accuracy of replication ($q_{min}$) at which the ordered quasi-species is converted into a changing population of sequences with finite life times is plotted as a function of reciprocal population size, $1/N$ (lower curve). The upper curve is the theoretical prediction according to eq 36.

tribution with mean $W$ and variance $v$, a similar result holding for other distributions):

$$\ln \sigma_{eff} = \frac{1}{2}\left(\frac{v}{W}\right)^2 + \ln \frac{W_0}{W} = \frac{1}{2}\left(\frac{v}{W}\right)^2 + \ln \frac{v}{W} + \frac{1}{2} \ln (\ln n) \qquad (35)$$

which makes the difference for the error threshold between evolutionary history on earth ($n \leq 10^{50}$) and a typical laboratory experiment in a test tube ($n \approx 10^{15}$) less than 50%. This further predicts absolute values of $\ln \sigma_{eff}$ that are of order 1—typically 3–5.

The error threshold relation was derived from kinetic equations and holds, strictly speaking, for infinite populations only. How sensitive is quasi-species localization to changes in population size? How small can a population be in which threshold phenomena are still observable? These and related questions are rather difficult to answer since their solution calls for an application of stochastic techniques. Such techniques are not routine for the description of nonlinear dynamical phenomena. A useful approximation to the population size dependence of the error threshold[31] has, however, been based on such a stochastic analysis:

$$\frac{\sigma Q_{min}(N)}{\sigma Q_{min}(N) - 1} = \alpha N \frac{Q_{min}(N) - Q_{min}(\infty)}{1 - Q_{min}(\infty)} \qquad (36)$$

where $\alpha$ is a small positive stability constant—$\alpha \approx 0.1$—and $N$ is the population size. Here, $Q_{min} = Q_{min}(N,\sigma,\nu)$ is the finite population size error rate at the threshold and $Q_{min}(\infty) = Q_{min}(\infty,\sigma,\nu)$ is the deterministic error threshold. Evaluation of the quadratic equation (36) yields an expression for the critical value $Q_{min}(N)$ and hence also for the critical single-digit accuracy for replication:

$$q_{min}(N) = (Q_{min}(N))^{1/\nu}$$

which is plotted and compared with computer simulation data[33] in Figure 5. In smaller populations, higher replication accuracy is required for quasi-species localization. This is to be expected, since stationary mutant distributions are *endangered* not only by accumulation of errors but also by natural fluctuations, which commonly fulfil a $N^{1/2}$ law and therefore become important in small populations. In addition, the computer simulations gave an answer to the second question. The smaller the population size, the more difficult it becomes to trace the position of the error threshold. Ultimately, around $N \approx 500$, variances becomes so large that the otherwise sharp transitions are smeared out by the scatter of data.

(33) Nowak, M.; Schuster, P. *Error Thresholds of Replication in Small Populations*, preprint.

So far, we have assumed that every part of the sequence to be replicated is equally important for the selective value of its carrier. There is ample evidence that this is not always so in nature. To give an example, the genes in some archaebacteria and in higher organisms contain coding and noncoding regions—so-called *exons* and *introns*—which are certainly under very different evolutionary constraints. The genomes of higher organisms are also full of DNA segments for which no evident usage is known at present. Can we apply the error threshold concept to such cases too? Clearly, the accuracy requirements have to be relaxed if parts of the genome contribute less or do not contribute at all to the selective values. Let us consider, for example, a sequence that consists of two different segments $A$ and $B$. Both segments exist in various versions, $I_1^A, I_2^A, ..., I_r^A$ and $I_1^B, I_2^B$, respectively. Every individual sequence can be characterized therefore as a particular combination of two segments $I_i^A$–$I_j^B$. The two segments need not be contiguous and might well appear in several disjoint parts. The concept of a segmented genome is easily incorporated into the kinetic differential equations (13). Under the assumption that the mutations during one replication occur either in segment $A$ or in segment $B$—but not simultaneously in both parts—a system of nonlinearly coupled differential equations was derived that allows one to handle the error propagation problem for individual segments.[34]

Here we consider only the simplest example: if segment $B$ has no specifically sequence-dependent influence on the replication rate of the polymer, the kinetic equations can be written as

$$\dot{x}_{i0} = f_B((A_iQ_{ii}^A - \bar{E}'(t))x_{i0} + \sum_{k \neq i} A_kQ_{ik}^A x_{k0});$$

$$i, k = 1, 2, ..., n \quad (37)$$

The variables $x_{i0}$ represent relative concentrations of all sequences carrying the segment $I_i^A$, irrespectively of the nature of segment $B$. The mean excess production $\bar{E}'(t)$—in this simplest case only—fulfills an expression identical with eq 12. The elements of the mutation matrix $Q^A = \{Q_{ik}^A\}$ refer exclusively to part $A$ and describe the frequencies of the mutations $I_k^A \rightarrow I_i^A$ with

$$\sum_{k=1}^{n} Q_{ik}^A = 1$$

Since mutations in the segment $B$ have no influence on the fitness of the carrier, this part can vary freely and only segment $A$ falls under the error threshold restriction. Consequently, the maximum chain length $\nu_{max}$ applies to part $A$ only.

In general, chain lengths $\nu$ are not conserved in polynucleotide replication: insertions and deletions as well as duplications of longer parts of sequences occur. The role of the factor $f_B$ in (37) is more subtle when sequences with different chain lengths compete. Whether or not a given polynucleotide is superior to its competitors of different length is determined largely by the chain-length dependence of the rate of replication. This dependence is rather complicated[8-10] and varies with environmental conditions. We expect an inequality $f_B \leq 1$ to hold on average above a certain length: longer sequences are replicated more slowly than shorter ones, but the greater variety in longer sequences may lead to occasional improvements.

The factor $f_B$ is an appropriate measure of the fitness of polynucleotides that carry unused or redundant parts. If the rate of replication does not depend on the chain length, $f_B$ will be close to 1 and longer sequences may compete successfully with shorter ones. Rates of replication that are largely insensitive to chain elongation are particularly important in the process of gene duplication. A coding sequence may be duplicated accidentally. The longer polynucleotide formed has to compete with shorter sequences for some intermediate period of time during which its fitness remains practically unchanged. If variation of the part that is not required for replication leads to a sequence that has a positive effect on the fitness of the whole polynucleotide, it will proliferate and eventually become dominant in the population.

The case when the unused part is translated and the translation products gains catalytic activity is particularly interesting. Successive gene duplications may lead to whole families of proteins with different substrate specificities.

The new catalytic function may eventually improve the accuracy of replication. Then gene duplication represents a powerful evolutionary mechanism to increase the catalytic capacity of a replicating system, which can in turn make replication more precise.

Present day structures of proteins provide some hints on the role of gene duplications in evolution. Proteins commonly occur in families. These are structurally related enzymes catalyzing reactions of the same class with different substrate specificities. The best studied examples are the families of proteases or dehydrogenases. One also observes interesting regularities in the structures of many globular proteins: substructures, so-called motifs, are often repeated exactly or with minor modifications only. Such repetitions occur within the same protein molecule as well as in different protein molecules. Both the modular structure of polymers and the existence of protein families may be readily explained by the gene duplication mechanism.

Last but not least, we may mention the exon–intron structure of eukaryotic genomes. Though error rates are generally adapted to the full genome size and throughout found to be smaller than $10^{-10}$, fidelity requirements may differ for exons and introns and even vary within both. For instance, part of the intron regions may have signal character and therefore ought to be highly conserved. The concept of the relaxed error threshold, introduced above, will allow a quantitative treatment of such cases.

## 6. Replication Dynamics and Equilibrium Statistical Mechanics

The threshold behavior of quasi-species localization, as reported in the previous section, is reminiscent of cooperative transitions observed with conformational equilibria in biopolymers. Equilibrium statistical mechanics of lattice models, in which the states of the system are combinations of local states defined at the points of a lattice, and the dynamics of Markov processes can indeed be described within the same mathematical discipline of *Markov random fields*,[35] and therefore, similarities in global behavior such as the existence of cooperative phenomena or phase transitions are not completely unexpected. Spin lattice models, in which the local states are one of a finite number of alternatives (as in the discrete spin states or magnetic moments in quantum theory), are particularly well-suited candidates, because several of them are sufficiently simple to allow derivations of analytical expressions for thermodynamic functions.

The basis of two recent attempts to work out analogies between replication dynamics[36,37] and spin lattice models[24,38] is a discrete dynamical system—modeled by a multitype branching process as shown in Figure 6. The distribution of sequences in the discrete system converges asymptotically toward the stationary solution of the differential equations (13) for replication–mutation dynamics.[39] Expected stationary distributions are derived from a proper statistics of *genealogies*. A genealogy is an individual time-ordered series of sequences that represents one particular recording of successive descendants: each member of a given genealogy was synthesized through a copying of its precursor as template (Figure 6).

The individual genealogy is identified with a spin–lattice. A one-dimensional array of *generalized spins* may be used for this purpose. The individual *spin* has $\kappa^\nu$ different states which represent different sequences $I_i$, $i = 1, ..., \kappa^\nu$. The discrete time of the dynamical system corresponds to the spatial dimension of the lattice. A genealogy is the analogue of one particular array of

(34) Eigen, M.; Schuster, P. *Relaxed Error Thresholds of Polynucleotide Replication*, preprint.

(35) Kindermann, R.; Snell, J. L. In *Contemporary Mathematics*; American Mathematical Society: Providence, RI, 1980; Vol. 1.
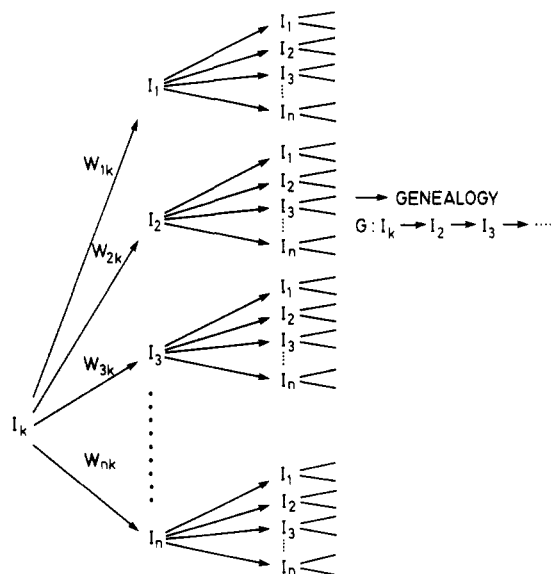
(36) Leuthäusser, I. *J. Chem. Phys.* **1986**, *84*, 1884.

(37) Leuthäusser, I. *J. Stat. Phys.* **1987**, *48*, 343.

(38) Demetrius, L. *J. Chem. Phys.* **1987**, *87*, 6939.

(39) Demetrius, L.; Schuster, P.; Sigmund, K. *Bull. Math. Biol.* **1985**, *47*, 239.

**Figure 6.** Polynucleotide replication as a multitype branching process. The probabilities of the different pathways are given by the elements of the value matrix **W**. The genealogy (G) of the branching process may be considered as an analogue of a particular one-dimensional array of generalized spins. Every generalized spin has $n = 2^\nu$ different states.

spins. Stationary mutant distributions are related to distributions of spin states as derived from the equilibrium statistics of spin lattices. The equivalence is proved by means of an optimization principle[24] which can be derived by application of a thermodynamic Gibbs measure to the statistics of genealogies.
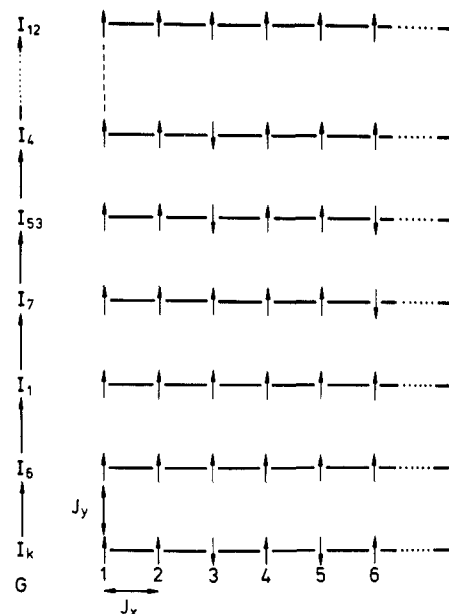
If we consider binary sequences replicating within the uniform error approximation, the model can be made more specific: individual sequences are now identified with the rows of two-dimensional spin lattices of size $\nu \times N$ (Figure 7). Since there are two classes of digits (0,1) in binary sequences, ordinary electronic spins ($m_s \pm 1/2$) are the proper magnetic equivalents. The partition function of the spin lattice, within the approximations of the Ising model, can be factorized by means of a ($2^\nu \times 2^\nu$) transfer matrix **P**, where $2^\nu$ refers to the number of different one-dimensional spin arrays. Identification of the transfer matrix **P** with the value matrix **W** from replication dynamics provides a relation between temperature $T$ and single-digit accuracy $q$:

$$\frac{1}{T} \Leftrightarrow |\log \epsilon| = \left|\log \frac{1-q}{q}\right| \qquad (38)$$

Localization thresholds appear now as an analogue of magnetic order–disorder transitions in spin lattices. Replication degenerates indeed to random production of sequences in the limit $q \to 1/2$ and corresponds to the limit $T \to \infty$, the state of maximum disorder. The localization of the sequence distribution and quasi-species formation are analogous to the appearance of a ferromagnetic phase in the spin system. An increase of the single-digit accuracy $q$ in the range $1/2 < q < 1$ is tantamount to the lowering of temperature $T$ until at absolute zero, the state of perfect order is reached. Zero kelvin thus corresponds to the limit of correct replication, $\lim q \to 1$ where the quasi-species converges toward a pure master sequence: the ultimate limit of localization is reached. The interaction in the vertical ($y$) direction of Figure 7 is given by the regularities of the mutation matrix, **Q**. It is well described by an Ising-type model. Interaction energies in the horizontal ($x$) direction represent the rate constant of replication of the corresponding sequences: $A_i$ ($i = 1, 2, ..., 2^\nu$). To approach realistic value landscapes, the restriction to nearest-neighbor spin interactions—as used in the Ising model—has to be abandoned in favor of more general spin glass Hamiltonians.[37]

**7. Conclusion**

The deterministic treatment of the replication–mutation model reveals three major regularities that may best be encompassed by the following headings:



**Figure 7.** Polynucleotide replication and two-dimensional spin lattices. In the two-dimensional model a genealogy is represented by a two-dimensional spin lattice. The individual spins exist in two states ($m_s \pm 1/2$) corresponding to the two digits in binary sequences. Every row of the lattice consists of $\nu$ digits and corresponds to a polynucleotide sequence. The *in row* interaction, described by the spin–spin coupling constants, $J_x$, is a property of the individual sequence and contributes to the rate constant of replication. The *vertical* coupling constant $J_y$, on the other hand, is a measure of the mutation frequency.

*1. Selection*: local stabilization of a quasi-species distribution in sequence space around one—or several degenerate—master species.

*2. Evolution*: destabilization of the local quasi-species upon arrival of an advantageous mutant that establishes a new quasi-species.

*3. Optimization*: tendency toward global stabilization guided by the population of nearly neutral mutants and their nonrandom distribution in sequence space.

To highlight the changes in the understanding of *Darwinian* systems introduced by this physicochemical model, in the chemical limit of large populations and sufficiently small genomes, we confront the current interpretation with the classical view under these three headings:

**Selection** may be characterized as a condensation phenomenon in sequence space that shows analogy to order–disorder transitions in condensed matter. The role of the critical temperature is played by the error threshold. Far below error threshold the quasi-species contracts and finally—at a state corresponding to $T = 0$ K—only the master sequence, the sequence corresponding to the highest peak in the value landscape, is populated. Surpassing the threshold means *melting* of the quasi-species due to accumulation of errors. Such an error catastrophy means a sharp loss of genetic information.

The *wild-type* is not a single individual but rather a distribution having a defined consensus sequence that usually coincides with the master sequence. Despite the fact that the whole system is often summarized by a single sequence—it is obtained, for example, in sequence determination by conventional methods—its fraction of the total population is usually very small, in fact undetectably small in the cases where it has been tested. Likewise, the term *fittest* is not related to any individual but rather to the complete quasi-species distribution acting as the target of selection. The *population of mutant states* in the quasi-species is strongly modulated by the fitness distribution. The effect is particularly strong for those mutants that are almost neutral—these are the mutants whose selective values are only slightly smaller than that of the master sequence. Population numbers depend not only on the selective values of individuals but also on those of neighboring mutants. The quasi-species typically has long asymmetric pro-

trusions that are of central importance to its evolution (see below).

The error threshold has been tested for several virus populations. This is not an easy task because both single-digit accuracies and superiority parameters are hard to determine experimentally. To measure error rates requires special care, and superiorities can be computed readily only if some features of the stationary mutant distribution are known. The difficulties in interpreting experimental data were discussed recently.[1] Here, we give only the results: the RNA Coli-phage $Q\beta$[11] and the animal RNA viruses foot-and-mouth disease virus,[40] influenza A virus,[13] and vesicular stomatitis virus[41] seem to operate under conditions close to the threshold. The actual lengths of their genomes are only slightly smaller than the maximum lengths computed by means of (33). We can think of two factors driving these small viruses toward high error rates. The parasite has to cope efficiently with the evolution of the host's protective measures, and this requires fast adaptation or as many mutants as possible. Second, the small genome of the virus does not allow much evolutionary flexibility, and therefore elongation of the RNA sequence would be of advantage. This can be done only below the error threshold so that the length available for coding protein is fixed near the maximum value and has to be used as well as possible. Indeed, we find multiple usage of the available genetic information by coding two proteins on the same piece of RNA by different reading frames or single stop-condons with a certain *read through* frequency.

**Evolution** may be viewed as a series of stabilizations and destabilizations of *localized* quasi-species that, in a constant environment, are associated with an increase in fitness. The evolutionary route avoids the vast regions of low fitness in sequence space by guidance along the ridges of high selective values. This process may occur in *jumps* if advantageous mutants are very rare or, in the other extreme, selection and evolution may coalesce into one process approaching the global fitness maximum if most mutations lead to advantageous variants. Evolutionary progress is expected to be greatest near the error threshold. Occasionally surpassing this limit for a short time may assist in the escape from metastable distributions in order to find higher fitness peaks, just as *simulated annealing* is used for efficient optimization.

Neutral or nearly neutral *mutants* appear in a new light. They are of utmost importance in determining the highly populated states at the periphery of the mutant spectrum and hence fix the route of evolution. As part of the quasi-species, they are rated not only by their selective values relative to the master but also with respect to their own mutant environment. The uncertainty of classical theory—how closely must the fitness of a mutant resemble that of the wild-type to be classified as neutral—is now replaced by new quantitative expressions. The results of the conventional *neutral theory* are valid exclusively for systems of low population numbers and large genomes: if the population number $N$ is smaller than the number of different one-error mutants—which amount to $3\nu$—not even the nearest neighbors in sequence can be populated and then the results of the neutral theory are expected to be representative. Otherwise the assumption of *blind* production of mutants has to be modified along the lines of the quasi-species concept.

The first test-tube experiments on evolutionary phenomena were done by Spiegelman.[4] He isolated the genome of the Coli-phage $Q\beta$—an RNA molecule 4220 bases long—and succeeded to replicate it in a cell-free medium that consists of an enzyme, $Q\beta$-replicase, and activated monomers, the nucleoside triphosphates, in excess. A suitable constraint is imposed on the system by serial transfer consisting of repetitions of up-growth and dilution phases over many generations. The viral RNA adapts to the new environment—consisting of ideal conditions for replication and release of all sequence constraints of coding protein—by *throwing away* about 90% of its chain. Small RNA molecules—a few hundred bases long only—are selected out of the serial transfer mixture that indeed grow much faster than the wild-type but are no longer infectious to *Escherichia coli* cells. One might well call this process degenerative or *retrograde* evolution, but it shows that adaptation to the environment is not restricted to living organisms. It occurs inevitably when the initially mentioned conditions are fulfilled.

The discovery of de novo synthesis of RNA by $Q\beta$-replicase[6,7,42,43] opened a new way of studying evolution processes in the laboratory under more favorable conditions. During de novo synthesis many different templates are being formed by individual enzyme molecules and then compete for growth—leading finally to the out-growth of a *fittest* sequence. Those experiments can be carried out in presence of conditions that exert specific selection pressures. Ethidium bromide, for example, provides such a selection pressure because it interferes with the replication process by intercalation between base pairs.

**Optimization** is not just a blind stochastic trial-and-error search for a better adapted mutant, but rather follows an inherent logic:

Selective advantage usually is to be expected in a distant mutant—at the periphery of the distribution.

Mutants far distant from an established master sequence arise from those that are less distant—from precursors en route.

The probability of producing a distant mutant depends critically on the population numbers of its precursors.

Population numbers of mutants are high relative to their equidistant competitors if they are situated in a domain of high fitness.

A fractal clustering of the fitness distribution in domains is likely to occur.

Precursors of distant advantageous mutants therefore have a higher chance to be populated than precursors of deleterious mutants. This results in a tremendous amplification of population numbers along high-value ridges.

The high dimensionality of sequence space aids the connectivity between nearly neutral mutants.

Evolutionary optimization proceeds along defined pathways in sequence space. There are alternative routes, but their number is so highly restricted that one has the impression of a phenomenon of automatic guidance to higher fitness. Optimization routes are to this degree deterministically ordained.

This logic, rather than the simplifying interpretation of Darwinian behavior as an interplay of *chance*—random mutation—and *necessity*—selective fixation of the advantageous mutant once produced—is the basis of evolutionary optimization or adaptation in the simplest living systems. The further study of these phenomena will be greatly facilitated by a serial transfer machine, designed to allow massively parallel evolution experiments, recently completed in our laboratory.[44]

(40) Domingo, E.; Davilla, M.; Ortin, J. *Gene* **1980**, *11*, 333.
(41) Spindler, K. R.; Horodyski, F. M.; Holland, J. J. *Virology* **1982**, *119*, 96.

(42) Sumper, M.; Luce, R. *Proc. Natl. Acad. Sci. U.S.A.* **1975**, *72*, 162.
(43) Biebricher, C. K.; Eigen, M.; Luce, R. *Nature (London)* **1986**, *321*, 89.
(44) Otten, H. Dissertation, Technical University, Braunschweig, 1988.