# Evolution on *"realistic"* fitness landscapes[*]

## Phase transitions, strong quasispecies, and neutrality

By

Peter Schuster[†]

**Abstract:** Evolution through mutation and selection in populations of asexually replicating entities is modeled by ordinary differential equations (ODEs) that are derived from chemical kinetics of replication. The solutions of the mutation-selection equation are obtained in terms of the eigenvectors of the value matrix $W = Q \cdot F$ with Q being the matrix of mutation frequencies and F the diagonal matrix of fitness values. The stationary mutant distribution of the population is given by the largest eigenvector of W called quasispecies $\bar{\Upsilon}$. In absence of neutrality a single variant, the master sequence $\mathbf{X}_m$, is present at highest concentration. The stationary frequency of mutants is determined by their Hamming distance from the master, by their fitness values, and by the fitness of their neighbors in sequence space. The quasispecies as a function of the mutation rate, $\bar{\Upsilon}(p)$, may show a sharp transition from an ordered regime into the uniform distribution ($\Pi$) at $p = p_{\mathrm{cr}}$ that is called error threshold. Three phenomena that are separable on model landscapes coincide at $p = p_{\mathrm{cr}}$: (i) steep decay in the concentration of the master sequence, (ii) phase transition like behavior, and (iii) a wide range of random replication where $\bar{\Upsilon}$ is the uniform distribution. "Realistic" model landscapes based on current knowledge of nucleic acid structures and functions show error thresholds but also other sharp transitions, where one quasispecies distribution is replaced by another quasispecies with a different master sequence at critical mutation rates $p = p_{\mathrm{tr}}$. Groups of nearest neighbors of high fitness are strongly coupled by mutation, behave like a single entity and are unlikely to be replaced in phase transitions. Such "strong quasispecies" – consisting of a master sequence and its most frequent mutants – maintain their identity over the entire range of mutation frequencies from $p = 0$ to the error threshold at $p = p_{\mathrm{cr}}$. Neutrality in the sense of identical fitness values for two or more sequences with Hamming distances $d_{\mathrm{H}} < 3$ leads to strongly coupled clusters of variants, which remain stable in the limit $\lim p \to 0$. Nearest neighbor or next nearest neighbor sequences appear at fixed ratios in the stationary distributions. Random selection of sequences by random drift occurs only at Hamming distances $d_{\mathrm{H}} \geq 3$.

**Key words:** error threshold – fitness landscapes – molecular evolution – mutation rates – neutral evolution – quasispecies – virus populations.

---

[†]Address: Institut für Theoretische Chemie der Universität Wien
   Währingerstraße 17, A-1090 Wien, Austria
   E-Mail: pks @ tbi.univie.ac.at , and
   Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA
   E-Mail: pks @ santafe.edu .

# 1 Introduction

About forty years ago Manfred Eigen [1] conceived a molecular theory of evolution based on chemical kinetics of replication and mutation as well as other knowledge from molecular biology, which has been extended, worked out in detail and presented in a number of publications [2–6]. Population dynamics is described by means of ordinary differential equations (ODEs) as in chemical kinetics. The replicating molecular entities are RNA or DNA molecules in test tube evolution experiments or genomes in the evolution of viroids, viruses, bacteria or higher organisms. In the original formulation the theory was thought to provide the proper frame for understanding pre-biotic evolution [1, 4, 7]. The theory is focused around two new concepts: (i) the *quasispecies*, which is the stationary mutant distribution in an asexually reproducing population of genotypes related by mutation, and (ii) the *hyper-cycle*, which is the simplest functional complex that suppresses competition of reproducing entities and allows for formation of an ensemble through cyclic mutual dependence. The theory of quasispecies has provided new insight into populations of viruses and was and is successfully applied to virus evolution and the development of novel antiviral therapies.

The quasispecies consists of a fittest genotype, the *master sequence*, surrounded by a cloud of frequent mutants and it replaces the notion of a single wild type genome in populations operating at high mutations rates like, for example RNA virus populations [8], as used in conventional population genetics. The most important result of quasispecies theory is the existence of an error threshold: Error accumulation leads to a complete breakdown of inheritance, if the (single point) mutation rate parameter $p$ is larger than some critical value, $p > p_{\mathrm{cr}}$. Accordingly, Darwinian evolution requiring ordered reproduction or at least partial conservation of genomes is possible only the range $0 \leq p \leq p_{\mathrm{cr}}$, which sets an upper bound to acceptable mutation rates, $p_{\max} = p_{\mathrm{cr}}$. Values for $p_{\mathrm{cr}}$ as a function of genome length ($\nu$) and fitness values ($f_j$) were derived by an approximation neglecting mutational backflow consisting of mutations from mutants back to the master sequence [1]. The first combined analytical and numerical study on a quasispecies with explicit consideration of all sequences in sequence space was published in 1982 [9] and gave insights into the structure of stationary mutant distributions as a function of the mutation rate parameter ($p$). Both approaches used, indirectly or directly, a simple distribution of fitness values: The master sequence has higher fitness than the rest of the population, which is assumed to have identical fitness values. Later the name *single peak fitness landscape* has been
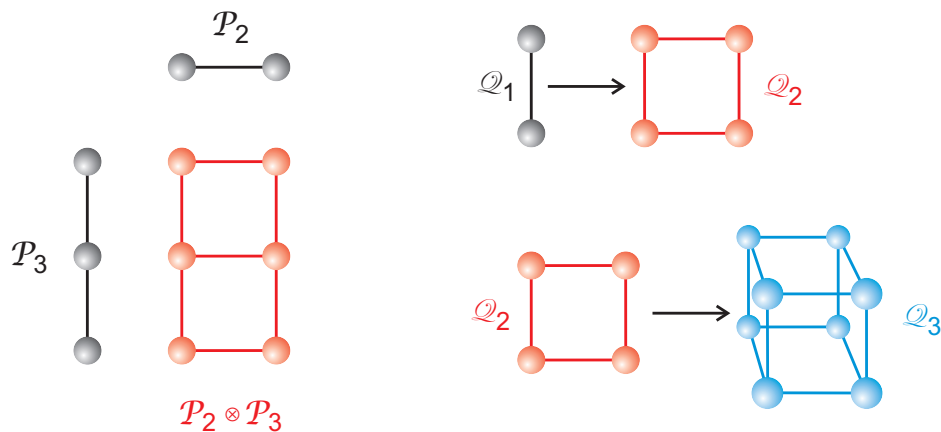
2

given to this simple fitness distribution. The transition from the ordered to the *random replication*[2] regime is sharp and the sharpness increases with the chain length of polynucleotides, $\nu$. At constant replication accuracy the error threshold defines a maximum chain length ($\nu < \nu_{\max}$) and this provides a limit for the length of genomes or polynucleotide sequences that can be reproduced faithfully. The error threshold has been applied to prebiotic scenarios in order to estimate the limitation of chain elongation in enzyme free replication [7] and by the same token the threshold defines an upper bound for the genome lengths of viruses, which are replicated by low accuracy replicases without proof reading [10]. Alternatively, the error threshold defines a maximal mutation rate, $p < p_{\mathrm{cr}} = p_{\max}$, for replication of genomes of the same chain lengths and has been exploited for the development of new antiviral drugs that increase the mutation rate through interfering with the replication of virus genomes [11].

This review consists of two parts: (i) an updated overview of the conventional quasispecies theory that focusses of the landscape concept (sections 2-5) and (ii) insights into quasispecies dynamics derived from the use of conventional and new classes of landscapes (sections 6-8). After a brief introduction into the landscape concept, chemical kinetics of replication is introduced by means of the flow reactor as as appropriate open system. Then, Charles Darwin's natural selection is modeled by means of the mathematics that was known at Darwin's lifetime already. Quasispecies theory is reviewed and the concept of the error threshold is introduced. Subsection 5.6 contains some new results on the nature of the error threshold and the approximations applied in the derivations. The second part starts by distinguishing simple and realistic landscapes. "Realistic"[3] landscapes are constructed in a way that mimics current knowledge on landscapes derived from biopolymer structure and function as well as from virus evolution. In essence, three new results were obtained. First, the error threshold phenomenon consists of a superposition of three features that coincide on some model landscapes and on the "realistic" landscapes studied here: (i) a steep decay of the stationary concentration of the master sequence, $\bar{x}_m(p)$, at small mutation rate parameters $p$, (ii) a sharp transition of the quasispecies $\bar{\Upsilon}(p)$ at the critical value $p = p_{\mathrm{cr}}$, and (iii) an extension of the domain of random replication from the point $p = \tilde{p} = \kappa^{-1}$ to a broad plateau covering the whole range $p_{\mathrm{cr}} < p \leq \tilde{p}$.

---

[2]Random replication expresses the fact that error accumulation destroys the relation between template and copy and inheritance is no longer possible.

[3]Realistic is put here between quotation marks, because we want to indicate that the knowledge on detailed shapes of landscapes is still incomplete.

**Figure 1: Path graphs and binary sequence spaces**. The structure of binary sequence spaces $\mathcal{Q}_l$ follows straightforwardly from the graph Cartesian product of the path graph $\mathcal{P}_2$. The definition of the graph Cartesian product is shown on the lhs of the figure. The sketch on the rhs presents the construction of binary sequence spaces: $\mathcal{Q}_1 = \mathcal{P}_2$, $\mathcal{Q}_2 = \mathcal{P}_2 \otimes \mathcal{P}_2$, $\mathcal{Q}_3 = \mathcal{P}_2 \otimes \mathcal{Q}_2 = \mathcal{P}_2 \otimes \mathcal{P}_2 \otimes \mathcal{P}_2$, and so on, and $\mathcal{Q}_l = (\mathcal{P}_2)^l$. $\mathcal{Q}_\nu$ is a hypercube of dimension $\nu$.

The choice of suitable model landscapes allows for a separation of the three features. Second, some "realistic" landscapes sustain quasispecies that are not dominated by a single master sequence but by a cluster of commonly four neighboring sequences in sequence space with high fitness values. These clusters show unusual stability, are not replaced in phase transitions, and were called *strong quasispecies* therefore. Third, pairs and groups of neutral sequences with Hamming distances $d_\mathrm{H} = 1$ and $d_\mathrm{H} = 2$ form strongly coupled clusters that are stable in the entire range of mutation rates up to the error thresholds. No such coupling exists for more distant pairs of sequences $(d_\mathrm{H} \geq 3)$, which are subjected to Kimura's random selection therefore. A digression on lethal mutants (section 9) and an account on limitations and perspectives of the quasispecies concept (section 10) finish the review.

## 2   Combinatorial spaces and landscapes

The replicating molecular entities are understood as strings **X** over some alphabet $\mathcal{A}_\kappa$ with $\kappa$ digits. Two alphabets are frequently used: (i) the binary alphabet $\mathcal{A}_2 = \{\mathbf{0}, \mathbf{1}\}$ with $\kappa = 2$ and the natural nucleobase alphabet $\mathcal{A}_4 = \{\mathbf{A}, \mathbf{U}, \mathbf{G}, \mathbf{C}\}$ with $\kappa = 4$.[4] Binary sequences are used here mainly for

---

[4]Here we shall be mainly dealing with RNA molecules. For DNA sequences **U** is automatically replaced by **T**.

simplicity since they show many features of four letter sequences but can be handled much more easily. It is also worth noticing that RNA molecules with reduced alphabets make perfect structures with ribozyme[5] functions (see [12] for $\kappa = 3$ and [13] for $\kappa = 2$). Without simplification binary sequences can be used to code for sequences over four letter alphabets (see subsection 5.1). Polynucleotide sequences or genomes are viewed as elements of the entire set of possible sequences, which constitutes a formal metric space called *sequence space*. All possible sequences of a given length $\nu$ form a combinatorial manifold and hence sequence space belongs to the class of combinatorial spaces [14].

The idea of sequence space without the explicit name has been used to order strings in informatics for quite some time (see, e.g., [15]) before the word has been coined for proteins [16] and nucleic acids [1]. The sequence space for binary sequences is a hypercube $Q_\nu$ of dimension $\nu$ where $\nu$ is the length of the string. The building principle of sequence spaces by means of the graph Cartesian product is illustrative and can be used for sequences over arbitrary alphabets and, in particular, also for the natural **AUGC** alphabet. The Cartesian product of two graphs is illustrated in figure 1 by means of two path graphs:[6] The product graph, $\mathcal{P}^{(1)} \otimes \mathcal{P}^{(1)}$ is two-dimensional and carries $\mathcal{P}^{(1)}$ on its horizontal and $\mathcal{P}^{(2)}$ on its vertical margin, respectively. There are many ways to visualize binary sequence spaces as hypercubes – one, the consecutive product of $\mathcal{P}_2$ graphs is illustrated in figure 1:

$$\mathcal{Q}_\nu \;=\; \mathcal{P}_2 \otimes \mathcal{P}_2 \otimes \ldots \otimes \mathcal{P}_2 \;=\; \left(\mathcal{P}_2\right)^\nu . \tag{1}$$

The construction of sequence spaces a graph Cartesian products has the advantage of being generalizable. If we choose a complete graph $\mathcal{K}_\kappa$ as unit the consecutive Cartesian product yields the corresponding sequence space for sequences of chain length $\nu$:

$$\mathcal{Q}_\nu^{(\kappa)} \;=\; \mathcal{K}(l,\kappa) \;=\; \mathcal{K}_\kappa \otimes \mathcal{K}_\kappa \otimes \ldots \otimes \mathcal{K}_\kappa \;=\; \left(\mathcal{K}_\kappa\right)^\nu . \tag{2}$$
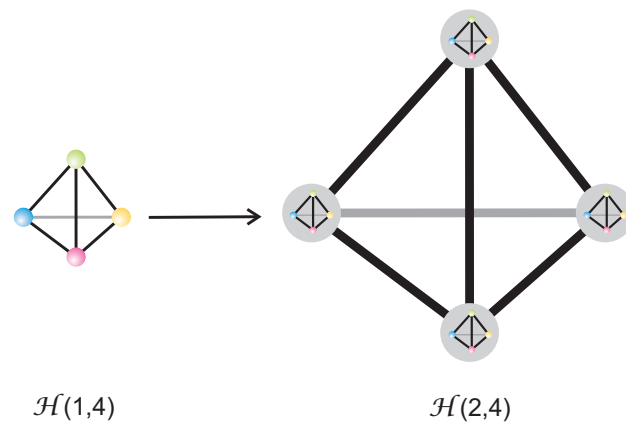
The most important case is the natural alphabet with $\kappa = 4$ (figure 2).

Sequence spaces are characterized by high dimensionality, and this makes it difficult to imagine distances. Considering, for example, the binary sequence space for strings of chain length $\nu = 10$ that contains $2^{10} = 1024$ sequences. Were sequence space a (one dimensional) path graph, the longest

---

[5]A ribozyme is a catalyst built from RNA with functions like a protein enzyme. The name is derived from **ribo**(nucleic acid en)**zyme**.

[6]A path graph $\mathcal{P}_n$ is a one-dimensional graph with $n$ nodes. Two nodes at the ends have vertex degree one and all other $n - 2$ nodes have vertex degree two.

**Figure 2: Four letter sequence spaces**. The sequence space derived from the four letter alphabet (**AUGC**; $\kappa = 4$) are the Hamming graphs $\mathcal{H}(l,4)$. The Hamming graph for a single nucleotide is the complete graph $\mathcal{H}(1,4) = \mathcal{K}_4$ (lhs) and for the 16 two letter sequences the space is $\mathcal{H}(2,4) = \mathcal{K}_4 \otimes \mathcal{K}_4$ (rhs). The general case, the space of sequences of chain length $\nu$, $\mathcal{H}(\nu,4)$ is the graph Cartesian product with $\nu$ factors $\mathcal{K}_4$.

distance would be 1023. In two dimensions is would be 62 and on the hypercube $\mathcal{Q}_{10}^{(2)}$ it is shrunk to only 10. The most natural metric for sequence spaces is the Hamming distance $d_{\mathrm{H}}$.[7]

The concept of a fitness landscape goes back to a metaphor used by the American population geneticist Sewall Wright [18]: Fitness is plotted on a support that is given by a (discrete) *recombination space* where each possible allele combination of the genotype is represented by a node [19, 20]. Populations and subpopulations of species are assumed to migrate on such fitness landscapes until they find a local fitness optimum where they stay until a change in the environmental conditions eliminates the peak. Wright assumed fitness landscapes to be rugged with a multitude of peaks, but in the days of his publication fitness was a rather abstract quantity that was inaccessible to precise measurement and moreover very little was known on the allele composition of genomes. The situation has changed completely over the years, the fitness values of individual variants – molecules, viruses or bacteria – can now be measured (see, e.g. [21–24]), and a rapidly growing body of knowledge in particular in virology sets the stage for a quantitative population biology at the molecular level. A recent example of a large scale

---

[7]The Hamming distance named after Richard Hamming counts the number of positions in which two aligned sequences differ. The appropriate alignment of two sequences requires knowledge of their functions [17]. Here, we shall be concerned only with the simplest case: end-to-end alignment of sequences of equal lengths.

study on HIV-1 aims at an understanding of the fitness landscape and the changes of viral fitness on medication [25]. The size of sequence spaces, however, is (still) prohibitive for an exhaustive experimental determination of fitness landscapes even in the simplest cases. Consequently the application of model landscapes is unavoidable (see, e.g., [26]). Simple examples of such landscapes are the *additive fitness landscape* (see e.g. [27]) where each mutation away from the fittest sequence causes the same additive detrimental contribution and, the *multiplicative fitness landscape* (see, e.g., [28] or [29]) where the detrimental contribution is a factor, or the *single peak fitness landscape* [9]. More complex landscapes are based on a deterministic as well as a random components. The so-called *Nk-model* proposed and developed by Stuart Kauffman [30, 31] may serve as a prominent example.

Here we intend to combine knowledge from experimentally measured fitness values, mutation studies on biomolecular structures as well as plausible inter- and extrapolations in order to get the required information for computing and analyzing quasispecies dynamics. The name *"realistic" rugged fitness landscapes* (RRL) is suggested to be used for landscapes that are based on three parameters: (i) the fitness of the fittest sequence, the master sequence $\mathbf{X}_m$ with fitness $f_m$, (ii) the mean fitness of all sequences with the exception of the master sequence, $\bar{f} = \bar{f}_{-m}$, and (iii) the band width $d$ of a random scatter of fitness values around the mean value $f$. The details of the scatter can be defined individually, for example, by the seeds of the pseudorandom number generator. A major problem of the random approach concerns the relation to experimental data: How can we reduce the number of required empirical parameters in such a way that a direct comparison between theoretical predictions and experimental measurements becomes possible? Here we make an attempt in this direction, which is based upon a combination of mathematical analysis and numerical computation.

Population dynamics of sequences that have the same fitness is, in principle, not accessible through ODEs because the time development of an ensemble of *neutral sequences* is a purely stochastic process without a deterministic drift term. Ever since Motoo Kimura's approach [32, 33], which allowed for handling *neutrality* within population genetics, neutral evolution is understood as a diffusion like process in sequence space that requires a stochastic approach. The major result of the *neutral theory* is that the loss of variants and the approach towards a homogeneous population does not require fitness differences. Since the outcome of the stochastic process is completely undetermined, the notion of *random selection* characterizes the phenomenon well. It is worth noticing that the time span required for the approach towards ho-

mogeneity in the population increases with decreasing fitness difference until it reaches a finite maximal value for differential fitness zero. It is also important to recognize that slightly smaller fitness does not lead to elimination of the less fitter variant in finite time as described by the *nearly neutral theory* of evolution developed by Tomoko Ohta [34–36]. Closely related neutral sequences behave differently from the predictions of Kimura's neutral theory (see section 8).

## 3   Replication in the flow reactor

Selection and evolution cannot take place at thermodynamic equilibrium. Required is an open system that exchanges matter and energy with an environment. The flow reactor (figure 3) is sufficiently simple and suitable for both theoretical modeling and experimental work particularly in chemical engineering [37]. Reactors with similar functions but extensive external control were and are used in microbiology and have different names characterizing the special conditions of the cell culture like chemostat [38, 39], turbidostat [40, 41] or cellstat [42]. Dispensing from all technical details two features of a continuously stirred tank reactor (CSTR) as shown in figure 3 are important: (i) Continuous influx of a stock solution provides the materials that are consumed by the reactions in the reactor and the increase in volume is compensated by an equal outflux of reactor solution, and (ii) spatial homogeneity is assumed to be achieved (almost) instantaneously by efficient stirring. In addition, temperature control is assumed to be provided by a heat bath, in other words heat produced or consumed by the reactions in the reactor is considered to be compensated by an isothermal environment.
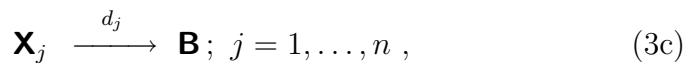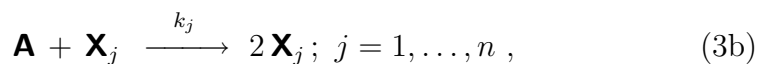
The population dynamics of $n$ classes of replicating entities or replicators, $(\mathbf{X}_1, \mathbf{X}_2, \ldots, \mathbf{X}_n)$ with the concentrations $\mathbf{x}^{\mathrm{t}} = (x_1, x_2, \ldots, x_n)^{\mathrm{t}}$ [8] is modeled comprehensively by a system of irreversible chemical reactions embedded in the flow reactor, which as said provides an energy and material flow in order

---

[8]In general vectors are understood as $(n \times 1)$ matrices. In order to save printing space we shall commonly write them in transposed form indicated by 't'.

to sustain non-equilibrium conditions

$$\star \xrightarrow{\;k\;} \mathbf{A} \;, \tag{3a}$$

$$\mathbf{A} + \mathbf{X}_j \xrightarrow{\;k_j\;} 2\,\mathbf{X}_j \;;\; j = 1, \ldots, n \;, \tag{3b}$$

$$\mathbf{X}_j \xrightarrow{\;d_j\;} \mathbf{B} \;;\; j = 1, \ldots, n \;, \tag{3c}$$

$$\mathbf{B} \xrightarrow{\;h\;} \mathbf{A} \;, \quad \text{and} \tag{3d}$$

$$\mathbf{A} \,,\; \mathbf{X}_{(i)} \,,\; \mathbf{B} \xrightarrow{\;d\;} \oslash \;. \tag{3e}$$

The molecular species $\mathbf{A}$ stands for the material required to synthesize a molecule $\mathbf{X}_{(j)}$ on a molecule $\mathbf{X}_{(j)}$ acting as template in the sense of a copying process. The species $\mathbf{A}$ enters the system in a zeroth order reaction with a constant rate $k$ and is either used in the synthesis of one of the $\mathbf{X}_{(j)}$ molecules, or degraded with a first order rate $d\,[\mathbf{A}]$.[9] The replicators $\mathbf{X}_{(j)}$ are either degraded with a rate $d_j\,[\mathbf{X}_{(j)}]$ to yield compounds $\mathbf{B}$, which can be recycled with rate parameter $h$ in order to return $\mathbf{A}$ into the system, or they are removed or degraded like $\mathbf{A}$ and $\mathbf{B}$ with rate parameter $d$ in an unspecific process. Real experimental studies will be carried out under conditions were some reactions can be neglected, and mathematical analysis will not consider all possible reactions but use a suitable subset. If, for example, recycling of material is negligible the reactions (3c) and (3d) as well as the molecular species $\mathbf{B}$ can be omitted.

The implementation of mechanism (3) requires some physical setup, for example a flow reactor as shown in figure 3. The supply of material is provided by an influx of stock solution with concentration $[\mathbf{A}] = a_0$ resulting in $k = r\,a_0$, the recycling reaction is neglected ($h = 0$), and the unspecific degradation is replaced by the outflux of the reactor, $d = r$. With these assumptions the kinetic ordinary differential equations are the form
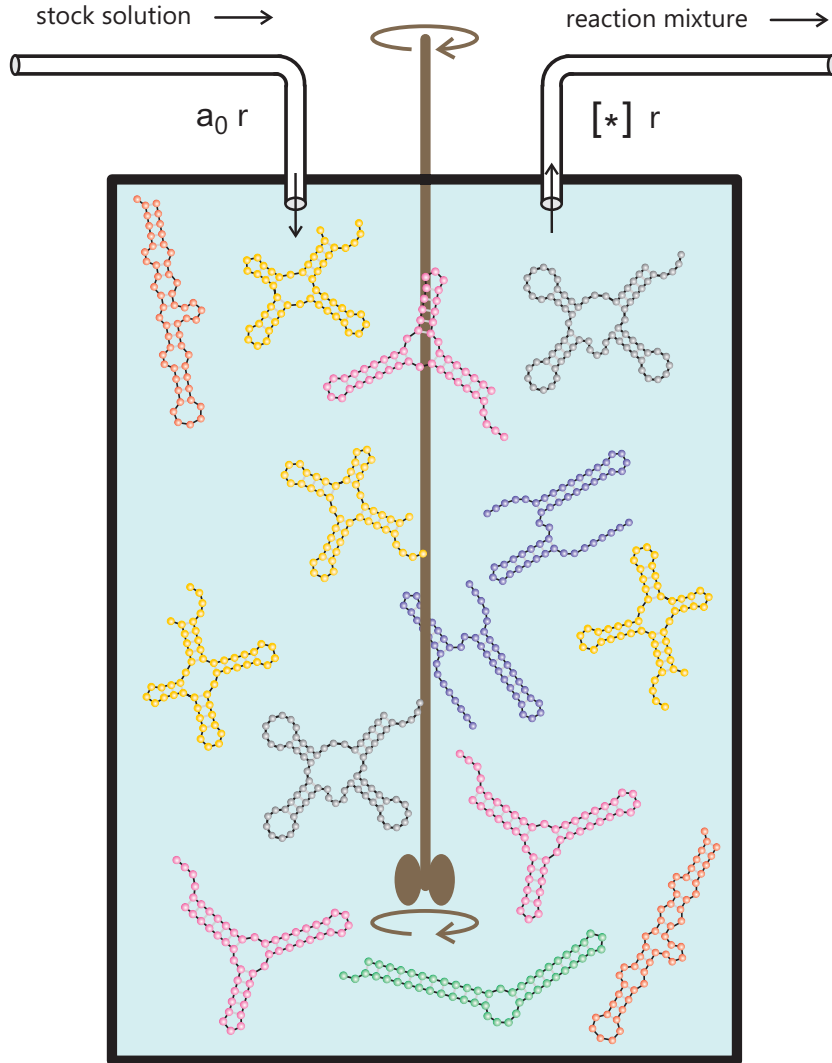
$$\frac{\mathrm{d}a}{\mathrm{d}t} = -a\left(r + \sum_{i=1}^{n} k_i\,c_i\right) + r\,a_0 \;, \tag{4a}$$

$$\frac{\mathrm{d}c_j}{\mathrm{d}t} = \left(k_j\,a - (d_j + r)\right)c_j \;;\; j = 1, \ldots, n \;, \quad \text{and} \tag{4b}$$

$$\frac{\mathrm{d}b}{\mathrm{d}t} = \sum_{i=1}^{n} d_j\,c_j - r\,b \;. \tag{4c}$$

---

[9]Concentrations, $a$, will be indicated by square brackets $[\mathbf{A}]$. For replicators we use particle numbers, $N_j$, or concentrations, $[\mathbf{X}_j] = c_j$, respectively. See also the notations in the appendix.

**Figure 3: The flow reactor for the evolution of RNA molecules.** A stock solution containing all materials for RNA replication ($[\mathbf{A}] = a_0$) including an RNA polymerase flows continuously at a flow rate $r$ into a well stirred tank reactor (CSTR) and an equal volume compensating for the influx and containing a fraction of the reaction mixture ($[\star] = \{a, b, c_i\}$) leaves the reactor (For different experimental setups see, e.g., Watts [43]). The population of RNA molecules $(\mathbf{X}_1, \mathbf{X}_2, \ldots, \mathbf{X}_n)$ is present in numbers $N_1, N_2, \ldots, N_n$ with $N = \sum_{i=1}^{n} N_i$) in the reactor and fluctuates around a mean value, $N \pm \sqrt{N}$. RNA molecules replicate and mutate in the reactor, and the fastest replicators are selected (see [44, pp.21-60]). The RNA flow reactor has been used also as an appropriate model for computer simulations [45–47]. There, other criteria for selection than fast replication can be applied. For example, fitness functions are defined that measure the distance to a predefined target structure and mean fitness increases during the approach towards the target [47].

The flow rate $r$ is the mean reciprocal residence time of a volume element in the reactor, $r = \tau_R^{-1}$. Equation (4) sustains $(n+1)$ stationary states, which are fulfilling the conditions $\dot{a} = 0, \dot{b} = 0, \dot{c}_j = 0$ for $j = 1, 2, \ldots, n$. Every stationarity conditions for one particular class of replicating molecules $\mathbf{X}_j$

$$\bar{c}_j \left( k_j \, \bar{a} - (d_j + r) \right) = 0$$

has two solutions (i) $\bar{c}_j = 0$ and (ii) $\bar{a} = (d_j + r)/k_j$. Since any pair of type (ii) conditions is incompatible,[10] only two types of solutions remain: (i) $\bar{c}_j = 0 \ \forall \ j = 1, 2, \ldots, n$, the *state of extinction*, because no replicating molecule survives and (ii) $n$ states with $\bar{c}_j = (a_0/(d_j + r) - 1/k_j) \, r$ and $\bar{c}_k = 0 \ \forall \ k \neq j$. Steady state analysis through linearization and diagonalization of the Jacobian matrix at the stationary points yields the result that only one of the $n$ states is asymptotically stable, and this is the one referring to species $\mathbf{X}_m$ that is defined by

$$k_m \, a_0 - d_m = \max\{a_j k_j - d_j, \ j = 1, 2, \ldots, n\} \, . \tag{5}$$

Accordingly, species $\mathbf{X}_m$ is selected and we denote this state as *state of selection*. The proof is straightforward and yields simple expressions for the eigenvalues $\lambda_k$ $(k = 0, 1, \ldots, n)$ of the Jacobian matrix when degradation is neglected, $d_j = 0$ $(j = 1, 2, \ldots, n)$. For the state of extinction we find

$$\lambda_0 = -r \quad \text{and} \quad \lambda_j = k_j \, a_0 - r \, . \tag{6}$$

It is asymptotically stable as long as $r > k_m \, a_0$ is fulfilled. If $r > k_m a_0$ then $r > k_j \, a_0 \ \forall \ j \neq m$ is valid by definition because of the selection criterion (5) for $d_j = 0$. For all other $n$ pure states, $\{\bar{c}_j = a_0 - r/k_j, \ \bar{c}_j = 0, \ j \neq i\}$ the eigenvalues of the Jacobian are:

$$\begin{aligned}
\lambda_0 &= -r \, , \\
\lambda_j &= -k_j \, a_0 + r \, , \quad \text{and} \\
\lambda_i &= -\frac{r}{k_j} \, (k_i - k_j) \ \forall \ i \neq j.
\end{aligned} \tag{7}$$

All pure states except the state at which $\mathbf{X}_m$ is selected (state of selection: $c_m = 0, \ c_j = 0, \ j = 1, \ldots, n, j \neq m$) have at least one positive eigenvalue and are unstable. Therefore we have proven the occurrence of selection of the molecular species with the largest value of $k_j$ (or $k_j \, a_0 - d_j$, respectively), because only at $\bar{c}_m \neq 0$ all eigenvalues of the Jacobian matrix are negative.

---

[10]In this section degenerate or neutral cases with $d_i = d_j$ and $k_i = k_j$ are excluded (see also section 8).

It is worth indicating that the dynamical system (4) has a invariant manifold $\mathbb{W}_{n+2}^{(a_0)}$: $\bar{w} = \bar{a} + \bar{b} + \sum_{i=1}^{n} \bar{c}_i = a_0$. From $\dot{w} = \dot{a} + \dot{b} + \sum_{i=1}^{n} \dot{c}_i = (a_0 - w) r$ we find by straightforward integration that the sum of all concentrations, $w(t)$, follows a simple exponential relaxation process towards the asymptotically stable steady state $\bar{w} = a_0$:

$$w(t) = a_0 - \big(a_0 - w(0)\big) \exp(-r\,t)\,,$$

with the flow rate $r$ being the relaxation constant.

## 4   The selection equation

As illustrated and analyzed in section 3 the basis of Darwinian selection is reproduction, which may be reduced to an overall autocatalytic reaction step, $\mathbf{A} + \mathbf{X}_j \to 2\mathbf{X}_j$. The system is simplified further by assuming that the material consumed in the reproduction process, $\mathbf{A}$, is present in excess: $[\mathbf{A}] = a_0$, and we indicate the buffering of the concentration by means of parentheses, $(\mathbf{A}) + \mathbf{X}_j \to 2\mathbf{X}_j$. The concentration of $\mathbf{A}$ is constant and can be absorbed in the rate constant: $f_j = k_j\,[\mathbf{A}] = k_j\,a_0$. In addition we neglect the degradation terms by putting $d_j = 0 \ \forall \ j$. In terms of chemical reaction kinetics selection based on pure reproduction is described by the dynamical system

$$\frac{\mathrm{d}c_j}{\mathrm{d}t} = f_j\,c_j - \frac{c_i}{c_0} \sum_{i=1}^{n} f_i\,c_i = c_i \left( f_i - \frac{\sum_{i=1}^{n} c_i}{c_0}\,\phi(t) \right) ; j = 1, 2, \ldots, n\,,$$

$$\phi(t) = \frac{1}{\sum_{i=1}^{n} c_i} \sum_{i=1}^{n} f_i\,c_i = \overline{f}\,.$$

$(8)$

The variables $c_j(t)$ are the concentrations of the genotypes $\mathbf{X}_j$ as before, $c_0$ is the total concentration of the invariant manifold $\mathbb{S}_n^{(c_0)}$: $\bar{c} = \sum_{j=1}^{n} c_j(t) = c_0$. The quantities $f_j$ are reproduction rate parameters corresponding to overall replication rate constants in molecular systems or, in general, the fitness values of the genotypes. A global flux $\phi(t)$ has been introduced to compensate for the net growth of the system. In the particular case here it is identical to the mean fitness of the population.

Transformation to relative concentrations, $x_j = c_j/c$, $\sum_{i=1}^{n} x_i(t) = 1$ and

$c_0 = 1$ simplifies further analysis:[11]

$$\frac{\mathrm{d}x_j}{\mathrm{d}t} = f_j\, x_j - x_j \sum_{i=1}^{n} f_i\, x_i = x_j(f_j - \phi) \quad \text{with}$$

$$\phi = \sum_{i=1}^{n} f_i\, x_i = \overline{f} \quad \text{and} \quad i = 1, 2, \ldots, n \ . \tag{9}$$

Because of this conservation relation, $\sum_{i=1}^{n} x_i(t) = 1$, only $n - 1$ variables $x_j$ are independent. In the space of $n$ Cartesian variables, $\mathbb{R}^n$, the $x$-variables represent a projection of the positive orthant onto the unit simplex

$$\mathbb{S}_n^{(1)} = \left\{ x_i \geq 0\, \forall\, i = 1, 2, \ldots, n \,\wedge\, \sum_{i=1}^{n} x_i = 1 \right\} \ . \tag{10}$$

The simplex $\mathbb{S}_n^{(1)}$ is an invariant manifold of the differential equation (9). This means that every solution curve $\mathbf{x}(t) = \big(x_1(t), x_2(t), \ldots, x_n(t)\big)$ that starts in one point of the simplex will stay on the simplex forever. Moreover, the boundary of the simplex consists of invariant subsimplices [48] and no trajectory can cross it neither from inside to outside – contradicting non-negativeness of particle numbers or concentrations – nor from outside to inside.

In order to analyze the stability of $\mathbb{S}_n^{(1)}$ we relax the conservation relation $\sum_{i=1}^{n} x_i(t) = c(t) \neq 1$ and assume that only the conditions

$$\{ f_j > 0 \,\wedge\, 0 \leq x_j(0) < \infty \}\, \forall\, i = 1, 2, \ldots, n \ ,$$

are fulfilled. According to this assumption all replication rate parameters are strictly positive – a condition that will be replaced by the weaker condition $f_i \geq 0\, \forall\, i \neq k \wedge f_k > 0$ below – and the concentration variables are non-negative quantities. Asymptotic stability of the unit simplex requires that all solution curves converge to the simplex from every initial condition, $\lim_{t \to \infty} \big( \sum_{i=1}^{n} x_i(t) \big) = 1$.

This conjecture is readily proved: From $\sum_{i=1}^{n} x_i(t) = c(t)$ follows

$$\frac{\mathrm{d}c}{\mathrm{d}t} = c\,(1 - c)\,\phi(t) \quad \text{with} \quad \phi(t) > 0 \ . \tag{11}$$

For $\mathrm{d}c/\mathrm{d}t = 0$ we find the two stationary states: a saddle point at $\bar{c} = 0$ and an asymptotically stable state at $\bar{c} = 1$. There are several possibilities to verify its asymptotic stability, we choose to solve the differential equation and find:

$$c(t) = 1 - \big(1 - c(0)\big)\, \exp\left( -\int_0^t \phi(\tau)d\tau \right) \ .$$

---

[11]Care is needed for the application of relative concentrations, because the total concentration $c(t)$ might vanish and then relative concentrations become spurious quantities.

Starting with any initial value $c(0)$ the population approaches the unit simplex. When it starts on $\mathbb{S}_n$ it stays there and returns to it also in presence of fluctuations.[12] Therefore, we can restrict population dynamics to the simplex without loosing generality and characterize the state of a population at time $t$ by the vector $\mathbf{x}(t)$ which fulfils the $\mathbb{L}^{(1)}$ norm $\sum_{i=1}^n x_i(t) = 1$.

The necessary and sufficient condition for the stability of the simplex, $\phi(t) > 0$, enables us to relax the condition for the rate parameters $f_i$. In order to have a positive flux it is sufficient that one rate parameter is strictly positive provided the corresponding variable is non-zero:

$$\phi(t) > 0 \implies \exists k \in \{1, 2, \ldots, n\} \text{ such that } f_k > 0 \wedge x_k > 0 .$$

For the variable $x_k$ it is sufficient that $x_k(0) > 0$ holds because $x_k(t) \geq x_k(0)$ when all other products $f_j x_j$ were zero at $t = 0$. This relaxed condition for the flux is important for the handling of lethal mutants with $f_j = 0$.

The time dependence of the mean fitness or flux $\phi$ is given by

$$\begin{aligned}
\frac{\mathrm{d}\phi}{\mathrm{dt}} &= \sum_{i=1}^n f_i \frac{\mathrm{d}x_i}{\mathrm{dt}} = \sum_{i=1}^n f_i \left( f_i x_i - x_i \sum_{j=1}^n f_j x_j \right) = \\
&= \sum_{i=1}^n f_i^2 x_i - \sum_{i=1}^n f_i x_i \sum_{j=1}^n f_j x_j = \\
&= \overline{f^2} - \left(\overline{f}\right)^2 = \mathrm{var}\{f\} \geq 0 .
\end{aligned} \tag{12}$$

Since a variance is always nonnegative, equation (12) implies that $\phi(t)$ is a non-decreasing function of time. The value $\mathrm{var}\{f\} = 0$ refers to a homogeneous population of the fittest variant, then $\phi(t)$ cannot increase any further and hence, it has been optimized during selection.

It is also possible to derive analytical solutions for equation (9) by a transform called integrating factors ([49], p.322ff.):
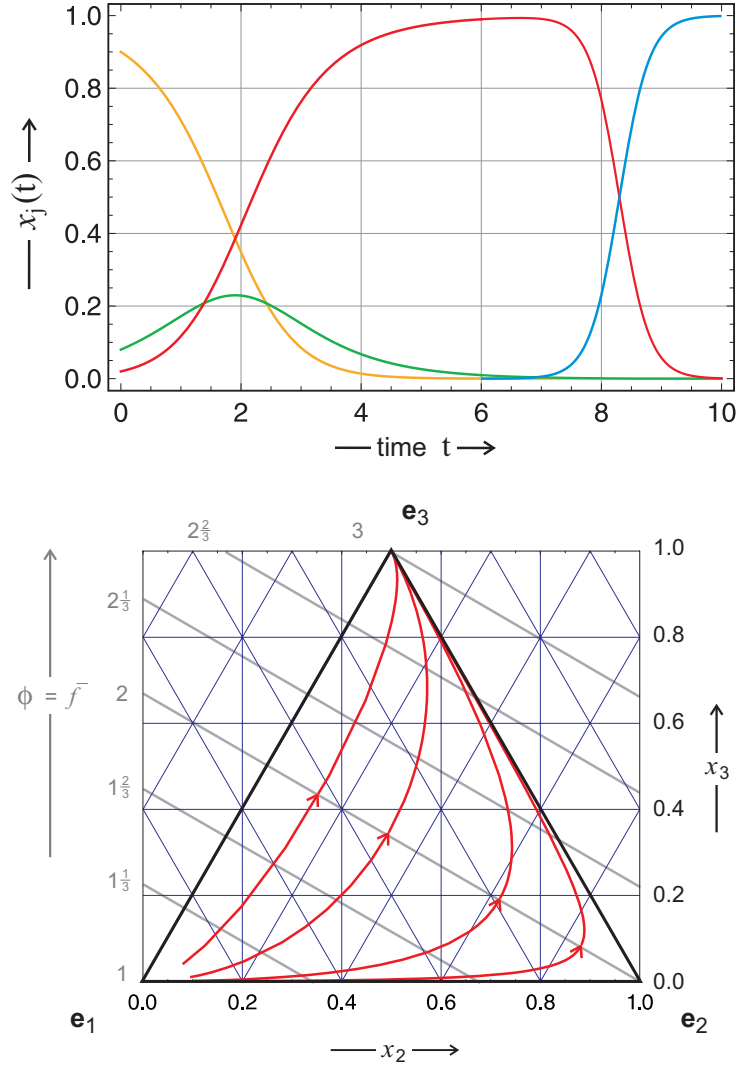
$$z_j(t) = x_j(t) \exp\left(\int_0^t \phi(\tau)d\tau\right) . \tag{13}$$

Insertion into (9) yields

$$\frac{\mathrm{d}z_j}{\mathrm{dt}} = f_j z_j \text{ and } z_j(t) = z_j(0) \exp(f_j t) ,$$

$$x_j(t) = x_j(0) \exp(f_j t) \exp\left(-\int_0^t \phi(\tau)d\tau\right) \text{ with}$$

$$\exp\left(\int_0^t \phi(\tau)d\tau\right) = \sum_{i=1}^n x_i(0) \exp(f_i t) ,$$

---

[12]Generalization to arbitrary but finite population sizes $c \neq 1$ is straightforward: For $\sum_{i=1}^n x_i(0) = c_0$ the equation $\mathrm{d}x_j / \mathrm{dt} = f_j x_j - (x_j/c_0) \sum_{i=1}^n f_i x_i$, $j = 1, 2, \ldots, n$ plays the same role as equation (9) did for $\sum_{i=1}^n x_i(0) = 1$.
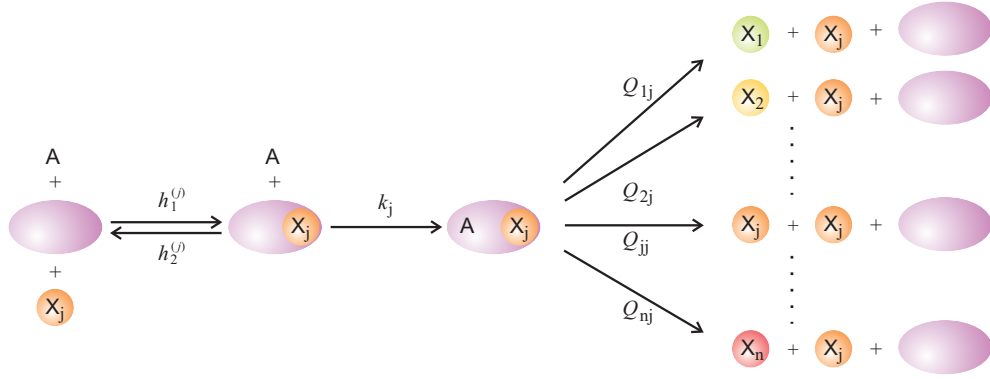
**Figure 4: Selection on the unit simplex $\mathbb{S}_3^{(1)}$.** In the upper part of the figure we show solution curves $\mathbf{x}(t)$ of equation (9) with $n = 3(4)$. The parameter values are: $f_1 = 1\,[\mathrm{t}^{-1}]$, $f_2 = 2\,[\mathrm{t}^{-1}]$, $f_3 = 3\,[\mathrm{t}^{-1}]$, and $f_7 = 7\,[\mathrm{t}^{-1}]$ where $[\mathrm{t}^{-1}]$ is an arbitrary reciprocal time unit. Initial conditions: $\mathbf{x}(0) = (0.90, 0.08, 0.02, 0)$ and $x_4(6) = 0.0001$. Color code: $x_1(t)$ yellow, $x_2(t)$ green, $x_3(t)$ red, and $x_4(t)$ blue. $\mathbf{X}_4$ is injected into the previously equilibrated system at time $t = 6$ as a fitter variant, which takes over in rather short time.

The lower part of the figure shows parametric plots $\mathbf{x}(t)$ on the simplex $\mathbb{S}_3^{(1)}$. Constant level sets of $\phi$ are straight lines (grey). Choice of parameters: $f_1 = 1\,[\mathrm{t}^{-1}]$, $f_2 = 2\,[\mathrm{t}^{-1}]$, and $f_3 = 3\,[\mathrm{t}^{-1}]$.

where we have used $z_j(0) = x_j(0)$ and the condition $\sum_{i=1}^{n} x_i = 1$. The solution finally is of the form

$$x_j(t) = \frac{x_j(0)\,\exp(f_j t)}{\sum_{i=1}^{n} x_i(0)\,\exp(f_i t)}\,;\ j = 1, 2, \ldots, n\,. \tag{14}$$

Under the assumption that the largest fitness parameter is non-degenerate, $\max\{f_j;\ j = 1, 2, \ldots, n\} = f_m > f_j\,\forall\,j \neq m$, every solution curve fulfill-

**Figure 5: A molecular view of replication and mutation**. The replication device **E**, commonly a replicase molecule or a multi-enzyme complex (violet) binds the template DNA or RNA molecule ($\mathbf{X}_j$, orange) in order to form a replication complex $\mathbf{E}\cdot\mathbf{X}_j$ with a binding constant $H_j = h_1^{(j)}[\mathbf{E}][\mathbf{X}_j] \,/\, h_2^{(j)}[\mathbf{E}\cdot\mathbf{X}_j]$ and replicates with a rate parameter $f_j$. During the template copying process reaction channels leading to mutations are opened through replication errors. The reaction leads to a correct copy with frequency $Q_{jj}$ and to a mutant $\mathbf{X}_k$ with frequency $Q_{kj}$ commonly with $Q_{jj} \gg Q_{kj} \,\forall\, k \neq j$. Stoichiometry of replication requires $\sum_{i=1}^{n} Q_{ij} = 1$, since the product has to be either correct or incorrect. The reaction is terminated by full dissociation of the replication complex. The sum of all activated monomers is denoted by **A**.

ing the initial condition $x_i(0) > 0$ approaches a homogenous population: $\lim_{t\to\infty} x_m(t) = \bar{x}_m = 1$ and $\lim_{t\to\infty} x_j(t) = \bar{x}_j = 0 \,\forall\, j \neq m$, and the flux approaches the largest fitness parameter monotonously, $\phi(t) \to f_m$ (Examples are shown in figure 4).

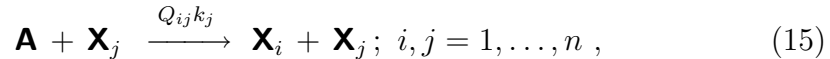Qualitative analysis of stationary points and their stability yields the following results:

(i) The only stationary points of equation (9) are the corners of the simplex, represented by the unit vectors $\mathbf{e}_k = \{x_k = 1, x_i = 0 \,\forall\, i \neq k\}$,

(ii) only one of these stationary points is asymptotically stable, the corner where the mean fitness $\phi$ adopts its maximal value on the simplex ($\mathbf{e}_m$: $\bar{x}_m = 1$ defined by $\max\{f_i;\, i = 1, 2, \ldots, n\} = f_m > f_i \,\forall\, i \neq m$), one corner is unstable in all directions, a source where the value of $\phi$ is minimal ($\mathbf{e}_s$: $\bar{x}_s = 1$ defined by $\min\{f_i;\, i = 1, 2, \ldots, n\} = f_s < f_i \,\forall\, i \neq s$), and all other $n - 2$ equilibria are saddle points, and

(iii) since $x_i(0) = 0$ implies $x_i(t) = 0 \,\forall\, t > 0$, every subsimplex of $\mathbb{S}_n^{(1)}$ is an invariant set, and thus the whole boundary of the simplex consists of invariant sets and subsets down the corners [48] (which represent members of class $\mathbb{S}_1^{(1)}$).

## 5  The mutation-selection equation

Correct replication and mutation are visualized as parallel reaction channels. An example of a simple molecular mechanism meeting this concept is shown in figure 5: The template $\mathbf{X}_j$ binds to the replicase $\mathbf{E}$ forming a complex $\mathbf{E} \cdot \mathbf{X}_j$ with a binding constant $H_j = h_1^{(j)}[\mathbf{E}][\mathbf{X}_j] \,/\, h_2^{(j)}[\mathbf{E} \cdot \mathbf{X}_j]$, replication is initiated with a rate parameter $k_j$, and then during the replication process the various reaction channels are opened through replication errors. At the end of the replication process template and correct copy or mutant – formed with the probabilities $Q_{jj}$ or $Q_{ij}$, respectively – dissociate from the enzyme. This mechanism is a rough description of what happens in viral RNA replication by virus specific replicases.

In formal terms mutation is readily introduced into the mechanism (3) by replacing equation (3b) by

$$\mathbf{A} \,+\, \mathbf{X}_j \xrightarrow{Q_{ij}k_j} \mathbf{X}_i \,+\, \mathbf{X}_j \,;\; i,j = 1, \ldots, n \,, \tag{15}$$

the rate parameter $k_j$ determines the rate at which some new molecule $\mathbf{X}_{(i)}$ is synthesized on the template $\mathbf{X}_j$, which is $k_j\,[\mathbf{X}_j]\,[\mathbf{A}]$, and $Q_{ij}$ describes the fraction of replication events that lead to the synthesis of precisely $\mathbf{X}_i$. Hence $Q_{jj}$ is the fraction of correctly copied molecules and $Q_{ij}$ is the fraction of mutations $\mathbf{X}_j \to \mathbf{X}_i$. Since we assume here that the enumeration $i = 1, \ldots, n$ is exhaustive – all possible molecules are *preexisting* even if most of them are not (yet) in the system $x_{(i)}(t) = 0$ at time $t$ – we have the conservation rule $\sum_{i=1}^{n} Q_{ij} = 1$ or, in other words, every copy is either correct or incorrect, and the mutation matrix

$$Q \;=\; \begin{pmatrix} Q_{11} & Q_{12} & \cdots & Q_{1n} \\ Q_{21} & Q_{22} & \cdots & Q_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ Q_{n1} & Q_{n2} & \cdots & Q_{nn} \end{pmatrix} ,$$

is a stochastic matrix with the elements of a given column summing up to one. In case one makes the assumption of equal probabilities for $\mathbf{X}_j \to \mathbf{X}_i$ and $\mathbf{X}_i \to \mathbf{X}_j$, as it is made for example in the uniform error rate model (see below and [2,6]), Q is symmetric and hence a bistochastic matrix where summation over all elements in the rows yields one as well.

### 5.1  Mutation

The uniform error rate model assumes that the mutation rate parameter per site and replication event, $p$, is independent of the position along the

polynucleotide chain.[13] Under this assumption all elements of the mutation matrix Q can be expressed in terms of three parameters only no matter how long the genomes are:

$$Q_{ji} \;=\; Q_{ij} \;=\; (1-p)^{\nu-d_{ij}}\,p^{d_{ij}} \;=\; (1-p)^{\nu}\,\varepsilon^{d_{ij}} \;\; \text{with} \;\; \varepsilon = \frac{p}{1-p}\,. \quad (16)$$

Apart from chain length $\nu$ and mutation rate parameter $p$ equation (16) contains the Hamming distance between two sequences $\mathbf{X}_j$ and $\mathbf{X}_i$, which is the number of positions in which the two aligned sequences differ [50]: $d_{ij} = d_{\mathrm{H}}(\mathbf{X}_i, \mathbf{X}_j)$.

Individual replicators $\mathbf{X}_{(j)}$ – molecules, viruses bacteria, or higher organisms – have genomes, which are polynucleotides, RNA or DNA, and every reproduction events is necessarily accompanied by genome replication. In a two letter alphabet we represent the genomes as binary strings of chain length $\nu$, which are individual points in sequence space (see figure 6),

$$\mathbf{X}_i = (\mathbf{0010010110011}\cdots\mathbf{10})\,,$$
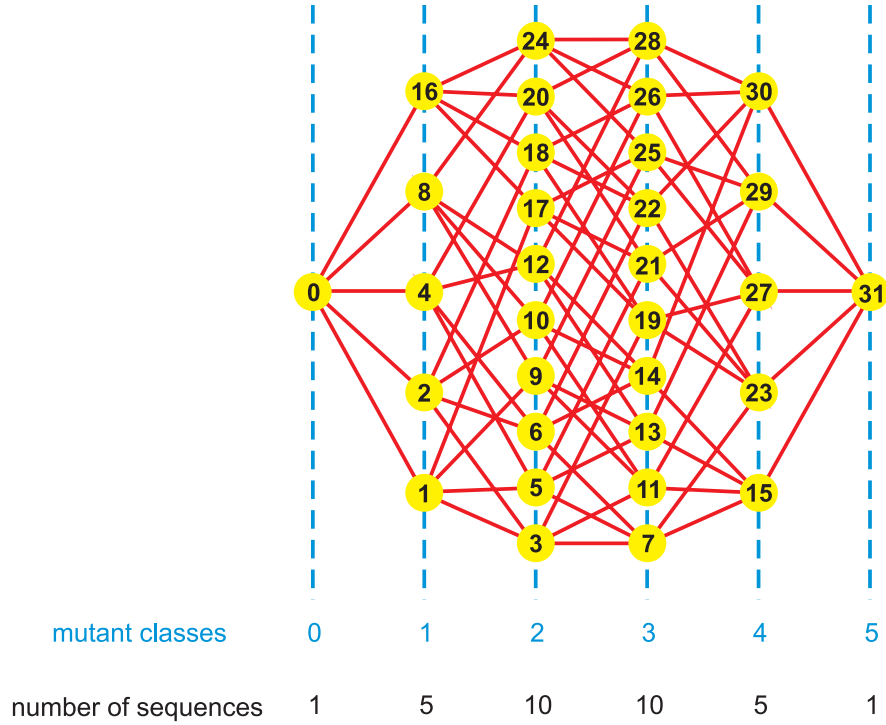$$\mathbf{X}_j = (\mathbf{0010110110001}\cdots\mathbf{10})\,,$$

and the distance between them is expressed as Hamming distance – $d_{\mathrm{H}}(\mathbf{X}_i, \mathbf{X}_j) = d_{ij} = 2$ in the example shown above. Natural polynucleotides of chain length $\nu$ can be represented by binary sequences of chain length $2\nu$ where the individual nucleobases of the four letter alphabet are encoded by two digits each, for example $\mathbf{C} \equiv \mathbf{00}$, $\mathbf{U} \equiv \mathbf{01}$, $\mathbf{A} \equiv \mathbf{10}$, and $\mathbf{G} \equiv \mathbf{11}$. A particular sequence and its binary encoding look, for example

$$\mathbf{X}_k = (\,\mathbf{C}\,|\,\mathbf{A}\,|\,\mathbf{U}\,|\,\mathbf{U}\,|\,\mathbf{A}\,|\,\mathbf{G}\,|\,\mathbf{A}\,|\cdots|\,\mathbf{A}\,)\,,$$
$$\mathbf{X}_k = (\mathbf{00}|\mathbf{10}|\mathbf{01}|\mathbf{01}|\mathbf{10}|\mathbf{11}|\mathbf{10}|\cdots|\mathbf{10})\,.$$

Complementarity in base pairs, $\mathbf{C} \equiv \mathbf{G}$ and $\mathbf{U} = \mathbf{A}$ is thereby retained provided the grouping in doublets of binary symbols is respected. Computation of Hamming distances between four letter sequences in binary encoding required some care: $\mathbf{00} \rightarrow \mathbf{11}$ corresponding to $\mathbf{C} \rightarrow \mathbf{G}$ is a single point mutation of Hamming distance $d_{\mathrm{H}} = 1$ whereas two digits are changed in the binary encoded sequences, $\mathbf{10} \rightarrow \mathbf{11}$ is only a change in a single digit, it corresponds to $\mathbf{A} \rightarrow \mathbf{G}$ and the Hamming distance is $d_{\mathrm{H}} = 1$ for binary and four letter sequences in this case.

---

[13]Uniform error rates or symmetry in the direction of mutations is commonly not fulfilled in nature. It is here introduced as a simplification, which facilitates the derivation of exact solutions of the differential equation (18). Moreover, neither the assumption of uniform error rates nor the condition of a symmetric mutation matrix Q are essential for the forthcoming analysis.

**Figure 6: Mutant classes in sequence space.** The sketch shows the sequence space for binary sequences of chain length $\nu = 5$, which are given in terms of their decadic encodings: "**0**" $\equiv$ **00000**, "**1**" $\equiv$ **00001**, . . . , "**31**" $\equiv$ **11111**. All pairs of sequences with Hamming distance $d_\mathrm{H} = 1$ are connected by red lines. The number of sequences in mutant class $k$ is $\binom{\nu}{k}$.

### 5.2  Replication-mutation dynamics

The implementation of mechanism (3) with reaction (15) instead of (3b) in the flow reactor (figure 3), whereby degradation is replaced by the outflux, $d = d_{(j)} = r$, leads to the kinetic differential equations

$$\frac{\mathrm{d}a}{\mathrm{d}t} = -a \left( r + \sum_{i=1}^{n} k_i\, c_i \right) + r\, a_0\, , \tag{17a}$$

$$\frac{\mathrm{d}c_j}{\mathrm{d}t} = a \left( \sum_{i=1}^{n} Q_{ji}\, k_i\, c_i \right) - r\, c_j\, ;\; j = 1, \ldots, n\, . \tag{17b}$$

Equation (17) can be analyzed straightforwardly [44, pp.21-60] but the results are more difficult to interpret than those derived from an idealized model, which is of almost general validity.

The simplifications were introduced by Eigen [1], they are essentially the same as those used in the previous section 4 or in population genetics [51, 52], and they concern two assumptions: (i) the material required to synthesize replicators is assumed to be present in excess, $a(t) = a_0$, and

can be absorbed in the rate parameter, $f_j = k_j a_0$, and (ii) the increase in replicator concentrations is compensated by means of an unspecific (dilution) flux $\phi(t)$. Condition (ii) leads to a constant total concentration of replicators, $c_0 = c(t) = \sum_{i=1}^{n} c_i(t)$, which allows for normalization of variables, $x_j(t) = c_j(t)/c_0$ with $\sum_{i=1}^{n} x_i(t) = 0$, and eventually the dynamics in populations of replicating and mutation molecules is cast into the differential equation

$$\frac{\mathrm{d}x_j}{\mathrm{dt}} = \sum_{i=1}^{n} Q_{ij} f_j x_i - x_j \cdot \phi\,;\; j = 1, 2, \ldots, n \;\; \text{with} \;\; \phi = \sum_{i=1}^{n} f_i x_i = \overline{f}\,. \quad (18)$$

It has been proven that the solution curves $x_j(t)$ are – up to a transformation of the time axis – independent of the total concentration as long as $c(t)$ stays finite and does not vanish, $c(t) \neq \{0, \infty\}$ [4].[14]

The mean fitness or flux $\phi(t)$ does not contain mutation terms because of the conservation relation $\sum_{i=1}^{n} Q_{ji} = 1$. A straightforward summation yields an equation for the time dependence of the total concentration, which is identical to equation (11)

$$\frac{\mathrm{d}c}{\mathrm{dt}} = c \left( 1 - \frac{c}{c_0} \right) \phi(t)\,, \quad (19)$$

which sustains three stationary solutions, (i) $\mathbf{P}_1$ : $\bar{c} = c_0$ and (ii) $\mathbf{P}_2$ : $\bar{c} = 0$. The third stationary state $\mathbf{P}_3$ is defined by $\phi(t) = 0$ and it is identical to $\mathbf{P}_2$.[15] The stability of the steady states can be analyzed by differentiation with respect to the total concentration $c$:

$$\frac{\partial}{\partial c} \left( \frac{\mathrm{d}c}{\mathrm{dt}} \right) = \phi(t) - \frac{c}{c_0} \left( 2\phi(t) - c_0 \frac{\partial \phi}{\partial c} \right) - \frac{c^2}{c_0} \frac{\partial \phi}{\partial c} = -\lambda\,.$$

Insertion of the stationary solutions yields

$$\lambda^{(1)} = -\phi|_{c=c_0} < 0 \quad \text{and} \quad \lambda^{(2)} = \phi|_{c=0} = 0$$

The first steady state $\mathbf{P}_1$ is asymptotically stable: $c(t)$ decreases for $c > c_0$ and increases for $0 < c < c_0$. The second state, $\mathbf{P}_2$ is marginally stable, any fluctuation $\delta c > 0$ leads to $\lambda^{(2)} > 0$, and implies progression towards increasing values of $c$ until $\mathbf{P}_1$ has been reached. The concentrations at state $\mathbf{P}_2$ are confined to the simplex $\mathbb{S}_n^{(c_0)}$ and after normalization the accessible

---

[14]Equation 18, accordingly, does not apply to situations of vanishing populations like lethal mutagenesis in virology where is has wrongly be used for comparison [53, 54].

[15]For strictly positive fitness values, $f_i > 0 \,\forall\, i = 1, 2, \ldots, n$, the condition $\phi = 0$ can only be fulfilled by $x_i = c_i = 0 \,\forall\, i = 1, 2, \ldots, n$, which is identical to state $\mathbf{P}_2$. If some $f_i$ values are zero – corresponding to lethal variants – the respective variables vanish in the infinite time limit because of $dc_i/dt = -c_i\,\phi(t)$ with $\phi(t) > 0$.

space for the variables $x_{(j)}$ is the unit simplex $\mathbb{S}_n^{(1)}$ defined in equation (10). Asymptotic stability of $\mathbf{P}_2$ implies that the system converges to the unit simplex, as it did without mutations. For initial values of the variables chosen on the simplex, $\sum_{i=1}^{n} x_i(0) = 1$, it remains there.

There is one important difference between the replication and the replication-mutation system: In the latter the boundary of the unit simplex, $\mathbb{S}_n^{(1)}$, is not invariant. Although no orbit starting on the simplex will leave it, which is a *conditio sine qua non* for chemical reactions requiring non-negative concentrations, trajectories flow from outside into $\mathbb{S}_n^{(1)}$. In other words, the condition $x_j(0) = 0$ does not lead to $x_j(t) = 0 \, \forall \, t > 0$ (figure 7). The chemical interpretation is straightforward: If a variant $\mathbf{X}_j$ is not present initially, it can be formed through a mutation event.

## 5.3   Numerical solution

Before discussing the role of the flux $\phi$ in the selection-mutation system with respect to optimization, we shall derive exact solutions of equation (18) by means of the integrating factor transformation as in the mutation-free case (see [49, p.322ff.] and [55, 56]). At first the variables $x_j(t)$ are transformed:

$$z_j(t) \;=\; x_j(t) \cdot \exp\left( \int_0^t \phi(\tau)d\tau \right) \; .$$

From $\sum_{i=1}^{n} x_i(t) = 1$ follows straightforwardly, again as in the selection-only case,

$$\exp\left( \int_0^t \phi(\tau)d\tau \right) \;=\; \sum_{i=1}^{n} z_i(t) \; .$$

What remains to be solved is a linear first order differential equation

$$\frac{dz_j}{dt} \;=\; \sum_{i=1}^{n} Q_{ji} \, f_i \, z_i \, ; \; j = 1, 2, \ldots, n \quad \text{or} \quad \frac{d\mathbf{z}}{dt} \;=\; Q \cdot F \, \mathbf{z} \; . \qquad (20)$$

which is readily achieved by means of standard linear algebra. We define a matrix $W = \{W_{ji} = Q_{ji} \, f_i\} = Q \cdot F$ where $F = \{F_{ii} = f_i \delta_{ij}\}$ is a diagonal matrix, and obtain $d\mathbf{z}/dt = W \cdot \mathbf{z}$. Provided matrix $W$ is diagonalizable, which will always be the case when the mutation matrix $Q$ is based on real chemical reaction mechanisms, we can transform variables by means of the two $n \times n$ matrices $B = \{b_{ij}\}$ and $B^{-1} = H = \{h_{ij}\}$ $(i, j = 1, \ldots, n)$,

$$\mathbf{z}(t) \;=\; B \cdot \boldsymbol{\zeta}(t) \quad \text{and} \quad \boldsymbol{\zeta}(t) \;=\; B^{-1} \cdot \mathbf{z}(t) \; ,$$

such that $B^{-1} \cdot W \cdot B = \Lambda$ is diagonal and its elements, $\lambda_k$, are the eigenvalues of the matrix $W$. The right-hand eigenvectors of $W$ are given by the columns

of B, $\mathbf{b}_j = (b_{i,j}; i = 1, \ldots, n)$, and the left-hand eigenvectors by the rows of $B^{-1} = H$, $\mathbf{h}_k = (h_{k,i}; i = 1, \ldots, n)$, respectively. These eigenvectors are the *normal modes* of selection-mutation kinetics. For strictly positive off-diagonal elements of W, implying the same for Q which says nothing more than every mutation $\mathbf{X}_i \rightarrow \mathbf{X}_j$ is possible although the probability might be extremely small, Perron-Frobenius theorem holds (see, for example, [57] and next paragraph) and we are dealing with a non-degenerate largest eigenvalue $\lambda_0$,

$$\lambda_0 > |\lambda_1| \geq |\lambda_2| \geq |\lambda_3| \geq \ldots \geq |\lambda_n| , \tag{21}$$

and a corresponding dominant eigenvector $\mathbf{b}_0$ with strictly positive components, $b_{i0} > 0 \, \forall \, i = 1, \ldots, n$.[16] In terms of components the differential equation in $\boldsymbol{\zeta}$ has the solutions

$$\zeta_k(t) = \zeta_k(0) \exp(\lambda_k t) . \tag{22}$$

Transformation back into the variables $\mathbf{z}$ yields

$$z_j(t) = \sum_{k=0}^{n-1} b_{jk} \beta_k(0) \exp(\lambda_k t) , \tag{23}$$

with the initial conditions encapsulated in the equation

$$\beta_k(0) = \sum_{i=1}^{n} h_{ki} z_i(0) = \sum_{i=1}^{n} h_{ki} x_i(0) . \tag{24}$$

From here we obtain eventually the solutions in the original variables $x_j$ through normalization

$$x_j(t) = \frac{\sum_{k=0}^{n-1} b_{jk} \beta_k(0) \exp(\lambda_k t)}{\sum_{i=1}^{n} \sum_{k=0}^{n-1} b_{ik} \beta_k(0) \exp(\lambda_k t)} . \tag{25}$$

For sufficiently long times the contribution of the largest eigenvalue dominates the summations and we find for the stationary solutions

$$\bar{x}_j(t) = \frac{b_{j0} \beta_0(0) \exp(\lambda_0 t)}{\sum_{i=1}^{n} b_{i0} \beta_0(0) \exp(\lambda_0 t)} , \tag{26}$$

which represent the components of the quasispecies.

Perron-Frobenius theorem comes in two versions [57] which we shall now consider and apply to the selection-mutation problem. The stronger version provides a proof for six properties of the largest eigenvector of non-negative primitive matrices[17] T:

---

[16]We introduce here an asymmetry in numbering rows and columns in order to point at the special properties of the largest eigenvalue $\lambda_0$ and the dominant eigenvector $\mathbf{b}_0$.

[17]A square non-negative matrix $T = \{t_{ij}; i, j = 1, \ldots, n; t_{ij} \geq 0\}$ is called *primitive* if there exists a positive integer $m$ such that $T^m$ is strictly positive: $T^m > 0$ which implies $T^m = \{t_{ij}^{(m)}; i, j = 1, \ldots, n; t_{ij}^{(m)} > 0\}$.

(i) The largest eigenvalue is real and positive, $\lambda_0 > 0$,

(ii) a strictly positive right eigenvector $\mathbf{b}_0$ and a strictly positive left eigenvector $\mathbf{h}_0$ are associated with $\lambda_0$,

(iii) $\lambda_0 > |\lambda_k|$ holds for all eigenvalues $\lambda_k \neq \lambda_0$,

(iv) the eigenvectors associated with $\lambda_0$ are unique up to constant factors,

(v) if $0 \leq B \leq T$ is fulfilled and $\beta$ is an eigenvalue of B, then $|\beta| \leq \lambda_0$, and, moreover, $|\beta| = \lambda_0$ implies $B = T$,

(vi) $\lambda_0$ is a simple root of the characteristic equation of T.

The weaker version of the theorem holds for irreducible matrices[18] T. All the above given assertions hold except (iii) has to be replaced by the weaker statement

(iii) $\lambda_0 \geq |\lambda_k|$ holds for all eigenvalues $\lambda_k \neq \lambda_0$.

Irreducible cyclic matrices can be used straightforwardly as examples in order to demonstrate the existence of conjugate complex eigenvalues (An example is discussed below). Perron-Frobenius theorem, in its strict or weaker form, holds not only for strictly positive matrices $T > 0$ but also for large classes of mutation or value matrices ($W \equiv T$ being a primitive or an irreducible non-negative matrix) with off-diagonal zero entries corresponding to zero mutation rates. The occurrence of a non-zero element $t_{ij}^{(m)}$ in $T^m$ implies the existence of a mutation path $\mathbf{X}_j \to \mathbf{X}_k \to \ldots \to \mathbf{X}_l \to \mathbf{X}_i$ with non-zero mutation frequencies for every individual step. This condition is almost always fulfilled in real systems (for exceptions see section 9).
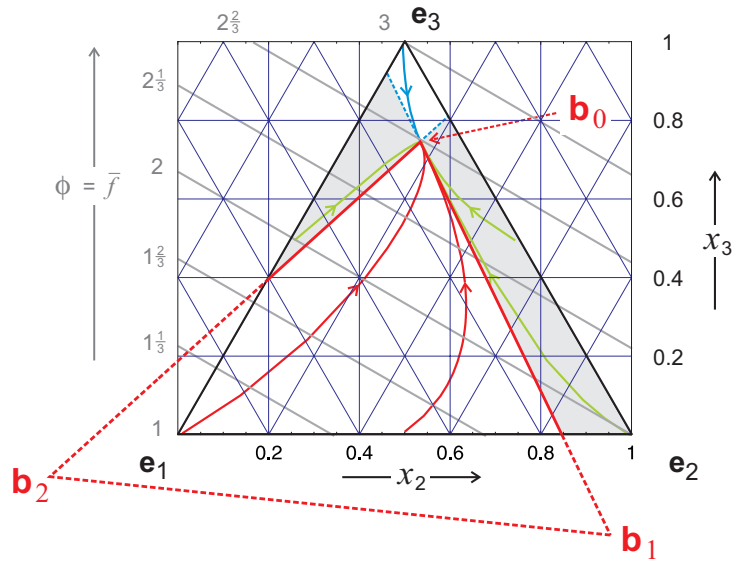
### 5.4   Complex eigenvalues

In order to address the existence of complex eigenvalues of the value matrix W we start by considering the straightforward case of a symmetric mutation matrix Q. Replication rate parameters, $f_i$ are subsumed in a diagonal matrix: $F = \{f_i \cdot \delta_{i,j}; i, j = 1, \ldots, n\}$, the value matrix is the product $W = Q \cdot F$, and, in general, W is not symmetric. A similarity transformation,

$$F^{\frac{1}{2}} \cdot W \cdot F^{-\frac{1}{2}} \;=\; F^{\frac{1}{2}} \cdot Q \cdot F \cdot F^{-\frac{1}{2}} \;=\; F^{\frac{1}{2}} \cdot Q \cdot F^{\frac{1}{2}} \;=\; W' \,.$$

---

[18] A square non-negative matrix $T = \{t_{ij}; i, j = 1, \ldots, n; t_{ij} \geq 0\}$ is called *irreducible* if for every pair $(i, j)$ of its index set there exists a positive integer $m_{ij} \equiv m(i, j)$ such that $t_{ij}^{m_{ij}} > 0$. An irreducible matrix is called *cyclic* with period $d$, if the period of (all) its indices satisfies $d > 1$, and it is said to be *acyclic* if $d = 1$.

**Figure 7: The quasispecies on the unit simplex.** Shown is the case of three variables $(x_1, x_2, x_3)$ on $\mathbb{S}_3^{(1)}$. The dominant eigenvector, the quasispecies denoted by $\mathbf{b}_0$, is shown together with the two other eigenvectors, $\mathbf{b}_1$ and $\mathbf{b}_2$. The simplex is partitioned into an *optimization cone* (white; red trajectories) where the mean replication rate $\bar{f}(t)$ is optimized, a second zone, the *master cone* where $\bar{f}(t)$ always decreases (white; blue trajectory), and two other zones where may increase, decrease or change nonmonotonously (grey; green trajectories). In this illustration $\mathbf{X}_3$ is chosen to be the master sequence. Solution curves are presented as parametric plots $\mathbf{x}(t)$. In particular, the parameter values are: $f_1 = 1.9\,[\mathrm{t}^{-1}]$, $f_2 = 2.0\,[\mathrm{t}^{-1}]$, and $f_3 = 2.1\,[\mathrm{t}^{-1}]$, the Q-matrix was assumed to be bistochastic with the elements $Q_{ii} = 0.98$ and $Q_{ij} = 0.01$ for $i, j = \{1, 2, 3\}$. Then the eigenvalues and eigenvectors of W are:

| k | $\lambda_k$ | $b_{1k}$ | $b_{2k}$ | $b_{3k}$ |
|---|---|---|---|---|
| 1 | 2.065 | 0.093 | 0.165 | 0.742 |
| 2 | 1.958 | 0.170 | 1.078 | -0.248 |
| 3 | 1.857 | 1.327 | -0.224 | -0.103 |

The mean replication rate $\bar{f}(t)$ is monotonously increasing along red trajectories, monotonously decreasing along the blue trajectory, and not necessarily monotonous along green trajectories.

yields a symmetric matrix [58], since $\mathrm{F}^{\frac{1}{2}} \cdot \mathrm{Q} \cdot \mathrm{F}^{\frac{1}{2}}$ is symmetric if Q is. Symmetric matrices have real eigenvalues, a similarity transformation does not change the eigenvalues and hence W has only real eigenvalues if Q is symmetric.

The simplest way to yield complex eigenvalues is introduction of cyclic symmetry into the matrix Q in such a way that the symmetry with respect

to the main diagonal is destroyed. An example is the matrix

$$Q \;=\; \begin{pmatrix} Q_{11} & Q_{12} & Q_{13} & \cdots & Q_{1n} \\ Q_{1n} & Q_{11} & Q_{12} & \cdots & Q_{1,n-1} \\ Q_{1,n-1} & Q_{1n} & Q_{11} & \cdots & Q_{1,n-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ Q_{12} & Q_{13} & Q_{14} & \cdots & Q_{11} \end{pmatrix} \;,$$

with different entries $Q_{ij}$. For equal replication parameters the eigenvalues contain complex $n$-th roots of one, $\gamma_k^n = 1$ or $\gamma_k = \exp(2\pi i k/n), i = 1, \ldots, n$, and for $n \geq 3$ most eigenvalues come in complex conjugate pairs. As mentioned earlier symmetry in mutation frequencies is commonly not fulfilled in nature. In case of point mutations the replacement of one particular base by another one does usually not occur with the same frequency as the inverse replacement, **G**→**A** versus **A**→**G** for example. Needless to stress, cyclic symmetry in mutation matrices is also highly improbable in real systems. The validity of Perron-Frobenius theorem, however, is not effected by the occurrence of complex conjugate pairs of eigenvectors. In addition, it is unimportant for most purposes whether a replication-mutation system approaches the stationary state monotonously or through damped oscillations (see next paragraph).

### 5.5  Optimization

In order to consider the optimization problem in the selection-mutation case, we choose the eigenvectors of W as the basis of a new coordinate system (see, for example, figure 7):

$$\mathbf{x}(t) \;=\; \sum_{i=1}^{n} x_k(t)\,\mathbf{e}_i \;=\; \sum_{k=0}^{n-1} \xi_k(t)\cdot\mathbf{b}_k \;,$$

wherein the vectors $\mathbf{e}_i$ are the unit vectors of the conventional Cartesian coordinate system and $\mathbf{b}_k$ the eigenvectors of W. The unit vectors represent the corners of $\mathbb{S}_n^{(1)}$ and in complete analogy we denote the space defined by the vectors $\mathbf{b}_k$ as $\tilde{\mathbb{S}}_n^{(1)}$. Formally, the transformed differential equation

$$\frac{d\xi_k}{dt} \;=\; \xi_k\,(\lambda_k - \phi)\,, \;\; k = 0, 1, \ldots, n-1 \;\; \text{with} \;\; \phi = \sum_{k=0}^{n-1} \lambda_k \xi_k = \overline{\lambda}$$

is identical to equation (9) and hence the solutions are the same,

$$\xi_k(t) \;=\; \xi_k(0)\,\exp\left(\lambda_k\,t - \int_0^t \phi(\tau)\,d\tau\right)\,, \;\; k = 0, 1, \ldots, n-1\;,$$

as well as the maximum principle on the simplex $\tilde{\mathbb{S}}_n^{(1)}$

$$\frac{d\phi}{dt} = \sum_{k=0}^{n-1} \frac{d\xi_k}{dt} \lambda_k = \sum_{k=0}^{n-1} \xi_k \lambda_k (\lambda_k - \phi) = <\lambda^2> - <\lambda>^2 \geq 0 . \quad (12a)$$

The difference between the representation of selection and selection-mutation comes from the fact that the simplex $\tilde{\mathbb{S}}_n$ does not coincide with the physically defined space $\mathbb{S}_n$ (see figure 7 for a three-dimensional example). Indeed, only the dominant eigenvector $\mathbf{b}_0$ lies in the interior of $\mathbb{S}_n^{(1)}$: It represent the stable stationary distribution of genotypes, the quasispecies $\bar{\Upsilon}$ towards which the solutions of the differential equation (18) converge. All other $n-1$ eigenvectors, $\mathbf{b}_1, \ldots, \mathbf{b}_{n-1}$ lie outside $\mathbb{S}_n^{(1)}$ in the *physically inaccessible* range where one or more variables $x_i$ are negative. The quasispecies $\bar{\Upsilon}$ represented by $\mathbf{b}_0$ is commonly dominated by a single genotype, the master sequence $\mathbf{X}_m$, having the largest stationary relative concentration: $\bar{x}_m \gg \bar{x}_i \forall i \neq m$, reflecting $f_m > f_i \forall i \neq m$ – and the same sequence in the elements of the matrix W: $W_{mm} > W_{ii} \forall i \neq m$ – and for sufficiently small mutation rates $p$: $W_{im} \ll \{W_{mm}, W_{ii}\} \forall i \neq m$. As sketched in figure 7 the quasispecies is then situated close to the unit vector $\mathbf{e}_m$ in the interior of $\mathbb{S}_n^{(1)}$.

For the discussion of the optimization behavior the simplex is partitioned into three zones: (i) The zone of maximization of $\phi(t)$, the (large) lower white area in figure 7 where equation (12a) holds and which we shall denote as *optimization cone*,[19] (ii) the zone that includes the unit vector of the master sequence, $\mathbf{e}_m$, and the quasispecies, $\mathbf{b}_0$, as corners, and that we shall characterize as *master cone*,[19] and (iii) the remaining part of the simplex $\mathbb{S}_n^{(1)}$ (two zones colored grey in figure 7). It is straightforward to proof that increase of $\phi(t)$ and monotonous convergence towards the quasispecies is restricted to the optimization cone [59]. From the properties of the selection equation (9) we recall and conclude that the boundaries of the simplex $\tilde{\mathbb{S}}_n^{(1)}$ are invariant sets. This implies that no orbit of the differential equation (18) can cross these boundaries. The boundaries of $\mathbb{S}_n^{(1)}$, on the other hand, are not invariant but have the restriction that they can be crossed exclusively in one direction: from outside to inside.[20] Therefore, a solution curve starting

---

[19]The exact geometry of the optimization cone or the master cone is a polyhedron that can be approximated by a pyramid rather than a cone. Nevertheless, we prefer the inexact notion *cone* because it is easier to memorize and to imagine in high-dimensional space.

[20]This is shown easily by analyzing the differential equation, but follows also from the physical background: No acceptable process can lead to negative particle numbers or concentrations. The process, however, can start at zero concentrations and this means the orbit begins at the boundary and goes into the interior of the physical concentration space, here the simplex $\mathbb{S}_n^{(1)}$.

in the optimization cone or in the master cone will stay inside the cone where it started and eventually converge towards the quasispecies, $\mathbf{b}_0$.

In zone (ii), the master cone, all variables $\xi_k$ except $\xi_0$ are negative and $\xi_0$ is larger than one in order to fulfill the $L^{(1)}$-norm condition $\sum_{k=0}^{n-1} \xi_k = 1$. In order to analyze the behavior of $\phi(t)$ we split the variables into two groups, $\xi_0$ the frequency of the quasispecies and the rest [59], $\{\xi_k; \, k = 1, \ldots, n-1\}$ with $\sum_{k=1}^{n-1} \xi_k = 1 - \xi_0$:

$$\frac{d\phi}{dt} = \lambda_0^2 \xi_0 + \sum_{k=1}^{n-1} \lambda_k^2 \xi_k - \left( \lambda_0 \xi_0 + \sum_{k=1}^{n-1} \lambda_k \xi_k \right)^2 .$$

Next we replace the distribution of $\lambda_k$ values in the second group by a single $\lambda$-value, $\tilde{\lambda}$ and find:

$$\frac{d\phi}{dt} = \lambda_0^2 \xi_0 + \tilde{\lambda}^2 (1 - \xi_0) - \left( \lambda_0 \xi_0 + \tilde{\lambda}(1 - \xi_0) \right)^2 .$$

After a view simple algebraic operations we find eventually

$$\frac{d\phi}{dt} = \xi_0 (1 - \xi_0) (\lambda_0 - \tilde{\lambda})^2 . \tag{27}$$

For the master cone with $\xi_0 \geq 1$, this implies $d\phi(t)/dt \leq 0$, the flux is a non-increasing function of time. Since we are only interested in the sign of $d\phi/dt$, the result is exact, because we could use the mean value $\tilde{\lambda} = \bar{\lambda} = (\sum_{k=1}^{n-1} \lambda_k \xi_k)/(1 - \xi_0)$, the largest possible value $\lambda_1$ or the smallest possible value $\lambda_{n-1}$ without changing the conclusion. Clearly, the distribution of $\lambda_k$-values matters for quantitative results. As it has to be, equation (27) applies also to the optimization cone and gives the correct result that $\phi(t)$ is non-decreasing. Decrease of mean fitness or flux $\phi(t)$ in the master cone is readily illustrated: Consider, for example, a homogeneous population of the master sequence as initial condition: $x_m(0) = 1$ and $\phi(0) = f_m$. The population becomes inhomogeneous because mutants are formed. Since all mutants have lower replication rate constants by definition, $(f_i < f_m \, \forall \, i \neq m)$, $\phi$ becomes smaller. Finally, the distribution approaches the quasispecies $\mathbf{b}_0$ and $\lim_{t \to \infty} \phi(t) = \lambda_0 < f_m$.

An extension of the analysis from the master cone to the grey zones, where not all $\xi_k$ values with $k \neq 0$ are negative, is not possible. It has been shown by means of numerical examples that $d\phi(t)/dt$ may show nonmonotonous behavior and can go through a maximum or a minimum at finite time [59].

## 5.6 Mutation rates and error threshold

In order to illustrate the influence of mutation rates on the selection process we apply (i) binary sequences, (ii) the uniform error rate approximation (16)

$$Q_{ij} = p^{d_{ij}} (1-p)^{\nu - d_{ij}} = (1-p)^{\nu} \varepsilon^{d_{ij}} \quad \text{with} \quad \varepsilon = \frac{p}{1-p}$$

with $d_{ij}$ being the Hamming distance between the two sequences $\mathbf{X}_i$ and $\mathbf{X}_j$, $\nu$ the chain length and $p$ the single point mutation or error rate parameter per site and replication, and (iii) a simple model for the distribution of fitness values known as *single-peak fitness landscape* [9],

$$f(\mathbf{Y}_k) = \begin{cases} f_0 & \text{if} \quad k = 0 = (m), \\ f & \text{if} \quad k = 1, \ldots, \nu \ (k \neq m). \end{cases} \quad (28)$$

All sequences are ordered in mutant classes $\mathbf{Y}_k$ with respect to the Hamming distance from the master sequence. In absence of neutrality the zero-error class contains only the master sequence ($\mathbf{Y}_0 : \{\mathbf{X}_m \equiv \mathbf{X}_0\}$), the one-error class comprises all single point mutations, the two-error class all double point mutations, etc.[21] Since the error rate $p$ is independent of the sequence and because of the assumption of a single-peak fitness landscape all molecules belonging to the same mutant class have identical fitness valued $f_{(k)} = \{f_0, f\}$, it is possible to introduce variables for entire mutant classes $\Gamma_k$ (figure 6):

$$y_k = \sum_{j, \mathbf{X}_j \in \Gamma_k} x_j, \quad k = 0, 1, \ldots, \nu, \quad \sum_{k=0}^{\nu} y_k = 1 . \quad (29)$$

The mutation matrix Q has to be adjusted to transitions between classes [9,60]. For mutations from class $\Gamma_l$ into $\Gamma_k$ we calculate:
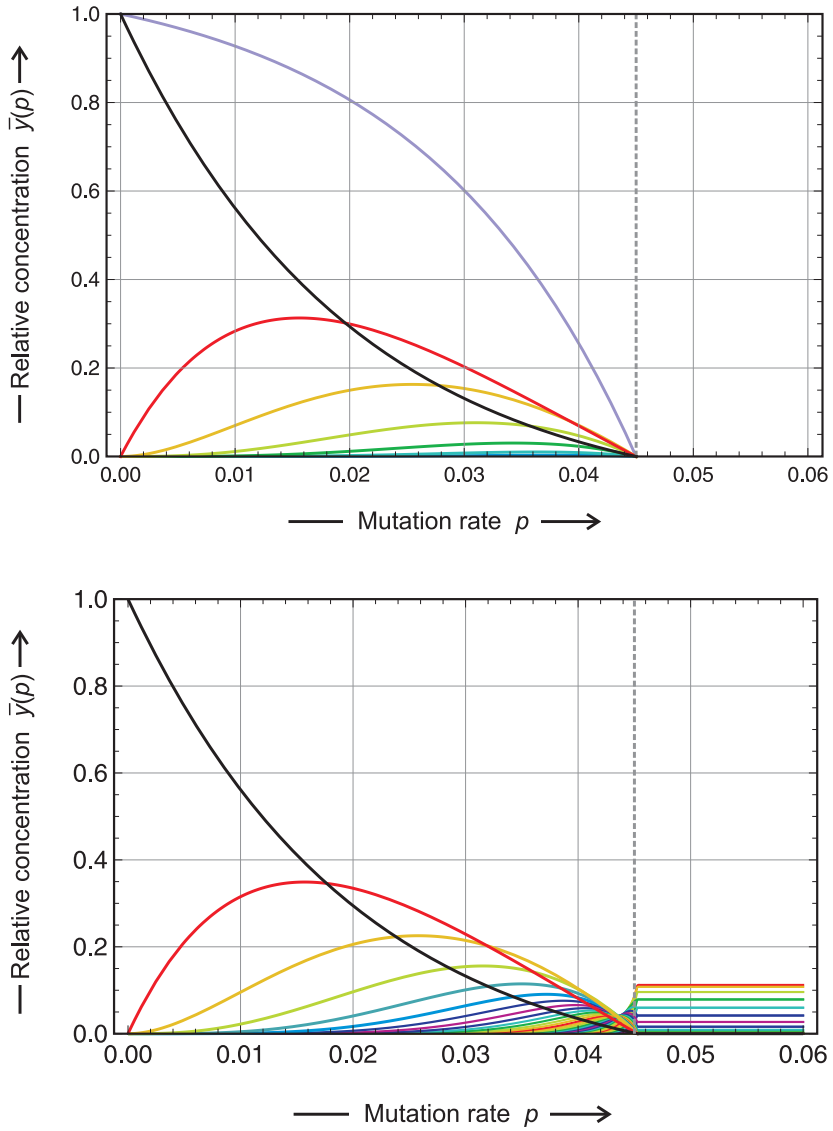
$$Q_{kl} = \sum_{i=l+k-\nu}^{\min(k,l)} \binom{k}{i} \binom{\nu-k}{l-i} p^{k+l-2i} (1-p)^{\nu-(k+l-2i)} . \quad (30)$$

The mutation matrix Q for error classes is not symmetric, $Q_{kl} \neq Q_{lk}$ as follows from equation (30).

The approaches to calculate error thresholds – numerical solutions and analytical approximations – are discussed here in rather fussy detail, because we shall make use of them later in the analysis of mutationally coupled clusters. In the following we present the zero mutational backflow approximation, an approach by Raleigh-Schrödinger perturbation theory and the numerical results computed for the (full) mutation-selection equation.

---

[21] Centering the quasispecies around the master sequence suggests to use decadic equivalents for the individual (binary) sequences: $\mathbf{X}_0 = \mathbf{000}\ldots\mathbf{00}$, $\mathbf{X}_1 = \mathbf{000}\ldots\mathbf{01}$, $\mathbf{X}_2 = \mathbf{000}\ldots\mathbf{10}, \cdots, \mathbf{X}_{2^{\nu}-1} = \mathbf{111}\ldots\mathbf{11}$.

**Figure 8: The quasispecies as a function of the point mutation rate $p$.**
The plot shows the stationary mutant distribution of sequences of chain length $\nu = 50$ on a single-peak fitness landscape as a function of the point mutation rate $p$. The upper part contains the approximate curves obtained through neglect of mutational backflow according to equation (32) and is compared with the numerical results presented in the lower part of the figure. Plotted are the relative concentration of entire mutant classes (figure 6): $\bar{y}_0$ (black) is the master sequence $\mathbf{X}_m \equiv \mathbf{X}_0$, $\bar{y}_1$ (red) is the sum of the concentrations of all one-error mutants of the master sequence, $\bar{y}_2$ (yellow) that of all two-error mutants, $\bar{y}_3$ (green) that of all three-error mutants, and so on. In the perturbation approach the entire population vanishes at a critical mutation rate $p_{cr}$ called the error threshold (which is indicated by a broken gray line at $p_{cr} = 0.04501$; the total concentration $\bar{c}^{(0)}(p)$ is shown in violet) whereas a sharp transition to the uniform distribution ($\Pi$)is observed with the numerical solutions. In the uniform distribution the concentration of class $k$ is given by $\binom{\nu}{k}/2^{\nu}$ with a largest value of $\bar{y}_{25} = 0.1123$ and a smallest value of $\bar{y}_0 = \bar{y}_{50} = 8.8818 \times 10^{-16}$. Choice of parameters: $f_m = 10$, $f = f_j = 1 \forall j = 1, \ldots, \nu$; $j \neq m$, and hence $\overline{f}_{-m} = 1$.

29

*Zero mutational backflow.*  Neglect of mutational backflow from mutants to the master sequence allows for the derivation of analytical approximations for the quasispecies [1,5]. The backflow is of the form

$$\Phi_{m\leftarrow(i)} \;=\; \sum_{i=1}^{n} Q_{mi} f_i \bar{x}_i \;=\; \sum_{i=1}^{n} W_{mi}\, \bar{x}_i \;,$$

and if $W_{mi} << |W_{mm} - W_{ii}|\, (i \neq m)$ is fulfilled $Q_{mi} = 0\, \forall\, i = 1, \ldots, n;\, i \neq m$ is a valid approximation for small mutation rates [9]. Insertion into equation (18) yields the following ODE for the master sequence[22]

$$\frac{\mathrm{d}x_m^{(0)}}{\mathrm{d}t} \;=\; (W_{mm} - \phi)\, x_m^{(0)} \;=\; (Q_{mm} f_m - \phi)\, x_m^{(0)} \;. \tag{31a}$$

$$\frac{\mathrm{d}x_j^{(0)}}{\mathrm{d}t} \;=\; W_{jm}\, x_m^{(0)} - \phi\, x_j^{(0)} \;=\; Q_{jm} f_m\, x_m^{(0)} - \phi\, x_j^{(0)} \;. \tag{31b}$$

The differential equation (31a) sustains two stationary states: (i) $\bar{x}_m^{(0)} = 0$, the state of extinction, and (ii) $Q_{mm} f_m - \phi = 0$. In the latter case we split $\phi$ as we did previously:

$$\phi \;=\; \sum_{i=1}^{n} f_i\, x_i^{(0)} \;=\; f_m\, x_m^{(0)} + \sum_{i=1, i\neq m}^{n} f_i\, x_i^{(0)} \;=\; f_m\, x_m^{(0)} + (1 - x_m^{(0)})\, \overline{f}_{-m}$$

$$\text{with} \quad \overline{f}_{-m} \;=\; \frac{1}{1 - x_m} \sum_{i=1, i\neq m}^{n} f_i\, x_i \;.$$

Insertion into condition (ii) yields

$$Q_{mm} f_m \;-\; f_m\, \bar{x}_m^{(0)} \;-\; (1 - \bar{x}_m^{(0)})\, \overline{f}_{-m} \;,$$

which can be evaluated to yield an expression for $\bar{x}_m^{(0)}$. For known concentrations of the master sequence we obtain the concentration of the mutants from equation (31b). For simplicity we introduce the assumption of the single peak landscape leading to $\overline{f}_{-m} = f$:

$$\bar{x}_j^{(0)} \;=\; \frac{Q_{jm} f_m\, \bar{x}_m^{(0)}}{(f_m - f)\, \bar{x}_m^{(0)} + f} \;.$$

which after some algebraic operations leads to an equation for the stationary concentrations of all members of the quasispecies

$$\bar{x}_m^{(0)} \;=\; \frac{W_{mm} - \overline{f}_{-m}}{f_m - \overline{f}_{-m}} \;=\; \frac{Q_{mm} - \sigma_m^{-1}}{1 - \sigma_m^{-1}} \;, \tag{32a}$$

$$\bar{x}_j^{(0)} \;=\; \varepsilon^{d_{i0}}\, \bar{x}_m^{(0)} \;;\; j = 1, \ldots, n\,,\, j \neq m \tag{32b}$$

$$\text{with} \;\; \sigma_m = f_m \,/\, \overline{f}_{-m} = f_m \,/\, f \;.$$

---

[22]The superscript '(0)' stands for zeroth order perturbation theory and means (total) neglect of mutational backflow, although the approach does not correspond to a defined order of perturbation theory (see next paragraph).

The superiority $\sigma_m$ is a measure of the advantage in fitness the master has over the rest of the population, and $\overline{f}_{-m}$ is the mean fitness of this rest.[23] In case of the single peak fitness landscape we have the trivial result: $\overline{f}_{-m} = f$. The superiority of the master sequence can also be understood as an empirical quantity that can be determined through direct measurements of replication efficiencies of cloned sequences.

The stationary concentration of the master sequence $\bar{x}_m$ as a function of the error rate $p$, and according to equation (32b) also the concentration of all other sequences $\bar{x}_j$ belonging to the quasispecies vanish at the critical error rate $p_{cr}$. An illustrative example is shown in figure 8. The stationary concentrations for all sequences except the master sequence pass through a maximum before the vanish at $p = p_{cr}$ and the position of this maximum, $p = p_{max}$, for sequence $\mathbf{X}_i$ can be obtained from the implicit equation

$$Q(p_{max}\left(d_{i0} - \nu p_{max}\right) = (1 - p_{max})^{\nu}\left(d_{i0} - \nu p_{max}\right) = d_{i0}\,\sigma_m^{-1}\,.$$

The maximum is shifted towards the error threshold with increasing distance from the maser sequence. The value of $p_{cr}$ is readily calculated from equation (32):

$$
\begin{aligned}
\bar{x}_m^{(0)} &= 0 \implies Q_{mm} - \sigma_m^{-1} = (1 - p_{cr})^{\nu} - \sigma_m^{-1} = 0\,, \\
p_{cr} &= 1 - \sigma_m^{-1/\nu} \ \text{ or } \ \nu\ln(1 - p_{cr}) = -\ln\sigma_m \ \text{ and} \\
p_{cr} &= p_{max} \approx \frac{\ln\sigma_m}{\nu} \ \text{ and } \ \nu_{max} \approx \frac{\ln\sigma}{p}\,,
\end{aligned}
\tag{33}
$$

where the relation $\ln(1 - z) \approx 1 - (z + \ldots) \forall -1 \leq z < 1$ has been used. Truncation after the first term in $z$ is valid for sufficiently small error rates.

The error threshold has two meanings: (i) for constant chain length $\nu$ the genetically tolerable error rate is limited, $p < p_{max} = p_{cr}$ and (ii) for a given replication accuracy the length of polynucleotide that can be faithfully replicated is limited by $\nu < \nu_{max}$. Relation (i) is important in virology and relation (ii) has been used in models of enzyme-free prebiotic reproduction of oligo- and polynucleotides [7]. RNA viruses commonly have mutation rates close to the error threshold [10]. The error rates can be increased by pharmaceutical drugs interfering with virus replication and accordingly, a new antiviral strategy has been developed, which drives virus replication into an error catastrophe [11, 62]. Recently, a direct extinction mechanism of lethal mutagenesis in virus infections has been extensively discussed [53, 54].

---

[23]An exact calculation of $\overline{f}_{-m}$ is difficult because it requires knowledge of the stationary concentrations of all variants in the population: $\bar{x}_i$; $x = 1, \ldots, n$. For computational details see [3, 5, 6, 61].

Zero mutational backflow fails to account for the quasispecies at mutation rates above the threshold $p_{\text{cr}}$: Perron-Frobenius theorem states that the concentrations of all members of the quasispecies are positive definite: $\bar{x}_i > 0 \, \forall \, i = 1, \ldots, n$ but zero mutational backflow yields $\bar{x}_i = 0 \, \forall \, i = 1, \ldots, n$ at $p = p_{\text{cr}}$. Considering the problem more closely this is no surprise since the zero mutational backflow assumption violates the conditions for the validity of the theorem: The requirement for matrix W was irreducibility and this implies that every sequence can be reached from every other sequence in a finite number of mutation steps – zero mutational backflow implies that the master cannot be reached from the mutants. Beyond the error threshold we have to consider either full first order perturbation theory or the numerical solutions. The manipulation of the elements of the matrix Q has also the consequence that the stationary total concentration $\bar{c}^{(0)} = \sum_{i=1}^{n} \bar{x}_i^{(0)}$ is not constant but vanishes at the error threshold

$$\bar{c}^{(0)}(p) = \frac{1}{Q} \frac{Q - \sigma_m^{-1}}{1 - \sigma_m^{-1}} .$$

Clearly, the excellent agreement between the zero mutational backflow approximation and the exact solution is fortuitous but as many examples have shown it is quite general and it would be worth to search for the reason.

*Perturbation theory.* Application of first and second order Rayleigh-Schrödinger perturbation theory to calculate the quasispecies $\bar{\Upsilon}$ as a function of the mutation rate has been performed in the past [9]. Here we present the full analytical first order expressions $\bar{x}_i^{(1)}(p)$. The second expressions are rather clumsy and bring only limited improvement for small mutation rates $p$.

The largest eigenvalue is the same by zero mutational backflow and in first order perturbation theory:

$$\lambda_0^{(0)} = \lambda_0^{(1)} = W_{mm} = Q_{mm} f_m .$$

For the computation of the largest eigenvector we make use of the first order expression from perturbation theory of the matrix W (32b):[24]

$$\bar{x}_j^{(1)} = \frac{W_{jm}}{W_{mm} - W_{jj}} \, \bar{x}_m^{(1)} \, ; \; j = 1, \ldots, n \, , \, j \neq m \, .$$
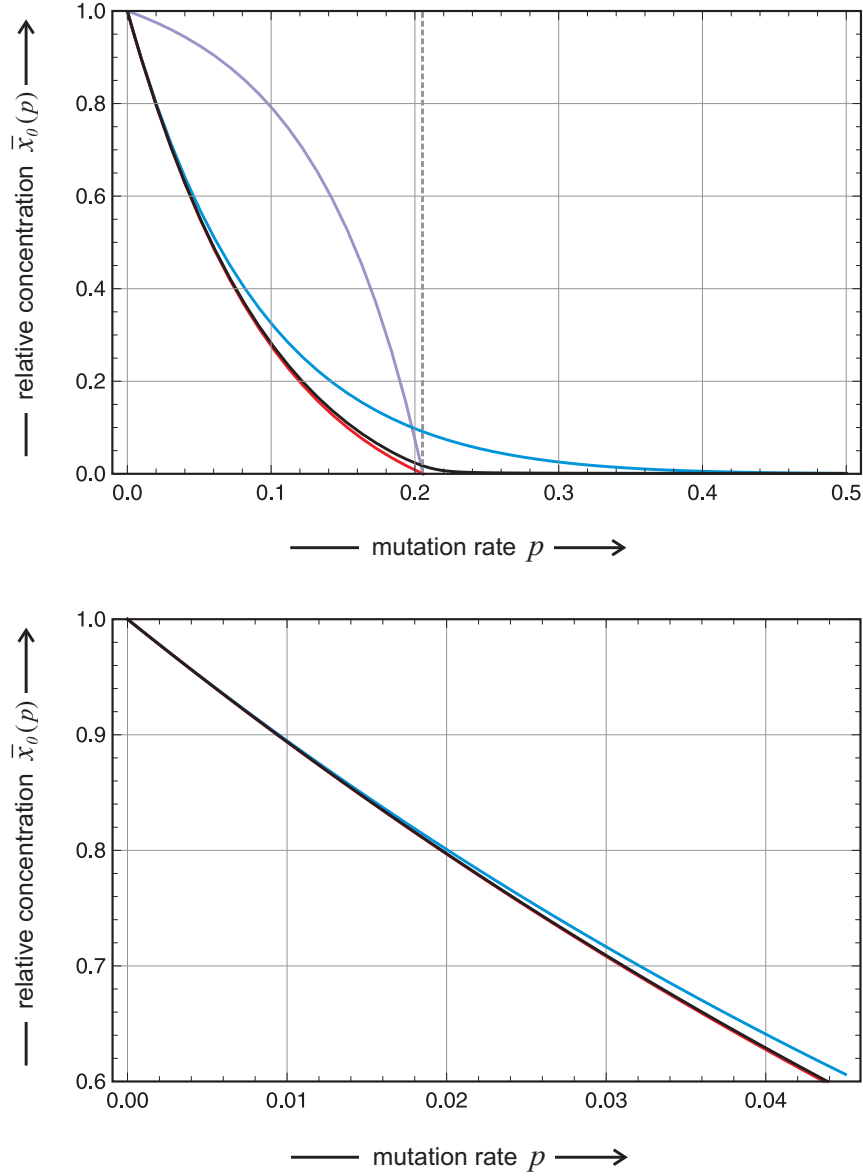
Making use of the normalization condition $\sum_{i=1}^{n} \bar{x}_i^{(1)} = 1$ we obtain for the master sequence

$$\bar{x}_m^{(1)} = \frac{1}{1 + \sum_{i=1, i \neq m}^{n} \frac{W_{im}}{W_{mm} - W_{ii}}} .$$

---

[24] As a matter of fact, the first order perturbation expressions are used in the zero mutational backflow approximation for the calculation of the concentrations of the mutants, because

**Figure 9: Quasispecies calculated by perturbation theory.** The upper plot presents a comparison of the stationary concentration of the master sequence in the zero mutational backflow approximation, $\bar{x}_m^{(0)}(p)$ (red), with first order perturbation theory, $\bar{x}_m^{(1)}(p)$ (blue), and numerical solution $\bar{x}_m$ (black). The violet curve is the total concentration in the zero mutational backflow approximation, $\bar{c}^{(0)} = \sum_{i=1}^{n} \bar{x}_i^{(0)}$. The lower plot is an enlargement of the curves at low mutation rates demonstrating that both approximations almost coincide with the exact solution curve. The error threshold indicated by a broken vertical line occurs at $p = p_{\mathrm{cr}} = 0.2057$. Choice of parameters: $n = 10$, $f_0 = 10\,[\mathrm{t}^{-1}]$, and $f = 1\,[\mathrm{t}^{-1}]$.

Straightforward calculations yield for the stationary concentrations

$$\bar{x}_m^{(1)}(p) \;=\; \frac{Q\,(1 - \sigma_m^{-1})}{1 - Q\,\sigma_m^{-1}} \;, \tag{34a}$$

$$\bar{x}_j^{(1)}(p) \;=\; \frac{Q_{im}}{Q\,(1 - \sigma_m^{-1})}\,\bar{x}_m^{(0)}\;; \;\; j = 1, \ldots, n\,, \; j \neq m \tag{34b}$$

$$\text{with } \; Q = (1 - p)^\nu \;.$$

33

As shown in figure 9 the curve $\bar{x}_m^{(1)}(p)$ extends to the point $p = \tilde{p} = \frac{1}{2}$ and further, but it does not pass precisely through the uniform distribution,

$$\bar{x}_m^{(1)}\left(\tfrac{1}{2}\right) \;=\; \left(\tfrac{1}{2}\right)^\nu \frac{1 - \sigma_m^{-1}}{1 - (\tfrac{1}{2})^\nu \sigma_m^{-1}} \;.$$

The deviation of $\bar{x}_m^{(1)}$ from numerical solution is much larger than in the zero mutational backflow approximation approximation and the error threshold phenomenon is not detectable. In summary, first order perturbation theory provides a consistent approximation to the eigenvalues and eigenvectors of the value matrix W. The results, however, are not nearly as good as those of the zero back mutation approach. Improvements by second order are possible at very small error rates but the calculations are rather tedious and the solutions for the eigenvalue $\lambda_0^{(2)}$ become unstable for larger error rates [9]. A combination of zero mutation flux approximation and first order perturbation theory in the sense of equations (31a) and (34b) [1, 9] leads to slightly better results than the zero mutation flux approach alone but is not recommended because of the lack of consistency.

*Numerical solutions.* Full solutions can be computed numerically through solving the eigenvalue problem of matrix W for different values of of the mutation rate $p$, and for a typical example the normalized concentrations of error classes $\bar{y}_{(k)}(p)$ are shown in figure 8. The agreement between the numerical results and the zero mutational backflow curve for the master class, $\bar{y}_0(p)$ in the region above the error threshold is remarkable indeed. The other solution curves $\bar{y}_{(k)}(p)$ $(k \neq 0)$ agree well too but the deviations become larger with increasing $k$.

The numerical solution for the master sequence (black curve) decreases monotonously from $p = 0$ to $p = \tilde{p} = 1/2$, this is between two points for which analytical solutions exist. At vanishing error rates, $\lim p \to 0$, the master sequence is selected, $\lim_{t \to \infty} x_0(t) = \lim_{t \to \infty} y_0(t) = \bar{y}_0 = \bar{x}_0 = 1$, and all other error classes vanish in the long time limit. Increasing error rates are reflected by a decrease in the stationary relative concentration of the master sequence and a corresponding increase in the concentration of all mutant classes. Except $\bar{y}_0(p)$ all concentrations $\bar{y}_k(p)$ with $k < \nu/2$ go through a maximum at values of $p$ that increases with $k$ – as in case of zero mutational backflow where we had an implicit analytical expression for the maximum, and approach the curves for $\bar{y}_{\nu-k}$ – whereas the zero mutational backflow curves still go through a maximum because they vanish at $p = p_{\mathrm{cr}}$. At $p = \tilde{p} = 1/2$ we have $\tilde{p} = 1 - \tilde{p}$ for binary sequences, and again the eigenvalue problem can be solved exactly. The value matrix $W = Q \cdot F$ is of

the form

$$
W \;=\; \frac{1}{2^\nu}
\begin{pmatrix}
1 & 1 & \cdots & 1 \\
1 & 1 & \cdots & 1 \\
\vdots & \vdots & \ddots & \vdots \\
1 & 1 & \cdots & 1
\end{pmatrix}
\cdot
\begin{pmatrix}
f_1 & 0 & \cdots & 0 \\
0 & f_2 & \cdots & 0 \\
\vdots & \vdots & \ddots & \vdots \\
0 & 0 & \cdots & f
\end{pmatrix}
\;=\;
\frac{1}{2^\nu}
\begin{pmatrix}
f_1 & f_2 & \cdots & f_n \\
f_1 & f_2 & \cdots & f_n \\
\vdots & \vdots & \ddots & \vdots \\
f_1 & f_2 & \cdots & f_n
\end{pmatrix} .
$$

The matrix W consists of $n$ identical columns and hence has $n-1$ singularities corresponding to eigenvalues $\lambda_i = 0 \,\forall\, i = 1, \ldots, n-1$. The $n$-th eigenvalue follows from the trace of the matrix W since the trace is invariant under the similarity transformation $\Lambda = B^{-1} \cdot W \cdot B$ and hence

$$
\sum_{i=0}^{n-1} \lambda_i \;=\; \frac{1}{2^\nu} \sum_{i=1}^{n} f_i \quad \text{and hence} \quad \lambda_0 \;=\; \frac{1}{2^\nu} \sum_{i=1}^{n} f_i \, .
$$

The largest eigenvalue $\lambda_0$ is strictly positive and the corresponding largest eigenvector $\mathbf{b}_0$ is given by the uniform distribution:

$$
\mathbf{b}_0 \;=\; \frac{1}{2^\nu} \, (1, 1, \ldots, 1)^{\mathrm t} \, .
$$

All normalized stationary concentrations, $\bar{x}_1 = \bar{x}_2 = \ldots = \bar{x}_n = 1/2^\nu$, are the same. The individual class variables are given by $\bar{y}_k = \binom{\nu}{k} \big/ 2^\nu$. The uniform distribution ($\Pi$) is a result of the fact that correct digit incorporation and point mutation are equally probable for binary sequences at $\tilde{p} = 1/2 = 1 - \tilde{p}$ and therefore we may characterize this scenario as random replication.[25] It is worth mentioning that the range of high mutation rates $\tilde{p} \leq p \leq 1$ is also meaningful: At $p = 1$ the complementary digit is incorporated with ultimate accuracy, $\mathbf{0} \to \mathbf{1}$ and $\mathbf{1} \to \mathbf{0}$, and accordingly, the range at high $p$-values describes error-prone complementary or plus-minus replication [9].

The special situation at the error threshold is the occurrence of an (almost) uniform distribution far away from the point $p = \tilde{p}$ – in figure 8 the critical mutation rate is $p_{\mathrm{cr}} = 0.045 \ll \tilde{p} = 0.5$ (figure 15 illustrates the extension of the domain of the uniform distribution). As we shall see in the next section 6, the error threshold on the single peak landscape is characterized by the coincidence of three phenomena: (i) the concentration of the master sequence becomes very small and this is expressed in term of level crossing values $\bar{y}_0(p)|_{p=p_{(1/M)}} = 1/M$ where $M$ is 100, 1000 or higher depending on the size of $2^\nu$, (ii) a sharp change in the quasispecies distribution within a narrow band of $p$-values that reminds of a phase transition [63–67], and (iii) a transition to the uniform distribution, which implies that the domain within

---

[25]The extension to sequences over an alphabet with $\kappa$ classes of digits is straightforward. In the frame of uniform errors random replication occurs at $\tilde{p} = 1/\kappa = (1 - \tilde{p})/(\kappa - 1)$. For the natural four letter alphabet we have $\tilde{p} = 1/4$.

**Figure 10: Existence of the error threshold.** The plots represent the exact solution (black) together with the zero mutational backflow approximation (green), the uniform backflow approximation (red) and the error-class one backflow approximation (red). The numerically exact solution is entrapped between the uniform and the one error-class approximation. Since both approximations converge to zero in the limit of long chain lengths ($\nu \to \infty$) the exact curve does as well. The error threshold as indicated by a broken vertical line occurs at $p = p_{cr} = 0.2057$. Choice of parameters: $n = 10$, $f_0 = 10\,[\mathrm{t}^{-1}]$, and $f = 1\,[\mathrm{t}^{-1}]$.

which the uniform distribution is fulfilled to a high degree of accuracy has the form of a broad plateau ($p_{cr} = 0.045 < p < \tilde{p} = 0.5$ in figure 8). It is worth considering the numerical data from the computations shown in the figure: $p_{cr} = 0.04501$ from zero mutational backflow versus the level crossing values $p_{(1/100)} = 0.04360$, $p_{(1/1000)} = 0.04509$, and $p_{(1/10000)} = 0.04525$.

*Proof for the existence of an error threshold.* In this paragraph we present a rigorous proof for the existence of an error threshold on the single peak landscape in the sense that the exact solution converges to the zero mutational back flow result in the limit of infinite chain length $\nu$. Previously we stated that the agreement between the (exact) numerical solution for the stationary quasispecies and the zero mutational backflow in surprisingly good and here we shall give a rigorous basis for this agreement. The proof proceeds in three steps: (i) Models are derived that provide upper and lower bounds for the exact solution, (ii) the models are evaluated analytically in order to yield expressions for the relative stationary concentration of the master sequence $\mathbf{X}_m$ at the position of the error threshold, $\bar{x}_m^{(\text{flow})}(p_{\text{cr}})$, and (iii) we show that the values for the upper and the lower bound coincide in the limit $\nu \to \infty$.

The zero mutational backflow approximation neglects backflow completely; now we introduce two other approximations that are based on model backflows that represent lower and upper bounds for the exact backflow. Computation of the mutational backflow requires either knowledge of the distribution of concentrations of all sequence of an assumption about it. In order to be able to handle the problem analytically the backflow must lead to an autonomous equation for the master concentration $x_0$. The minimal backflow can be estimated by the assumption of a uniform distribution (II) for all sequences except the master. In this case all sequences contribute equally no matter whether a particular sequence is close to the master sequences or far apart. For the concentrations $x_i = (1 - x_0^{(\text{II})})/(n-1) \, \forall \, i = 1, \ldots, n$ with $n = \kappa^\nu$ we obtain under the further assumption of a single peak landscape and the uniform error rate model the ODE for $x_0^{(\text{II})}$:

$$
\frac{\mathrm{d}x_0^{(\text{II})}}{\mathrm{d}t} = x_0^{(\text{II})}(Q_{00}f_0 - \phi) + f\frac{1 - x_0^{(\text{II})}}{n-1} \tag{35}
$$
$$
\text{with} \quad \phi = f + (f_0 - f)\,x_0^{(\text{II})}
$$

The stationary concentration is obtained as the solution of a quadratic equation

$$
\bar{x}_0^{(\text{II})} = \frac{Qf_0 - f - f\gamma(1-Q) + \sqrt{\left(Qf_0 - f - f\gamma(1-Q)\right)^2 + 4(f_0 - f)(1-Q)f\gamma}}{2(f_0 - f)}
$$
$$
\text{with} \quad Q = Q_{00} = (1-\mathrm{p})^\nu \quad \text{and} \quad \gamma = \frac{1}{\mathrm{n}-1}
$$

Insertion of the value of the mutation rate parameter at the error threshold, $p = p_{\text{cr}} = 1 - \sigma^{-1/\nu}$ or $Q = (1 - p_{\text{cr}})^\nu = \sigma^{-1}$ leads to the result

$$
\bar{x}_0^{(\text{II})}(p_{\text{cr}}) = \frac{1}{2}\frac{\sqrt{1 + 4\sigma(n-1)} - 1}{\sigma(n-1)}, \tag{36}
$$

which yields in the limit of long chains or large $\nu$-values

$$\bar{x}_0^{(\mathrm{II})}(p_{\mathrm{cr}}) \approx \frac{1}{\sqrt{\sigma\,n}} = \frac{1}{\sqrt{\sigma\,\kappa^\nu}} . \qquad (36')$$

Ultimately the value of the stationary concentration of the master sequence decays exponentially with one halt of the chain lengths as exponent: $\bar{x}_0^{(\mathrm{II})}(p_{\mathrm{cr}}) \propto \kappa^{-\nu/2}$. It is straightforward to show that the uniform mutational backflow approximation becomes exact at the point $p = \tilde{p}$ and insertion in the quadratic equation yields: $\bar{x}_0^{(\mathrm{II})}(\frac{1}{2}) = (\frac{1}{2})^\nu$. Generalization, of course, is straightforward and yields $p = \tilde{p}$ and insertion in the quadratic equation yields: $\bar{x}_0^{(\mathrm{II})}(1/\kappa) = (1/\kappa)^\nu$

In order to find an upper bound for the stationary solution of the master sequence we assume that mutational backflow comes only form the sequences in the one error class, $\Gamma_1^{(0)}$, which can be assumed to be present at equal concentrations, $x_i = (1 - x_0^{(\mathrm{I})})/\nu$, and $\sum_{i=1}^{\nu} x_i = 1 - x_0$. All other sequences except the master sequence and the one error class are absent. Pointing at the fact that $\Gamma_1^{(0)}$ represents the entire mutant cloud we shall denote this distribution by **I**. For the corresponding elements of the mutation matrix Q we use the usual expressions, which are all equal: $Q_{0i} = Q_{0(1)} = Q_{01}\,\forall\,i = 1, \ldots, \nu$. The ODE for the master sequence is then again autonomous and can be readily solved for the stationary state:

$$\frac{\mathrm{d}x_0^{(\mathrm{I})}}{\mathrm{d}t} = x_0^{(\mathrm{I})}(Q_{00}f_0 - \phi) + Q_{01}f\,(1 - x_0^{(\mathrm{I})})$$

$$\text{with } \phi = f + (f_0 - f)\,x_0^{(\mathrm{I})} . \qquad (37)$$

The stationary concentration is again obtained from a quadratic equation of similar structure as before

$$\bar{x}_0^{(\mathrm{I})} = \frac{Q_{00}f_0 - Q_{01}f - f + \sqrt{\left(Q_{00}f_0 - Q_{01}f - f\right)^2 + 4(f_0 - f)Q_{01}f}}{2(f_0 - f)}$$

$$\text{with } Q_{00} = (1 - p)^\nu \text{ and } Q_{01} = (1 - p)^{\nu-1}\,p .$$

It is shown straightforwardly that the curve for the class one backflow passes through the point $p = \tilde{p} = \kappa^{-\nu}$. For the stationary concentration of the master sequence at the error threshold we find

$$\bar{x}_0^{(\mathrm{I})}(p_{\mathrm{cr}}) = -\frac{Q_{01}f}{2\,(f_0 - f)} + \sqrt{\frac{Q_{01}f}{f_0 - f}} \cdot \sqrt{1 + \frac{Q_{01}f}{4\,(f_0 - f)}} , \qquad (38)$$

with three components. Before we can discuss the individual terms we have to examine the asymptotic dependence of the mutation rate $p$ on the chain length $\nu$, which is encapsulated in the series expansion

$$Q_{01} = (1 - p)^{\nu-1}\,p = p - (\nu - 1)\,p^2 + \frac{(\nu - 1)(\nu - 2)}{2}\,p^3 - \cdots$$

with the first term being $p$. The critical mutation rate can be approximated by $p_{\mathrm{cr}} \approx \ln\sigma/\nu$ and we can consider equation (38). The negative term in equation

38

shows an asymptotic dependence on the chain length of $\nu^{-1}$, the first factor behaves asymptotically like $1/\sqrt{\nu}$ whereas the second factor converges to unity. What remains in the limit of long chains or large $\nu$-values is

$$\bar{x}_0^{(\mathrm{I})}(p_{\mathrm{cr}}) \approx \sqrt{\frac{f \ln \sigma}{f_0 - f}} \cdot \frac{1}{\sqrt{\nu}} = \sqrt{\frac{\ln \sigma}{\sigma - 1}} \cdot \frac{1}{\sqrt{\nu}} . \tag{38'}$$
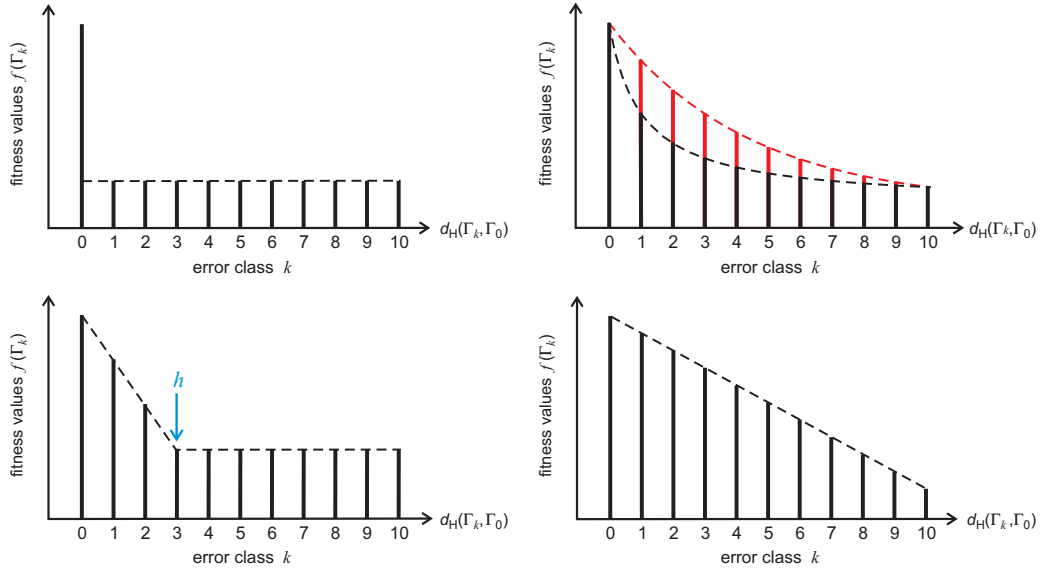
The value of the stationary concentration of the master sequence decays with the reciprocal square root of the chain length: $\bar{x}_0^{(\mathrm{I})}(p_{\mathrm{cr}}) \propto 1/\sqrt{\nu}$. Although the class one uniform distribution is not an impressively good upper bound for the exact solution curve, it is sufficient for our purpose here because $\bar{x}_0^{(\mathrm{I})}(p_{\mathrm{cr}})$ vanishes in the limit $\nu \to \infty$.

In summary the four values of the solution of the mutation-selection equation (18) and the three approximation appear in the following order at the critical mutation rate $p_{\mathrm{cr}}$:

| | | $\bar{x}_0$-value at the mutation rate | | |
|---|---|---|---|---|
| mutational backflow | notation | $p = 0$ | $p = p_{\mathrm{cr}}$ | $\tilde{p} = \frac{1}{\kappa}$ |
| class one uniform | $\bar{x}_0^{(\mathrm{I})}(p)$ | 1 | $\sqrt{\ln \sigma/(\sigma - 1)} \cdot 1/\sqrt{\nu}$ | $\kappa^{-\nu}$ |
| exact | $\bar{x}_0(p)$ | 1 | computed | $\kappa^{-\nu}$ |
| all uniform | $\bar{x}_0^{(\mathrm{II})}(p)$ | 1 | $1/\sqrt{\sigma \kappa^{\nu}}$ | $\kappa^{-\nu}$ |
| zero | $\bar{x}_0^{(0)}(p)$ | 1 | 0 | negative |

The exact solution is indeed entrapped between the two approximations for the mutational backflow and, since both converge asymptotically to zero the exact curve approaches the zero mutational backflow approximation in the limit of long chains. All four curves (figure 10) start at the point $\bar{x}_m(0)$ and all except the zero backflow approximation end at the correct value $\tilde{p} = \kappa^{-\nu}$. The stationary master concentrations at the critical mutation rate $p = p_{\mathrm{cr}}$ illustrate the relative importance of the mutational backflow (figure 38): Zero backflow assumption causes the stationary concentration $\bar{x}_0^{(0)}$ to vanish. The all uniform backflow is a little more than one half of the computed exact value, and this approximation is a excellent lower bound for the exact solution. The error one class backflow is about three time as large as the exact solution. Nevertheless it is an upper bound for the real mutation flow and serves the purpose for which it is intended here. If one is interested in an approximation apart from this proof the zero mutational back flow approximation $\bar{x}_0^{(\mathrm{II})}(p)$ in the region $0 \le p < p_{\mathrm{cr}}$ and the uniform backflow approximation in the entire range are suitable approximations. In particular the uniform backflow approximation is well suited because it is exact for $p = 0$ and $p = \tilde{p} = 1/\kappa$ and it has also a correct asymptotic behavior in the long chain limit at the error threshold.

The error threshold has been put in relation to phase transitions [63–65]. Here, we are in a position to prove the behavior of the exact solution curve $\bar{x}_0(p)$ in the

**Figure 11: Some examples of model fitness landscapes**. The figure shows five model landscapes with identical fitness values for all sequences in a given error class: (i) the single peak landscape (upper left drawing), (ii) the hyperbolic landscape (upper right drawing, black curve), (iii) the step-linear landscape (lower left drawing), (iv) the multiplicative landscape (upper right drawing, red curve), and (v) the additive or linear landscape (lower right drawing). Mathematical expressions are given in the text.

limit $\nu \to \infty$. The critical mutation rate converges to the value zero: $\lim_{\nu \to \infty} p_{\mathrm{cr}} = \lim_{\nu \to \infty} \ln \sigma / \nu = 0$. At the same time we have $\lim_{\nu \to \infty} \bar{x}_0 = 0$ for $p > 0$ and thus the quasispecies degenerates to an "L"-shaped distribution, $\bar{x}_0(0) = 1$ and $\bar{x}_0(p) = 0 \,\forall\, p > 0$, and we are left with a pathological phase transition at $p_{\mathrm{cr}} = 0$.

## 6 "Simple" landscapes

The existence and form of the error threshold turned out to be dependent on the nature of the fitness landscape [68]. In particular, two types of landscapes that are commonly used in population genetics, the additive and the multiplicative landscape, don't show error thresholds at all. In addition to the single peak landscape four examples of simple model landscapes are compared, in which identical fitness values are assigned to all members of the same mutant class:

(i) the additive or linear landscape
$$f(\mathbf{Y}_k) = f_k = f_0 - (f_0 - f)\,k/\nu\,; \ \ k = 0, 1, \ldots, \nu\,,$$

(ii) the multiplicative landscape
$$f(\mathbf{Y}_k) = f_k = f_0 \left(\tfrac{f}{f_0}\right)^{k/\nu}\,; \ \ k = 0, 1, \ldots, \nu\,,$$

(iii) the hyperbolic landscape
$$f(\mathbf{Y}_k) = f_k = f_0 - (f_0 - f)\left(\tfrac{\nu+1}{\nu}\right)\left(\tfrac{k}{k+1}\right)\,; \ \ k = 0, 1, \ldots, \nu\,, \text{ and}$$

40

(iv) the step-linear landscape

$$f(\mathbf{Y}_k) = f_k = \begin{cases} f_0 - (f_0 - f)\, k/h & \text{if } k = 0, 1, \ldots, h-1\,, \\ f & \text{if } k = h, \ldots, \nu\,. \end{cases}$$

In order to be able to compare the different landscapes the values of $f$ were chosen such that all landscapes are characterized by the same superiority of the master sequence: $\sigma_m = \sigma_0 = f_0 \,/\, \overline{f}_{-0}$ with $\overline{f}_{-0} = \sum_{i=1}^{n} y_i f_i \,/\, (1 - y_0)$. Since the distribution of concentrations is not known *a priori* we have to make an assumption. As said in the previous subsection 5.6 the uniform distribution extends down to the error threshold for decreasing mutation rates and hence, the assumption of the uniform distribution in the calculation of $f$ is well justified,

$$f = \left( \overline{f}_{-0}\,(2^\nu - 1) - f_0\,(2^{\nu-1} - 1) \right) \Big/ 2^{\nu-1}\,, \tag{39a}$$

$$f = \left( \left( \overline{f}_{-0}\,(2^\nu - 1) + f_0 \right)^{1/\nu} - f_0^{1/\nu} \right)^\nu\,, \tag{39b}$$

$$f = \left( \overline{f}_{-0}\,\nu\,(2^\nu - 1) - f_0\,(2^\nu - \nu + 1) \right) \Big/ \left( 2^\nu(\nu - 1) + 1 \right)\,, \tag{39c}$$

$$f = \frac{\overline{f}_{-0}\,(2^\nu - 1) - f_0 \left( \sum_{k=0}^{h-1} \binom{\nu}{k} \frac{h-k}{h} - 1 \right)}{\sum_{k=0}^{h-1} \binom{\nu}{k} \frac{k}{h} + \sum_{k=h}^{\nu} \binom{\nu}{k}}\,. \tag{39d}$$

In figures 12 and 13 the solution curves $\bar{y}_k(p)$ are compared for the three landscapes showing error thresholds, single-peak, hyperbolic and step-linear. The superiority was adjusted to $\sigma_0 = 10$ be means of equation (39). As already seen in figure 8 the calculated value for $p_{\mathrm{cr}}$ coincides perfectly with the position of the transition on the single-peak landscape, and the $p$-values for concentration level crossing lie close together and near $p_{\mathrm{cr}}$ (see table 1) indicating a rather steep decrease of $\bar{y}_0$ in the range on the left-hand side of the transition to the uniform distribution. The decrease of $\bar{y}_0$ is flatter on the hyperbolic landscape, the actual transition occurs slightly above the error threshold on the single-peak landscape and does not result in the uniform distribution. Instead we observe a mutant distribution above the error threshold, which changes slightly with $p$. On the step-linear landscape, eventually, the curve of $\bar{y}_0$ is even flatter, the transition is shifted further to higher $p$-values, and the transition leads to the uniform distribution as in the single-peak case (for an explanation see the discussion of the additive landscape).

Next we consider two examples of landscapes, which do not sustain error thresholds, the additive or linear landscape and the multiplicative landscape. As shown in figure 14 the change from the homogeneous distribution at $p = 0$

**Figure 12: Error thresholds on different model landscapes.** The figures show stationary concentrations of mutant classes as functions of the error rate, $\bar{y}_k(p)$, for sequences of chain length $\nu = 100$ with $f_0 = 10$ and $\overline{f}_{-0} = 1$ on three different model landscapes: the single peak landscape (upper part, $f = 1$), the hyperbolic landscape (middle part, $f = 10/11$), and the step-linear landscape (lower part, $f = 1$). The dashed line indicates the value of the error threshold calculated by zero mutational backflow, $p_{\mathrm{cr}} = 0.022768$.

**Figure 13: Error thresholds on different model landscapes.** The three figures are enlargements of the plots from in figure 12. Stationary concentrations of mutant classes, $\bar{y}_k(p)$, are shown for the single peak landscape (upper part), the hyperbolic landscape (middle part), and the step-linear landscape (lower part; see the caption figure 12 for details).

**Table 1: Concentration level crossing near the error threshold.** The decline of the master class, $\bar{y}_0 = \bar{x}_0$, at $p$-values below the error threshold $p_{\mathrm{cr}}$ is illustrated by means of the points $p_{(1/M)}$ where $\bar{y}_0(p)$ crosses the level $1/M$ for the three fitness landscapes that sustain error thresholds. Parameters: $\nu = 100$, $f_0 = 10$, and $\overline{f}_{-0} = 1$.

| Landscape | Level crossing | | | Error threshold |
|---|---|---|---|---|
| | $p_{(1/100)}$ | $p_{(1/1000)}$ | $p_{(1/10000)}$ | $p_{\mathrm{cr}}$ |
| Single-peak | 0.02198 | 0.02274 | 0.02282 | 0.02277 |
| Hyperbolic | 0.01450 | 0.01810 | 0.02036 | 0.02277 |
| Step-linear | 0.01067 | 0.01774 | 0.02330 | 0.02277 |

to the uniform distribution at $p = \tilde{p}$ is gradual without any abrupt change. The stationary concentrations of all error classes $\Gamma_k$ in the range $1, \ldots, \lfloor \nu/2 \rfloor$ pass through a maximum, whereas the others with higher class indices change monotonously until they reach the value $\bar{y}_k = \binom{\nu}{k}/2^\nu = \bar{y}_k = \binom{\nu}{\nu-k}/2^\nu$ at $\tilde{p}$. Fast decay of the concentration of the master class $\bar{y}_0(p)$ is observed in both cases but no abrupt change in the stationary mutant distribution occurs and the domain of the uniform distribution is restricted to $p = \tilde{p}$ (figure 15). Knowing now the behavior of the quasispecies on the single-peak and the linear landscape an interpretation of the observed plots for the step-linear landscape is straightforward: In the range of small Hamming distances from the master sequence the fitness landscape has the same shape as the linear landscape and for small mutation rates the quasispecies is dominated by sequences, which are near the master in sequence space, at higher mutation rates $p$ sequences that are further away from the master gain importance, and indeed we observe a similarity of the quasispecies with that on the linear landscape at small $p$-values whereas an error threshold and the uniform distribution beyond it are observed at higher mutation rates $p$. In the step-linear landscape the position of the step, $h$ can be varied as well. For the parameters $f_0 = 10$ and $f = 1$ we observe error thresholds in the range $0 \leq h \leq 35$, at higher $h$-values it becomes softer and eventually around $h = 45$ it has completely disappeared.[26] A useful indicator for the existence of an error threshold is the upper envelope of all individual curves $\bar{y}_k(p)$: The absence of a threshold leads to a monotonous decrease off the envelope (fig-
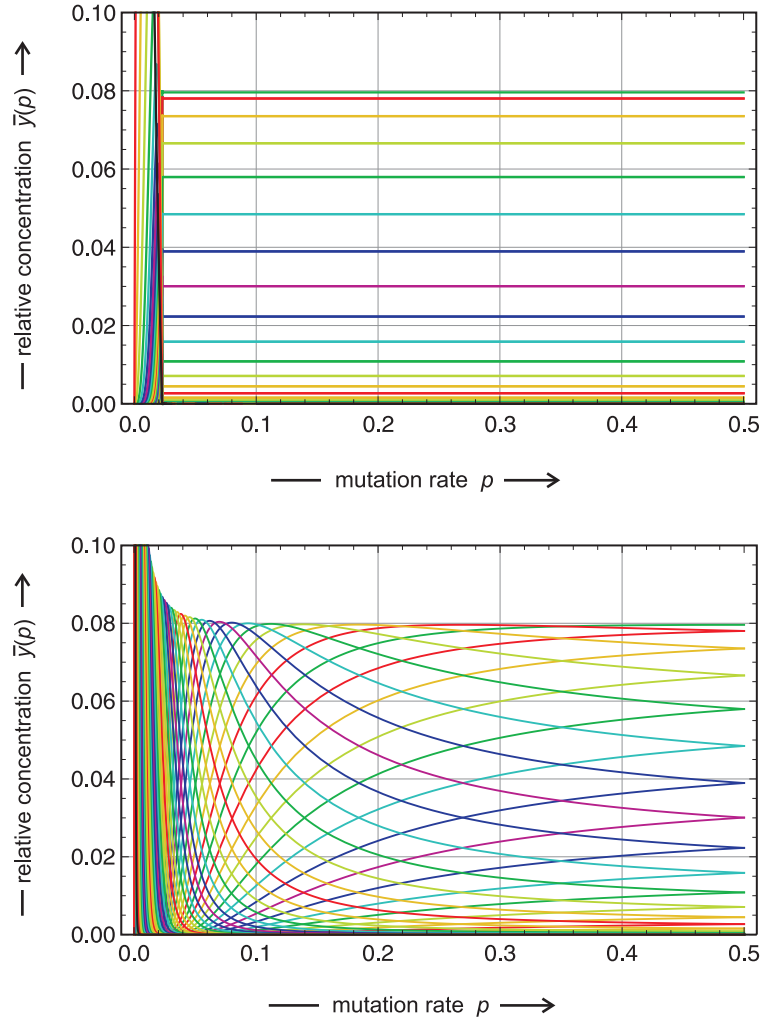
---

[26]Like in physics we distinguish hard and soft transitions. A hard transition is confined to a very narrow range of the order parameter – here the error rate $p$ – and becomes steeper and steeper as the system grows to infinity.

**Figure 14: Model landscapes not sustaining error thresholds.** The figures show stationary concentrations of mutant classes as functions of the error rate, $\bar{y}_k(p)$, for sequences of chain length $\nu = 100$ with $f_0 = 10$ and $\overline{f}_{-0} = 1$ on two different model landscapes: the multiplicative landscape ($f = 0.09472$; upper part: range of the error threshold; middle part: full range), and the additive or linear landscape ($f = 0.001$; lower part). The fitness value $f$ leading to a superiority $\sigma_0 = 10$ is negative for the linear landscape and accordingly a small positive value has been chosen instead.

ure 14 whereas an error threshold manifests itself in a pronounced minimum of the envelope just below $p_{\mathrm{cr}}$ (figure 13).

The comparison of the quasispecies development as a function of the mutation rate $p$ between the single-peak and the multiplicative landscape shown in figure 15 illustrates the interpretation that the uniform distribution being the exact solution at $p = \tilde{p}$ is extended to lower $p$-values down to the error threshold. The stationary concentrations of error classes, $\bar{y}_k(p)$ are almost perfect straight lines over the whole range for landscapes that sustain error thresholds, whereas on the additive and the multiplicative landscape the curves for $\bar{y}_k(p)$ and $\bar{y}_{\nu-k}(p)$, which have identical values in the uniform

**Figure 15: Comparison of landscapes with and without error thresholds.**
The figures show stationary concentrations of mutant classes as functions of the
error rate, $\bar{y}_k(p)$, for sequences of chain length $\nu = 100$ with $f_0 = 10$ and $\overline{f}_{-0} = 1$ on
two different model landscapes: the single-peak landscape ($f = 1$; upper part), and
the multiplicative landscape ($f = 0.09472$; lower part). The uniform distribution
($\Pi$) representing the exact solution at $p = \tilde{p} = 1/2$ extends over the whole range
from the threshold at $p = p_{\mathrm{cr}}$ to $p = \tilde{p}$ in case of the single-peak landscape.

distribution don't meet visually before the point $p = \tilde{p}$ on the multiplicative
landscape.

The fact that the behavior of quasispecies depends strongly on the nature
of the fitness landscape is not surprising. Fitness values after all play the
same role as rate parameters in chemical kinetics and the behavior of a sys-
tem can be changed completely by a different choice of rate parameters. The
most relevant but also most difficult question concerns the relation between
rate parameters and observed stationary distributions: Can we predict the
quasispecies from a knowledge of the fitness landscape? Or the even more
difficult inverse problem [69]: Does the observed behavior of the quasispecies
allow for conclusions about the distribution of fitness values? A few regu-

larities can be recognized already from the simple model landscapes: (i) fast decay of the master concentration, $\bar{y}_0(p)$ may occur without the appearance of a sharp transition, (ii) a sharp transition may occur on fitness landscapes with gradually changing fitness values provided the decay of $f(\mathbf{Y}_k)$ with $k$ is sufficiently steep, (iii) a sharp transition may occur without leading to the uniform distribution, and (iv) the appearance of the uniform distribution at $p_{\mathrm{cr}}$-values lower than $\tilde{p}$ requires a flat part of the fitness landscapes in the sense that fitness values of neighboring classes are (almost) the same.

# 7 "Realistic" rugged landscapes (RRL)

The majority of data on the relation between sequences and molecular properties comes from structural biology of biopolymers, in particular RNA and protein. RNA secondary structures provide a simple and mathematically accessible example of a *realistic* mapping of biopolymer sequences onto structures [70]. The RNA model is commonly restricted to the assignment of a single structure to every sequence but the explicit consideration of suboptimal conformations is possible as well [71]. Two features are characteristic for landscapes derived from RNA molecules: (i) *ruggedness* – pairs of sequences situated nearby in sequence space, i.e., having Hamming distance $d_{\mathrm{H}} = 1$, may give rise to very similar or entirely different structures and properties – and (ii) *neutrality* – two or more sequences may have identical structures and properties.[27] Both properties are easily illustrated by means of RNA secondary structures, which are defined in terms of Watson-Crick and $\mathbf{G} - \mathbf{U}$ base pairs: Exchange of one nucleobase in a base pair, e.g., $\mathbf{C} \to \mathbf{G}$ in $\mathbf{G} \equiv \mathbf{C}$, may open the base pair, destroy a stack and eventually lead to an entirely different structure with different properties, or leave the structure unchanged, e.g., $\mathbf{A} \to \mathbf{G}$ in $\mathbf{A} = \mathbf{U}$. Neutrality is equally well demonstrated: Exchanging both bases in a base pair may leave structure and (most) properties unchanged, $\mathbf{G} \equiv \mathbf{C} \to \mathbf{C} \equiv \mathbf{G}$ may serve as an example. Evolutionary dynamics is clearly influenced by the shape of fitness landscapes and the interplay of the two characteristic features was found to be essential for the success of evolutionary searches [46,47,73]. Methods were developed recently that allow for efficient construction of fitness landscapes for catalytically active RNA molecules [74]. Ruggedness and neutrality are not restricted to RNA-molecules, similar results providing direct evidence were found with

---

[27] *Identical* in the context of neutrality does not mean identical in strict mathematical sense but indistinguishable for the experimental setup or for natural selection [33,72].

proteins [75] and simple organisms like viruses [25]. Apart from a few exceptions experimental comprehensive information on fitness landscapes or conformational free energy surfaces is still rare but the amount of reliable data is rapidly growing. It seems to be appropriate therefore to conceive and construct model landscapes that account for the known features and to study evolutionary dynamics on them.

Rugged fitness landscapes, which are more elaborate than the simple ones discussed in section 6, have been proposed. The most popular example is the $Nk$-model conceived by Stuart Kauffman [30, 31, 76] that is based on individual loci on a genome and interactions between them: $N$ is the number of loci and $k$ is the number of interactions. A random element, which is drawn from a predefined probability distribution – commonly the normal distribution – and which defines the interaction network, is added to the otherwise deterministic model: $N$ and $k$ are fixed and not subjected to variation. Here a different approach is proposed that starts out from the nucleotide sequence of a genome rather than from genes and alleles, and consequently it is based on the notion of sequence space. Ruggedness (this section 7) and neutrality (see section 8) are introduced by means of tunable parameters, $d$ and $\lambda$, and pseudorandom numbers are used to introduce random scatter, which reflects the current ignorance with respect to detailed fitness values and which is thought to be replaced by real data when they become available in the near future.

A new type of landscapes, the *realistic rugged landscape* (RRL), is introduced and analyzed here. Ruggedness is modeled by assigning fitness differences at random within a predefined band of fitness values with adjustable width $d$. The highest fitness value is attributed to the master sequence $\mathbf{X}_m \doteq \mathbf{X}_0$, $f_m = f_0$, and the fitness values of all other sequences are obtained by means of the equation

$$
f(\mathbf{X}_j) \; = \; f_j \; = \; \begin{cases} f_0 & \text{if } j = 0 \;, \\ f + 2d(f_0 - f)\left(\eta_j^{(s)} - 0.5\right) & \text{if } j = 1, \ldots, \kappa^l \;, \end{cases} \tag{40}
$$

where $\eta_j^{(s)}$ is the $j$-th output random number from a pseudorandom number generator with a uniform distribution of numbers in the range $0 \leq \eta_j^{(s)} \leq 1$. The random number generator is assumed to have been started with the seed $s$,[28] which will be used to characterize a particular distribution of fitness values (figure 16). The parameter $d$ determines the amount of scatter around

---

[28]The seed $s$ indeed determines all details of the landscape, which is completely defined by $s$ and the particular type of the pseudorandom number generator as well as by $f_0$, $f$, and $d$.

the mean value $\bar{f}_{-0} = f$, which is independent of $d$: $d = 0$ yields the single peak landscape, and $d = 1$ leads to fully developed or maximal scatter where individual fitness values $f_j$ can reach the value $f_0$. A given landscape can be characterized by

$$\mathcal{L} = \mathcal{L}(\lambda, d, s; \nu, f_0, f) , \tag{41}$$

where $\lambda$ is the degree of neutrality (see section 8; here we have $\lambda = 0$). The parameters $\nu$, $f_0$ and $f$ have the same meaning as for the single peak landscape (28).

Two properties of realistic rugged landscapes fulfilled by fitness values relative to the mean except the master, $\varphi_j = f_j - f \,\forall\, j = 0, \ldots, \kappa^\nu - 1$, are important: (i) the ratio of two relative fitness values of sequences within the mutant cloud is independent of the scatter $d$ and (ii) the ratio of the relative fitness values of a sequence from the cloud and the master sequence is proportional to the scatter $d$:

$$\frac{\varphi_j}{\varphi_k} = \frac{\eta_j^{(s)} - 0.5}{\eta_k^{(s)} - 0.5} ; \; j, k = 1, \ldots, \kappa^\nu - 1 \; \text{ and} \tag{42a}$$

$$\frac{\varphi_j}{\varphi_0} = 2\,d\left(\eta_j^{(s)} - 0.5\right); \; j = 1, \ldots, \kappa^\nu - 1 . \tag{42b}$$

The second equation immediately shows that $\sum_{j=1}^{\kappa^\nu - 1} \varphi_j = 0$.

## 7.1 Single master quasispecies

We are now in a position to explore whether or not the results derived from simple model landscapes are representative for mutation-selection dynamics in real populations. At first the influence of random scatter on quasispecies and error thresholds will be studied. The chain length for which diagonalization of the value matrix W can be routinely performed lies at rather small values around $\nu = 10$ giving rise to a matrix size of $1\,000 \times 1\,000$. Accordingly, it has to be confirmed first whether or not such a short chain length is sufficient to yield representative results. In figure 17 the stationary concentrations of mutant classes, $\bar{y}_k$ $(k = 0, 1, \ldots, 10)$ are shown for different band widths $d$ of random scatter: The purely deterministic case $d = 0$ representing the single-peak landscape, $d = 0.5$, and $d = 0.95$, the maximal scatter that sustains a single quasispecies over the entire range, $0 \leq p < p_{\mathrm{cr}}$.[29] Despite the short chain length of $\nu = 10$ the plots reflect the threshold phenomena rather well, the width of the transition to the uniform distribution is hardly
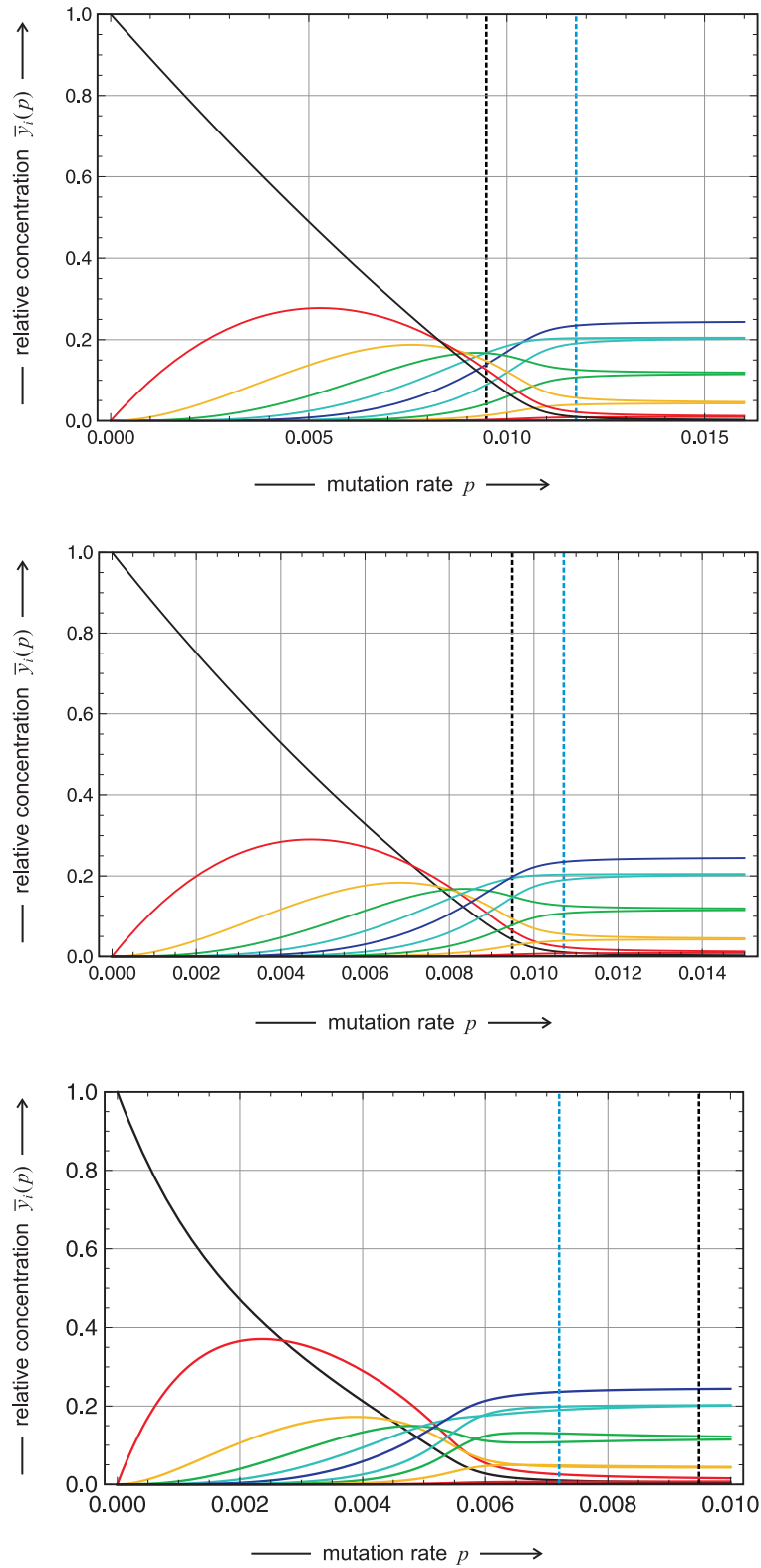
---

[29]As shown below in detail (figure 22) individual quasispecies may be replaced by others at certain critical $p$-values, $p_{\mathrm{tr}}$. For a given scatter $s$ the number of such transitions becomes larger with increasing values of $d$.
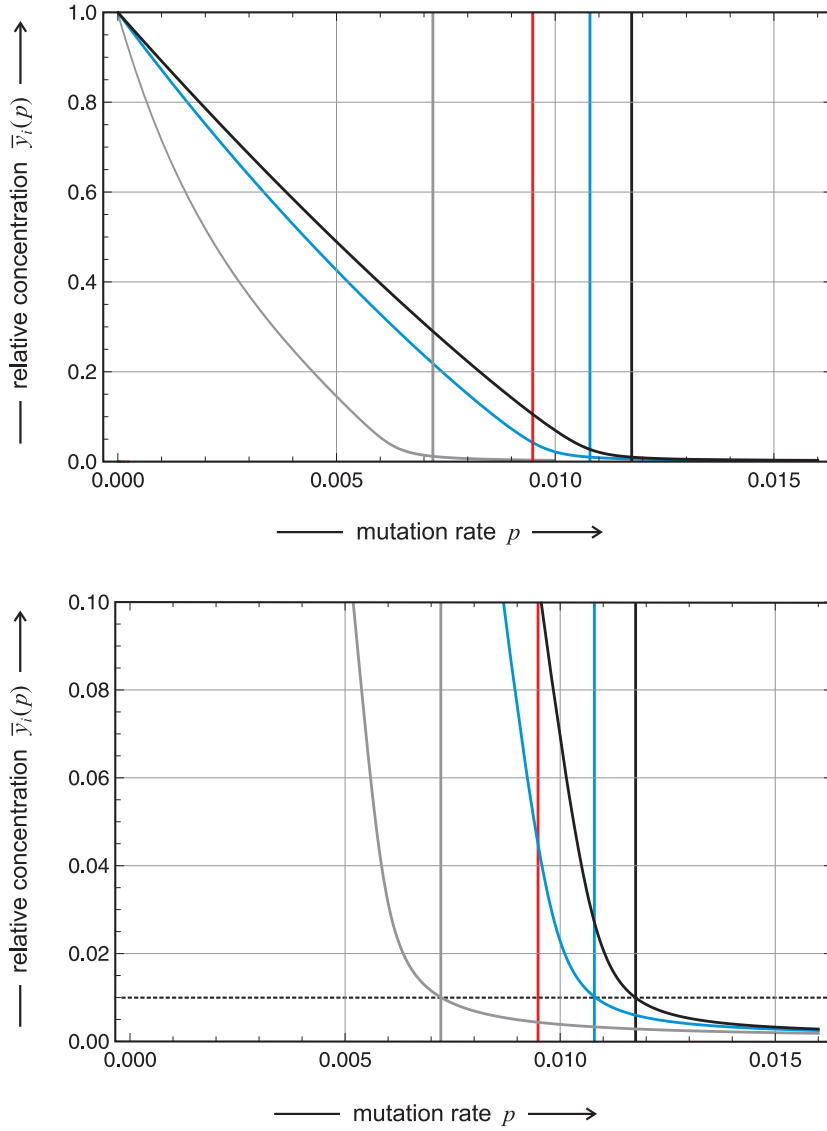
**Figure 16:** *Realistic* **rugged fitness landscapes.** The landscapes for binary sequences with chain length $\nu = 10$ are constructed according to equation (40). In the upper plot the band width of random scatter was chosen to be $d = 0.5$ and a seed of $s = 919$ was used for the random number generator. For the lower plot showing maximal random scatter $d = 1$ and $s = 637$ was applied. Careful inspection allows for the detection of individual differences. The broken blue lines separate different mutant classes.

changing, and values for level crossing (section 6 and table 1) are shifted towards smaller $p_{(1/M)}$-values with increasing $d$. Answering the initial question, computer studies with $\nu = 10$ are suitable for investigations on quasispecies behavior.

For $d > 0$ the fitness values for individual sequences within one class are no longer the same and hence the curves $\bar{x}_j(p)$ differ from each other and form a band for each class that increases in width with the amplitude $d$ of the ran-

50

**Figure 17: Error thresholds on a *realistic* model landscape with different random scatter $d$.** Shown are the stationary concentrations of classes $\bar{y}_j(p)$ on the *realistic* landscape with $s = 023$ for $d = 0$ (upper plot), $d = 0.5$ (middle plot), and $d = 0.95$ (lower plot). The error threshold calculated by zero mutational backflow lies at $p_{\mathrm{cr}} = 0.009486$ (black dotted line), the values for level crossing decrease with the width of random scatter $d$ (blue dotted lines). Other parameters: $\nu = 10$, $f_0 = 1.1$, and $f = 1.0$.

**Figure 18: Error threshold and decay of the master sequence $\mathbf{X}_0$.** Shown are the stationary concentrations of the master sequence $\bar{x}_0(p)$ and the level cross-ing values $p_{(1/100)}$ (vertical lines) on a landscape with $s = 023$ for $d = 0$ (black), $d = 0.5$ (blue), and $d = 0.950$ (grey). The error threshold lies at $p_{\mathrm{cr}} = 0.094857$ (red). The lower plot enlarges the upper plot and shows the level $\bar{x}_0 = 0.01$ (dotted horizontal line, black). Other parameters: $\nu = 10$, $f_0 = 1.1$ and $f = 1.0$.

dom component (figure 19). The separation of the bands formed by curves belonging to different error classes is always recognizable at sufficiently small mutation rates $p$ but the bands overlap and merge at higher $p$-values. As expected the zone where the bands begin to mix moves in the direction $p = 0$ with increasing scatter $d$. Interestingly, the error threshold phenomenon is fully retained thereby, only the level-crossing value $p_{(1/100)}$ is shifted towards lower error rates (figures 18, 19, and 20). Indeed, the approaches towards the uniform distribution on the landscape without a random component ($d = 0$) and on the landscape with $d = 0.5$ are very similar apart from the rela-
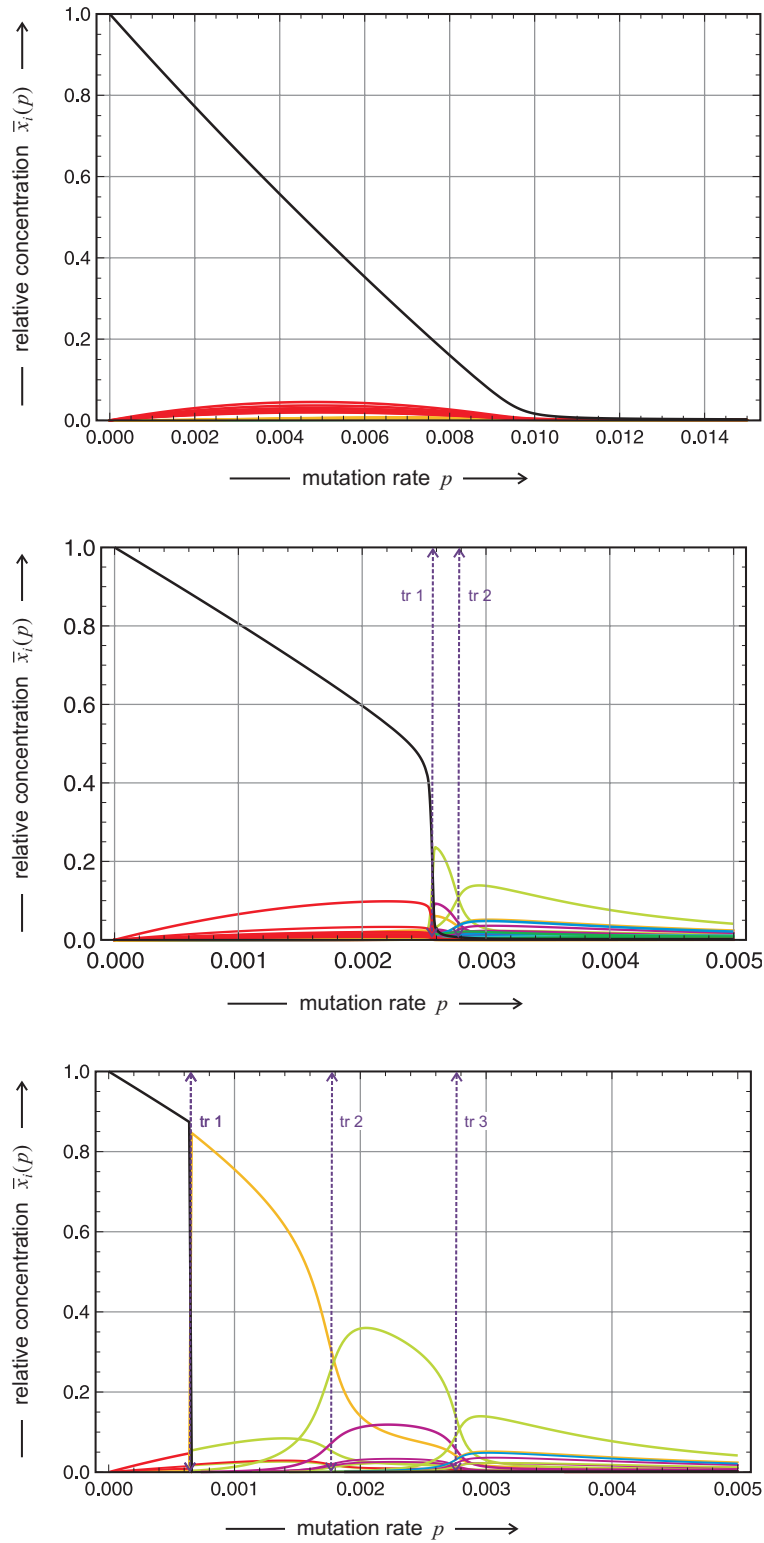
**Figure 19: Quasispecies on a *realistic* model landscape with different random scatter $d$.** Shown are the stationary concentrations $\bar{x}_j(p)$ on a landscape with $s = 491$ for $d = 0$ (upper plot), $d = 0.5$ (middle plot), and $d = 0.9375$ (lower plot) for the classes $\Gamma_0$, $\Gamma_1$, and $\Gamma_2$. In the topmost plot the curves for all sequences in $\Gamma_1$ (single point mutations, $d_H(\mathbf{X}_0, \mathbf{X}_{(1)}) = 1$) coincide, and so do the curves in $\Gamma_2$ (double point mutations, $d_H(\mathbf{X}_0, \mathbf{X}_{(2)}) = 2$) since zero scatter, $d = 0$, has been chosen. The error threshold calculated by zero mutational backflow lies at $p_{cr} = 0.066967$. Other parameters: $\nu = 10$, $f_0 = 2.0$, and $f = 1.0$.

53

**Figure 20: Level-crossing values for the master sequence of different model landscape with different random scatter** $d$. Shown are the level crossing values for $M = 100$ as functions of the random scatter $p_{(1/100)}(d)$. The error threshold calculated by zero mutational backflow lies at $p_{\mathrm{cr}} = 0.0117481$. Color code for different seeds $s$: $023 =$ orange, $229 =$ red, $367 =$ green, $491 =$ black, $577 =$ chartreuse, $637 =$ blue, $673 =$ yellow, $877 =$ magenta, $887 =$ turquoise, $919 =$ blue violet, and $953 =$ hot pink. Other parameters: $\nu = 10$, $f_0 = 1.1$, and $f = 1.0$.

**Figure 21: A *realistic* model landscape with a transition between quasispecies.** Shown are the stationary concentrations $\bar{x}_j(p)$ on a landscape with $s = 023$ for $d = 0.5$ (upper plot), $d = 0.999$ (middle plot), and fully developed scatter $d = 1.0$ (lower plot). Other parameters: $\nu = 10$, $f_0 = 1.1$, and $f = 1.0$.

**Figure 22: A *realistic* model landscape with multiple transitions between quasispecies.** Shown are the stationary concentrations $\bar{x}_j(p)$ on a landscape with $s = 637$ for $d = 0.5$ (upper plot), $d = 0.995$ (middle plot), and fully developed scatter $d = 1.0$ (lower plot). Other parameters: $\nu = 10$, $f_0 = 1.1$, and $f = 1.0$.

tively small shift towards lower $p$-values, whereas the shift for $d = 0.95$ is substantially larger and the solution curve $\bar{x}_0(p)$ is curved upwards more strongly. Closer inspection of the shift of the level-crossing value shows non-monotonous behavior for some landscapes: The level crossing value is shifted towards larger $p$-values at first, passes a maximum value and then follows the general shift towards lower values of $p$ with increasing scatter $d$ (figure 20).

In addition, transitions between quasispecies may be observed at critical mutation rates $p = p_{\mathrm{tr}}$: One quasispecies, $\Upsilon_0$, which is the stationary solution of the mutation-selection equation (18) in the range $0 \leq p < p_{\mathrm{tr}}$, is replaced by another quasispecies, $\Upsilon_k$, which represents the stationary solution above the critical value up to the error threshold $p_{\mathrm{tr}} < p < p_{\mathrm{cr}}$, or up to a second transition, $(p_{\mathrm{tr}})_1 < p < (p_{\mathrm{tr}})_2$. More than two transitions are also possible, an example is shown in figure 22 (lower plot). The mechanism by which quasispecies replace each other is easily interpreted [59]:[30] The stationary mutational backflow from the sequences $\mathbf{X}_i$ $(i = 1, \ldots, n)$ to the master sequence $\mathbf{X}_0$ is determined by the sum of the product terms $\psi_0 = \sum_{i=1}^n W_{0i} = Q_{0i} f_i$ and likewise for a potential master sequence $\mathbf{X}_k$, $\psi_k = \sum_{i=0, i \neq k}^n W_{ki} = Q_{ki} f_i$. The necessary – but not sufficient – condition for the existence of a transition is $\Delta\psi = \psi_0 - \psi_k < 0$. Since the fitness value $f_0$ is the largest by definition we have $f_0 > f_i$ $(i = 1, \ldots, n)$ and at sufficiently small mutation rates $p$ the differences in the values, $\Delta\omega = \omega_0 - \omega_k = W_{00} - W_{kk} = Q_{00} f_0 - Q_{kk} f_k > 0$, will always outweigh the difference in the backflow, $\Delta\omega > |\Delta\psi|$. With increasing values of $p$, however, the replication accuracy and $\Delta\omega$ will decrease because of the term $Q_{00} = Q_{kk} \approx (1 - p)^\nu$ in the uniform error approximation. At the same time $\Delta\psi$ will increase in absolute value and provided $\Delta\psi < 0$ there might exist a mutation rate $p = p_{\mathrm{tr}}$ smaller than the threshold value $p_{\mathrm{tr}} < p_{\mathrm{cr}}$ at which the condition $\Delta\omega + \Delta\psi = 0$ is fulfilled and consequently, the quasispecies $\Upsilon_k$ is the stable stationary solution of equation (18) at $p > p_{\mathrm{tr}}$. The influence of a distribution of fitness values instead of the single value of the single-peak landscapes can be predicted straightforwardly: Since $f_0$ is independent of the fitness scatter $d$ the difference $f_0 - f_k$ will decrease with increasing scatter $d$. Accordingly, the condition for a transition between quasispecies can be fulfilled at lower $p$-values and we expect to find one or more transitions preceding the error threshold $p_{\mathrm{cr}}$. No transition can occur on the single peak landscape but as $d$ increases and the difference $\Delta\omega$ becomes smaller and it becomes more likely that the difference in backflow becomes sufficiently strong for a replacement of $\Upsilon_0$ by $\Upsilon_k$ below $p_{\mathrm{cr}}$. Fig-

---

[30]Thirteen years after this publication the phenomenon has been observed in quasispecies of digital organisms [77] and was called *survival of the flattest*.

ure 21 presents a typical example: No quasispecies transition is found up to a random scatter of $d = 0.95$. Then, a soft transition becomes observable at $d = 0.975$ and eventually dominates the plot of the quasispecies against the mutation rate $p$ at random scatter close to the maximum ($d = 0.995$ and $d = 1.0$). An example with multiple transitions increasing in number with increasing $d$ is shown in figure 22.

An explicit computation of the transition point $p = p_{tr}$ has been performed some time ago [59]. A simple model is used for the calculation of the critical value that is based on a zero mutational flow assumption between the two quasispecies. The value matrix W corresponding to all $2^\nu$ sequences of chain length $\nu$ is partitioned into two diagonal blocks and the rest of the matrix: (i) Block $\bar{\Upsilon}_0$ contains sequence $\mathbf{X}_0$ with the highest fitness value $f_0$, which is the master sequence in the range $0 \leq p < p_{tr}$ and all its one-error mutants $\mathbf{X}_{(1)}^{(0)} = \{\mathbf{X}_j \in \Gamma_1^{(0)}\}$ with a fitness value $f_1^{(0)}$, (ii) block $\bar{\Upsilon}_m$ contains sequence $\mathbf{X}_m$ with the fitness value $f_m$, which is the master sequence in the range $p_{tr} \leq p < p_{cr}$, and all its one-error mutants $\mathbf{X}_{(1)}^{(m)} = \{\mathbf{X}_j \in \Gamma_1^{(m)}\}$ with a fitness value $f_1^{(m)}$, and (iii) the rest of the matrix $W$ is neglected completely as all entries are set equal zero:

$$
W = q^\nu \begin{pmatrix}
f_0 & f_1^{(0)}\varepsilon & \cdots & f_1^{(0)}\varepsilon & \cdots & 0 & 0 & \cdots & 0 \\
f_0\varepsilon & f_1^{(0)} & \cdots & f_1^{(0)}\varepsilon^2 & \cdots & 0 & 0 & \cdots & 0 \\
\vdots & \vdots & \ddots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
f_0\varepsilon & f_1^{(0)}\varepsilon^2 & \cdots & f_1^{(0)} & \cdots & 0 & 0 & \cdots & 0 \\
0 & 0 & \cdots & 0 & \cdots & f_m & f_1^{(m)}\varepsilon & \cdots & f_1^{(m)}\varepsilon \\
0 & 0 & \cdots & 0 & \cdots & f_m\varepsilon & f_1^{(m)} & \cdots & f_1^{(m)}\varepsilon^2 \\
\vdots & \vdots & \ddots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
0 & 0 & \cdots & 0 & \cdots & f_m\varepsilon & f_1^{(m)}\varepsilon^2 & \cdots & f_1^{(m)}
\end{pmatrix}.
$$

Each block is now represented by a $2 \times 2$ matrix

$$
W_0 = q^\nu \begin{pmatrix} f_0 & \nu f_1^{(0)}\varepsilon \\ f_0\varepsilon & f_1^{(0)}\left(1 + (\nu-1)\varepsilon^2\right) \end{pmatrix} \quad \text{and}
$$

$$
W_m = q^\nu \begin{pmatrix} f_m & \nu f_1^{(m)}\varepsilon \\ f_m\varepsilon & f_1^{(m)}\left(1 + (\nu-1)\varepsilon^2\right) \end{pmatrix}.
$$

Calculation of the two largest eigenvalues $\lambda_0$ and $\lambda_m$ yields the condition for the occurrence of the transitions: $\lambda_0 = \lambda_m$. The result is

$$
p_{tr} = 1 - \sqrt{1 - \frac{\bar{\vartheta}_{tr}}{\nu - 1}}, \tag{43}
$$

**Table 2: Computed and numerically observed quasispecies transitions.**
In the table we compare the numerically observed values of $p$ at transitions between quasispecies, $p_{tr}$, with the values calculated from equation (43), $(p_{tr})_{eval}$, and the error thresholds, $p_{cr}$. The table is adopted from [59].

| Chain length | Qsp. $\bar{\Upsilon}_0$ | | Qsp. $\bar{\Upsilon}_m$ | | Critical mutation rates | | |
|---|---|---|---|---|---|---|---|
| $\nu$ | $f_0$ | $f_1^{(0)}$ | $f_m$ | $f_1^{(m)}$ | $p_{tr}$ | $(p_{tr})_{eval}$ | $p_{cr}$ |
| 20 | 10 | 1 | 9.9 | 2 | 0.0520 | 0.0567 | 0.1130 |
| 50 | 10 | 1 | 9.9 | 2 | 0.0362 | 0.0366 | 0.0454 |
| 50 | 10 | 1 | 9.9 | 5 | 0.0148 | 0.0147 | 0.0470 |
| 50 | 10 | 1 | 9.0 | 5 | 0.0445 | 0.0456 | 0.0453 |

with $\vartheta_{tr}$ being the result of the equation

$$\vartheta_{tr} = \frac{1}{2}\left(\alpha + \beta - \gamma + \sqrt{(\alpha + \beta - \gamma)^2 - 4\alpha\beta}\right) \quad \text{with} \qquad (43')$$

$$\alpha = \nu - \frac{f_0 - f_m}{f_1^{(m)} - f_1^{(0)}} \ ,$$

$$\beta = \nu - \frac{f_0 f_m (f_1^{(m)} - f_1^{(0)})}{f_1^{(0)} f_1^{(m)} (f_0 - f_m)} \ , \quad \text{and}$$

$$\gamma = \frac{(f_0 f_1^{(0)} - f_m f_1^{(m)})^2}{(\nu - 1) f_1^{(0)} f_1^{(m)} (f_0 - f_m)(f_1^{(m)} - f_1^{(0)})} \ .$$
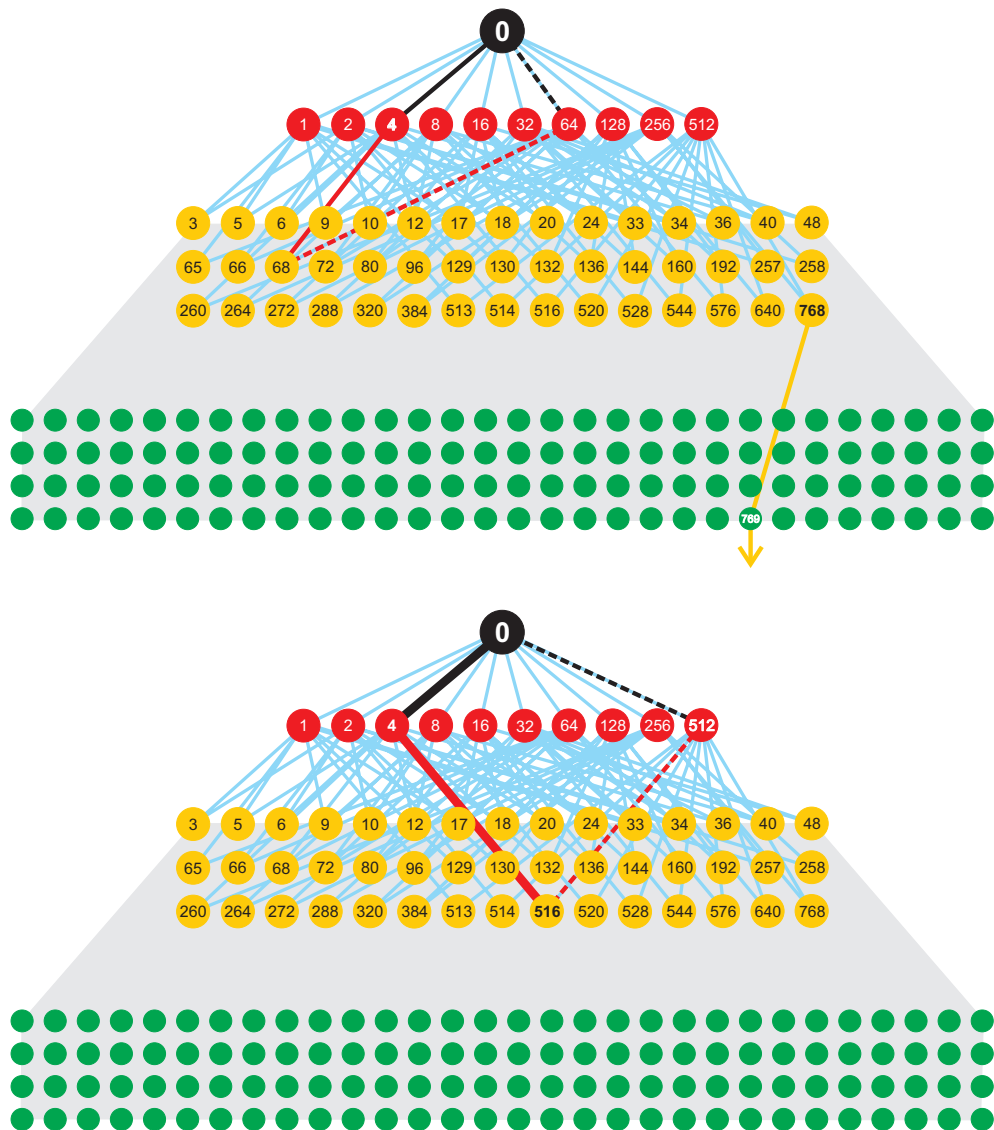
Although the complexity of these equations is prohibitive for further manipulations, the accuracy of the zero mutational flow approximations is relevant for the next subsection 7.2 were we shall apply a similar approximation. The corresponding table 2 is reproduced from [59]. The agreement is very good indeed but is cases where $p_{tr}$ and $p_{cr}$ are very close it can nevertheless happen that the calculated value lies above the error threshold.

## 7.2 Clusters of coupled sequences

A certain fraction of landscapes gives rise to characteristic quasispecies distributions as a function of the mutation rate $p$ that is substantially different from the one shown in figure 22 and discussed above. No transitions are observed, not even at fully developed scatter $d = 1$ (figure 23). Another feature concerns the classes to which the most frequent sequences belong. On the landscape defined by $s = 919$ these sequences are the master sequence ($\mathbf{X}_0$; black curve), one one-error mutant ($\mathbf{X}_4$; red curve), and one two-error

**Figure 23: Error thresholds on a *realistic* model landscape with different random scatter $d$ and transitions between quasispecies.** The landscape characteristic is $s = 919$. Shown are the stationary concentrations $\bar{x}_j(p)$ for $d = 0.5$ (upper plot), $d = 0.995$ (middle plot), and fully developed scatter $d = 1.0$ (lower plot). Other parameters: $\nu = 10$, $f_0 = 1.1$, and $f = 1.0$.

**Figure 24: Mutation flow in quasispecies.** The sketch shows two typical situations in the distribution of fitness values in sequence space. In the upper diagram $(s = 637)$ the fittest two-error mutant, $\mathbf{X}_{768}$, has its fittest nearest neighbor, $\mathbf{X}_{769}$, in the three-error class $\Gamma_3$, and the fittest sequence in the one-error neighborhood of $\mathbf{X}_4$ (being the fittest sequence in the one-error class), $\mathbf{X}_{68}$, is different from $\mathbf{X}_{768}$, the mutational flow is not sufficiently strong to couple $\mathbf{X}_0$, $\mathbf{X}_4$, and $\mathbf{X}_{68}$, and transitions between different quasispecies are observed (figure 22). The lower diagram $(s = 919)$ shows the typical fitness distribution for a strong quasispecies: The fittest two-error mutant, $\mathbf{X}_{516}$, has its fittest nearest neighbor, $\mathbf{X}_4$, in the one-error class $\Gamma_1$ and it coincides with the fittest one-error mutant. Accordingly, the three sequences ($\mathbf{X}_0$, $\mathbf{X}_4$, and $\mathbf{X}_{516}$) are strongly coupled by mutational flow and a strong quasispecies is observed (figure 23).

mutant ($\mathbf{X}_{516}$; yellow curve).[31] The three sequences are situated close-by in

---

[31]Naïvely we would expect a band of one-error sequences at higher concentration than the two-error sequence.

sequence space – Hamming distances $d_H(\mathbf{X}_1, \mathbf{X}_4) = d_H(\mathbf{X}_4, \mathbf{X}_{516}) = 1$ and $d_H(\mathbf{X}_1, \mathbf{X}_{516}) = 2$)– form a cluster, which is dynamically coupled by means of strong mutational flow (figure 24). Apparently, such a quasispecies is not likely to be replaced in a transition by another one that is centered around a single master sequence and accordingly, we call such clusters *strong quasispecies*. The problem that ought to be solved now is the prediction of the occurrence of strong quasispecies from know fitness values.

First, a heuristic is mentioned that serves as an (almost perfect) diagnostic tool for detecting whether or not a given fitness landscape gives rise to a strong quasispecies: (i) For every mutant class we identify the sequence with the highest fitness value, $f_0$, $(f_{(1)})_{\max} = f(\mathbf{X}_{m(1)})$, $(f_{(2)})_{\max} = f(\mathbf{X}_{m(2)})$, $\ldots$ , and call them *class-fittest* sequences. Next we determine the fittest sequences in the one-error neighborhood of the class-fittest sequences. Clearly, for the class $k$-fittest sequence $\mathbf{X}_{m(k)}$ this sequence lies either in class $k-1$ or in class $k+1$.[32] Simple combinatorics is favoring classes closer to the middle of sequence space. Any sequence in the two-error class, for example, has two nearest neighbors in the one-error class but $n-2$ nearest neighbors in the three-error class. To be a candidate for a strong quasispecies requires that – against probabilities – the fittest sequence in the one-error neighborhood of $\mathbf{X}_{m(2)}$ lies in the one-error class: $(f_{(\mathbf{X}_{m(2)})_{m(1)}})_{\max}$ with $(\mathbf{X}_{m(2)})_{m(1)} \in \Gamma_1$ and preferentially, this is the fittest one-error sequence, $(\mathbf{X}_{m(2)})_{m(1)} \equiv \mathbf{X}_{m(1)}$. Since all mutation rates between nearest neighbor sequences in neighboring classes are the same – $(1-p)^{n-1}p$ within the uniform error model – the strength of mutational flow is dependent only on the fitness values, and the way in which the three sequences were determined guarantees optimality of the flow: If such a three-membered cluster was found it is the one with the highest internal mutational flow for a given landscape. Figure 24 (lower picture, $s = 919$) shows an example where such three sequences form a strongly coupled cluster. There is always a fourth sequence – here $\mathbf{X}_{512}$ – belonging to the cluster but it may play no major role because of low fitness. The heuristic presented was applied to 21 fitness landscapes with different random scatter and three strong quasispecies ($s =$401, 577, and 919) were observed. How many would be expected by combinatorial arguments? The probability for a sequence in $\Gamma_2$ to have a neighbor in $\Gamma_1$ is $2/10 = 0.2$ and, since the sequence $\mathbf{X}_{m(1)}$ is fittest in $\Gamma_1$ and hence also in the one-error neighborhood of $\mathbf{X}_{m(2)}$, this is also the probability for finding a strong quasispecies. The sample that has been investigated in this study comprised

---

[32] For class $k = 1$ we omit the master sequence $\mathbf{X}_0$, which trivially is the fittest sequence in the one-error neighborhood, and search only in class $k = 2$.

21 landscapes an hence we expect to encounter $21/5 = 4.2$ cases, which is – with respect to the small sample size – in agreement with the three cases that we found.

The suggestion put forward in the heuristic mentioned above – a cluster of sequences coupled by mutational flow that is stronger within the group than to the rest of sequence space because of frequent mutations and high fitness values – will now be analyzed and tested through the application of perturbation theory. Instead of a single master sequence we consider a *master cluster* of sequences and then proceed in full analogy to subsection 5.6 by applying zero mutational backflow from the rest of sequence space to the cluster. In order to be able to deal with a cluster of sequences we rearrange the value matrix W:

$$
W \;=\;
\begin{pmatrix}
W_{11} & W_{12} & \cdots & W_{1k} & W_{1,k+1} & \cdots & W_{1n} \\
W_{21} & W_{22} & \cdots & W_{2k} & W_{2,k+1} & \cdots & W_{2n} \\
\vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
W_{k1} & W_{k2} & \cdots & W_{kk} & W_{k,k+1} & \cdots & W_{kn} \\
& & & & & & \\
W_{k+1,1} & W_{k+1,2} & \cdots & W_{k+1,k} & W_{k+1,k+1} & \cdots & W_{k+1,n} \\
\vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
W_{n1} & W_{n2} & \cdots & W_{nk} & W_{n,k+1} & \cdots & W_{nn}
\end{pmatrix} . \tag{44}
$$

The upper left square part of the matrix W will be denoted by $w_m$. It represents the core of the quasispecies in the sense of a mutationally coupled master cluster, $\mathcal{C}_m = \{\mathbf{X}_{m1}, \ldots, \mathbf{X}_{mk}\}$, and after neglect of mutational backflow from sequences outside the core we are left with the eigenvalue problem

$$
w_m \, \boldsymbol{\zeta}_{mj} \;=\; \lambda_{mj} \, \boldsymbol{\zeta}_{mj}; \; j = 0, \ldots, k-1 \; . \tag{45}
$$

In the uniform error rate model the elements of the mutation matrix Q are of the form

$$
Q_{mi,mj} \;=\; (1-p)^{n-d_{mi,mj}} \, p^{d_{mi,mj}} \;=\; (1-p)^{n-k} \, q_{mi,mj} \;\; \text{with}
$$

$$
q_{mi,mj} \;=\; (1-p)^{k-d_{mi,mj}} \, p^{d_{mi,mj}}
$$

Apart from the reduced dimension the eigenvalue problem (45) is in complete analogy to the eigenvalue problem in subsection 5.3. The common factor $(1-p)^{n-k}$ leaves the eigenvectors unchanged and is a multiplier for the eigenvalues: $\lambda_{mj} \Rightarrow (1-p)^{n-k} \lambda_{mj} \, \forall \, j = 0, \ldots, k-1$. Only the largest eigenvalue $\lambda_{m0}$ and the corresponding eigenvector $\boldsymbol{\zeta}_{m0}$ – with the components $\zeta_i^{(m0)}$ and $\sum_{i=0}^{k} \zeta_i^{(m0)} = 1$ – are important for the discussion of the quasispecies. By the same tokens as in subsection 5.6, equation (32a), we obtain
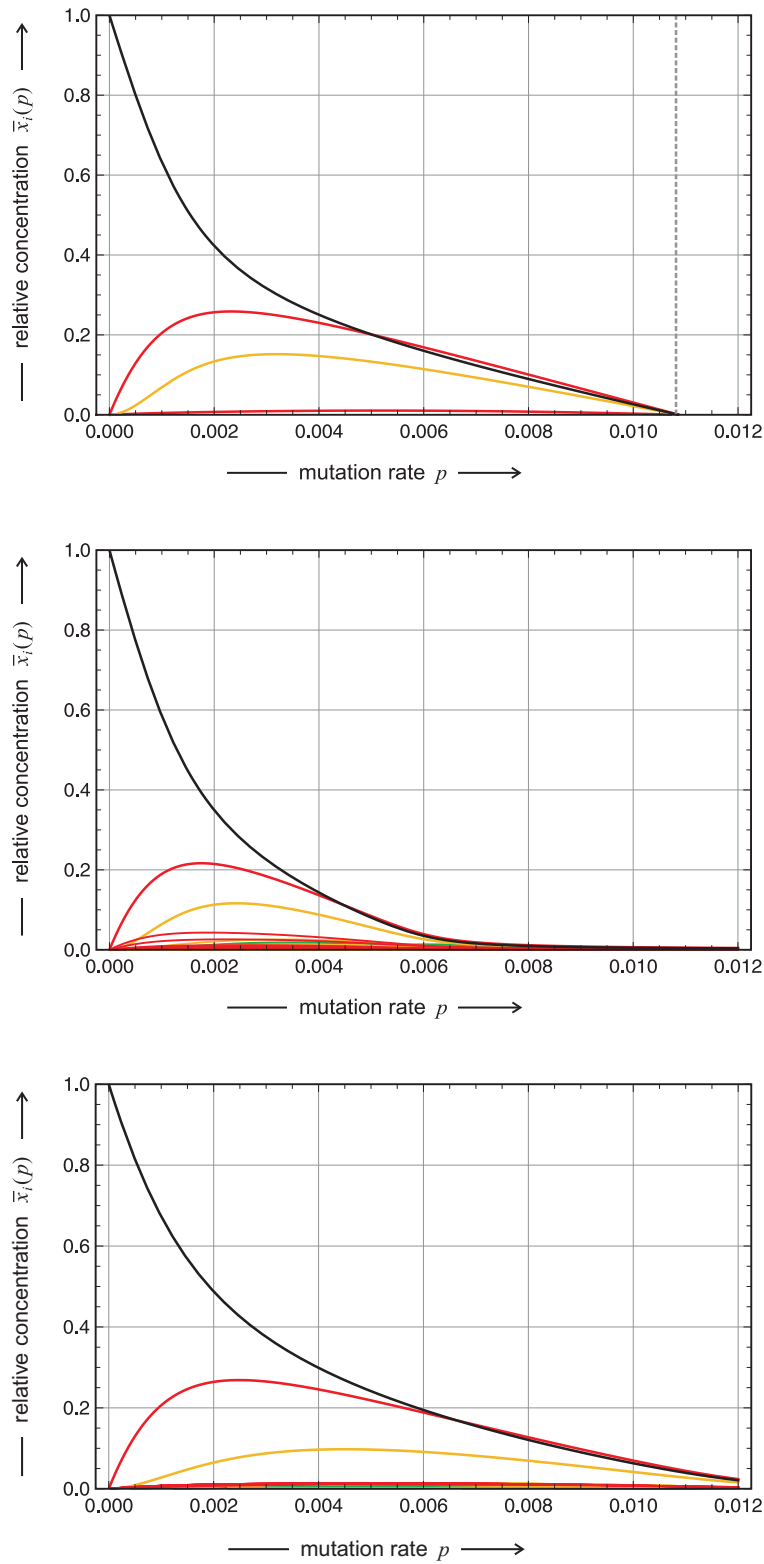
the stationary solution

$$\bar{c}_m^{(0)} = \frac{\lambda_{m0}(1-p)^{n-k} - \overline{f}_{-m}}{\overline{f}_m - \overline{f}_{-m}} \quad \text{with}$$

$$\bar{x}_{mj}^{(0)} = \zeta_j^{(m0)} \bar{c}_m^{(0)} ; \quad j = 1, \dots, k, \quad \text{and} \qquad (46)$$

$$\overline{f}_m = \sum_{i=1}^{k} \zeta_i^{(m0)} f_i \quad \text{and} \quad \overline{f}_{-m} = \sum_{i=k+1}^{n} \bar{x}_i f_i \left/ \sum_{i=k+1}^{n} \bar{x}_i \right. .$$

The calculations of the concentrations of the sequences not belonging to the master core is straightforward but more involved than in the case of a single master sequence. We dispense here from details because we shall not make use of the corresponding expressions. In the forthcoming examples we shall apply a modified single peak landscape where all sequences except those in the master core have the same fitness values $f$ and then the equation $\overline{f}_{-m} = f$ is trivially fulfilled.

For the purpose of illustration of the analysis of sequence clustering in strong quasispecies full numerical computations are compared with the zero mutational backflow approximation for the four membered cluster on the fitness landscape $\mathcal{L}(\lambda = 0, s = 919, d = 1.00)$ in figure 25. Although differences are readily recognized and the agreement between the full calculation and the approximation is not as good as in the case of a single master sequence, the appearance of the cluster is very well reproduced by the zero mutational backflow approximation. In particular, the relative frequencies of the four sequences forming the cluster are reproduced well. In comparison to the full calculation, the critical mutation rate at the error threshold, the point $p = p_{cr}$ at which the entire quasispecies $\overline{v}^{(0)}$ vanishes, appears at a higher $p$-value than the level crossings of the full calculation. The difference in the critical mutation rates is readily interpreted: The full calculation is based on a landscape with fully developed random scatter whereas the zero mutational backflow calculation compares best with a *four peak landscape* where the four peaks correspond to the members of the cluster ($\mathbf{X}_0$, $\mathbf{X}_4$, $\mathbf{X}_{516}$, $\mathbf{X}_{512}$) and all other sequences have identical fitness values. In order to show that this interpretation is correct the cluster has been implemented on a single peak landscape ($d = 0$) with the same fitness values ($f_0 = 1.1$ and $f = 1.0$) and the error threshold on this landscape is shifted slightly to higher values of the mutation rate parameter $p$. The agreement with the zero mutational backflow approximation is remarkably good. This agreement can be taken as a strong indication that the interpretation of strong quasispecies being a result of the formation of mutationally linked clusters of sequences within the population.

**Figure 25: Zero backflow approximation for a quasispecies on a *realistic* model landscape.** The landscape characteristic is $\mathcal{L}(\lambda = 0, d = 1.00, s = 919)$. Shown are the stationary concentrations $\bar{x}_j(p)^{(0)}$ $(j = 1, 2, 3, 4)$ for the cluster obtained through zero mutational backflow (upper plot), the results of the full numerical computation (middle plot), and of a full numerical computation where the cluster was implemented on the single peak landscape (lower plot, $d = 0$). Other parameters: $\nu = 10$, $f_0 = 1.1$, and $f = 1.0$.

**Table 3: Strong quasispecies.** Shown are error thresholds and level crossing values for three cases of strong quasispecies.

| Random seeds | $s = 401$ | | $s = 577$ | | $s = 919$ | |
|---|---|---|---|---|---|---|
| | $j$ | fitness | $j$ | fitness | $j$ | fitness |
| Core sequences | **0** | 1.1000 | **0** | 1.1000 | **0** | 1.1000 |
| | **64** | 1.0981 | **64** | 1.0951 | **4** | 1.0966 |
| | **16** | 1.0772 | **256** | 1.0894 | **512** | 0.9296 |
| | **80** | 1.0987 | **320** | 1.0999 | **516** | 1.0970 |
| | $j$ | $p_{(1/100)}$ | $j$ | $p_{(1/100)}$ | $j$ | $p_{(1/100)}$ |
| Level crossing, $d = 0$ | **0** | 0.01396 | **0** | **0.01410** | **0** | 0.01320 |
| | **64** | **0.01406** | **64** | 0.01402 | **4** | **0.01348** |
| | **16** | 0.01318 | **256** | 0.01377 | **512** | 0.00828 |
| | **80** | 0.01389 | **320** | 0.01410 | **516** | 0.01304 |
| Level crossing, $d = 1$ | **0** | **0.008443** | **0** | **0.006921** [a] | **0** | 0.007876 |
| | **64** | 0.008359 | **64** | 0.006481 | **4** | **0.008385** |
| | **16** | 0.007003 | **256** | 0.006440 | **512** | - - - [b] |
| | **80** | 0.007876 | **320** | 0.006733 | **516** | 0.007476 |
| Error threshold ($p_{\text{cr}}$) | 0.01134 | | 0.01145 | | 0.01087 | |

[a] The quasispecies with $s = 577$ shows a small smooth transition just above the error threshold. The following three sequences have the same or higher level crossing values: $p_{(1/100)}(\mathbf{X}_{899}) = 0.008026$, $p_{(1/100)}(\mathbf{X}_{931}) = 0.007186$, and $p_{(1/100)}(\mathbf{X}_{962}) = 0.006842$.

[b] The stationary concentration $\bar{x}_{512}(p)$ never exceeds nor reaches the value 0.01.

Three strong quasispecies with the values $s = 401$, 577, and 919 were found among the 21 landscapes studied here. The most important computed data are summarized in table 3. Like in the single master case on the single peak landscape the level crossing values $p_{(1/100)}$ occur at higher mutation rates than the error threshold (see figure 18). The shift in the strong quasispecies is about $\Delta p = 0.0265$ somewhat larger than that for the single master lying at $\Delta p = 0.00226$. In the single master case the error threshold was calculated

to be $p_{\rm cr} = 0.094875$ whereas here it is shifted to higher $p$-values by about $\Delta p = 0.0046$. The interpretation is straightforward: The core taken together has a higher effective fitness than a single master and this is reflected by the shift to higher mutation rates. This shift is smallest in case of the core with $s = 919$ being the in agreement with a particularly small fitness value of one of the two class 1 mutant ($f_{512} = 0.9296$). In fact, the core in this case consists practically of three sequences only: $\mathbf{X}_0$, $\mathbf{X}_4$, and $\mathbf{X}_{516}$. In the computation with fully developed scatter ($d = 1$) for the strong quasispecies with $s = 577$ we observe $p_{(1/100)}$-values that are smaller than in the other two cases. Again the explanation is straightforward: There is a small and smooth transition at a $p_{\rm tr}$ value just below the error threshold and the stationary concentration of the master beyond the transition is higher than that of the dominating sequence in the core, $\bar{x}_{899} > \bar{x}_0$, and the $p_{(1/100)}$-value for the sequence $\mathbf{X}_{899}$ is indeed higher, $p_{(1/100)}(\mathbf{X}_{899}) = 0.008026$.

The mutation rate at which the last stationary concentration crosses the value 1/100 shows some scatter: For the twenty one random landscapes that were investigated here it amounts to $p_{(1/100)}(\mathbf{X}_{\rm last}) = 0.00812\pm0.00071$. Interestingly the values for strong quasispecies lie close together at $p_{(1/100)}(\mathbf{X}_{\rm last}) = 0.0084$. The observed scatter in the level crossing of the concentration of $\mathbf{X}_{\rm last}$ is definitely smaller than that found for the master sequence ($p_{(1/100)}(\mathbf{X}_0)$ in figure 20), which is an obvious result since $\bar{x}_0$ decays to small values at transitions that occur before the error threshold, at values $p_{\rm tr} < p_{\rm cr}$.

## 8 "Realistic" rugged and neutral landscapes (RNL)

The second property of realistic fitness landscapes mentioned in section 7 is *neutrality* and the challenge is to implement it together with ruggedness. In order to be able to handle both features together we conceived a two parameter landscape: (i) the random scatter is denoted by $d$ as before and (ii) a degree of neutrality $\lambda$ is introduced. The value $\lambda = 0$ means absence of neutrality and $\lambda = 1$ describes the completely flat landscape in the sense of Motoo Kimura's *neutral evolution* [33]. The result of the theory of neutral evolution that is most relevant here concerns *random selection*: Although fitness differences are absent, one randomly chosen sequence is selected by means of the stochastic replication mechanism, $\mathbf{X} \to 2\,\mathbf{X}$ and $\mathbf{X} \to \oslash$. For most of the time the randomly replicating population consists of a dominant genotype and a number of neutral variants at low concentration.

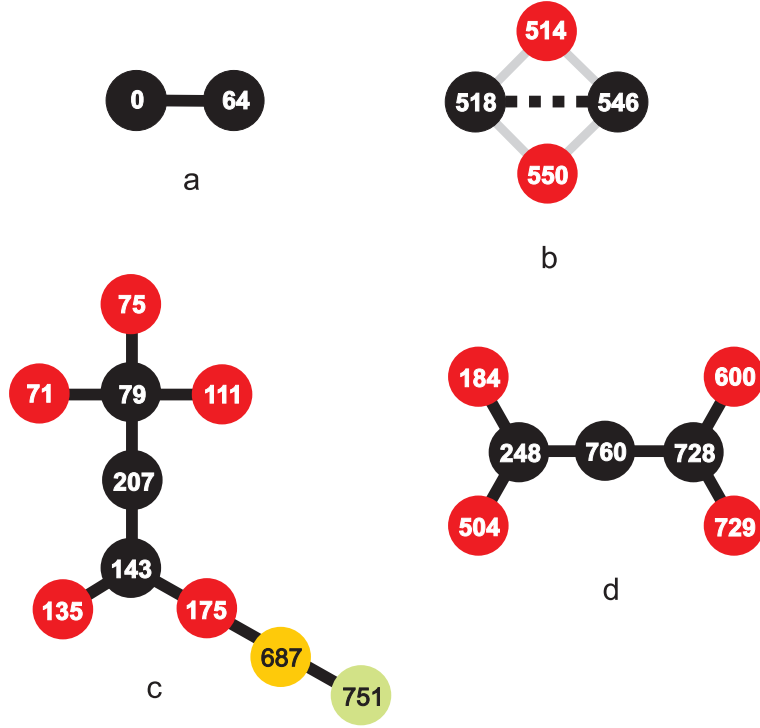An important issue of the landscape approach is the random positioning

of neutral master sequences in sequence space, which is achieved by means of the same random number generator that is used to compute the random scatter of the other fitness values obtained from pseudorandom numbers with a uniform distribution in the interval $0 \leq \eta \leq 1$:

$$
f(\mathbf{X}_j) = \begin{cases} f_0 & \text{if } j = 0 , \\[2ex] f_0 & \text{if } \eta_j^{(s)} \geq 1 - \lambda , \\[2ex] f + \frac{2d}{1-\lambda}(f_0 - f)\left(\eta_j^{(s)} - 0.5\right) & \text{if } \eta_j^{(s)} < 1 - \lambda , \\[1ex] & j = 1, \ldots, \kappa^l;\ j \neq m . \end{cases}
\tag{47}
$$

The rugged and neutral fitness landscape (47) is the complete analogue to the rugged fitness landscape (40) under the condition that several master sequences exist, which have the same highest fitness values $f_0$. The fraction of neutral mutants is determined by the fraction random numbers, which fall into the range $1 - \lambda < \eta \leq 1$, apart from statistical fluctuations this fraction is $\lambda$. At small values of the degree of neutrality $\lambda$ isolated peaks of highest fitness $f_0$ will appear in sequence space. Increasing $\lambda$ will result in the formation of clusters of sequences of highest fitness. Connecting all fittest sequences of Hamming distance $d_{\mathrm{H}} = 1$ by an edge results in a graph that has been characterized as *neutral network* [14, 78]. Neutral networks were originally conceived as a tool to model, analyze, and understand the mapping of RNA sequences into secondary structures [70, 79, 80]. The neutral network in RNA sequence-structure mappings is the preimage of a given structure in sequence space and these networks can be approximated in zeroth order by random graphs [81, 82]. Whereas neutral networks in RNA sequence-structure mappings are characterized by a relatively high degree of neutrality around $\lambda \approx 0.3$ and sequence space percolation is an important phenomenon, we shall be dealing here with much lower $\lambda$-values.

## 8.1 Small neutral clusters

The two smallest clusters of fittest sequences have Hamming distances $d_{\mathrm{H}} = 1$ and $d_{\mathrm{H}} = 2$ (figure 26). In the former case we are dealing with a minimal neutral neutral network, in the latter case the Hamming distance two sequences are coupled through two intermediate sequences similarly as in the core of strong quasispecies. An exact mathematical analysis for both cases is possible in the limit of vanishing mutation rates, $\lim p \to 0$ [59], led to

**Figure 26: Neutral networks in quasispecies.** The sketch presents four special cases that were observed on rugged neutral landscapes defined in equation (47). Part **a** shows the smallest possible network consisting of two sequences of Hamming distance $d_{\mathrm{H}} = 1$ observed with $s = 367$ and $\lambda = 0.01$. Part **b** contains two sequences of Hamming distance $d_{\mathrm{H}} = 2$, which are coupled through two $d_{\mathrm{H}} = 1$ sequences; it was found with $s = 877$ and $\lambda = 0.01$. The neutral network in part **c** has a core of three sequences, surrounded by five one-error mutants, one of them having a chain of two further mutants attached to it; the parameters of the landscape are $s = 367$ and $\lambda = 0.1$. Part **d** eventually shows a symmetric network with three core sequences and four one-error mutants attached to it, observed with $s = 229$ and $\lambda = 0.1$. Choice of further parameters: $n = 10$, $f_0 = 1.1$, $f = 1.0$, and $d = 0.5$. Color code: core sequences in black, one-error mutants in red, two-error mutants in yellow, and three-error mutants in green.

results that differ from Kimura's neutral theory:

$$\lim_{p \to 0} \bar{x}_{\mathrm{I}} = \frac{1}{2}, \; \lim_{p \to 0} \bar{x}_{\mathrm{II}} = \frac{1}{2} \qquad \text{for} \;\; d_{\mathrm{H}}(\mathbf{X}_{\mathrm{I}}, \mathbf{X}_{\mathrm{II}}) = 1 , \qquad (48\text{a})$$

$$\lim_{p \to 0} \bar{x}_{\mathrm{I}} = \frac{\alpha}{1 + \alpha}, \; \lim_{p \to 0} \bar{x}_{\mathrm{II}} = \frac{1}{1 + \alpha} \qquad \text{for} \;\; d_{\mathrm{H}}(\mathbf{X}_{\mathrm{I}}, \mathbf{X}_{\mathrm{II}}) = 2 , \qquad (48\text{b})$$

$$\lim_{p \to 0} \bar{x}_{\mathrm{I}} = 1 , \; \lim_{p \to 0} \bar{x}_{\mathrm{II}} = 0 \qquad \text{or} \qquad \lim_{p \to 0} \bar{x}_{\mathrm{I}} = 0 , \; \lim_{p \to 0} \bar{x}_{\mathrm{II}} = 1 ,$$

$$\text{for} \;\; d_{\mathrm{H}}(\mathbf{X}_{\mathrm{I}}, \mathbf{X}_{\mathrm{II}}) \geq 3 . \qquad (48\text{c})$$

If the two neutral fittest sequences, $\mathbf{X}_{\mathrm{I}}$ and $\mathbf{X}_{\mathrm{II}}$, are nearest neighbors in sequence space, $d_{\mathrm{H}}(\mathbf{X}_{\mathrm{I}}, \mathbf{X}_{\mathrm{II}}) = 1$, they are present at equal concentrations in the quasispecies in the low mutation rate limit, in case they are next
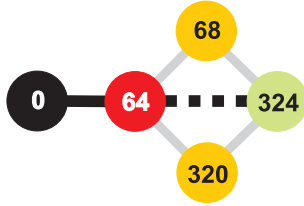
nearest neighbors in sequence space, $d_H(\mathbf{X}_I, \mathbf{X}_{II}) = 2$, they are observed at some ratio $\alpha$, and in both cases none of the two sequences vanishes. Only for Hamming distances $d_H(\mathbf{X}_I, \mathbf{X}_{II}) \geq 3$ Kimura's scenario of random selection occurs. It is important to stress a difference between the two scenarios, the deterministic ODE approach leading to clusters of neutral sequences and the random selection phenomenon of Motoo Kimura: In the quasispecies we have strong mutational flow within the cluster of neutral sequences – which is not substantially different from the flow within the non-neutral clusters in subsection 7.2 – and this flow outweighs fluctuations. In the random replication scenario mutations don't occur and the only drive for change in particle numbers is random fluctuations. For Hamming distances $d_H$ of three and more the mutational flow is too weak to counteract random drift.

The question now is whether or not the exact results derived for $\lim p \to 0$ are of more general validity. In order to find an answer numerical computations of quasispecies as functions of the mutation rate $p$ were performed. Random landscapes with a degree of neutrality of $\lambda = 0.01$ were searched and indeed the desired small networks with distances $d_H = 1$ and one for $d_H = 2$ between the master sequences were found for $s = 367$ and $s = 877$ (figure 26, parts **a** and **b**, respectively). Figures 27 and 29 show the solutions curves $\bar{x}_j(p)$ for the two examples of small neutral clusters. The concentration ratios of the two fittest sequences fulfil the predictions of the analytical approach in the limit of small mutation rates, $\lim p \to 0$: The ratio for $\mathbf{X}_0$ and $\mathbf{X}_{64}$ in the Hamming distance one case, $\bar{x}_0(0)/\bar{x}_{64}(0) = 1$, and some finite ratio $\bar{x}_{518}(0)/\bar{x}_{546}(0) = \alpha = 1.2259$ in the Hamming distance two case, respectively.

Figure 27 illustrates the dependence of the quasispecies formed by two master sequences of Hamming distance $d_H = 1$ on the mutation rate $p$. The extrapolation of the exact result, $\bar{x}_0/\bar{x}_{64} = 1$ to nonzero mutation rates turns out to be successful: Indeed, the red curve behind the black curve is hardly to be seen in the topmost plot as well as in the enlargement (middle plot). A precise calculation of this ratio shows a slight increase until at $\hat{p} = 0.009405$ a maximum of $\bar{x}_0(\hat{p})/\bar{x}_{64}(\hat{p}) = 1.0610$ is reached. Then the ratio decreases and apparently becomes unity again at $\tilde{p} = 0.5$. The plot of all stationary concentrations in the quasispecies belonging to the network **a** in figure 26 shows an interesting detail: The non-master sequence with the highest concentration, $\mathbf{X}_{324}$, does not belong to the combined one-error neighborhood of the two master sequences but lies at Hamming distance $d_H = 2$ and $d_H = 3$ from the two masters, $\mathbf{X}_0$ and $\mathbf{X}_{64}$, respectively. The explanation follows straightforwardly from an inspection of the fitness landscape. Sequence $\mathbf{X}_{324}$

**Figure 27: Cluster on a weakly neutral rugged model landscape with** $s = 367$. The plot in the middle is an enlargement of the topmost plot. in the bottom plot only the curves of the dominant cluster, consisting of the two master sequences, $\mathbf{X}_0$ and $\mathbf{X}_{64}$, their one-error neighborhoods, and the third fittest neutral sequence $\mathbf{X}_{324}$, are shown. Further parameters: $n = 10$, $f_0 = 1.1$, $f = 1.0$, $\lambda = 0.01$, $d = 0.5$. Color code see appendix.
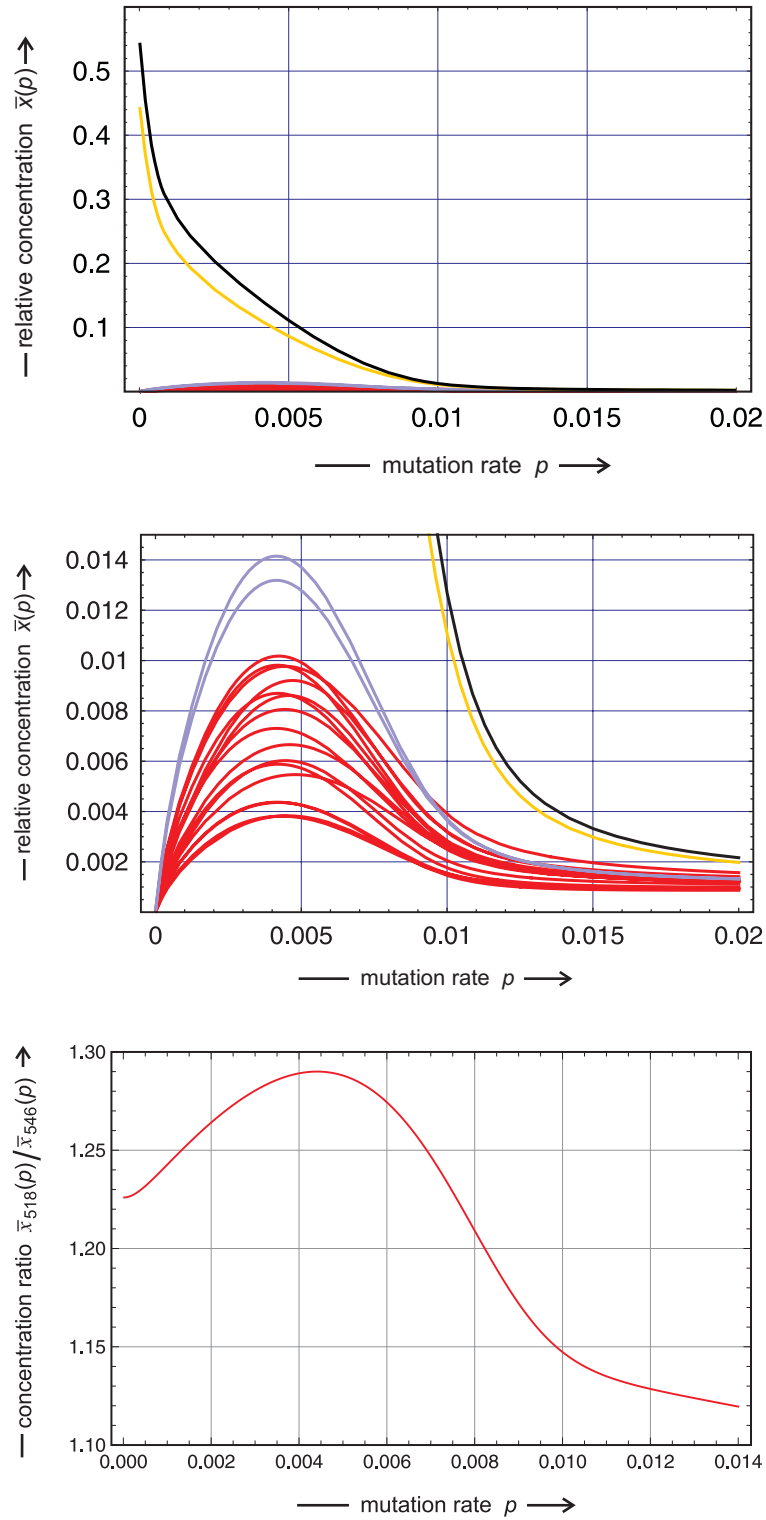
**Figure 28: A small neutral cluster in a quasispecies.** The color code from the appendix that is different from figure 26 is used: $\mathbf{X}_0$ black, the one-error mutant $\mathbf{X}_{64}$ red, the two-error mutants in yellow, and the fittest three-error mutant $\mathbf{X}_{324}$ in green.

belongs to the class of fittest neutral sequences but it is not coupled by an edge to the dominant network (figure 28). Instead it forms an Hamming distance two cluster together with $\mathbf{X}_{64}$ with $\mathbf{X}_{68}$ and $\mathbf{X}_{320}$ being the intermediates. The plot at the bottom of figure 27 contains the curves of the stationary concentrations of the Hamming distance one master pair (black and red) and their complete one-error neighborhood (red and yellow, respectively) together with that of the neutral sequence $\mathbf{X}_{324}$. It is interesting that all curves of the neighborhood sequences and the curve of $\mathbf{X}_{324}$ have their maxima at almost the same position near $p \approx 0.05$ whereas the maxima of all other curves (comparison with the middle plot of figure 28) are shifted towards higher $p$-values. An error threshold expressed by means of $p_{1/100}$-values occurs at somewhat higher mutation rates than in the case of a single master sequence: $p_{1/100}(\mathbf{X}_0) = 0.01073$ and $p_{1/100}(\mathbf{X}_{\text{last}}) = 0.01082$ with $\mathbf{X}_{\text{last}} \equiv \mathbf{X}_{64}$ compared to $p_{1/100}(\mathbf{X}_0) = 0.01065$ in the non-neutral case. As expected the Hamming distance one pair is equivalent to a master that is slightly stronger than a single sequence. Indeed, the fitness value of $\mathbf{X}_{64}$ is raised from $f_{64} = 1.04923$ to $f_{64} = 1.1$ on the landscape with neutrality.

An isolated cluster with a distance $d_{\text{H}} = 2$ between the two master sequences, $\{\mathbf{X}_{518}, \mathbf{X}_{546}\}$, has been observed on the rugged neutral landscape with $\lambda = 0.01$ and $s = 877$. In the limit $p \to 0$ the two fittest neutral sequences $\mathbf{X}_{518}$ and $\mathbf{X}_{546}$ are present at the stationary concentrations $\bar{x}_{518} = 0.5507$ and $\bar{x}_{546} = 0.4493$, respectively, and their ratio is $\bar{x}_{518}(0) \,/\, \bar{x}_{546}(0) = \alpha(0) = 1.2259$. Both stationary concentrations decrease with increasing $p$-values, the ratio increases at first but then decreases and approaches the value one corresponding to the uniform distribution: $\lim_{p \to \tilde{p}=1/2} \alpha(p) = 1$. The function $\alpha(p)$ passes a (local) maximum of $\alpha(p) = 1.29$ at $p = 0.00441$ (see figure 29, plot at the bottom). The plot in the middle of the figure demonstrates that the two sequences lying in between the master pair, $\mathbf{X}_{514}$ and $\mathbf{X}_{550}$, appear at higher concentrations than the rest of the one-error

**Figure 29: Quasispecies on weakly neutral rugged model landscape with $s = 877$.** The topmost part of the figure refers to the landscape with $s = 877$ and presents the solution curves for the master pair, $\{\mathbf{X}_{518}, \mathbf{X}_{546}\}$ and their one-error mutants. The plot in the middle is an enlargement and highlights the curves for the two intermediate sequences $\mathbf{X}_{514}$ and $\mathbf{X}_{550}$ in pastel blue. The plot at the bottom shows the ratio between the stationary concentrations of the two master sequences, $\alpha(p)$. Further parameters: $n = 10$, $f_0 = 1.1$, $f = 1.0$, $\lambda = 0.01$, $d = 0.5$. Color code see appendix.

**Figure 30: Quasispecies on weakly neutral rugged model landscape with $s = 877$.** The topmost part of the figure refers to the landscape with $s = 877$ and presents all solution curves Further parameters: $n = 10$, $f_0 = 1.1$, $f = 1.0$, $\lambda = 0.01$, $d = 0.5$. Color code see appendix.

cloud.[33] It is interesting to note that the landscape $\mathcal{L}(\lambda = 0.01, s = 877)$ sustains another Hamming distance two pair of fittest sequences $\{\mathbf{X}_0, \mathbf{X}_{132}\}$ with the intermediates $\mathbf{X}_4$ and $\mathbf{X}_{128}$. This second cluster is in competition with the first cluster as shown in figure 30 and gains in concentration with increasing mutation rates $p$, passes through a maximum and then decays through an error threshold to the uniform distribution at $p = \tilde{p}$. The position of the error threshold again is estimated by means of the $p_{1/100}$-values and one finds $p_{1/100}(\mathbf{X}_{518}) = 0.01053$ and $p_{1/100}(\mathbf{X}_{546}) = 0.01022$ with $\mathbf{X}_{518} \equiv \mathbf{X}_{\text{last}}$. On this neutral landscape the corresponding non-neutral master sequence is
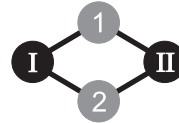
---

[33]In the case shown here, the two intermediate sequences have very similar fitness values, $f_{514} = 1.017$ and $f_{550} = 1.012$. Large fitness differences can outweigh the advantage caused by the mutation flow from both master sequences.

······· ACAU**G**CGAA ······· master sequence I
······· A**UAUA**CGAA ·······
······· ACAUGC**GC**A ·······
······· **G**CAUACGAA ·······
······· ACAUGC**U**AA ·······
······· ACAUGCGA**G** ·······
······· ACA**C**GCGAA ·······
······· AC**GU**ACGAA ·······
······· ACAU**A**GGAA ·······
······· ACAU**A**CGAA ······· master sequence II

······· ACAU$^\text{G}_\text{A}$CGAA ······· consensus sequence

······· ACA**GU**C**A**GAA ······· master sequence I
······· ACAGUC**C**GAA ······· intermediate 1
······· A**UA**AUCCGAA ·······
······· ACA**GU**C**A**GCA ·······
······· **G**CA**GU**C**A**GAA ·······
······· ACA**GU**C**AU**AA ·······
······· ACA**GU**C**A**GA**G** ·······
······· ACA**A**CCGAA ·······
······· AC**G**GUCAGAA ·······
······· ACA**GU**G**A**GAA ·······
······· ACA**A**UCAGAA ······· intermediate 2
······· ACA**A**UCCGAA ······· master sequence II

······· ACA$^\text{G}_\text{A}$UC$^\text{A}_\text{C}$GAA ······· consensus sequence
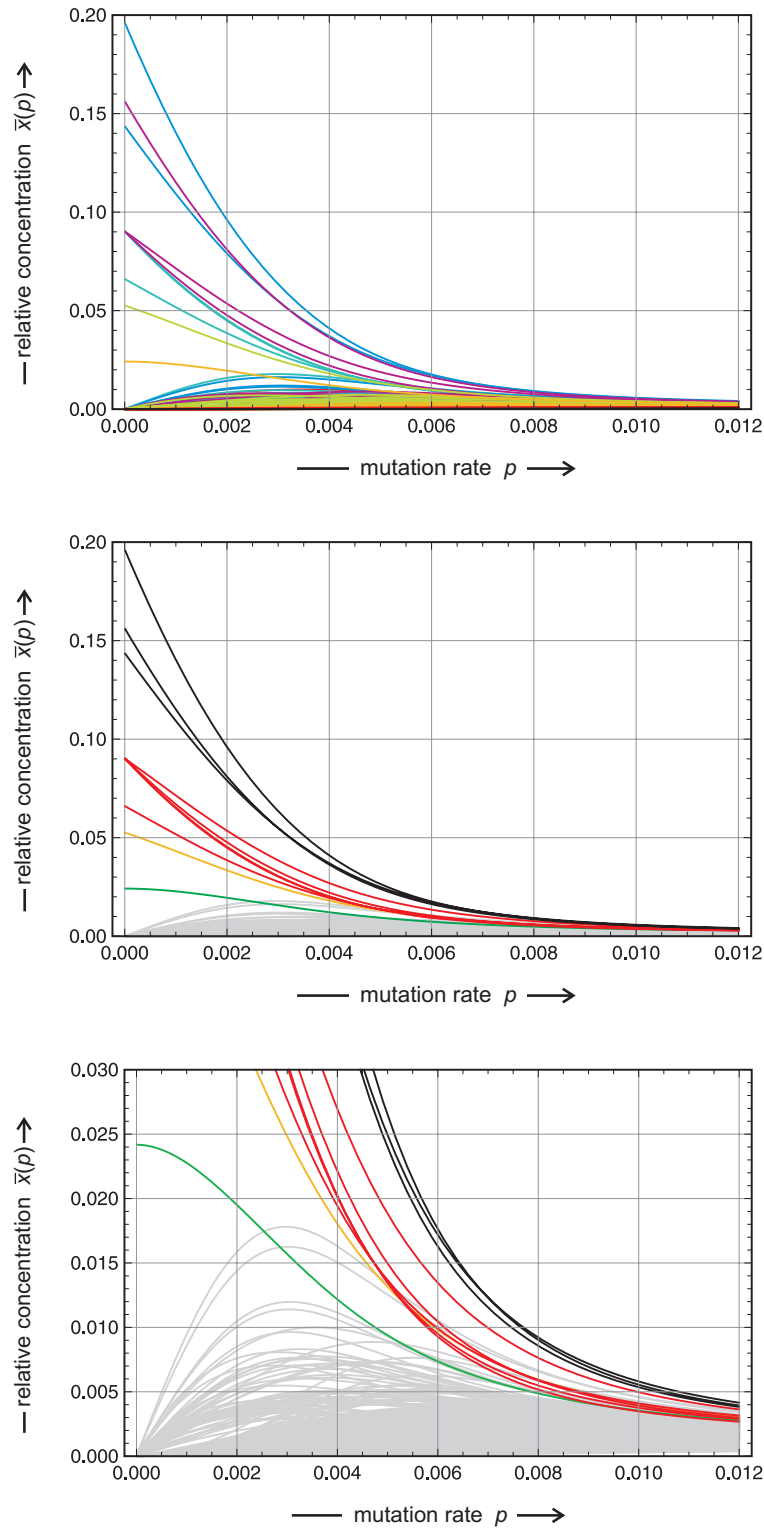
**Figure 31: Quasispecies with two neutral master sequences.** The sketch contains a set of sequences in order to demonstrate the role of neutrality in the determination of the consensus sequence of a population. Two fittest sequences with Hamming distance $d_\text{H} = 1$ (upper picture) lead to an ambiguity at one position. Mutations at other positions are wiped out by statistics whereas the one-to-one ratio of the two mater sequences leads to a 50/50 ratio of two nucleobases at the position of the mutation. The lower picture refers to two master sequences with Hamming distance $d_\text{H} = 2$: Ambiguities are obtained at two positions and the ratios of the nucleobases are given approximately by the value of $\alpha$ in equation(48b). The two intermediates are present in small stationary concentrations only but they are, nevertheless, more frequent than the other one-error mutants (see figure 29).

more stable than the cluster as expressed by $p_{1/100}(\mathbf{X}_0) = 0.01075$. This fact is difficult to interpret, because the original master $\mathbf{X}_0$ is not member of the cluster, which accordingly is situated in another part of sequences space with different fitness values of the neighboring sequences.

Eventually we consider a simple practical consequence of the existence of fittest neutral pairs of Hamming distance $d_\text{H} = 1$ and $d_\text{H} = 2$ for the sequence analysis in populations. Despite vast sequence heterogeneity [83], in particular of virus populations, average or consensus sequences are fairly insensitive to individual mutations provided the population size is sufficiently

**Figure 32: Quasispecies on rugged neutral model landscapes I.** Shown are the stationary concentrations for the landscape $\mathcal{L}(\lambda = 0.1, d = 0.5, s = 229)$. The topmost plot is drawn with the color code of the appendix, the plot in the middle applies the color code of the neutral network in figure 26 **d** with the curves of sequences not belonging to the net in grey. The plot at the bottom is an enlargement of the plot in the middle. Other parameters: $\nu = 10$, $f_0 = 1.1$, and $f = 1.0$.

**Figure 33: Quasispecies on rugged neutral model landscapes II.** Shown are the stationary concentrations for the landscape $\mathcal{L}(\lambda = 0.1, d = 0.5, s = 367)$. The topmost plot is drawn with the color code of the appendix, the plot in the middle applies the color code of the neutral network in figure 26 **c** with the curves of sequences not belonging to the net in grey. The plot at the bottom is an enlargement of the plot in the middle. Other parameters: $\nu = 10$, $f_0 = 1.1$, and $f = 1.0$.

large. Individual deviations in mutant sequences cancel through averaging in population with single master sequences. This will not be the case in the presence of neutral variants. In the 50/50 mixture of two master sequences with mutant clouds surrounding both the sequence difference between the masters is not going to cancel by averaging and ambiguities remain. Considering now the two cases discussed here: (i) two master sequences at Hamming distance one present at equal concentrations and (ii) two master sequences at Hamming distance two present at a concentration ratio $\alpha$, we expect to find sequence averages as sketched in figure 31. In the former case a 50/50 mixtures of two nucleotides is expected to occur at one position on the sequence, and in the latter case two positions will show nucleobase ambiguities with the ratio $\alpha$.
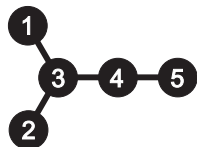
## 8.2 Medium size neutral clusters

An increase in the degree of neutrality $\lambda$ will results in the appearance of larger neutral networks that are scattered all over sequence space. We start here by the introduction of the adjacency matrix as an appropriate reference state of neutral networks and the discuss two examples of more complex neutral networks that are observed in form of quasispecies on landscapes $\mathcal{L}_n(\lambda = 0.1, d = 0.5, s)$.

The adjacency matrix of a graph contains an entry one at every off diagonal element that corresponds to an edge in the graph. We have, for example, the adjacency matrix A

$$
A = \begin{pmatrix}
0 & 0 & 1 & 0 & 0 \\
0 & 0 & 1 & 0 & 0 \\
1 & 1 & 0 & 1 & 0 \\
0 & 0 & 1 & 0 & 1 \\
0 & 0 & 0 & 1 & 0
\end{pmatrix}
$$

for the graph:



Now we consider a neutral network corresponding to this graph and obtain for the mutation matrix Q:

$$
Q = \begin{pmatrix}
(1-p)^n & (1-p)^{n-2}p^2 & (1-p)^{n-1}p & (1-p)^{n-2}p^2 & (1-p)^{n-3}p^3 \\
(1-p)^{n-2}p^2 & (1-p)^n & (1-p)^{n-1}p & (1-p)^{n-2}p^2 & (1-p)^{n-3}p^3 \\
(1-p)^{n-1}p & (1-p)^{n-1}p & (1-p)^n & (1-p)^{n-1}p & (1-p)^{n-2}p^2 \\
(1-p)^{n-2}p^2 & (1-p)^{n-2}p^2 & (1-p)^{n-1}p & (1-p)^n & (1-p)^{n-1}p \\
(1-p)^{n-3}p^3 & (1-p)^{n-3}p^3 & (1-p)^{n-2}p^2 & (1-p)^{n-1}p & (1-p)^n
\end{pmatrix}
.
$$

In the limit of small mutation rates we neglect all powers $f(p) \in o(p)$ and after multiplication with the fitness matrix $F = f_0 \cdot \mathbb{I}$, where $\mathbb{I}$ is the identity or unit matrix, the result is the value matrix in the zeroth order approximation

$$W^{(0)} = f_0 \, p^{n-1} \cdot \begin{pmatrix} 1-p & 0 & p & 0 & 0 \\ 0 & 1-p & p & 0 & 0 \\ p & p & 1-p & p & 0 \\ 0 & 0 & p & 1-p & p \\ 0 & 0 & 0 & p & 1-p \end{pmatrix} .$$

Since neither the addition of a constant to the diagonal elements nor the multiplication by a common factor changes eigenvectors, the adjacency matrix and the matrix $W^{(0)}$ have identical eigenvectors. Accordingly, the adjacency matrix of the neutral network is the appropriate reference for quasispecies rugged neutral landscapes. Clearly, this has been the case for the small cluster shown in figure 27 where the dominant eigenvector of the trivial adjacency matrix

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad \text{simply is} \quad \zeta_0^{(0)} = (\tfrac{1}{2}, \tfrac{1}{2})^{\mathrm{t}} \,,$$

and represents also the solution of the mutation selection equation (18) for $p = 0$.

In figure 32 the quasispecies as a function of the mutation rate, $\bar{\Upsilon}(p)$ is shown for the landscape $\mathcal{L}(\lambda = 0.1, d = 0.5, s = 229)$. The neutral network consists of seven sequences, three of them form a linear inner core and four are attached to it on the periphery (figure 26 d). The dominant eigenvector of the adjacency matrix is of the form

$$\zeta_0^{(0)} = (0.1, 0.1, 0.2, 0.2, 0.2, 0.1, 0.1) = (\bar{x}_{184}, \bar{x}_{504}, \bar{x}_{248}, \bar{x}_{760}, \bar{x}_{728}, \bar{x}_{600}, \bar{x}_{729})^{\mathrm{t}} \,.$$

The figure shows that the relative concentration within the quasispecies in the sense of three more frequent and four less frequent sequences are perfectly maintained almost up to the error threshold. The level crossing of the three core sequences occurs at: $p_{1/100}(\mathbf{X}_{248}) = 0.007712$, $p_{1/100}(\mathbf{X}_{760}) = 0.007307$, and $p_{1/100}(\mathbf{X}_{728}) = 0.007413$.

The neutral network on the landscape $\mathcal{L}(\lambda = 0.1, d = 0.5, s = 367)$ has a more complicated structure than the symmetric seven membered neutral cluster discussed in the previous paragraph. It contains ten individual sequences and has the form shown in figure 26 c: A core of three sequences is surrounded by five nearest neighbors and has a tail consisting of one Hamming distance two and one Hamming distance three sequence. Also

**Table 4: Neutral networks and adjacency matrix.** The largest eigenvector of the adjacency matrix of the neutral network in figure 26 c is compared with the quasispecies calculated for a small value of the mutation rate $p$.

| Class | No. | $j$ | $\zeta_0(\mathbf{X}_j)$ [a] | $\zeta_0^{(0)}(\mathbf{X}_j)$ [b] |
|---|---|---|---|---|
| Core | 1 | **79** | 0.196220 | 0.196281 |
| | 2 | **207** | 0.156240 | 0.156284 |
| | 3 | **143** | 0.143652 | 0.143688 |
| Class 1 | 4 | **175** | 0.090212 | 0.090231 |
| | 5 | **71** | 0.090206 | 0.090231 |
| | 6 | **75** | 0.090212 | 0.090231 |
| | 7 | **111** | 0.090206 | 0.090231 |
| | 8 | **135** | 0.066039 | 0.066053 |
| Class 2 | 9 | **687** | 0.052585 | 0.052934 |
| Class 3 | 10 | **751** | 0.024177 | 0.024177 |

[a] The entries in this column are the components of the quasispecies expressed as the elements of the largest eigenvector of the value matrix W computed with a mutation rate $p = 1 \times 10^{-6}$.

[b] The entries in this column are the components of the largest eigenvector $\zeta_0^{(0)}$ of the adjacency matrix of the graph representing the fittest neutral network on the landscape $\mathcal{L}(\lambda = 0.05, d = 0.5, s = 367)$.

for this more involved topology the low mutation rate limit of the quasispecies, $\lim_{p \to 0} \bar{\Upsilon}(p)$, converges exactly to the largest eigenvector of the adjacency matrix (table 4). The three core sequences stay together for the whole range of $p$-values up to the error threshold that is reached in terms of level crossing at: $p_{1/100}(\mathbf{X}_{79}) = 0.007649$, $p_{1/100}(\mathbf{X}_{207}) = 0.007462$, and $p_{1/100}(\mathbf{X}_{143}) = 0.007704$. It is not easy to guess that four out of the five nearest neighbor sequences have identical values in the eigenvector of the adjacency matrix – and it is the tail-free sequence, $\mathbf{X}_{135}$, and not the sequence carrying the tail, $\mathbf{X}_{175}$, which is different from the other four.

Further increase in the degree of neutrality, $\lambda$, gives rise to extended neutral networks, which eventually percolate whole sequence space. Whether or not such large clusters of neutral sequences play a role in real biology cannot

**Figure 34: Lethal mutants and replication errors.** The model for lethal mutants corresponding to a *single peak landscape* with $k_1 = 1$ and $k_2 = \ldots = k_n = 0$ is studied in the flow reactor. The concentrations of the master sequence (black) and the mutant classes (red, dark orange, light orange, etc.; full lines) are shown as functions of the error rate $p$. For the purpose of comparison the parameters were chosen with $\nu = 20$, $r = 1$, $a_0 = 2$, and $\eta = 2$. The plots are compared to the curves for the master sequence (grey; broken curve) and the one error class (light red; broken curve) of a quasispecies on the single peak landscape with $\nu = 20$, $f_0 = 2$, $f = 1$, and hence $\sigma = 2$.

be said in the moment but more empirical knowledge on fitness landscapes will help to decide this question.

## 9  Lethal mutations

Many antiviral drugs are powerful because they increase the mutation rate and drive virus populations to extinction but for a satisfactory molecular explanation of the mechanism the required information on the fitness landscape is still missing. Nevertheless, *lethal mutagenesis* is an important phenomenon and simplified models providing phenomenological explanations have been developed.[34] It is important to note that a quasispecies can exist also in cases where the Perron-Frobenius theorem is not fulfilled. As an example we consider an extreme case of lethality that allows for analytical solutions: Only the master genotype $\mathbf{X}_0$ has a positive fitness value, $f_0 > 0$ and $f_1 = \ldots = f_{n-1} = 0$, and hence only the entries $W_{k0} = Q_{k1}f_0$ of matrix

---

[34]An early paper [84] claimed that zero fitness values are incompatible with the existence of quasispecies and error threshold. The result, however, turned out to be an artifact of a rather naïve linear sequence space, since later works demonstrated that selection and mutation on realistic sequence spaces sustains error thresholds also in presence of lethal variants [85, 86].

W are nonzero:

$$W = \begin{pmatrix} W_{00} & 0 & \cdots & 0 \\ W_{10} & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ W_{n-1,0} & 0 & \cdots & 0 \end{pmatrix} \quad \text{and} \quad W^k = W_{00}^k \begin{pmatrix} 1 & 0 & \cdots & 0 \\ \frac{W_{10}}{W_{00}} & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \frac{W_{n-1,0}}{W_{00}} & 0 & \cdots & 0 \end{pmatrix}.$$

Accordingly, W is not primitive but under suitable conditions $\bar{\Upsilon} = \bar{\mathbf{x}} = (Q_{00}, Q_{10}, \ldots, Q_{n-1,0})$ is a stable stationary mutant distribution and for $Q_{00} > Q_{j0} \, \forall \, j = 1, \ldots, n-1$ – correct replication occurs more frequently than a particular mutation – genotype $\mathbf{X}_0$ is the master sequence. On the basis of a rather idiosyncratic mutation model consisting in a one-dimensional chain of genotypes the claim was raised that no quasispecies can be stable in presence of lethal mutants and accordingly, no error thresholds could exist [84]. Recent papers [85, 86], however, used a realistic high-dimensional mutation model and presented analytical results as well as numerically computed examples for error thresholds in the presence of lethal mutations.

In order to be able to handle the case of lethal mutants properly we have to go back to absolute concentrations in a realistic physical setup, the flow reactor applied in section 3 and shown in figure 3. We neglect degradation and find for $\mathbf{X}_0$ being the only viable genotype:[35]

$$\begin{aligned} \frac{da}{dt} &= -\left( \sum_{i=0}^{n-1} Q_{i0} k_0 \, c_1 \right) a + r \, (a_0 - a) \\ \frac{dc_i}{dt} &= Q_{i0} k_0 \, a \, c_0 - r \, c_i \, , \quad i = 0, 1, \ldots, n-1 \, . \end{aligned} \tag{49}$$

Computation of stationary states is straightforward and yields two solutions, (i) the state of extinction with $\bar{a} = a_0$ and $\bar{c}_i = 0 \, \forall \, i = 0, 1, \ldots, n-1$, and (ii) a state of selection of a quasispecies $\bar{\Upsilon}$ that consists of the master sequence $\mathbf{X}_0$ and its mutant cloud at the stationary concentrations $\bar{a} = r/(Q_{00} k_0)$, $\bar{c}_0 = Q_{00} a_0 - r/k_0$, and $\bar{c}_i = \bar{c}_1 (Q_{i0}/Q_{00})$ for $i = 1, 2, \ldots, n-1$. As an example we compute a maximum error rate for constant flow, $r = r_0$, again applying the uniform error rate model (16):

$$\begin{aligned} Q_{00} &= (1-p)^\nu \quad \text{and} \\ Q_{i0} &= p^{d_{i0}} \cdot (1-p)^{\nu - d_{i0}} \, , \end{aligned}$$

where $d_{i0}$ again is the Hamming distance between the two sequences $\mathbf{X}_i$ and $\mathbf{X}_0$. Instead of the superiority $\sigma$ of the master sequence – that diverges since $\bar{f}_{-m} = 0$ because of $f_1 = \ldots = f_{n-1} = 0$ – we define a dimensionless quantity

---

[35]We use $k_i$ for the rate constants as in section 3, since $a(t)$ is a variable here.

$\eta_0$, the *carrying surplus* of the reactor for the master sequence $\mathbf{X}_0$,[36] which can be defined to be

$$\eta = \frac{k_0 \, a_0}{r_0}$$

for the flow reactor. The value of $p$, at which the stationary concentration of the master sequence $\bar{c}_1(p)$ and those of all other mutants vanishes, represents the analogue of the error threshold (33), and for the sake of clearness it is called the *extinction threshold*. Using $\ln(1 - p) \approx -p$ again we obtain:
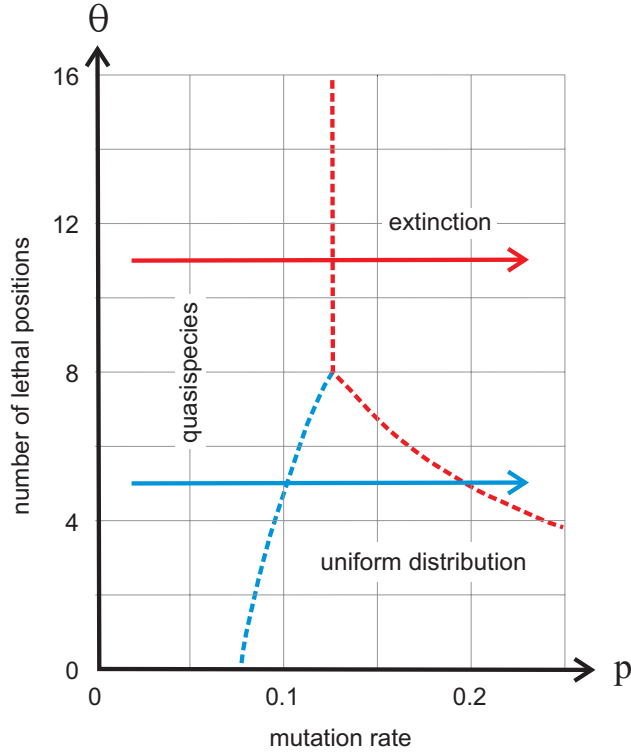
$$p_{\max} \approx \frac{\ln \eta}{\nu} \quad \text{for small } p \, . \tag{50}$$

The major difference between the error threshold (33) and the extinction threshold (50) concerns the state of the population at values $p > p_{\max}$: Replication with non-zero fitness of mutants leads to the uniform distribution, whereas the population goes extinct in the lethal mutant case. Accordingly, the transformation to relative concentrations fails and equation (9) is not applicable. In figure 34 we show an example for the extinction threshold with $\nu = 20$ and $\eta = 2$. For this case the extinction threshold is calculated from (50) to occur at $p_{\max} = 0.03466$ compared to a value of $0.03406$ observed in computer simulations. In the figure we see also a comparison of the curves for the master sequence and the one error class in the lethality model and on single peak landscape. The agreement of the two curves for the master sequences is not surprising since the models were adjusted to coincide in the values $\bar{c}_0(0) = 1$ and $p_{\max} = \ln 2/20$. The curves for the one error classes show some difference that is due to the lack of mutational backflow in case of lethal variants.

Manfred Eigen and Esteban Domingo originally explained lethal mutagenesis caused by pharmaceutical compounds increasing the mutation rate through driving populations beyond the error threshold [11,62]: At mutation rates above the error threshold replication becomes random, the result is a complete loss of the genetic information and eventually the viral life cycle breaks down. Later on James Bull and Claus Wilke studied the dynamics of lethal mutagenesis in more detail and claimed correctly that population extinction is a phenomenon independently of catastrophic error accumulation [53, 54] (For a critical review on this subject see [87]). Recently, lethal mutagenesis and error threshold crossing has been modeled by means of

---

[36]The name carrying surplus for $\eta_0 \geq 1$ is chosen in order to indicate how far the reactor is away from the state of extinction at $\eta_0 = 1$.

**Figure 35: Quasispecies and lethal mutations.** The sketch shows the long time behavior of mutation-selection dynamics in presence of lethal mutations. For a sufficiently large number of lethal positions ($\theta < 8$) in the virus genotype the quasispecies goes extinct at the extinction threshold (red) whereas for a smaller fraction of lethal variants two thresholds are observed: (i) an error threshold (blue) and an extinction threshold (red). The picture is redrawn from Tejero et al. [86], figure 2d. Choice of parameters: chain length $\nu = 20$, $f_m = 15$, $f_k = 3$, $\vartheta = 1$.

three-species replication-mutation kinetics with neglect of back mutation [86]

$$
\begin{aligned}
\frac{\mathrm{d}dc_m}{\mathrm{dt}} &= \left(f_m(1-p)^\nu - \vartheta\right)c_m \;, \\
\frac{\mathrm{d}dc_k}{\mathrm{dt}} &= f_m(1-p)^\theta\left(1-(1-p)^{(\nu-\theta)}\right)c_m + \left(f_k(1-p)^\theta - \vartheta\right)c_k\;, \quad \text{and} \\
\frac{\mathrm{d}dc_j}{\mathrm{dt}} &= f_m\left(1-(1-p)^\theta\right)c_m + f_k\left(1-(1-p)^\theta\right)c_k - \vartheta\,c_j\;; \quad c = \sum_{i=1}^{n} c_i \;,
\end{aligned}
\tag{51}
$$

where $\mathbf{X}_m$ is the master sequence with the concentration $[\mathbf{X}_m] = c_m$ and the replication rate parameter $f_m$, $c_k$ and $f_k$ refer to the class of non-lethal mutants, and $c_j$ to the class of lethal mutants, $\vartheta$ is the uniform degradation rate parameter for all sequences, $\nu$ is the chain length, $\theta$ the number of positions at which mutation yields a lethal variant and, eventually $p$ the single point mutation rate. This model [86] is characterized by two features: (i) The concentration of the material consumed in the reproduction process is assumed to be constant, $[\mathbf{A}] = a_0$, and $a_0$ is absorbed in the fitness parameter $f_i$ $(i = 1, \ldots, n)$, and (ii) a degradation rate $\vartheta$ is introduced for all species.

In contrast to the selection-mutation equation (18) and the flow reactor discussed in section 3, the model system (51) does not approach a stationary state but the total concentration $c$ either grows infinitely or goes extinct. Figure 35 illustrates the result concerning lethal mutagenesis. There are two different scenarios of quasispecies development with increasing mutation rate $p$, which depend on the degree of lethality that is expressed by the number of lethal sites $d$: (i) At low lethality the quasispecies reaches first the error threshold at $p = p_{\mathrm{cr}}$, passes a range of $p$-values and then becomes extinct at $p = p_{\mathrm{ext}}$, and (ii) at sufficiently high degree of lethality the error threshold merges with the extinction threshold and the quasispecies dies out directly at $p = p_{\mathrm{ext}}$. It is worth noticing that the stability of the quasispecies against mutation increases with increasing degree of lethality corresponding to a shift of the error threshold towards higher mutation rate. Lethal mutagenesis is understood at the phenomenological level but when it comes to molecular details, more experimental data and a comprehensive molecular theory is required. Studies based on more realistic landscapes including lethal variants into model landscapes in the sense of (40) and (47) are still missing too.

## 10 Limitations and perspectives

An implicit assumption of the mathematical analysis of Darwinian selection presented here is the applicability of kinetic differential equations to describe selection and mutation in populations. In principle the ODE approach neglects finite size effects and hence is exact for an infinite population size only. Commonly, large numbers of particles, molecules, individuals or agents are sufficient to justify the use of differential equations. Biological populations, however, may be relatively small and low frequency mutants may be often present in a single copy or very few copies only. The uniform distribution at error rates above the threshold presents an illustrative example that is relevant for modeling evolutionary dynamics: It can never be achieved in reality because the numbers of possible polynucleotide sequences are huge compared to the largest accessible populations ranging from $10^6$ to $10^{15}$ individuals in replication experiments with bacteria, viruses, or RNA molecules. The human population size is approximately $7 \times 10^9$. Typical situations in biology differ from scenarios in chemistry where large populations are distributed upon a few chemical species. Are the results derived from the differential equations then representative for real systems? Two situations can be distinguished: (i) Individual mutations are rare events and it is ex-

tremely unlikely that the same mutation occurs twice or is precisely reversed after it has occurred, and (ii) mutations are sufficiently frequent and occur in both directions within the time of observation. The first scenario seems to be fulfilled with higher organisms. The second scenario is typical for virus evolution and *in vitro* evolution experiments with molecules. Bacteria may be in an intermediate situation. A typical example of the low mutation rate scenario (i) is Muller's ratchet named after the American biologist Hermann Muller [88, 89]. Lost mutants are not likely to be replaced, all variants starting with the fittest one will disappear sooner or later, and it is a only matter of time before a situation is reached where all genotypes have been replaced by others no matter what there fitness values were (for a comparison between the error threshold phenomenon and Muller's ratchet see [84]). The frequent mutation scenario (ii) allows for modeling and studying the kinetic equations of reproduction and selection as stochastic processes [60, 90–92] – examples are multitype branching or birth-and-death processes – chemical master equations [93], which are amenable to computer simulations [94] (for an overview of stochastic modeling see, e.g., [95]). The expression for the error threshold can be readily extended to finite populations [60]. Formation of stable quasispecies requires a replication fidelity that is the higher the smaller the population size is.

How relevant is the error threshold in realistic situations? According to the results presented in subsection 5.6 the question boils down to an exploration of natural fitness landscapes: Are biopolymer landscapes rugged or smooth? All evidence obtained so far points towards a rather bizarre structures of these landscapes. Single nucleotide exchanges may lead to large effects, small effects or no consequences at all as in the case of neutral mutations. Since biomolecules are usually optimized with respect to their functions within an organism, most mutations have deleterious effects or no effect. Biopolymer landscapes have three characteristic features, out of which all three are hard to visualize: (i) high dimensionality, (ii) ruggedness, and (iii) neutrality. In case equally fit genotypes are nearest or next nearest neighbors in sequence space they form joint quasispecies as described in [59]. When they are not closely related, however, neutral evolution in the sense of Motoo Kimura is observed [33]. Neutrality in genotype space still allows for the formulation of a selection model in phenotype space [96, 97]. Then, the variables are concentrations of phenotypes that are obtained through lumping together all concentrations of genotypes, which form the same phenotype and an analysis similar to the one presented here can be carried out. The genotypic error threshold is relaxed and the system gives rise to a phenotypic

error threshold below which the fittest or master phenotype is conserved in the population. The ODE model is readily supplemented by a theory of phenotype evolution based on a new concept of evolutionary nearness of phenotypes in sequence space [47, 73, 98], which is confirmed by computer simulations of RNA structure optimization in a flow reactor of the type discussed here [47, 73, 99]. Random drift of populations occurs on neutral subspaces of sequence space, which is visualized by a series of snapshots showing the spreading of a population that breaks up into individual clones [99] as it has been observed already earlier with model simulations [100, 101]. Computer simulations were also successful in providing evidence for the occurrence of error thresholds in stochastic replication-mutation systems [102].

The molecular approach to evolutionary phenomena does not only provide a firm basis that is rooted in chemical kinetics but at the same time it represents a frame that can be readily extended to complications that are very likely to be prohibitive for other theories. Complex reproduction mechanisms can be incorporated into the mutation-selection equations. One successful example is the detailed multistep reaction of $Q\beta$ RNA replication [103–106]: Different outcomes of the selection processes for different experimental conditions are correctly predicted. The integration of genetic regulation on the DNA or RNA level into the evolutionary model is not simple but straightforward, since gene regulation is based on biochemical kinetics. A still unsolved problem of high current actuality in population dynamics concerns the role of epigenetics in the evolution of phenotypes [107–110]. A well understood and illustrative example is provided by the early embryonic development of the fruit fly *drosophila*: The embryo is shaped by a cooperative interaction between maternal and zygotic genes [111, chapter 9] and it is meaningless to separate genetics and epigenetics in this case. A theoretical approach based on molecular biology and dealing simultaneously with several generations might bring the solution.

**Acknowledgements**

# References

[1] M. Eigen. Selforganization of matter and the evolution of biological macromolecules. *Naturwissenschaften*, 58:465–523, 1971.

[2] M. Eigen and P. Schuster. The hypercycle. A principle of natural self-organization. Part A: Emergence of the hypercycle. *Naturwissenschaften*, 64:541–565, 1977.

[3] M. Eigen and P. Schuster. The hypercycle. A principle of natural self-organization. Part B: The abstract hypercycle. *Naturwissenschaften*, 65:7–41, 1978.

[4] M. Eigen and P. Schuster. The hypercycle. A principle of natural self-organization. Part C: The realistic hypercycle. *Naturwissenschaften*, 65:341–369, 1978.

[5] M. Eigen, J. McCaskill, and P. Schuster. Molecular quasispecies. *J. Phys. Chem.*, 92:6881–6891, 1988.

[6] M. Eigen, J. McCaskill, and P. Schuster. The molecular quasispecies. *Adv. Chem. Phys.*, 75:149–263, 1989.

[7] M. Eigen and P. Schuster. Stages of emerging life - Five principles of early organization. *J. Mol. Evol.*, 19:47–61, 1982.

[8] E. Domingo, C. K. Biebricher, M. Eigen, and J. J. Holland. *Quasispecies and RNA Virus Evolution: Principles and Consequences*. Landes Bioscience, Austin, TX, 2001.

[9] J. Swetina and P. Schuster. Self-replication with errors - A model for polynucleotide replication. *Biophys. Chem.*, 16:329–345, 1982.

[10] J. W. Drake. Rates of spontaneous mutation among RNA viruses. *Proc. Natl. Acad. Sci. USA*, 90:4171–4175, 1993.

[11] E. Domingo, ed. Virus entry into error catastrophe as a new antiviral strategy. *Virus Research*, 107(2):115–228, 2005.

[12] J. Rogers and G. F. Joyce. A ribozyme that lacks cytidine. *Nature*, 402:323–325, 1999.

[13] J. S. Reader and G. F. Joyce. A ribozyme composed of only two different nucleotides. *Nature*, 420:841–844, 2002.

[14] C. M. Reidys and P. F. Stadler. Combinatorial landscapes. *SIAM Review*, 44:3–54, 2002.

[15] R. W. Hamming. Error detecting and error correcting codes. *Bell Syst. Tech. J.*, 29:147–160, 1950.

[16] J. Maynard-Smith. Natural selection and the concept of a protein space. *Nature*, 225:563–564, 1970.

[17] D. W. Mount. *Bioinformatics. Sequence and Genome Analysis*. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, second edition, 2004.

[18] S. Wright. The roles of mutation, inbreeding, crossbreeding and selection in evolution. In D. F. Jones, editor, *Int. Proceedings of the Sixth International Congress on Genetics*, volume 1, pages 356–366, Ithaca, NY, 1932. Brooklyn Botanic Garden.

[19] P. F. Stadler and G. P. Wagner. Algebraic theory of recombination spaces. *Evolutionary Computation*, 5:241–275, 1997.

[20] B. R. M. Stadler, P. F. Stadler, M. Shpak, and G. P. Wagner. Recombination spaces, metrics, and pretopologies. *Z. phys. Chem.*, 216:217–234, 2002.

[21] C. K. Biebricher. Quantitative analysis of mutation and selection in self-replicating RNA. *Adv. Space Res.*, 12(4):191–197, 1992.

[22] P. Carrasco, J. A. Darós, P. Agudelo-Romero, and S. F. Elena. A real-time RT-PCR assay for quantifying the fitness of tobacco etch virus in competition experiments. *J. Virological Methods*, 139:181–188, 2007.

[23] H. Ahn, H. La, and L. J. Forney. System for determining the relative fitness of multiple bacterial populations without using selecitve markers. *Appl. Environ. Microbiol.*, 72:7383–7385, 2006.

[24] C. F. Pope, T. D. McHugh, and S. H. Gillespie. Methods to determine fitness in bacteria. *Methods Mol. Biol.*, 642(3):113–121, 2010.

[25] R. D. Kouyos, G. E. Leventhal, T. Hinkley, M. Haddad, J. M. Whitcomb, C. J. Petropoulos, and S. Bonhoeffer. Exploring the complexity of the HIV-1 fitness landscape. *PLoS Genetics*, 8:e1002551, 2012.

[26] S. Gavrilets. *Fitness Landscapes and the Origin of Species*. Princeton University Press, Princeton, NJ, 2004.

[27] T. Aita and Y. Husimi. Fitness landscape of a biopolymer participating in a multi-step reaction. *J. Theor. Biol.*, 191:377–390, 1998.

[28] G. Woodcock and P. G. Higgs. Population evolution on a mutiplicative single-peak fitness landscape. *J. Theor. Biool.*, 179:61–73, 1996.

[29] L. P. Maia, D. F. Botelho, and J. F. Fontanari. Analytical solution of the evolution dynamics on a multiplicative-fitness landscape. *J. Math. Biol.*, 47:453–456, 2003.

[30] S. Kauffman and S. Levin. Towards a general theory of adaptive walks on rugged landscapes. *J. Theor. Biol.*, 128:11–45, 1987.

[31] S. A. Kauffman and E. D. Weinberger. The N-k model of rugged fitness landscapes and its application to the maturation of the immune response. *J. Theor. Biol.*, 141:211–245, 1989.

[32] M. Kimura. Evolutionary rate at the molecular level. *Nature*, 217:624–626, 1968.

[33] M. Kimura. *The Neutral Theory of Molecular Evolution*. Cambridge University Press, Cambridge, UK, 1983.

[34] T. Ohta. Slightly deleterious mutant substitutions in evolution. *Nature*, 246:96–98, 1973.

[35] T. Ohta and J. H. Gillepie. Development of neutral and nearly neutral theories. *Theor. Pop. Biol.*, 49:128–142, 1996.

[36] T. Ohta. Mechanisms of molecular evolution. *Porc. Natl. Acad. Sci. USA*, 99:16134–16137, 2002.

[37] L. D. Schmidt. *The Engineering of Chemical Reactions.* Oxford University Press, New York, 1998.

[38] A. Novick and L. Szilard. Description of the chemostat. *Science*, 112:715–716, 1950.

[39] V. Bryson and W. Szybalski. Microbial selection. *Science*, 116:45–51, 1952.

[40] P. Sorgeloos, E. Van Outryve, G. Persoone, and A. Cattoir-Reynaerts. New type of turbidostat with intermittent determination of cell density outside the culture vessel. *Applied and Environmental Microbiology*, 31:327–331, 1976.

[41] T. G. Watson. The present status and future prospects of the turbibostat. *J. Appl. Chem. Biotechnol.*, 22:5832–5838, 1971.

[42] Y. Husimi, K. Nishigaki, Y. Kinoshita, and T. Tanaka. Cellstat – A continuous culture system of a bacteriophage for the study of the mutation rate and the selection process at the DNA level. *Rev. Sci. Instrum.*, 53:517–522, 1982.

[43] A. Watts and G. Schwarz, editors. *Evolutionary Biotechnology – From Theory to Experiment*, volume 66/2-3 of *Biophysical Chemistry*, pages 67–284. Elesvier, Amsterdam, 1997.

[44] P. E. Phillipson and P. Schuster. *Modeling by Nonlinear Differential Equations. Dissipative and Conservative Processes*, volume 69 of *World Scientific Series on Nonlinear Science A*. World Scientific, Singapore, 2009.

[45] W. Fontana and P. Schuster. A computer model of evolutionary optimization. *Biophys. Chem.*, 26:123–147, 1987.

[46] M. A. Huynen, P. F. Stadler, and W. Fontana. Smoothness within ruggedness. The role of neutrality in adaptation. *Proc. Natl. Acad. Sci. USA*, 93:397–401, 1996.

[47] W. Fontana and P. Schuster. Continuity in evolution. On the nature of transitions. *Science*, 280:1451–1455, 1998.

[48] P. Schuster, K. Sigmund, and R. Wolff. Dynamical systems under constant organization I. Topological analysis of a family of non-linear differential equations – A model for catalytic hypercycles. *Bull. Math. Biol.*, 40:734–769, 1978.

[49] D. Zwillinger. *Handbook of Differential Equations.* Academic Press, San Diego, CA, third edition, 1998.

[50] R. W. Hamming. *Coding and Information Theory.* Prentice-Hall, Englewood Cliffs, NJ, 2nd edition, 1986.

[51] D. L. Hartl and A. G. Clark. *Principles of Population Genetics.* Sinauer Associates, Sunderland, MA, third edition, 1997.

[52] C. O. Wilke. Quasispecies theory in the context of population genetics. *BMC Evolutionary Bioloy*, 5:e44, 2005.

[53] J. J. Bull, L. Ancel Myers, and M. Lachmann. Quasispecies made simple. *PLoS Comp. Biol.*, 1:450–460, 2005.

[54] J. J. Bull, R. Sanjuan, and C. O. Wilke. Theory for lethal mutagenesis for viruses. *J. Virology*, 81:2930–2939, 2007.

[55] C. J. Thompson and J. L. McBride. On Eigen's theory of the self-organization of matter and the evolution of biological macromolecules. *Math. Biosci.*, 21:127–142, 1974.

[56] B. L. Jones, R. H. Enns, and S. S. Rangnekar. On the theory of selection of coupled macromolecular systems. *Bull. Math. Biol.*, 38:15–28, 1976.

[57] E. Seneta. *Non-negative Matrices and Markov Chains*. Springer-Verlag, New York, second edition, 1981.

[58] D. S. Rumschitzki. Spectral properties of Eigen evolution matrices. *J. Math. Biol.*, 24:667–680, 1987.

[59] P. Schuster and J. Swetina. Stationary mutant distribution and evolutionary optimization. *Bull. Math. Biol.*, 50:635–660, 1988.

[60] M. Nowak and P. Schuster. Error thresholds of replication in finite populations. Mutation frequencies and the onset of Muller's ratchet. *J. Theor. Biol.*, 137:375–395, 1989.

[61] P. Schuster. Mathematical modeling of evolution. Solved and open problems. *Theory in Biosciences*, 130:71–89, 2011.

[62] M. Eigen. Error catastrophe and antiviral strategy. *Proc. Natl. Acad. Sci. USA*, 99:13374–13376, 2002.

[63] I. Leuthäusser. An exact correspondence between Eigen's evolution model and a two-dimensional ising system. *J. Chem. Phys.*, 84:1884–1885, 1986.

[64] I. Leuthäusser. Statistical mechanics of Eigen's evolution model. *J. Stat. Phys.*, 48:343–360, 1987.

[65] P. Tarazona. Error thresholds for molecular quasispecies as phase transitions: From simple landscapes to spin glasses. *Phys. Rev. A*, 45:6038–6050, 1992.

[66] R. V. Solé, S. C. Manrubia, B. Luque, J. Delgado, and J. Bascompte. Phase transitions and complex systems – Simple, nonlinear models capture complex systems at the edge of chaos. *Complexity*, 1(4):13–26, 1996.

[67] N. Wagner, E. Tannenbaum, and G. Ashkenasy. Second order catalytic quasispecies yields discontinuous mean fitness at error threshold. *Phys. Rev. Letters*, 104:188101, 2010.

[68] T. Wiehe. Model dependency of error thresholds: The role of fitness functions and contrasts between the finite and infinite sites models. *Genet. Res. Camb.*, 69:127–136, 1997.

[69] H. W. Engl, C. Flamm, P. Kügler, J. Lu, S. Müller, and P. Schuster. Inverse problems in systems biology. *Inverse Problems*, 25:123014, 2009.

[70] P. Schuster, W. Fontana, P. F. Stadler, and I. L. Hofacker. From sequences to shapes and back: A case study in RNA secondary structures. *Proc. Roy. Soc. Lond. B*, 255:279–284, 1994.

[71] P. Schuster. Prediction of RNA secondary structures: From theory to models and real molecules. *Reports on Progress in Physics*, 69:1419–1477, 2006.

[72] T. Ohta. Mechanisms of molecular evolution. *Phil. Trans. Roy. Soc. London B*, 355:1623–1626, 2000.

[73] W. Fontana and P. Schuster. Shaping space. The possible and the attainable in RNA genotype-phenotype mapping. *J. Theor. Biol.*, 194:491–515, 1998.

[74] J. N. Pitt and A. R. Ferré-D'Amaré. Rapid construction of empirical RNA fitness landscapes. *Science*, 330:376–379, 2010.

[75] Y. Hayashi, T. Auta, H. Toyota, Y. Husimi, I. Urabe, and T. Yomo. Experimental rugged fiteness landscape in protein sequence space. *PLoS One*, 1:e96, 2006.

[76] J. J. Welch and D. Waxman. The *nk* model and population genetics. *J. Theor. Biol.*, 234:329–340, 2005.

[77] C. O. Wilke, J. L. Wang, and C. Ofria. Evolution of digital organisms at high mutation rates leads to survival of the flattest. *Nature*, 412:331–333, 2001.

[78] C. Reidys, P. F. Stadler, and P. Schuster. Generic properties of combinatory maps. Neutral networks of RNA secondary structure. *Bull. Math. Biol.*, 59:339–397, 1997.

[79] W. Grüner, R. Giegerich, D. Strothmann, C. Reidys, J. Weber, I. L. Hofacker, and P. Schuster. Analysis of RNA sequence structure maps by exhaustive enumeration. I. Neutral networks. *Mh. Chem.*, 127:355–374, 1996.

[80] W. Grüner, R. Giegerich, D. Strothmann, C. Reidys, J. Weber, I. L. Hofacker, and P. Schuster. Analysis of RNA sequence structure maps by exhaustive enumeration. II. Structures of neutral networks and shape space covering. *Mh. Chem.*, 127:375–389, 1996.

[81] P. Erdős and A. Rényi. On random graphs. I. *Publicationes Mathematicae*, 6:290–295, 1959.

[82] P. Erdős and A. Rényi. On the evolution of random graphs. *Publications of the Mathematical Institute of the Hungarian Academy of Sciences*, 5:17–61, 1960.

[83] E. Domingo, D. Sabo, T. Tangichi, and C. Weissmann. Nucleotide sequence heterogeneity of an RNA phage population. *Cell*, 13:735–744, 1978.

[84] G. P. Wagner and P. Krall. What is the difference between models of error thresholds and Muller's ratchet. *J. Math. Biol.*, 32:33–44, 1993.

[85] N. Takeuchi and P. Hogeweg. Error-thresholds exist in fitness landscapes with lethal mutants. *BMC Evolutionary Biology*, 7:15:1–11, 2007.

[86] H. Tejero, A. Marín, and F. Moran. Effect of lethality on the extinction and on the error threshold of quasispecies. *J. Theor. Biol.*, 262:733–741, 2010.

[87] J. Summers and S. Litwin. Examining the theory of error catastrophe. *J. Virology*, 80:20–26, 2006.

[88] H. J. Muller. Some genetic aspects of sex. *American Naturalist*, 66:118–138, 1932.

[89] H. J. Muller. The relation of recombination to mutational advance. *Mutat. Res.*, 106:2–9, 1964.

[90] P. Schuster and K. Sigmund. Random selection - A simple model based on linear birth and death processes. *Bull. Math. Biol.*, 46:11–17, 1984.

[91] J. S. McCaskill. A stochastic theory of macromolecular evolution. *Biol. Cybern.*, 50:63–73, 1984.

[92] L. Demetrius, P. Schuster, and K. Sigmund. Polynucleotide evolution and branching processes. *Bull. Math. Biol.*, 47:239–262, 1985.

[93] C. W. Gardiner. *Stochastic Methods. A Handbook for the Natural Sciences and Social Sciences*. Springer Series in Synergetics. Springer-Verlag, Berlin, fourth edition, 2009.

[94] D. T. Gillespie. Stochastic simulation of chemical kinetics. *Annu. Rev. Phys. Chem.*, 58:35–55, 2007.

[95] R. A. Blythe and A. McKane. Stochastic models of evolution in genetice, ecology and linguistics. *J. Stat. Mech.: Theor. Exp.*, page P07018, 2007.

[96] C. Reidys, C. Forst, and P. Schuster. Replication and mutation on neutral networks. *Bull. Math. Biol.*, 63:57–94, 2001.

[97] N. Takeuchi, P. H. Poorthuis, and P. Hogeweg. Phenotypic error threshold – Additivity and epistasis in RNA evolution. *BMC Evolutionary Biology*, 5 (Feb.), 2005. Art.No.9.

[98] B. R. M. Stadler, P. F. Stadler, G. P. Wagner, and W. Fontana. The topology of the possible: Formal spaces underlying patterns of evolutionary change. *J. Theor. Biol.*, 213:241–274, 2001.

[99] P. Schuster. Molecular insight into the evolution of phenotypes. In J. P. Crutchfield and P. Schuster, editors, *Evolutionary Dynamics – Exploring the Interplay of Accident, Selection, Neutrality, and Function*, pages 163–215. Oxford University Press, New York, 2003.

[100] P. G. Higgs and B. Derrida. Stochastic models for species formation in evolving populations. *J. Phys. A: Math. Gen.*, 24:L985–L991, 1991.

[101] B. Derrida and L. Peliti. Evolution in a flat fittness landscape. *Bull. Math. Biol.*, 53:355–382, 1991.

[102] A. Kupczok and P. Dittrich. Determinants of simulated RNA evolution. *J. Theor. Biol.*, 238:726–735, 2006.

[103] C. K. Biebricher. Darwinian selection of self-replicating RNA molecules. In M. K. Hecht, B. Wallace, and G. T. Prance, editors, *Evolutionary Biology, Vol. 16*, pages 1–52. Plenum Publishing Corporation, 1983.

[104] C. K. Biebricher, M. Eigen, and J. William C. Gardiner. Kinetics of RNA replication. *Biochemistry*, 22:2544–2559, 1983.

[105] C. K. Biebricher, M. Eigen, and J. William C. Gardiner. Kinetics of RNA replication: Plus-minus asymmetry and double-strand formation. *Biochemistry*, 23:3186–3194, 1984.

[106] C. K. Biebricher, M. Eigen, and J. William C. Gardiner. Kinetics of RNA replication: Competition and selection among self-replicating RNA species. *Biochemistry*, 24:6550–6560, 1985.

[107] E. J. Richards. Population epigenetics. *Curr. Op. Genet. Develop.*, 18:221–226, 2008.

[108] L. J. Johnson and P. J. Tricker. Epigenomic plasticity within populations: Its evolutionary significance and potential. *Heredity*, 105:113–121, 2010.

[109] C. L. Richards, O. Bossdorf, and M. Pigliucci. What role does heritable epigenetic variation play in phenotypic evolution? *BioScience*, 60:232–237, 2010.

[110] J. L. Geoghegan and H. G. Spencer. Population epigenetic models of selection. *Theor. Pop. Biol.*, 81:232–242, 2012.

[111] S. F. Gilbert. *Developmental Biology*. Sinauer Associates, Sunderland, MA, sixth edition, 2000.

**Contents**

## Notation

| | |
|---|---|
| building blocks and degradation products | $\mathbf{M}, \mathbf{A}, \mathbf{B}, \ldots,$ |
| numbers of particles of $\mathbf{M}, \mathbf{A}, \mathbf{B}, \ldots,$ | $N_{\mathbf{M}}, N_{\mathbf{A}}, N_{\mathbf{B}}, \ldots,$ |
| concentrations of $\mathbf{M}, \mathbf{A}, \mathbf{B}, \ldots,$ | $[\mathbf{M}] = m, [\mathbf{A}] = a, [\mathbf{B}] = b, \ldots,$ |
| | |
| replicating molecular species | $\mathbf{X}_0, \mathbf{X}_1, \mathbf{X}_2, \ldots,$ |
| validity for all individual species of type $i$ | $\mathbf{X}_{(i)},$ |
| numbers of particles of $\mathbf{X}_1, \mathbf{X}_2, \ldots,$ | $N_1, N_2, \ldots,$ |
| concentrations of $\mathbf{X}_1, \mathbf{X}_2, \ldots,$ | $[\mathbf{X}_1] = c_1, [\mathbf{X}_2] = c_2, \ldots,$ |
| relative concentrations of $\mathbf{X}_1, \mathbf{X}_2, \ldots,$ | $[\mathbf{X}_1] = x_1, [\mathbf{X}_2] = x_2, \ldots,$ |
| quasispecies | $\Upsilon = \{\mathbf{X}_0, \mathbf{X}_1, \ldots\},$ |
| fitness values | $f_0, \ f_1, \ \ldots, \ \bar{f}_{-m} = \frac{\sum_{i \neq m} f_i}{1 - x_m} = f,$ |
| master sequence | $\mathbf{X}_0 \text{ or } \mathbf{X}_m, \ f_m = \max\{f_i\} \, \forall \, i,$ |
| notation of classes | $\Gamma_k, \ \mathbf{X}_{(k)} \text{ sequences} \in \Gamma_k,$ |
| partial sums of relative concentrations | $y_k = \sum_{i \in \Gamma_k} x_i,$ |
| "realistic" landscape | $\mathcal{L}(\lambda, d, s; \nu, f_0, f),$ |
| | |
| flow rate, influx concentration in the CSTR | $r, a_0,$ |
| rate parameters | $d_i, k_i, f_i, \ldots \quad i = 1, 2, \ldots,$ |
| global regulation flux | $\phi(t),$ |
| | |
| chain length of polynucleotides | $\nu,$ |
| superiority of the master sequence $\mathbf{X}_m$ | $\sigma_m = f_m \, / \, \bar{f}_{-m},$ |
| population entropy | $S = \sum_i x_i \ln x_i,$ |
| | |
| fitness landscape | $\mathcal{L}(\nu, f_0, f, \lambda, d, s),$ |
| degree of neutrality | $\lambda,$ |
| width of random scatter | $0 \leq d \leq 1,$ |
| seeds for random number generator | $s.$ |

## Color code

The color code for error classes is used in figures showing concentration plots on the binary sequence space with $\nu = 10$.

| class | color | sequence | line | # seqeunces |
|:---:|:---:|:---:|:---:|:---:|
| 0 | black | ● | —— | 1 |
| 1 | red | ● | —— | 10 |
| 2 | yellow | ● | —— | 45 |
| 3 | green | ● | —— | 120 |
| 4 | cyan | ● | —— | 210 |
| 5 | blue | ● | —— | 252 |
| 6 | magenta | ● | —— | 210 |
| 7 | chartreuse | ● | —— | 120 |
| 8 | yellow | ● | —— | 45 |
| 9 | red | ● | —— | 10 |
| 10 | black | ● | —— | 1 |

# Index