

The effects of population structure and the genotype-phenotype map on evolutionary dynamics

Szöllősi Gergely János

Ph.D. Thesis

Advisor: **Dr. Imre Derényi**, D.Sc., associate professor

Department of Biological Physics
Eötvös Loránd University, Budapest



Graduate School of Physics

Doctoral Program for Statistical Physics, Biological Physics, and Quantum
Systems

Head of School: Dr. Horváth Zalán, regular member of HAS

Head of Program: Dr. Kürti Jenő, D.Sc., professor

2009

I think that science suggests to us (tentatively of course) a picture of a Universe that is inventive or even creative; of a Universe in which *new things* emerge on *new levels*.

There is on the first level the emergence of heavy atomic nuclei in the center of big stars, and, on a higher level, the evidence for the emergence somewhere in space of organic molecules.

On the next level there is the emergence of life. Even if the emergence of life should one day become reproducible in the laboratory, life creates something that is utterly new in the Universe: the peculiar activity of organisms; especially the often purposeful actions of animals. All organisms are constant problem solvers; even though they are not conscious of most of the problems they are trying to solve.

On the next level the great step is the emergence of conscious states. With the distinction between conscious states and unconscious states, again something of the utterly new and of the greatest importance enters the Universe. It is a new world a world of conscious experience.

On the next level this is followed by the emergence of the products of the human mind, such as the works of art; and also the works of science; especially scientific theories.

I think that scientists, however sceptical, are bound to admit that the Universe, or Nature, or whatever we may call it, is creative. For it has produced creative men: it has produced Shakespeare and Michelangelo and Mozart, and thus indirectly their works. It has produced Darwin, and so created the theory of natural selection. Natural selection has destroyed the proof for the miraculous specific intervention of the Creator. But it has left us with the marvel of the creativeness of the Universe, of life, and of the human mind. Although science has nothing to say about a personal Creator, the fact of the emergence of novelty, and of creativity, can hardly be denied.

Sir Karl Raimund Popper, The first Darwin Lecture at Darwin College, Cambridge University, Recorded by the BBC on the 8th of November 1977

Contents

1	Introduction	1
1.1	Logistic growth	2
1.2	Drift and mutation	3
1.2.1	The neutral theory of molecular evolution	7
2	The effects of population structure	10
2.1	Introduction	10
2.2	Evolutionary games	12
2.3	Finite population size	13
2.4	A minimal model of population structure	16
2.5	Cooperation in populations with hierarchical levels of mixing	22
2.6	Three strategies	25
2.6.1	The RPS game	25
2.6.2	The repeated prisoner's dilemma game	27
2.7	The maintenance of natural genetic transformation in bacteria	30
2.7.1	The dynamics of genome organization in bacteria	32
2.7.2	Self-Consistent Migration	36
2.7.3	A Phase Diagram for the survival NGT	36
2.7.4	The effects of NGT on genome organization	38
2.8	Discussion	41
3	The mapping between genotype and phenotype	46
3.1	Introduction	46
3.1.1	Genetic robustness	48
3.2	The neutral evolution of genetic robustness	50
3.2.1	The infinite population limit	52
3.2.2	Finite populations	54

3.3	The effects of recombination	55
3.3.1	The infinite population limit	58
3.3.2	Finite populations	59
3.4	Congruent evolution of robustness in microRNA	61
3.4.1	MicroRNA stem-loops as a molecular phenotype	61
3.4.2	Robustness of microRNA sequences	64
3.4.3	The case for congruent evolution of robustness	70
3.4.4	The temperature of mutations	70
3.5	Discussion	74
4	Methods	76
4.1	Genome organization dynamics	76
4.2	Random neutral networks	78
4.3	RNA secondary structure	79
4.3.1	The Vienna package	79
4.3.2	Scaled down microRNA neutral network	80
4.4	MicroRNA sequences	80
4.4.1	Sampling of sequences with a given MFE structure	83
4.4.2	Estimating the effective temperature of mutations	83

Acknowledgements

I would like to thank first of all my advisor, Imre Derényi, who afforded me the freedom to pursue my own directions; directions along which I would most certainly have gotten lost, without the guiding light of his insight, and the invaluable constraint of his patiently argued advice.

I would like to thank Tibor Vellai, who introduced us to the fascinating world of gene swapping bacteria.

I would like to thank Michael Lässig, Ville Mustonen and Stephan Schiffels, work and discussion with whom, while not presented here, has contributed tremendously.

I would like to thank my wife, Ági, for surrendering me to the laptop screen as often as was necessary, but not more. My parents for their relentless support and love, and my grandmother, who tried as hard as she could to teach me how to write properly.

Chapter 1

Introduction

"We will now discuss in a little more detail the Struggle for Existence."

Charles Darwin, *The Origin of Species by Means of Natural Selection, or the Preservation of Favoured Races in the Struggle for Life*. 1859, John Murray, London UK

“Survival of the fittest”, a phrase coined by Herbert Spencer [1] an English philosopher and contemporary of Charles Robert Darwin and Alfred Russel Wallace is merely a metaphor not a scientific description. It has, however, been a very influential metaphor, despite being a very incomplete picture of biological evolution. Evolution is not driven by natural selection (i.e., survival of the fittest) alone, all populations are influenced by the nonadaptive forces of mutation, recombination and random genetic drift, and constrained by population structure and the phenotypic response to genetic variation. A rich body of theoretical and experimental work has demonstrated that these are not mere embellishments around a primary axis of selection, but are quite the opposite they dictate what selection can and cannot do [2, 3].

But what then is evolution? Darwin in the *Origin of Species* [4] laid the foundations of modern evolutionary biology on two central ideas: (i) all species are related to one another through a history of common descent and (ii) the exquisite match between a species and its environment is explained by natural selection, a process in which individuals with beneficial mutations leave more offspring. These two ideas concern very different time-scales and have historically been the subject of different disciplines: (i) a group of scientists (mostly in natural history museums) developed the field of phylogenetics to infer the evolutionary history of life while (ii) another group of scientists, population geneticists, developed mathematical models to describe the behavior of mutations in populations. This dichotomy

of long and short time scales is a fundamental feature of evolution; a fundamental problem of evolutionary theories trying to understand the evolution and maintenance of sex [5] or altruism [6]; a fundamental theme of the evolution of life on Earth, as reflected in the punctuated tempo of the progression of form in the paleontological record [7] or the maintenance of extant biodiversity [8].

It is this dichotomy of long and short time-scale in evolution that we explore here from a statistical physics perspective. By a statistical physics perspective we mean the bottom up program of attempting to discern statistical laws, and understand emergent complexity, from the behaviour of large numbers (a population) of simple (or at least simplified) individuals subject to shared constraints. The program of trying to understand long time scale evolutionary phenomena as the emergent result of evolutionary dynamics taking place in structured populations and under the constraints of the mapping between genotype and phenotype.

To start down this road we first have to describe the evolutionary dynamics of unstructured populations composed of the simplest individuals with the least degree of interaction. These minimal models will be the building blocks that form the foundations of models taking into account more complex scenarios.

1.1 Logistic growth

The realization that a population free of constraints, a population of individuals reproducing in an environment with an inexhaustible supply of resources, would grow exponentially is commonly attributed to the Reverend Thomas Robert Malthus. Malthus was concerned that the “The power of population is indefinitely greater than the power in the earth to produce subsistence for man” [9], and advocated measures he perceived as promoting longer-term stability of the economy above short-term expediency. The first mathematical model of population growth limited by available resources, the logistic equation, was developed by Verhulst [10]. Let us, following in the steps of Verhulst, consider a homogeneous population of identical individuals that reproduce by using resources in their environment to construct perfect copies of themselves. A population of size $n(t)$ with a rate of growth r obeys the differential equation:

$$\frac{dn}{dt} = rn \tag{1.1}$$

A larger population will inevitably lead to a reduction in resources available for reproduction. Reduction in resources will lead in turn to a reduced growth rate. In the simplest case,

where the rate of growth decreases linearly with population size, we obtain the logistic equation:

$$\frac{dn}{dt} = rn\left(1 - \frac{n}{N}\right), \quad (1.2)$$

where N is called the carrying capacity (of the environment with respect to a given population) and is equal to the asymptotic value of the population size. We will primarily be interested in the long time limit where population size can, to good approximation, be considered constant.

Provided constant population size we can turn to the situation where we have several different types of individuals. Let us e.g. assume that the population is divided into d types $i \in \{1, \dots, d\}$ with frequencies $\{x_i\}$. In general the number of offspring (fitness) of an individual is linked to the type and frequency of its competitors. Denoting the fitness $\pi_i(x_1, \dots, x_d)$ of type i , in the limit of very large populations, we may express the relative rate of increase of type \dot{x}_i/x_i of type i as the difference between the fitness $\pi_i(x_1, \dots, x_d)$ and the average fitness $\bar{\pi}(x_1, \dots, x_d) = \sum_i x_i \pi_i(x_1, \dots, x_d)$ obtaining

$$\frac{\dot{x}_i}{x_i} = \text{fitness of type } i - \text{average fitness}. \quad (1.3)$$

This leads to the replicator equation:

$$\frac{\dot{x}_i}{x_i} = \pi_i(x_1, \dots, x_d) - \bar{\pi}(x_1, \dots, x_d). \quad (1.4)$$

1.2 Drift and mutation

The replicator equation is a very general description of evolutionary dynamics, it is, however, far from complete. In natural populations individuals are not immutable and finite population size inevitably introduces stochasticity. These two considerations, mutation and finite population size, are traditionally the subject of the field of population genetics.

Population genetics is concerned with the genetic basis of evolution, it views evolution as the change of the frequencies of genotypes through time. The basic model in population genetics, which we will use to introduce the effects of finite population size, closely following Ref. [11], is the Wright-Fisher model. It assumes fixed population size and non-overlapping generations, i.e., each individual is replaced in every generation. In the most simple case a population of constant size N is considered, where all individuals have equal fitness, but different genotypes¹. Each generation we pick N individuals randomly and

¹The Wright-Fisher model is traditionally developed in the context of diploid individuals and motivated

with replacement from the population to form the next generation. We will describe the state of the population by \mathcal{G} defined as the probability that two individuals differ by origin. In the absence of variation in the population $\mathcal{G} = 1$, while in contrast if all individuals have different genotypes $\mathcal{G} = 0$. The value after one generation can be expressed using the current value of \mathcal{G} as

$$\mathcal{G}' = \frac{1}{N} + \left(1 - \frac{1}{N}\right) \mathcal{G} \quad (1.5)$$

That is two individuals are related after one generation either because, with probability $\frac{1}{N}$, they are descendants of the same individual from the previous generation, or because, with probability $(1 - \frac{1}{N})\mathcal{G}$, they are the descendants of two distinct individuals from the previous generation (this occurs with probability $(1 - \frac{1}{N})$) who were of common descent (this occurs with probability \mathcal{G}).

The time course of \mathcal{G} is most readily expressed in terms of $\mathcal{H} = 1 - \mathcal{G}$, the probability of two individuals having distinct origins, using which Eq. (1.5) can be written as

$$\mathcal{H}' = 1 - \mathcal{G}' = \left(1 - \frac{1}{N}\right) \mathcal{H}, \quad (1.6)$$

which leads to

$$\mathcal{H}' - \mathcal{H} = -\frac{1}{N} \mathcal{H}. \quad (1.7)$$

That is the probability that two individuals are of different descent decays at a rate $1/N$ per generation as a result of sampling effects inherent to finite populations – as a result of *genetic drift*. By induction we can readily see that the decay is geometric, i.e., denoting the probability of two individuals being of distinct descent in generation t as \mathcal{H}_t we have

$$\mathcal{H}_t = \mathcal{H}_0 \left(1 - \frac{1}{N}\right)^t. \quad (1.8)$$

Genetic variation removed by drift is restored by mutation. Provided individuals suffer mutations at some rate μ per generation, variation is introduced into the population at a rate $N\mu$ if we assume each mutation to be unique². Taking mutations into account leads to the modified form of Eq. (1.5):

$$\mathcal{G}' = (1 - \mu)^2 \left[\frac{1}{N} + \left(1 - \frac{1}{N}\right) \mathcal{G} \right] \quad (1.9)$$

by an underlying picture of random mating, with new generations sampling their genes from a very large haplotype pool – from a formal point of view this only differs in replacing N by $2N$ throughout.

²This assumption is appropriate in the context of the very large number of alleles possible at the molecular level.

The term $(1 - \mu)^2$ being the probability that both individuals escape suffering a mutation. As current estimates of mutation rates give 5.0×10^{-10} /bp/generation and 5.4×10^{-8} /bp/generation among prokaryotes and vertebrates respectively [2], and population sizes are large (often well in excess of 10^4), we can approximate the above as:

$$\mathcal{G}' \approx (1 - 2\mu) \left[\frac{1}{N} + \left(1 - \frac{1}{N}\right) \mathcal{G} \right] \approx \frac{1}{N} + \left(1 - \frac{1}{N}\right) \mathcal{G} - 2\mu\mathcal{G}. \quad (1.10)$$

Introducing again $\mathcal{H} = 1 - \mathcal{G}$ we have

$$\mathcal{H}' \approx \left(1 - \frac{1}{N}\right) \mathcal{H} + 2\mu(1 - \mathcal{H}), \quad (1.11)$$

leading to

$$\mathcal{H}' - \mathcal{H} = -\frac{1}{N}\mathcal{H} + 2\mu(1 - \mathcal{H}), \quad (1.12)$$

Examining the first term of the right hand side of the above we can readily identify the right hand side of Eq. (1.7) and organize the terms as

$$\Delta\mathcal{H} = \Delta_N\mathcal{H} + \Delta_u\mathcal{H} \quad (1.13)$$

where $\Delta\mathcal{H} = \mathcal{H}' - \mathcal{H}$ and the terms corresponding to drift and mutation are:

$$\Delta_N\mathcal{H} = -\frac{1}{N}\mathcal{H} \quad \text{and} \quad \Delta_u\mathcal{H} = 2\mu(1 - \mathcal{H}). \quad (1.14)$$

The change due to drift is always negative (drift decreases genetic variation), whereas the change due to mutation is always positive (mutation increases genetic variation). These two forces are in equilibrium when the probability that two individuals in the population have identical genotypes is:

$$\mathcal{G}_{\text{eq.}} \approx \frac{1}{1 + 2N\mu}. \quad (1.15)$$

We can go further and derive the rate u at which (neutral) substitutions occur during the evolution of the population. Substitutions occur because new mutations enter the population each generation and genetic drift causes a small fraction of them to fix each generation. Of the μN new mutations that enter the population any given generation a fraction $\frac{1}{N}$ will fix, hence the (neutral) substitution rate – independent of N – is equal to the (neutral)

mutation rate³:

$$u = \mu. \quad (1.16)$$

Each allele is descended from an allele in the previous generation. Following this hierarchy of descent we may trace back the lineage of a pair of alleles in or population to the allele that was their most recent common ancestor. The probability that any two alleles shared a common ancestor in the previous generation is $1/N$. Looking back across generations the probability that two alleles first had a common ancestor t generations ago, i.e., they *coalesce* after t generations, is just the probability that they do not coalesce for $t - 1$ generations $(1 - 1/N)$, multiplied by the probability that they do $1/N$. The probability of coalescence at time in generation t is

$$P_C(t_C = t) = \left(1 - \frac{1}{N}\right)^{(t-1)} \left(\frac{1}{N}\right). \quad (1.17)$$

For sufficiently large N the above is well approximated by

$$P_C(t_C = t) = \frac{1}{N} e^{-t/N}, \quad (1.18)$$

which implies that the mean and the standard deviation of the time to coalescence t_C are both N .

To incorporate the effects of natural selection let us consider the most minimal genotype, a single genomic locus with two alleles with different fitness, and mutation causing each allele to change into the other. As in the neutral case each generation we pick N individuals randomly and *with replacement* from the population to form the next generation, this time, however, as individuals have different fitness, we pick each individual with a probability proportional to its fitness. In the limit of infinite population size and the absence of mutation the population dynamics is described according to the replicator equation (1.4)

$$\dot{x}/x = 1 + s - [x(1 + s) + (1 - x)1] \Rightarrow \dot{x} = sx(1 - x), \quad (1.19)$$

where x is the frequency of the fitter allele $1 + s$ is the fitness of the fitter allele (and 1 is the fitness of the less fit). Combining Eqs. (1.14) and (1.19) (keeping in mind that $\mathcal{H} = x(1 - x)$ in our case) we can motivate the Kimura equation, which we will not derive here. The Kimura equation is the Fokker-Planck equation describing time evolution of the

³This is a very general result that remains true even in the presence of selection at nearby (linked) sites in the genome.

probability $p(x, t)$ that a fraction x of the population at time t has the fitter allele [12]:

$$\frac{\partial}{\partial t} p(x, t) = \left[\frac{1}{2N} \frac{\partial^2}{\partial x^2} x(1-x) - \mu \frac{\partial}{\partial x} (1-2x) - s \frac{\partial}{\partial x} x(1-x) \right] p(x, t). \quad (1.20)$$

It has the stationary solution:

$$p_{\text{eq.}}(x | 2N\mu, 2Ns) = \frac{1}{Z_{\text{eq.}}} e^{2Nsx} (x(1-x))^{-1+2N\mu}, \quad (1.21)$$

where the normalization factor $Z_{\text{eq.}}$ can be expressed analytically. We can see that the dynamics is described by the scaled selection strength $2Ns$ and mutation rate $2N\mu$. Provided the product of the population size and mutation rate is small⁴ $\mu N \ll 1$ the rate of fixation of both the fitter u_+ and the less fit allele u_- can be derived from the conditional stationary probability densities of the trajectories that eventually reach $x = 0$ and $x = 1$ starting from $x = 1$ and $x = 0$, respectively, and give the classic Kimura-Ohta substitution rates:

$$u_+ = \mu \frac{2Ns}{1 - e^{-2Ns}} \quad \text{and} \quad u_- = \mu \frac{-2Ns}{1 - e^{2Ns}}. \quad (1.22)$$

The probability of being fixed in the fitter λ_+ or the less fit λ_- state readily follows as

$$\lambda_+ = \frac{u_+}{u_+ + u_-} = \frac{e^{2N(1+s)}}{e^{2N1} + e^{2N(1+s)}} \quad \text{and} \quad \lambda_- = \frac{u_-}{u_+ + u_-} = \frac{e^{2N1}}{e^{2N1} + e^{2Ns}}. \quad (1.23)$$

This suggest an analogy with statistical physics, with enthalpy corresponding to the negative of the (additive) fitness and the temperature scale kT corresponding to $1/2N$. This analogy, describing the equilibrium distribution of fixed genotypes, can in fact be carried further to describe more complex genotypes with arbitrary fitness schemes, provided no more than two alleles exist simultaneously in the population [13].

1.2.1 The neutral theory of molecular evolution

By the early 1960's there was growing evidence that the number of amino acid differences (substitutions) between different pairs of orthologous⁵ genes from different pairs of

⁴This condition is appropriate for multicellular eukaryotes with small effective population size and high recombination rates, but not necessarily for prokaryotes or unicellular eukaryotes, cf. section 3.1.1.

⁵Genes of common descent are in general homologous. If they were separated by a speciation event, i.e., both are derived from a single ancestral gene in the last common ancestor of the two extant species in which they now reside, they are orthologous. If, however, they were separated by a gene duplication, i.e., both are derived from the duplication of an ancestral gene in some lineage leading up to the last common ancestor of the two extant species in which they now reside, they are paralogous.

mammalian species are roughly proportional to the time since they had diverged from one another, as inferred from the fossil record. This led Zuckerkandl and Pauling [14] to propose the “molecular clock hypothesis” according to which the rate of substitutions is constant over time. This, however, posed the problem of proposing an evolutionary mechanism that could produce a constant rate of substitution. Coupled with the observation, derived from electrophoresis studies of proteins, that natural populations show unanticipated levels of variability, to a degree that is difficult to reconcile with strong selection the scene was set for Kimura [15] and King and Jukes [16] to propose that the overwhelming majority of divergence between and variability within species is neutral.

It is important to understand that this does not require that the overwhelming majority of mutations be neutral. Let us consider a gene as a sequence of L sites with two alleles, each described by its own selection strength s_i with contributions from different alleles contributing additively and a common mutation rate μ such that the condition $\mu NL \ll 1$ holds. From Eqs. (1.22) and (1.23) we can derive the equilibrium rate at which substitution event occurs at a site with selection strength $2Ns$:

$$u_{\text{tot.}}(2Ns) = \lambda_-(2Ns)u_+(2Ns) + \lambda_+(2Ns)u_-(2Ns) = \mu \frac{2Ns}{\sinh(2Ns)}, \quad (1.24)$$

that is the substitution rate drops exponentially as a function of increasing selection strength: $u_{\text{tot.}}(2Ns) \propto \mathcal{O}(\exp(-2Ns))$. High fitness sites (low energy states) are constrained by selection to remain in the fitter state and the equilibrium fitness scale over which genetic drift (thermal noise) dominates is determined by the inverse of the population size. The same can be shown for the probability of sampling two distinct alleles from a population using Eq. (1.21). This means that even if the individual selection coefficients $\{s_i\}$ are drawn from a distribution which implies that a significant fraction of sites have high $2Ns \gg 1$ selection coefficients, the overwhelming majority of substitutions observed will none the less have selection coefficients $2Ns_i \leq 1$ and will occur at a constant rate approximately equal to the mutation rate μ .

It is of course clear that the central assumption of mutation-selection-drift equilibrium cannot be true on arbitrarily long time scales. Fitness is in general time-dependent, reflecting environmental and coevolutionary change. As organisms adapt to changing natural-selection by evolving new phenotypes, new genomic opportunities open for beneficial mutations [17, 18]. It is not clear, however, what fraction of substitutions separating a given pair of species is due to drift. In the last decade the sequencing of entire genomes of multiple individuals within and across species has become the new standard for genomic data

sets [19]. Population genetic studies exploiting the flood of new data on genome wide polymorphism and substitution data using of statistical tests of increased sensitivity, have lead to an emerging consensus that the extent of adaptive evolution at the molecular level, at least for *Drosophila*, has been underestimated [20, 21, 22]. It is now widely accepted that a substantial fraction of the divergence between *Drosophila* species is driven by adaptive substitutions. How wide-spread this reevaluation turns out to be remains to be seen.

Chapter 2

The effects of population structure

“It may help to classify the various theories; first according to the time scale on which selection is supposed to act, and then according to the ‘unit of selection’ – population, individual, or gene.”

John Maynard Smith, *The Evolution of sex*. 1977, Cambridge University Press, Cambridge UK

2.1 Introduction

The long term advantages of both sexual reproduction (for e.g. bringing together favorable mutations), and altruism (for the community) are rather straightforward. Concerns over the short term stability of both these processes, however, have made both of these traits the focus of intensive theoretical investigation. Let us briefly consider each in turn.

A readily apparent problem with sex is the so called two fold cost of sex: a female that reproduces sexually only passes on half of its genes to any given offspring, compared to an asexually reproducing competitor that passes on all of them. More generally, reproducing sexually, that is, by exchanging genetic information with other members of the population, is costly due to the risk of incorporating deleterious alleles and the risk of breaking up favorable combinations of alleles. Moreover, unless adaptive mutations are much more frequent than thought, it is hard to understand what the evolutionary advantage is that prevents sexual recombination from being lost as a result of chance. The long term disadvantage of asexuality, or more generally of a reduced recombination rate are, however, quite clear from the typically short lifespan of asexual lineages or the poor state of affairs¹ on chromosomes

¹Ancient Y-chromosomes of various organisms contain few active genes and an abundance of repetitive DNA [23]. The human Y chromosome 24 Mb male-specific portion contains only 78 protein-coding genes,

suffering significantly reduced recombination rates, e.g. the Y chromosomes of humans or chromosome 4 of *Drosophila melanogaster*.

In the case of altruistic behaviour we must consider individuals that have a choice of distributing some limited resource between their own fitness and the fitness of the group. Individuals that distribute their resources solely toward their own fitness will outcompete those altruistic individuals that also contribute to group fitness even in cases where the benefit to the population is significantly greater than the cost to the individual. Selection will favour selfish individuals, the fraction of altruistic individuals will decline along with the mean fitness of the group.

As reflected in the replicator equation, the dynamics of Darwinian evolution is intrinsically frequency dependent. The fitness of individuals is tightly coupled to the type and number of competitors. Evolutionary dynamics acts, however, on populations, not individuals and as a consequence depends on not only population composition, but also population size and structure. Evolutionary game theory came about as the result of the realization that frequency dependent fitness introduces strategic aspects to evolution [24, 25, 26]. More recently the investigation of the evolutionary dynamics of structured populations, where individuals only compete with some subset of the population, e.g. their neighbors in space or more generally in some graph [27, 28], has lead to the recognition that the success of different strategies can be greatly influenced by the topology of interactions within the population. Fundamental differences were found – compared to well-mixed populations, where individuals interact with randomly chosen partners – in models that describe the evolution of cooperation (variants of the prisoner’s dilemma game [27, 29, 6, 30, 31]) or deal with the maintenance of biodiversity in the context of competitive cycles (variants of the rock-paper-scissors game [26, 32, 33, 34, 35, 36]).

In this chapter we first introduce evolutionary games and derive the stochastic replicator equations describing a finite population. Subsequently we develop a minimal model of population structure, described by two distinct hierarchical levels of interaction, by deriving the dynamics governing its evolution starting from fundamental individual level stochastic processes. Using our model we demonstrate the existence of a continuous transition leading to the dominance of cooperation as the benefit to cost ratio becomes smaller than the local population size. After examining the topologically more complex behaviour of games with three strategies we turn to Natural Genetic Transformation in bacteria. We demonstrate that the combination of minimal population structure, developed previously, and fluctuating

less than a third of the genomic average. The cause of the observed degradation is believed to be transportable element insertions and large deletions.

resource availability provides NGT – the bacterial equivalent of sex – with a short term advantage as a mechanism to reload locally lost, intermittently selected genes from the collective gene pool of the species through DNA uptake from migrants.

2.2 Evolutionary games

"If we repeated Noah's experiment – starting an ecosystem with one couple of each species – we would certainly not expect a restoration of the old régime. Numbers matter."

Josef Hofbauer and Karl Sigmund, *Evolutionary Games and Population Dynamics* 1998, Cambridge University Press, Cambridge, UK

The field of Game Theory arose around the end of the second World War with the work of Von Neumann and Morgenstern [37], and Nash [38]. Most of the early focus was on economic applications (and occasionally, in a more clandestine fashion, on analysing Cold War strategies). Evolutionary thinking in game theory – as well as game theoretical thinking in evolution – owns its origins to the seminal work of Maynard Smith and Price [25] who applied game theoretical concepts to the animal conflicts.

The question that motivated Maynard Smith and Price in applying game theory to evolutionary biology was the following: Why is it that animal conflicts so seldom involve full use of the available destructive power? After all stags would clearly have the ability and motivation to engage in lethal combat over females (a resource of obvious value), none the less the outcome of clashes between males seldom results in serious injury. Maynard-Smith and Price originally considered a series of strategies (hawk, mouse, bully, retaliator and prober-retaliator), but subsequent studies tended to focus on just two strategies hawk and mouse².

Hawks and Mice

Let us consider two types of individuals “hawks” and “mice”. Individuals interact in the context of contests over some resource, e.g. a morsel of food, the boundary line between territories or a potential mate. The prize for winning the contest corresponds to a morsel of fitness b , but also exposes the contender to the possibility of a fitness reducing injury c .

²The mouse strategy was later renamed dove, most probably after the metaphorical bird of peace rather than the actual animal, which is notoriously aggressive. Both names are inappropriate in the sense that individuals with different strategies are supposed to be of the same species.

If two mice meet they may attempt to intimidate each other, but eventually one of them will retreat. The winner obtains b , the loser gets nothing, hence the average increase in fitness for the participants if two mice meet is $b/2$. A mice meeting a hawk immediately retreats, obtaining nothing and always leaving the hawk b . In the unfortunate event that two hawks meet they will both escalate the fight until one of the contenders is knocked out, receiving a fitness penalty of c while the other obtains b , the average change in fitness for the participants is $(b - c)/2$. The above rules may be summarized in a *pay-off matrix*:

	meeting a hawk	meeting a mouse	
hawk gets	$\frac{b-c}{2}$	b	(2.1)
mouse gets	0	$\frac{b}{2}$	

In a population consisting mostly of mice, hawks will spread as they are likely to meet mostly mice. Conversely, in a population of mostly hawks a mice will on the average be better off provided $b < c$, as mice avoid all fights keeping their fitness unchanged, while that of hawks declines due to frequent fights with fellow hawks. Denoting frequency of hawks by x the change in fitness of hawks is $x(b - c)/2 + (1 - x)b$, while the change in fitness is $(1 - x)b/2 + x0$. Equality will hold if and only if $x = b/c$. If $x < b/c$ the frequency of hawks will increase, whereas if $x > b/c$ it will decrease.

We see that if the cost of injury c is large hawks, and as a consequence conflicts resulting in serious injury, will be rare. This is consistent with the observation that “gloved fist” conflicts prevail among heavily armed species (e.g. deer). Apparently harmless animals, such as doves, on the other hand seem not to have developed traits for avoiding escalation. While under natural conditions doves cannot inflict serious injuries upon each other, they will fight to the death when confined in a cage [26].

2.3 Finite population size

Traditionally, frequency dependent evolutionary dynamics is described by deterministic replicator dynamics (Eq. (1.4)) under the implicit assumption of infinite population size. In natural population of finite size stochasticity resulting from finite population size must be taken into account. Following Traulsen et al. [39, 40] we derive the stochastic replicator dynamics describing a finite population of competing individuals. We develop the model from the level of the individual by considering the Moran process as the microscopic process of population dynamics. The Moran process [41] consist of three basic processes:

1. *selection*: an individual is randomly selected for reproduction with a probability proportional to its fitness
2. *reproduction*: the selected individual produces a single identical offspring
3. *replacement*: the offspring replaces a randomly selected individual.

The Moran process is closely related to the Wright-Fisher process and allows the derivation of fixation probabilities of mutations. Considering overlapping generations, in place of the discrete nonoverlapping generations considered in the Wright-Fisher model, allows us to more easily connect with a continuous time replicator equation.

Let us first consider a population of two types of individuals which we will call **A** and **B**. The fitness (or payoff) of an individual of either type depends on the type of individuals with which it comes into interaction according to the payoff matrix:

$$\begin{array}{rcc}
 & \text{with } \mathbf{A} & \text{with } \mathbf{B} \\
 \mathbf{A} \text{ interacts} & a & b \\
 \mathbf{B} \text{ interacts} & c & d.
 \end{array} \tag{2.2}$$

If each individual can be considered to interact with a random sample of the population, the average payoff of individuals of type **A** and type **B** will be determined by the fraction of both types in the population. Excluding self-interactions, this implies the expressions

$$\pi_{\mathbf{A}} = \pi_{\text{base}} + \frac{a(n-1) + b(N-n)}{N-1} \quad \text{and} \quad \pi_{\mathbf{B}} = \pi_{\text{base}} + \frac{c(n-1) + d(N-n-1)}{N-1}, \tag{2.3}$$

for the fitness of type **A** and **B** respectively, where n is the number of individuals of type **A** and N is the population size and π_{base} is some universal baseline fitness. The probability that the number of individuals of type **A** increases from n to $n+1$ is given by the product of the probability of an individual of type **A** being born and an individual of type **B** being chosen for replacement

$$T^+(n) = \frac{\pi_{\mathbf{A}} n (N-n)}{\bar{\pi} N}, \tag{2.4}$$

while the probability that the number of individuals of type **A** decreases from n to $n-1$ can be expressed as the product of the probability of an individual of type **B** being born and an individual of type **A** being chosen for replacement

$$T^-(n) = \frac{\pi_{\mathbf{B}} (N-n) n}{\bar{\pi} N}, \tag{2.5}$$

where $\bar{\pi} = \pi_{\mathbf{A}}n/N + \pi_{\mathbf{A}}(N - n)/N$ is the average fitness in the population.

The stochastic process described by the transition probabilities (2.4) and (2.5) can be formulated in terms of the master equation [42]

$$\begin{aligned} P^{\tau+1}(n) - P^{\tau}(n) &= P^{\tau}(n-1)T^{+}(n-1) - P^{\tau}(n)T^{-}(n) \\ &+ P^{\tau}(n+1)T^{-}(n+1) - P^{\tau}(n)T^{+}(n), \end{aligned} \quad (2.6)$$

where $P^{\tau}(n)$ is the probability that the system is in state n at time τ . Introducing the notation $x = i/N$, $t = \tau/N$, and the probability density $\rho(x, t) = NP^{\tau}(n)$ yields

$$\begin{aligned} \rho(x, t + N^{-1}) - \rho(x, t) &= \rho(x - N^{-1}, t)T^{+}(x - N^{-1}) \\ &+ \rho(x + N^{-1}, t)T^{-}(x + N^{-1}) \\ &- \rho(x, t)T^{-}(x) - \rho(x, t)T^{+}(x). \end{aligned} \quad (2.7)$$

For $N \gg 1$ (i.e., large but finite N) we may expand in a Taylor series the transition probabilities and the probability densities around x and T (this is called the Kramers-Moyal expansion [42]) and neglecting higher order terms in $1/N$, we get the Fokker-Planck equation that approximates master equation:

$$\frac{d}{dt} = -\frac{d}{dx} [a(x)\rho(x, t)] + \frac{1}{2} \frac{d^2}{dx^2} [c^2(x, t)\rho(x, t)] \quad (2.8)$$

where

$$a(x) = T^{+}(x) - T^{-}(x) \quad (2.9)$$

and

$$c(x) = \sqrt{(1/N)[T^{+}(x) + T^{-}(x)]}. \quad (2.10)$$

Subsequent steps of the Moran process are independent, as a consequence the internal noise is not correlated in time. This implies that the Itô calculus [42] may be applied to derive the Langevin equation:

$$\dot{x} = a(x) + c(x)\xi, \quad (2.11)$$

where ξ is uncorrelated Gaussian noise and $c(x) = 0$ for $x = 0$ and $x = 1$.

For $N \rightarrow \infty$ the diffusion term $c(x)$ vanishes $\mathcal{O}(1/\sqrt{N})$, and a deterministic dynamics is recovered:

$$\dot{x} = x [\pi_{\mathbf{A}}(x) - \bar{\pi}(x)] \frac{1}{\bar{\pi}(x)}. \quad (2.12)$$

That is the Moran process leads to an adjusted replicator equation with a $\frac{1}{\bar{\pi}(x)}$ time rescaling factor dependent on the composition of the population. As we will show below, this term, first derived by Traulsen et al., has important implications for the evolutionary stability of cooperation in structured populations. Traulsen et al. extended the above derivation to an arbitrary number of types [40]. We will use these results below in our treatment of the coevolutionary dynamics of games on minimally structured populations.

2.4 A minimal model of population structure

In general, in order to investigate the coevolutionary dynamics of games on structured populations the full set of connections between a potentially very large number of individuals must be specified. This is only possible by reducing the number of degrees of freedom considered, either through postulating a highly symmetrical structure (such as lattices [27, 30, 43, 44, 45]), or a fundamentally random connection structure (such as some random graph ensemble [46, 47]). The question of how one goes about the task of reducing the number of degrees of freedom – of choosing the relevant parameters to describe the population structure constrained to which the evolutionary dynamic is played out – is not trivial. Both the explicit spatial as well as the random graph ensemble approach have clear precedents in condensed matter physics and network theory, respectively. It is not, however, clear which – if either – approach best describes natural populations of cyclically competing species or societies composed of individuals playing the prisoner’s dilemma game.

As an example let us consider colicin producing bacteria, that play the so called "rock-paper-scissors" (RPS) game (for details see below). This system has recently been the subject of two experimental studies aimed at demonstrating the role of structured populations in the maintenance of diversity. In the first study [32, 33] bacteria were cultured *in vitro* in Petri dishes, effectively restricting competition between bacteria to neighbors on the (2D) Petri dish surface (Fig.2.1 top left), while in the second experiment [34] *in vivo* bacterial colonies were established in co-caged mice and their development was subsequently followed. In the case of the first experiment the analogy with explicit 2D spatial embedding (present by construction) is clear (Fig.2.1 bottom left). The population structure of the second experiment is, however, clearly different. The bacteria in individual mice can be readily considered as locally well-mixed populations, the coevolutionary dynamics of which reduces in the standard mean field limit to a system of non-linear differential equations (the adjusted replicator equations derived in section 2.3 following Ref. [39]). As the

experiments show, however, migration of bacteria between mice may also occur – resulting in the observed cyclic presence of the three strains in individuals. There are two distinct scales of mixing present in the system. Bacteria within each mice compete with each other forming local populations – an unstructured neighborhood composed of individual bacteria, while also being exposed to migrants from mice with whom they share the cage, together forming a global population – an unstructured neighborhood composed of individual local populations (Fig.2.1 top and bottom right). This setup is referred to in the ecology literature – albeit in significantly different contexts – as a “structured metapopulation” [48, 49] where “structured” here refers to the detailed consideration of the population dynamics of the individual populations (often called “patches”) comprising the metapopulation and is also related to the finite island models of population genetics [50].

The above example of co-caged mice is not unique, we may readily think of other ecological or sociological examples where an approximation with hierarchical scales of mixing with no internal structure can be relevant (such as human societies with two distinct scales of mixing present, the first within individual nations the between them at an international level). Below in section 2.7 we use a similar approach to construct a model of genetic exchange among bacteria of the same species (the bacterial equivalent of sex) with which we were able to take into account the effects of spatial and temporal fluctuations in a manner that can explain the benefit of such genetic exchange at the level of the individual as described in our publication Ref. [51].

In this section we construct a hierarchical mean field theory where the two distinct (i.e., local and global) scales of mixing are each taken into account in terms of two separate *mean field* approximations and fluctuations resulting from finite population size on the local scale of mixing are also considered. We subsequently explore the similarities and differences between this and other models of structured populations in the case of the “rock-paper-scissors” and prisoner’s dilemma games in the following two sections. Through these examples we suggest that our approach allows the separation of the effects of structured populations on coevolutionary dynamics into effects which are highly sensitive to and dependent on the details of the topology, and effects which only require the minimal structure present in our approximation and can consequently (in terms of sensitivity to the details of the topology) be considered more robust.

Let us consider an evolutionary game between d types (strategies) described by the $d \times d$ payoff matrix A with elements α_{kj} . Assuming finite and constant population size, natural selection can be described at the level of the individual by the Moran process described in section 2.3, during which at each time step an individual is selected randomly from the

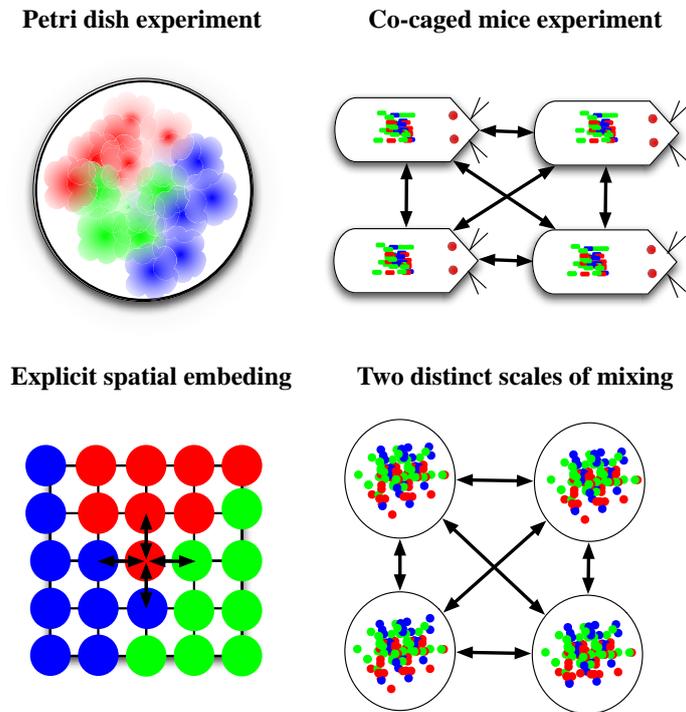


Figure 2.1: In the colicin version of the RPS game, strains that produce colicins (red) kill sensitive (green) strains, that outcompete resistant (blue) strains, that outcompete colicin producing strains (toxin production involves bacterial suicide). Experiments [32] show that colicin-producing strains cannot coexist with sensitive or resistant strains in a well-mixed culture, yet all three phenotypes are recovered in natural populations. Two recent experiments have examined the role of population structure in the maintenance of diversity among colicin-producing bacteria. In the first [32] *in vitro* colonies were established on an agar substrate in Petri dishes, a setup which effectively limits competition to neighbors on the petri dish in analogy with explicit spatial embedding in 2D. In the second [34] *in vivo* colonies were established in the intestines of co-caged mice, a setup which has two distinct scales of mixing, with no explicit structure on either scale.

population to be replaced (death) by the offspring of an individual that is chosen proportional to its fitness to reproduce (birth). This models a population in equilibrium, where the time scale of the population dynamics is set by the rate at which “vacancies” become available in the population. The fitness of each individual depends on the payoff received from playing the game described by A with competitors (an individual of type k receiving a payoff α_{kj} when playing with an individual of type j). In well-mixed populations, individuals can be considered to come into contact (compete) with equal probability with any member of the population excluding themselves – this allows one to calculate the fitness of an individual of type k in a mean field manner, yielding

$$\pi_k = \pi_{\text{base}} + \sum_{j=1}^d \frac{\alpha_{kj}(n_j - \delta_{kj})}{N - 1}, \quad (2.13)$$

where n_k is the number of individuals of type k in the population, $\sum_{k=1}^d n_k = N$ is the size of the population, π_{base} is some baseline fitness and the Kronecker delta symbol δ_{kj} is equal to unity if $k = j$ and is zero otherwise. From this we may calculate the transition probabilities of our stochastic process, i.e., the probability of an individual of type i being replaced by an offspring of an individual of type k is given by

$$T_{ik} = \frac{n_i \pi_k n_k}{N \bar{\pi} N}, \quad (2.14)$$

where $\bar{\pi} = \sum_{k=1}^d \pi_k n_k / N$. The state of any population is completely described by the frequency of the different strategies $x_k = n_k / N$. Due to the normalization $\sum_{k=1}^d x_k = 1$, the values of x_k are restricted to the unit simplex S_d [26]. For $d = 2$ this is the interval $[0, 1]$, S_3 is the triangle with vertices $\{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$ while S_4 is a tetrahedron etc.

As described in section 2.3 for sufficiently large, but finite populations the above stochastic process can be well approximated by a set of stochastic differential equations combining deterministic dynamics and diffusion (population drift) referred to as Langevin dynamics:

$$\dot{x}_k = a_k(\mathbf{x}) + \sum_{j=1}^{d-1} c_{kj}(\mathbf{x}) \xi_j(t), \quad (2.15)$$

where the effective deterministic terms $a_k(\mathbf{x})$ are given by

$$a_k(\mathbf{x}) = \sum_{j=1}^d (T_{jk} - T_{kj}) = x_k \frac{\pi_k(\mathbf{x}) - \bar{\pi}(\mathbf{x})}{\bar{\pi}(\mathbf{x})}, \quad (2.16)$$

$c_{kj}(\mathbf{x})$ are effective diffusion terms described in Ref. [40], and ξ_j are delta correlated $\langle \xi_k(t)\xi_j(t') \rangle = \delta_{kj}\delta(t-t')$ Gaussian white noise terms. As $N \rightarrow \infty$ the diffusion term tends to zero as $1/\sqrt{N}$ and we are left with the modified replicator equation.

In the context of our hierarchical mixing model the topology of connections can be described by two parameters, the populations size at the local scale of mixing N , and a second parameter μ , which tunes the strength of global mixing relative to the local dynamics. We take into account the second (global) scale of mixing – mixing among local populations – by introducing a modified version of the Moran process. In the modified process a random individual is replaced at each time step either with the offspring of an individual from the same population (local reproduction) or with an individual from the global population (global mixing). This is equivalent to considering the global population to be well-mixed at the scale of local populations.

Let us consider a global population that is composed of M local populations of size N . In each local population vacancies become available that local reproduction and global mixing compete to fill. In any local population l the probability of an individual of some type k filling a new vacancy due to local reproduction must be proportional to the number of individuals of type k multiplied by their fitness i.e., $\pi_k^l n_k^l$, where we consider π_k^l to be determined only by interactions with individuals in the same local population according to equation (2.13). To describe the tendency of individuals of some type k in local population l to contribute to global mixing we introduce the parameters σ_k^l . The choice of appropriate σ_k^l depends on the details of the global mixing mechanism, for systems where only the offspring of individuals mix globally it is proportional to the fitness of a given type, while for mechanisms such as physical mixing, by e.g. wind or ocean currents, it may be identical for each type. Irrespective of the details, however, the probability of an individual of some type k filling in a new vacancy due to global mixing should be proportional to the global average of the number of individuals of type k multiplied by their mixing tendency, which we denoted as $\langle \sigma_k n_k \rangle = \sum_{l=1}^M \sigma_k^l n_k^l / M$, and the strength of global mixing μ . These consideration lead to the new transition probabilities:

$$\hat{T}_{ik}^l = \frac{n_i^l}{N} \left(\frac{\pi_k^l n_k^l + \mu \langle \sigma_k n_k \rangle}{\sum_{k=1}^d (\pi_k^l n_k^l + \mu \langle \sigma_k n_k \rangle)} \right) = \frac{n_i^l}{N} \left(\frac{\pi_k^l n_k^l + \mu \langle \sigma_k n_k \rangle}{N(\bar{\pi}^l + \mu \langle \bar{\sigma} \rangle)} \right), \quad (2.17)$$

where $\bar{\pi}^l = \sum_{k=1}^d \pi_k^l n_k^l / N$ and $\langle \bar{\sigma} \rangle = \sum_{k=1}^d \langle \sigma_k n_k \rangle / N$.

We have found that the results presented below are qualitatively the same for both the *fitness dependent* choice of $\sigma_k^l = \pi_k^l$ and the *fitness independent* choice of $\sigma_k^l = 1$. Therefore, in the following we restrict ourselves to the somewhat simpler *fitness independent*

choice of $\sigma_k^l = 1$, which can be considered to correspond to some form of physical mixing mechanism. The transition probabilities (2.17) then reduce to:

$$\bar{T}_{ik}^l = \frac{n_i^l}{N} \left(\frac{\bar{\pi}^l}{\bar{\pi}^l + \mu} \frac{\pi_k^l n_k^l}{\bar{\pi}^l N} + \frac{\mu}{\bar{\pi}^l + \mu} \frac{\langle n_k \rangle}{N} \right). \quad (2.18)$$

We can see that after a vacancy appears either local reproduction occurs, with probability $\bar{\pi}^l/(\bar{\pi}^l + \mu)$, or global mixing, with probability $\mu/(\bar{\pi}^l + \mu)$. From equation (2.18) we may derive the Langevin equation describing the coevolutionary dynamics of population l from the

$$\dot{x}_k^l = \hat{a}_k(\mathbf{x}^l, \langle \mathbf{x} \rangle) + \sum_{j=1}^{d-1} \hat{c}_{kj}(\mathbf{x}^l, \langle \mathbf{x} \rangle) \xi_j(t), \quad (2.19)$$

with the modified deterministic terms given by

$$\hat{a}_k(\mathbf{x}^l, \langle \mathbf{x} \rangle) = \frac{x_k^l (\pi_k(\mathbf{x}^l) - \bar{\pi}(\mathbf{x}^l)) + \mu (\langle x_k \rangle - x_k^l)}{\bar{\pi}(\mathbf{x}^l) + \mu}, \quad (2.20)$$

where the vector $\langle \mathbf{x} \rangle = \sum_{l=1}^M \mathbf{x}^l / M$ with components $\langle x_k \rangle = \sum_{l=1}^M x_k^l / M$ describes the frequencies of the individual types in the global population and the diffusion terms $\hat{c}(x^l, \langle \mathbf{x} \rangle)$ can be expressed in terms of the modified transition probabilities \hat{T}_{ik}^l as above.

Equations (2.19) describe the coevolutionary dynamics of the global population through the coupled evolution of the $\{\mathbf{x}^1, \dots, \mathbf{x}^M\}$ local populations. In the limit of a large number of local populations ($M \rightarrow \infty$) the distribution of the local populations over the space of population states (the simplex S_d) is described by a density function $\rho(\mathbf{x})$ that is normalized over S_d , i.e., $\int_{S_d} \rho(\mathbf{x}) = 1$. The time evolution of $\rho(\mathbf{x})$ follows a $d - 1$ dimensional advection-diffusion equation – the Fokker-Planck equation corresponding to Eq. (2.19):

$$\dot{\rho}(\mathbf{x}) = -\nabla \cdot \{ \hat{\mathbf{a}}(\mathbf{x}, \langle \mathbf{x} \rangle) \rho(\mathbf{x}) \} + \frac{1}{2} \nabla^2 \cdot \{ \hat{\mathbf{b}}(\mathbf{x}, \langle \mathbf{x} \rangle) \rho(\mathbf{x}) \}, \quad (2.21)$$

with the global averages $\langle x_k \rangle = \int_{S_d} x_k \rho(\mathbf{x})$ coupled back in a *self-consistent* manner into the deterministic terms $\hat{a}_k(\mathbf{x}, \langle \mathbf{x} \rangle)$ and the diffusion matrix $\hat{b}_{kj}(\mathbf{x}, \langle \mathbf{x} \rangle) = \sum_{i=1}^{d-1} \hat{c}_{ki}(\mathbf{x}, \langle \mathbf{x} \rangle) \hat{c}_{ij}(\mathbf{x}, \langle \mathbf{x} \rangle)$. For large local populations ($N \rightarrow \infty$) the diffusion term vanishes as $1/N$.

The above advection-diffusion equation (2.21) presents an intuitive picture of the coevolutionary dynamics of the population at a global scale. We can see that local populations each attempt to follow the trajectories corresponding to the deterministic replicator dynamics, while under the influence of two additional opposing forces: (i) global mixing, which attempts to synchronize local dynamics and (ii) diffusion resulting from finite population

size effects, which attempts to smear them out over the simplex. The strength of these forces are tuned by two parameters μ and N , respectively.

If, further, the effects of synchronization are irrelevant, as for example in the case of populations where selection is externally driven by independent environmental fluctuations, we may replace the global population average with the time average of any single population. This is the approach we will use in our study of genetic mixing in bacteria in section 2.7 below.

Our approach readily generalizes for an arbitrary number of hierarchical mixing levels. For three levels of mixing we may consider the global population to be comprised of \mathcal{M} subpopulations each of which is in turn subdivided into M local populations. With $m \in \{1, \dots, \mathcal{M}\}$ running over subpopulations and $l \in \{1, \dots, M\}$ over local populations the transition probabilities can be written as:

$$\hat{T}_{ik}^{ml} = \frac{n_i^{ml}}{N} \left(\frac{\pi_k^{ml} n_k^{ml} + \mu^{(1)} \langle \sigma_k^{ml'} \rangle_{l'} + \mu^{(2)} \langle \langle \sigma_k^{m'l'} \rangle_{l'} \rangle_{m'}}{\sum_{k=1}^d (\pi_k^{ml} n_k^{ml} + \mu^{(1)} \langle \sigma_k^{ml'} \rangle_{l'} + \mu^{(2)} \langle \langle \sigma_k^{m'l'} \rangle_{l'} \rangle_{m'})} \right), \quad (2.22)$$

where primed indices indicate the scale of mixing over which the average is taken, $\mu^{(1)}$ describes the strength of mixing, and the $\sigma_k^{ml'}$ the tendencies of mixing among local populations within a subpopulation, while $\mu^{(2)}$ describes the strength of mixing, and the $\sigma_k^{m'l'}$ the tendencies of mixing among subpopulations in the global population.

During our numerical investigations presented in Ref. [52] and discuss in the following two sections ((2.5) and (2.6)) we found solving the advection-diffusion equation (2.21) numerically challenging, particularly in the $N \rightarrow \infty$ limit. We resorted instead to solving the coupled Langevin equations (2.19) for large $M = 10^4 - 10^5$ to simulate the time evolution of $\rho(\mathbf{x})$.

2.5 Cooperation in populations with hierarchical levels of mixing

The evolution of cooperation is a fundamental problem in biology, as natural selection under most conditions favors individuals who defect. Despite of this, cooperation is widespread in nature. A cooperator (**C**) is an individual who pays a cost c to provide another individual with some benefit b . A defector (**D**) pays no cost and does not distribute any benefits. This implies the payoff matrix

$$\begin{array}{rcc}
& \text{from C} & \text{from D} \\
\text{C gets} & b - c & -c \\
\text{D gets} & b & 0
\end{array} , \tag{2.23}$$

where b is the benefit derived from playing with a cooperator while c is the cost for cooperation. From the perspective of evolutionary game theory, which equates payoff with fitness, the apparent dominance of defection is simply the expression of the fact that natural selection *a priori* selects for fitness of individuals and not the fitness of groups.

In stark contrast to the case of hawks and mice, defection dominates cooperation in any well-mixed population [26]. Population structure induced by spatial structure [27, 45] and more general networks of interactions [46, 47, 53]) has, however, been found to facilitate the emergence and maintenance of cooperation. The mechanism responsible, termed spatial, or more generally, network reciprocity [54], depends strongly on the details of local topology. In particular, it seems that lattice like connectivity structures where three-site clique percolation occurs [44] and more general interaction graphs where the degree of nodes k does not exceed the ratio of benefit to cost (i.e., $k < b/c$) [47] are required for cooperation to be favored.

Examining the effects of hierarchical mixing on the evolutionary dynamics of cooperation, we found that a sharp, but continuous transition leads to the dominance of cooperation as the benefit to cost ratio becomes smaller than the local population size, i.e., $b/c < N$. If the benefit to cost ratio is larger than the local population size the global population is dominated by defectors. The mechanism leading to the dominance of cooperation arises due to the competition between local reproduction and global mixing. In local populations with lower average fitness – larger number of defectors – the influx of individuals from the global scale will be larger than in local populations with higher average fitness (cf. Eq. (2.18) where the relative strength of the two terms on the left hand side depends on the sum of the average fitness of population l and μ). The crucial ingredient for cooperation to be successful is population drift introduced by finite local population size. It is biased influx coupled with drift that can result in cooperation being favored in the global population (Fig 2.2.).

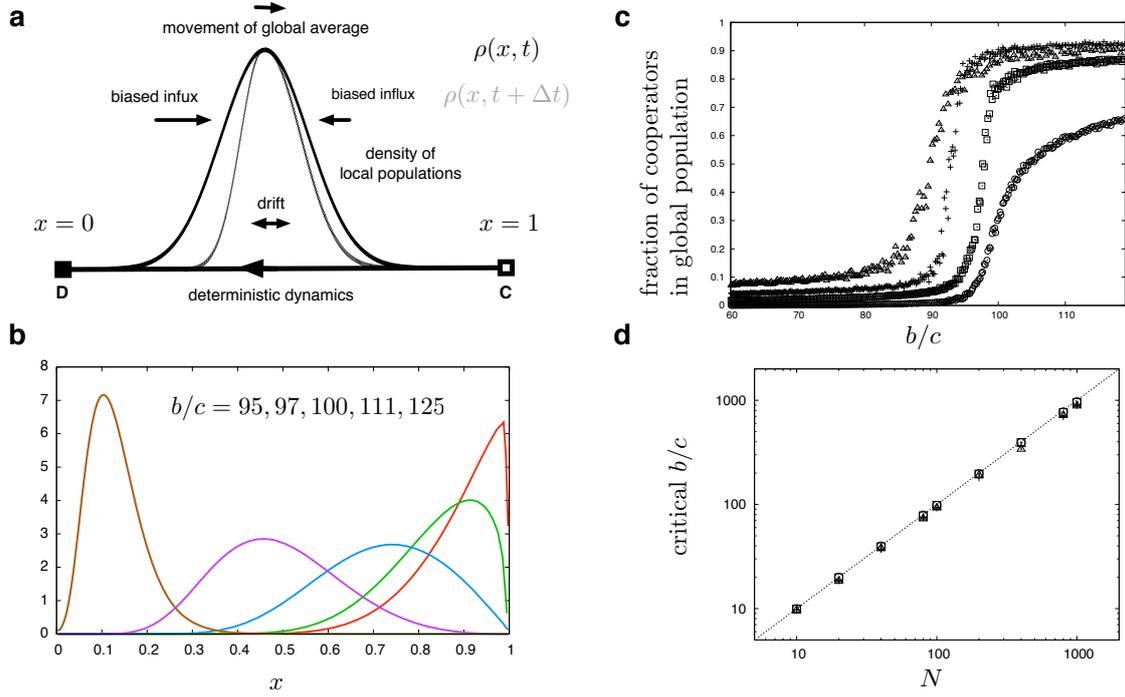


Figure 2.2: **(a)** In an infinitely large well-mixed population evolutionary dynamics is deterministic and leads to the extinction of cooperators as average fitness monotonically declines. The only stable fixed point corresponds to the point where the fraction of cooperators is zero ($x = 0$). To understand qualitatively the mechanism favoring cooperation in hierarchically mixed populations let us consider some density of local populations ($\rho(x, t)$) that is symmetric around its mean at time t . Due to global mixing all local populations are being driven toward the global average. Due to the influx bias, however, populations with a lower than average number of cooperators will be driven stronger (faster) than those on the other side of the average. Examining the density of local populations at some time $t + \Delta t$, this results in a net movement of the global average toward a larger fraction of cooperators. This is, of course, opposed by local reproduction that favors an increase in the number of defectors. For the global average to keep moving toward a higher number of cooperators and eventually to keep balance with local reproduction bias a density of local population with finite width is needed over which the effect of the influx bias can exert itself. It is drift caused by local population size that maintains this finite width, and this is the reason that the b/c threshold above which cooperation dominates depends on local population size. **(b)** Stationary density of local populations $\rho(x)$ for different values of b/c with $N = 100$, $\mu = 0.1$. **(c)** Transition toward a global dominance of cooperation for $\mu = 10$ (triangles), $\mu = 1$ (crosses), $\mu = 0.1$ (squares), $\mu = 0.01$ (circles) with $N = 100$. The critical value of b/c depends only weakly on μ changing by 20% over four orders of magnitude **(d)** Critical values of b/c as a function of N for different values of μ (notation as before). The dashed line corresponds to $b/c = N$. The critical b/c values (presented in our publication Ref. [52]) were determined by numerically finding the inflection point of the transition curves. $M = 10^3$ was used throughout.

2.6 Three strategies

2.6.1 The RPS game

To explore the effects of hierarchical mixing in the context of games with three strategies we first turn to the case of the so called “rock-paper-scissors” (RPS) game. In the original popular version of the game two players are afforded the chance to simultaneously display either rock (fist), paper (flat hand) or scissors (two fingers). If player one displays a flat hand while player two displays a fist, player one wins as paper wraps rock. Similarly scissors cut paper, and rocks smashes scissors. Several examples of this game have been found in nature (e.g. among lizards [55]), but it is bacteria that have received the most experimental and theoretical attention.

In ecology the often high diversity among microbial organisms in seemingly uniform environments, referred to as the “paradox of the plankton”, has been difficult to understand. Several models based on spatially explicit game theoretical models have been proposed to explain this diversity [35, 36, 32, 33]. These models are all variants of the RPS game played by colicin producing bacteria. Colicins are antibiotics produced by some strains of *Escherichia coli*. In experiments (see Fig.2.1) typically three strains are used: colicin producing (C), sensitive (S) and resistant (R). The coevolutionary dynamics of the three strains can be cast in terms of an RPS game, C strains kill S strains, but are outcompeted, by R strains, because toxin production involves the suicide of bacteria. The cycle is closed by S strains that outcompete R strains, because resistance requires mutant versions of certain membrane protein, which are less efficient than the wild type [32]. Despite the cyclic dynamics colicin-producing strains cannot coexist with sensitive or resistant strains in a well-mixed culture, yet all three phenotypes are recovered in natural populations. Local dispersal (modeled as explicit spatial embedding) has widely been credited with promoting the maintenance of diversity in this system [32, 33, 35, 36].

In its most symmetric form the RPS game is described by the payoff matrix

$$\begin{pmatrix} 0 & -\epsilon & \epsilon \\ \epsilon & 0 & -\epsilon \\ -\epsilon & \epsilon & 0 \end{pmatrix}, \quad (2.24)$$

and some $\pi_{\text{base}} > \epsilon$. The dynamics of this game in an infinitely large well mixed population consists of neutral orbits along which the product $x_R x_P x_S$ is conserved. For any finite N , however, fluctuations lead to the inevitable extinction of all but one of the strategies

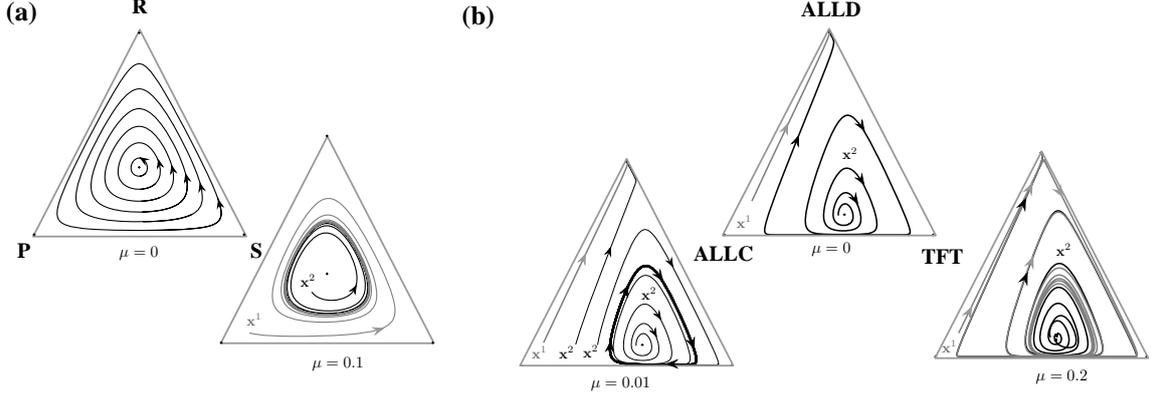


Figure 2.3: **(a)** Deterministic replicator dynamics (the $N \rightarrow \infty$ limit) of the symmetric RPS game consists of neutrally stable orbits along which the product of the strategy frequencies $x_R x_P x_S$ is conserved. If global mixing is present ($\mu > 0$) local populations deviate from these neutral orbits toward the global average $\langle \mathbf{x} \rangle$. Considering the simplest system with global mixing, that consisting of $M = 2$ local populations we see that in the presence of global mixing population \mathbf{x}^1 and population \mathbf{x}^2 move toward each other, respectively moving closer and further from the barycentre of the triangle until they become synchronized and subsequently pursue a common orbit. For deterministic local dynamics ($N \rightarrow \infty$) such synchronization invariably occurs for any M if $\mu > 0$ and typically converges to the barycentre of the simplex for sufficiently homogeneous initial conditions. **(b)** The deterministic replicator dynamics of the repeated PD game is markedly different from that of the RPS game in that the internal fixed point is unstable and in the absence of global mixing only ALLD survives. Again turning to the simplest scenario with $M = 2$ we see that if $\mu = 0$ any pair of populations \mathbf{x}^1 and \mathbf{x}^2 (gray and black lines) converge to the to the ALLD corner. As μ is increased above a critical value a second, stable configuration emerges: for a large subset of the possible initial conditions (all, but the left most \mathbf{x}^2) we see that one of the populations (\mathbf{x}^1) converges to ALLD, while the second (\mathbf{x}^1) approaches a limit cycle. If μ is increased further, the above configuration ceases to be stable, the population which initially converges to ALLD (\mathbf{x}^1) is subsequently “pulled out” by global mixing, following which the two populations synchronize and are finally absorbed together in ALLD. Simulations (presented in Ref. [52]), however, show that synchronization may be avoided for $M > 2$ if μ is not too large.

[56]. Spatial population structure can avert this reduction in diversity [35, 32] through the emergence of a stable fixed point at the barycentre of the simplex. The effect of the gradual randomization of different lattice topologies (where a small number of edges are randomly rewired) on the dynamics of the game has also been investigated. A Hopf bifurcation leading to global oscillations was observed [57, 58] as the fraction of rewired links was increased above some critical value.

Examining the dynamics of the symmetric RPS game in terms of our hierarchical mean field approximation we observed that an internal fixed point emerged for $N \rightarrow \infty$ (Fig.2.3a). More importantly, diversity was also maintained for finite local population sizes if global mixing was present. Simulations of the time evolution of $\rho(\mathbf{x})$ also revealed a

Hopf bifurcation leading to the oscillation of the global average as μ was increased above a critical value μ_c depending on N (Fig.2.4a). These results show that previous results obtained from simulations of populations constrained to different lattice topologies can be considered universal in the sense that not only lattices, but any population structure that can be approximated by two distinct internally unstructured scales of mixing are sufficient for their existence. In the context of the “paradox of the plankton” these results imply that aside of local dispersal (modeled as explicit spatial embedding) a minimal metapopulation structure (with local competition and global migration) can also facilitate the maintenance of diversity in cyclic competition systems.

2.6.2 The repeated prisoner’s dilemma game

In the general formulation of the prisoner’s dilemma (PD) game, two players have the choice to cooperate or to defect. Both obtain some payoff R for mutual cooperation and some lower payoff P for mutual defection. If only one of the players defects, while the other cooperates, the defector receives the highest payoff T and the cooperator receives the lowest payoff S . That is $T > R > P > S$ and defection dominates cooperation in any well-mixed population. New strategies become possible, however if the game is repeated, and players are allowed to choose whether to defect or cooperate based on the previous actions of the opponent. In the following we consider, similar to Refs. [59] and [60] that examined the role of finite population size and mutation and finite population size, respectively in terms of the repeated PD game with three strategies: always defect (ALLD), always cooperate (ALLC), and tit-for-tat (TFT). TFT cooperates in the first move and then does whatever the opponent did in the previous move. TFT has been a world champion in the repeated prisoner’s dilemma ever since Axelrod conducted his celebrated computer tournaments [6], although it does have weaknesses and may be defeated by other more complex strategies [61].

Previous results indicate that if only the two pure strategies are present (players who either always defect or ones who always cooperate) explicit spatial embedding [27] and some sufficiently sparse interaction graphs [47, 62] allow cooperation to survive and the behavior of populations is highly sensitive to the underlying topology of the embedding [44]. As we have shown in section 2.5 introducing global mixing into the PD game with only the two pure strategies present also allows cooperation to survive.

To investigate the effect of global mixing on the repeated PD game with three possible strategies: ALLD, ALLC and TFT following Imhof et al. [60] we considered the payoff matrix:

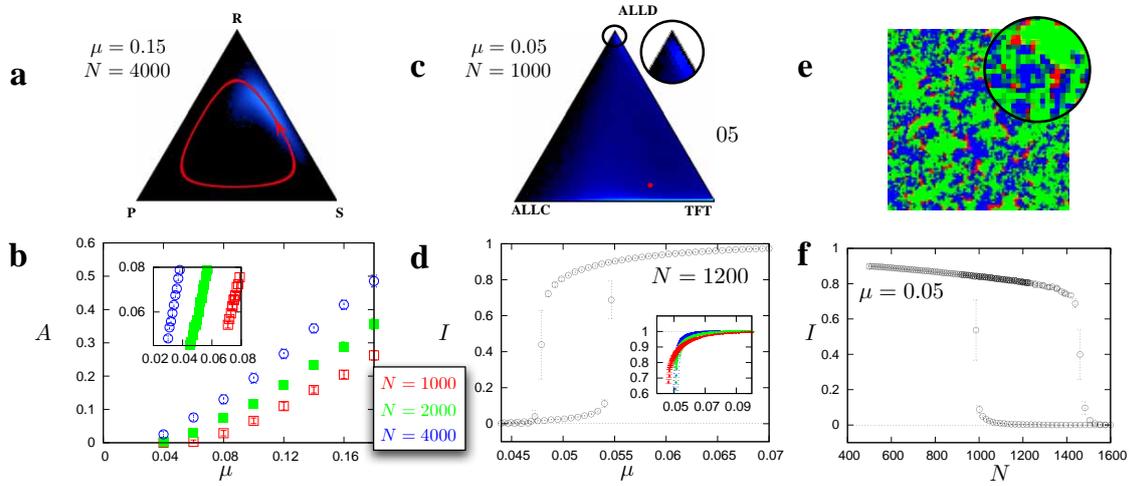


Figure 2.4: **(a)** In the case of the rock-paper-scissors game a Hopf bifurcation similar to that observed for populations evolving on gradually randomized lattices [57, 58] leads to the emergence of global oscillations (the red line indicates the trajectory of $\langle \mathbf{x} \rangle$) if μ is larger than a critical value $\mu_c(N)$ (see video S1 part of the supporting information of Ref. [52] and also available as part of the electronic supplement of the thesis at <http://angel.elte.hu/~ssolo/thesis.html>). The density $\rho(\mathbf{x})$ is indicated with a blue color scale. **(b)** The ratio A of the area of the global limit cycle and the area of the simplex is plotted as a function of μ for three different values of N . For the repeated prisoner's dilemma game the combination of finite local population size and global mixing $\mu > 0$ can lead to a stationary solution **(c)** qualitatively similar to that observed for explicit spatial embedding **(e)**. This state is characterized by a stable global average (large dot), just as the lattice system (data not shown) and sustained local cycles of cooperation, defection and reciprocity, also similar to the lattice case where groups of ALLD (red) individuals are chased by those playing TFT (blue, black), which are gradually outcompeted by ALLC (green). **(d)** As μ is decreased a discontinuous transition can be observed to the ALLD phase. The ratio I of populations on the internal cycle is plotted as a function of μ . The inset shows the transition for different values of N . **(f)** The same critical line in the μ - N plane can be approached by increasing N with μ fixed. A large hysteresis can be observed as N is decreased below the critical value indicating the discontinuous nature of the transition. We numerically simulated the time evolution of $\rho(\mathbf{x})$ by integrating the stochastic differential equation system defined by Eq. (2.19) for large M ($10^4 - 10^5$) throughout. For the RPS game we used $\pi_{\text{base}} = 1$ and $\epsilon = 0.5$, while in the case of the repeated PD game we followed Ref. [60], setting $T = 5, R = 3, P = 1, S = 0.1, m = 10$ and $c = 0.8$. Lattice simulations **(e)** were performed on 1000×1000 square lattice with an asynchronous local Moran process between neighbors and periodic boundary conditions.

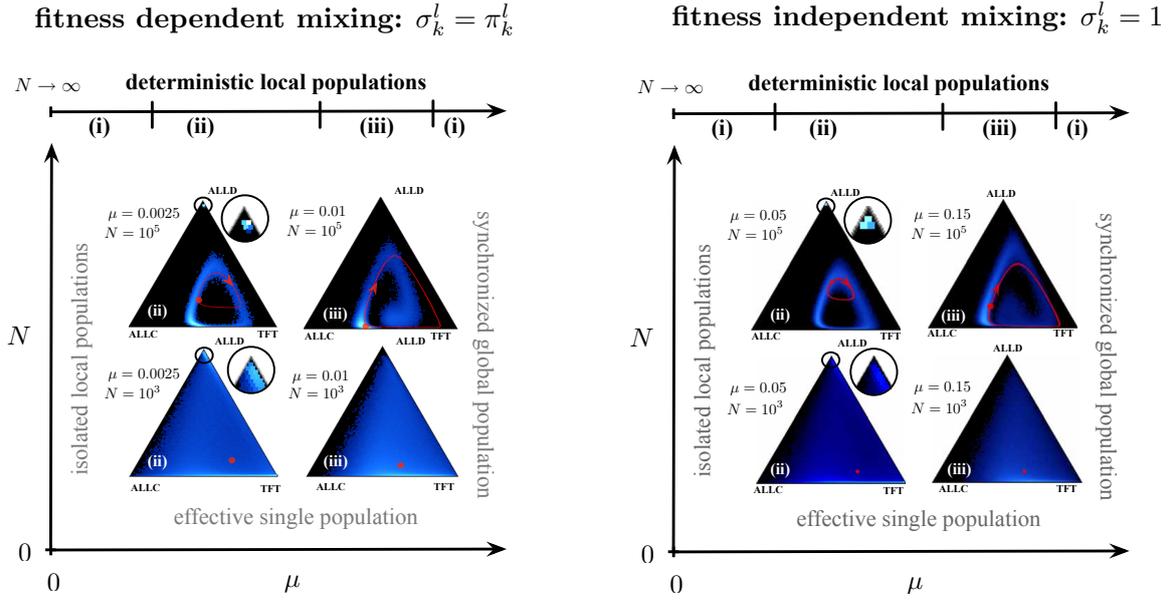


Figure 2.5: Phase space for the repeated prisoner’s dilemma game on a population structure with two distinct scales (see video S2 part of the supporting information of Ref. [52] and also available as part of the electronic supplement of the thesis at <http://angel.elte.hu/~ssolo/thesis.html>). Three different phases are possible depending on the values of μ and N : (i) only ALLD survives (ii) an internal limit cycles is maintained by global mixing due to a large density of local populations around the ALLD corner (iii) a globally oscillating self maintaining limit cycle is formed. For extreme values of μ the global dynamics reduces to that of some well-mixed population where only ALLD survives: As μ becomes negligible ($\mu \ll \pi_k$ for all k) we approach the limit of isolated local populations, while for $\mu \gg \pi_k$ we are left with a single synchronized population. Similarly for $N = 2$ – the smallest system with competition – the system can be described as a single well mixed population for any μ and ALLD again prevails. In the limit of deterministic local populations ($N \rightarrow \infty$) all three phases can be found depending on the value of μ . The density $\rho(\mathbf{x})$ is indicated with the color scale.

$$\begin{pmatrix} Rm & Sm & Rm \\ Tm & Pm & T + P(m - 1) \\ Rm - c & S + P(m - 1) - c & Rm - c \end{pmatrix}, \quad (2.25)$$

where the strategies are considered in the order ALLC, ALLD, TFT, m corresponds to the number of rounds played and c to the complexity cost associated with conditional strategies (TFT). The dynamics of this game has a single unstable internal fixed point and the state where each member of the population plays ALLD is the only nontrivial stable equilibrium (Fig.2.3b).

Introducing global mixing, between local well-mixed populations, however, causes new stationary states to emerge. Three phases can be identified: (i) ALLD wins (ii) large frac-

tion of local populations in the ALLD corner maintains local cycles of cooperation defection and reciprocity through providing an influx of defectors that prevent TFT players from being outcompeted by ALLC playing individuals (iii) a self maintaining internal globally oscillating cycle emerges. The simplest scenario of two ($M = 2$) deterministic ($N \rightarrow \infty$) local populations coupled by global mixing ($\mu > 0$) already leads to the emergence of phase (ii) as demonstrated in Fig.2.2b while phase (iii) only emerges for larger M . For larger M simulations show that in the limit of large local populations all global configurations with less than some maximum ratio of the populations I on the internal cycle are stable in phase (ii). A transition from phase (ii) to (i) happens as μ is decreased below a critical value $\mu_c^{ii \rightarrow i}$ and I approaches zero as $I = (1 - \mu_c^{ii \rightarrow i} / \mu)$ (data not shown). This can be understood if we considered that near the transition point a critical proportion $C = \mu(1 - I)$ of ALLD individuals needs to arrive to stabilize local cycles of cooperation defection and reciprocity. At the critical point $I = 0$ and $\mu = \mu_c^{ii \rightarrow i}$ which implies $C = \mu_c^{ii \rightarrow i}$ giving $I = (1 - \mu_c^{ii \rightarrow i} / \mu)$.

Exploring the $N - \mu$ phase space (Fig.2.5) we see that the transition from phase (i) to (ii) becomes discontinuous for finite N (Fig.2.4d,e). Further, for any given value of N and μ the global configuration is described by a unique I due to the presence of diffusion. For appropriate values of the parameters the global average converges to a stationary value in phase (ii) similarly to case of explicit spatial embedding (Fig.2.4c).

For very small ($\mu \ll \pi_k$ for all k) and very large ($\mu \gg \pi_k$) values of μ the global dynamics can be reduced to that of some well-mixed population where only ALLD persists (Fig.2.5.). For small N we again have an effective well-mixed population – the only limit where defectors do not dominate is $N \rightarrow \infty$. In comparison with previous results of Imhof et al. we can see that evolutionary cycles of cooperation defection and reciprocity can be maintained not only by mutation, but also by population structures with hierarchical levels of mixing.

2.7 The maintenance of natural genetic transformation in bacteria

Sexual reproduction is a process that brings genomes, or portions of genomes, from different individuals into a common cell, producing a new combination of genes: in eukaryotes, this occurs as a result of fertilization and meiotic recombination; in bacteria, it happens as a result of the acquisition of exogenous DNA. The ubiquity of genetic transfer in bacteria is reflected in the dynamic structure of their genomes, which are constantly being shaped

by two opposing forces: selection for shorter length (favoring DNA loss through deletion) and selection for gene function (driving genome loading by the acquisition of exogenous DNA) [63, 64]. The balance of these forces results in most bacteria having highly economized genomes with only a small fraction (around 10%) of noncoding sequences [65, 64]. DNA transfer into the bacterial cell can occur in three ways: (i) transduction by viruses, (ii) conjugation by plasmids, and (iii) natural genetic transformation (NGT) by developing competence, a regulated physiological state in which the bacterial cell is able to take up DNA fragments released by another cell [66, 67]. The genetic elements responsible for transduction and conjugation primarily survive as parasites, and are located on viral and plasmid DNA. The genes required for competence are, however, located on the bacterial chromosome, placing NGT under the direct control of the cell. While all three mechanisms play a role in rare gene transfer events between bacteria of different species, termed horizontal gene transfer, NGT is the most significant source of active and frequent genetic transfer within a species (for a comparative review of the three processes see Ref. [68]). In bacteria capable of NGT, alleles typically change more frequently by recombination (e.g. 5-10 fold in *Streptococcus pneumoniae* and *Neisseria meningitidis*) than by mutation [69, 70, 71, 68]. It is this combination of high throughput genetic mixing among members of the same species, and direct cellular control that is responsible for NGT often being referred to as the bacterial analogue of meiotic sex in eukaryotes [72, 73].

The persistence of NGT raises the same question as the prevalence of meiotic sex [74, 75, 76, 77]: What is the short-term advantage of genetic mixing to the individual? NGT is obviously costly, not just because the machinery of DNA uptake must be maintained, but also because bacteria undergoing transformation face the risk of incorporating defective alleles [78]. And while the long-term implications of NGT - for genome adaptation [79] and diversification [80, 81] - are clear, the short-term advantage to individuals (an advantage necessary to maintain NGT) remains elusive [82].

The ability to take up naked DNA through NGT has been detected across a wide phylogenetic spectrum, ranging from archaea through divergent subdivisions of bacteria, including representatives from Gram positive bacteria, cyanobacteria, *Thermus* spp., green sulphur bacteria and many other Gram negative bacteria [68]. The details of transformation, however, vary widely among bacteria of different species. With the exception of *Neisseria gonorrhoeae* most naturally transformable bacteria develop time-limited competence in response to specific environmental conditions such as altered growth conditions, nutrient access, cell density (by quorum sensing) or starvation. The conserved ability among a wide range of bacteria to acquire DNA molecules through NGT indicates that the genetic trait is

functionally important in the environment, enabling access to DNA as a source of nutrients or genetic information [68].

It has been convincingly demonstrated by Redfield et al. [78] that NGT's role in repairing deleterious mutations under constant selection is insufficient for its survival. The lack of other viable explanations has left no alternative except that the uptake of DNA provides nucleotides for food [82]. This, however, is difficult to reconcile with the facts that one of the strands is taken up intact (despite the apparent risks of degradation inside the cell), and that highly specific sequences are required for the binding and uptake of DNA in some bacterial species, e.g. *N. gonorrhoeae* or *Haemophilus influenzae* [67] (even though nucleotides from other sources confer the same nutritional benefits). In order to understand the prevalence of NGT as a vehicle of genetic information one must take into account not only a single population, but a collection of populations living under diverse and constantly changing ecological conditions, as only these together possess the complete set of genes common to the species [83, 84]. To find the short-term advantage that maintains NGT, we have to consider its role in allowing genetic mixing between populations facing variable selection. There is no reason to assume, however, that the details of population structure are important, hence the maintenance of NGT is an ideal application for the minimal model of population structure developed in section 2.4.

2.7.1 The dynamics of genome organization in bacteria

As multiple genomes from closely related bacteria have become available a surprising picture of genomic diversity has begun to emerge. While it has been clear for some time that only a handful of genes are shared across the entire tree of life (the current estimate is about 60 genes mostly related to translation [85]), it has come quite as a surprise that this diversity of genome content seems to be present – at least among prokaryotes – at every phylogenetic scale [86]: Lerat et al. [87] have estimated that the core genome of 13 γ -proteobacteria contains < 300 genes and, although all *Escherichia coli* and *Shigella* genomes sequenced to date have a total of > 4000 protein coding genes they only share < 3000 [88, 89]. This striking variability has raised the question of how to define the genome of a species, and the concept of “pan genome” was proposed as the set of all genes found in a species. Tettelin et al. [90] have demonstrated that both “closed” and “open” pan genomes can be found in Nature. In *Bacillus anthracis*, the pan genome is “closed” it is completely described by four genomes. The pan genome of group B *Streptococcus*, in contrast, remains “open”, the number of new genes contributed by every new genome sequence contributing an expected 30 genes independent of the number of genomes already present in the comparison – this

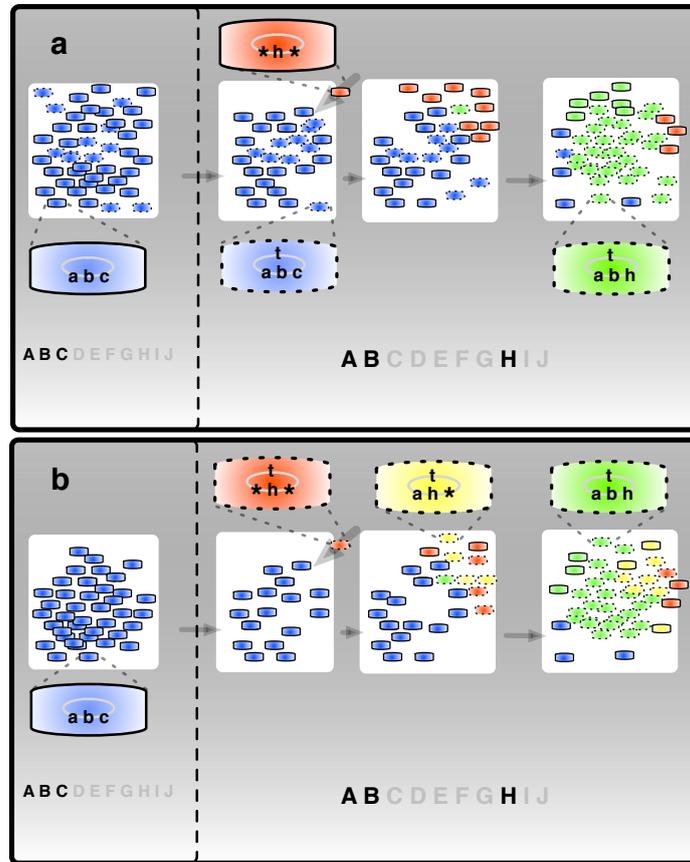


Figure 2.6: Schematic depiction of the two main processes responsible for the spread of NGT. In both scenarios local food types change from **ABC** to **ABH**. Bacteria with the most viable genotypes are shown with the relevant model genes indicated inside, those capable of NGT are emphasized by “perforated” borders. **(a)** Under rapidly changing conditions the number of bacteria with the **t** (transformation) model gene usually remains significant (cf. Fig.2.7b). Newly arrived bacteria possessing the model food gene **h** (red cells) spread due to a lack of competition, subsequently allowing local bacteria capable of NGT (blue cells with perforation) to successfully adapt by assembling the optimal combination of model genes **abh+t** (green cells with perforation), and to outcompete all the others, while also spreading the **t** model gene. **(b)** Under slowly changing conditions, gene **t** often disappears from the local population (cf. Fig.2.7b). It is, however, possible that a competent bacterium possessing gene **h** (red cells with perforation) arrives from another population. Progeny of this migrant subsequently adapt and outgrow the local bacteria by incrementally (indicated by yellow then green coloring) assembling the best adapted genome, **abh+t** (green cells with perforation), which also results in the spread of the model **t** gene.

indicates that the many more genome sequences are needed to “close” the pan genome of this species. Our model of a bacterial species subdivided into strains with locally specialized genomes presented below, while developed independently of these results, is fully consistent with the emerging picture of large *pan-genomes*.

We view a bacterial species as a *metapopulation* [49] that is composed of a large number of spatially distinct populations living under varying ecological conditions. Different populations experience selection for different combinations of numerous possible environmental factors (availability of particular metabolites, host recognition [91] and others [92]). Most of these factors fluctuate in time over a broad range of timescales (starting from daily weather changes, through seasonal alternations and decades-long host life-cycles, up to long-scale climate changes, to mention just a few). The populations are, on one hand, constantly adapting to their own locally changing environments and, on the other hand, connected by weak inter-population migration that can span large distances (even continents and oceans). What we aim to show is that under these conditions the advantage of NGT may lie in providing locally adapted populations with the ability to respond to environmental changes by importing genes from the collective gene pool of the species through taking up DNA from migrants. This way bacteria can economize their genomes by disposing of genes that are not currently in use in the local environment and picking up those that have just become useful. There is a growing body of experimental evidence showing that such genetic plasticity plays a central role in the adaptation of bacteria, the most well studied examples being virulence related genetic diversity (particularly of the genes responsible for capsule composition) in *S. pneumoniae* [91] and the hypervariable region of *Helicobacter pylori* [93] responsible for different pathophysiologies associated with chronic *H. pylori* infection in humans.

To test our hypothesis quantitatively, we consider a model of a bacterial species (Fig.2.6) that is capable of living on 10 different types of food (denoted by **A**, **B**, ..., **J**). Individual bacteria can utilize any of these foods only if they have the corresponding model food gene (**a**, **b**,..., **j**, respectively) present in their genomes. Each such model gene is understood to represent the complete group of genes necessary for the utilization of a particular resource type. While under natural conditions there may exist dozens of such resource–gene group pairs, numerical treatment of the model rapidly becomes intractable as this number is increased, forcing us to consider only a limited number of food types and gene groups. The relatively small model genomes used in our simulations (typically a few model genes long) intend to represent the much larger genomes of real bacteria. As a consequence, model gene numbers and model genome lengths have to be suitably rescaled when interpreting

the results of the model.

At any moment only 3 out of the 10 possible food sources are available, but their types are changing in time independently in each population of the metapopulation. This food change is characterized by the common rate R_{food} , at which one of the food sources is randomly replaced by another one not currently available in the population (e.g., **ABC** changes to **ABH**, and then to **ADH**, etc.). Bacteria reproduce at a rate r_r (maximized at $r_r^{\text{max}} = 2 \text{ hour}^{-1}$) that depends on the amount of available food they can utilize and decreases with the number of functional model food genes possessed (in order to impose genome economization), for details see Methods section 4.1. Bacteria can also be washed out of their populations at a fixed rate $r_w = 10^{-2} \text{ hour}^{-1}$, lose any of their functional genes by mutation at a rate $r_m = 10^{-4} \text{ hour}^{-1}$ per model gene, and those possessing a functional copy of the model gene for NGT (denoted by **t**) attempt to incorporate exogenous DNA into their genome from the surrounding medium at a rate $r_t = 10^{-3} \text{ hour}^{-1}$ per model gene. Transformation and mutation are considered at the level of model genes (**a - j** and **t**). That is to say, a single mutation event in the context of our model encompasses a series of events starting with a deleterious mutation (either a point mutation or a deletion) leading to loss of function and continuing with the subsequent gradual loss of the gene group responsible for the specific function through repeated deletion events. In other words, we only consider a given group of genes to be either present and fully functional or completely absent. We model transformation in a similar fashion, a bacterium may acquire or lose a complete functioning copy of any model gene as a result of a single transformation event. Note that the gene group responsible for transformation (the model **t** gene) is capable of eliminating even itself by taking up a defective copy from the environment. While the rates of a series of events leading to the acquisition or loss of a complete gene group (a single model gene) responsible for a given function are probably orders of magnitude smaller than those considered in our model, we argue in detail below that this does not effect the validity of the results obtained.

The frequency of model gene fragments in the surroundings is approximated by that in the living individuals of the population [79]. This is consistent with the assumption that the death of bacteria is largely independent of their gene composition, and experiments showing that DNA fragments from lysed bacteria persist for hours to days – as measured in natural transformation assays [68] – a time-scale that is too short to allow extracellular DNA to survive the time period between food changes in our model. For details of how the competition of individual bacteria with each possible genotype was modeled see the section 4.1.

Migration between the populations is a crucial element of the model, because without it NGT would just futilely reshuffle the existing genes inside the populations. To see this, let us suppose that in a fraction P of a population a certain gene has become defective. Then the frequency of repairing this defective gene by taking up a functional copy from a recently deceased member of the population, $r_t P(1 - P)$, is the same as that of accidentally replacing a functional gene by a defective one, $r_t(1 - P)P$. Thus, lacking any short-term advantage, NGT would rapidly disappear, in agreement with the results of [78]. If, however, we take into account migration, it is able to facilitate the spread of gene \mathbf{t} , as illustrated in the context of our model in Fig.2.6.

2.7.2 Self-Consistent Migration

Assuming that bacteria can migrate long distances (such that within this distance a large number of populations exist with independently changing food sources), migration can be taken into account very efficiently in terms of a *mean field* approach, commonly used in statistical physics. This means that the genotype distribution averaged over the populations within the migrational range can be well approximated by the time average of the genotype distribution of any single population. In short, spatial and temporal averages are interchangeable. Consequently, it is enough to consider only one population, and use its own history to compute the genotype distribution of the arriving migrants. Coupling back, through migration, an ever increasing fraction of a single population's past into its own dynamics results in the convergence of the genotype distribution of the metapopulation to its stationary value in a *self-consistent* manner, as described in more detail in the following section. The influx of migrants (the number of incoming bacteria per unit time) is defined in the model as the product of the migration rate R_{migr} and the average size of the population N (which was in the order of 10^8 in our simulations).

2.7.3 A Phase Diagram for the survival NGT

Performing extensive computer simulations for various values of the two main external parameters, the food change rate R_{food} and the migration rate R_{migr} , we have found that NGT (represented by model gene \mathbf{t}) indeed survives under a wide range of parameters, as can be seen in Fig.2.7a. In these stochastic population dynamics simulations we have numerically followed the time evolution of the number of bacteria in each of the possible 2^{10+1} genotypes (representing the presence or absence of the 10 model food genes and the model transformation gene) in a single population, as detailed in Methods section 4.1.

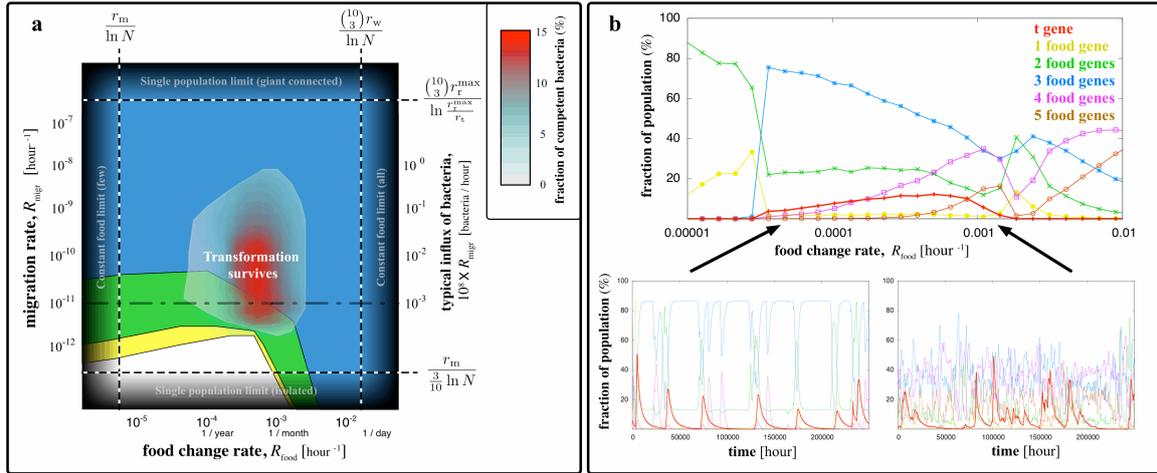


Figure 2.7: Phase diagram for the persistence of NGT. **(a)** The average percentage of bacteria capable of NGT is shown with color gradient (translucent from white to red) on the R_{food} - R_{migr} parameter space, with different points corresponding to different possible bacterial metapopulations. The maximum genome size (the number of functional model food genes) in systems where transformation has been turned off is underlain in color code (yellow regions only have bacteria with one functional food gene, green regions have bacteria with at most two, while the blue region corresponds to metapopulations where bacteria exist that can utilize all three food types). The dashed lines indicate the theoretical limits of the persistence of transformation as described in the text. **(b)** The upper panel shows a slice in the parameter space along the dashed-dotted black line in part (a), with data points corresponding to the fraction of bacteria with the t gene and the fraction of those having exactly 1, 2, 3, . . . functional food genes. The lower panel displays two illustrative time series from the simulations, corresponding to a low and a high food change rate, as indicated by the arrows.

The limits beyond which NGT cannot persist (dashed straight lines in Fig.2.7a), either because the temporal or the spatial fluctuations become irrelevant, can be estimated easily. If $R_{\text{food}} < r_m / \ln N$, then the food sources remain unchanged for such a long time that the **t** gene completely disappears by deletion before it could become beneficial at the next food change. At the other extreme, for $R_{\text{food}} > \binom{10}{3} r_w / \ln N$, i.e., when the rate at which any given food combination recurs ($R_{\text{food}} / \binom{10}{3}$) is larger than the rate at which bacteria carrying the corresponding combination of functional genes are completely washed out of the population ($\max. r_w / \ln N$), transformation cannot confer an advantage through assembling this combination as it is always present in the local population. In other words, the food changes so rapidly that populations effectively experience a constant feeding (with all possible food combinations), and NGT becomes useless.

There are similar constraints on the migration rate as well. Obviously, for very small migration rates, $R_{\text{migr}} N < r_m / (\frac{3}{10} \ln N)$, i.e., when the influx of bacteria with a newly required food gene ($\frac{3}{10} R_{\text{migr}} N$) is lower than the rate of the complete disappearance of the functional **t** gene ($r_m / \ln N$), migration loses its role in the propagation of NGT, thus, the metapopulation practically falls apart into isolated populations, in which transformation cannot survive. The other limit is $R_{\text{migr}} N > \binom{10}{3} r_r^{\max} / \ln \frac{r_r^{\max}}{r_t}$, i.e., when the migration is so intense that after a food change the rate of the arrival of a bacterium with the optimal combination of functional model genes ($R_{\text{migr}} N / \binom{10}{3}$) exceeds the rate at which the model gene that has just become beneficial proliferates and subsequently gets incorporated by a member of the original population ($r_r^{\max} / \ln \frac{r_r^{\max}}{r_t}$). Then the entire metapopulation effectively becomes a single giant population, in which transformation cannot survive either. Although these are rather crude estimates, which only give the limits outside of which transformation is certain not to survive, we have found none the less, that within these extremes NGT persists for most parameter values, indicating the robustness of our proposed mechanism.

2.7.4 The effects of NGT on genome organization

To demonstrate the dramatic effect of transformation on the genome composition, we take a cut through the parameter space (dashed-dotted line in Fig.2.7a) and plot the average fraction of bacteria possessing a functional copy of the **t** gene, as well as the fraction of those having exactly 1, 2, ... functional model food genes in Fig.2.7b. For low food change rate, but within the range where NGT survives, most bacteria only contain functional copies of the three necessary food genes, and the average fraction of functional copies of the model gene **t** is very low as it is used very rarely. Functional copies of the **t** gene can occasionally

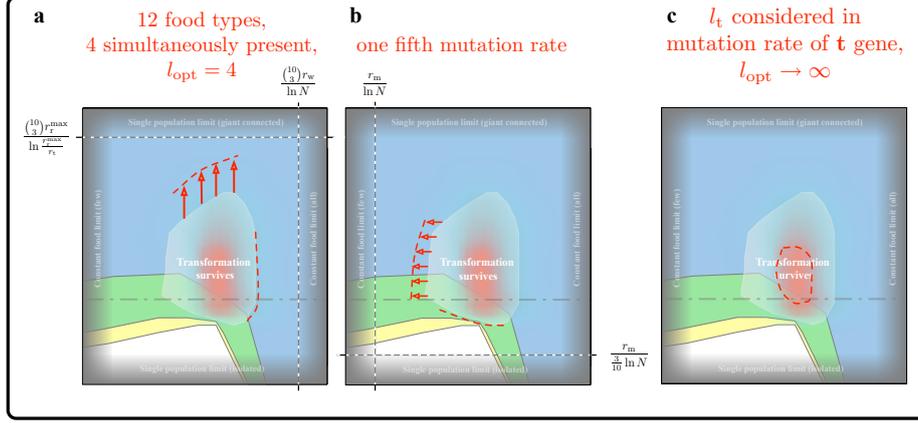


Figure 2.8: Robustness of the proposed mechanism for the survival of NGT. We investigated the robustness of the the survival range of NGT by varying several model parameters: In (a) we changed the number of possible food types from 10 to 12 and the number simultaneously present from 3 to 4, consequently also setting $l_{\text{opt}} = 4$. The survival range of NGT was found to extend to higher values of the migrational influx in accordance with our predictions, while it only extended to a much smaller degree toward faster food change rates. In (b) we have changed the mutation rate to $1/5$ of that considered in Fig.2.7. Under these conditions –similarly to (a)– the survival range of NGT extends in one direction – toward smaller food change rates, but does not significantly move in the other – toward smaller values of migrational influx. In the third panel (c) we demonstrate that the \mathbf{t} gene survives, although under a more limited set of parameters, even if selection for genome length is absent. During these simulations we gradually approached the limit $l_{\text{opt}} \rightarrow \infty$, where selection for genome length disappears. We also increased the mutation rate of food genes by an order of magnitude, while keeping that of the \mathbf{t} gene fixed. This way we were able to take into consideration, the length of the \mathbf{t} gene $l_{\mathbf{t}} = 0.1$ in terms of the relative mutation, while also keeping our simulations tractable.

disappear from the population and then subsequently be replanted by migrants as illustrated in Fig.2.6a. For larger food change rates, but still within the survival range of NGT, the \mathbf{t} gene becomes beneficial more often, its average fraction increases (except near the other end of the range), and the fraction of bacteria having four or more functional model food genes also increases. Two time series (at a low and a high food change rate) displaying the evolution of these fractions are shown in the lower panel of Fig.2.7b.

In the model we have chosen numerically tractable values for the rates ($r_{\mathbf{t}}^{\text{max}}$, r_w , r_m , $r_{\mathbf{t}}$) characterizing a given bacterial species. Some of these rates may be off by up to a few orders of magnitude for certain types of bacteria, this, however, should not fundamentally effect the width of the parameter range where NGT survives (neither in terms of R_{food} or R_{migr}), as this primarily depends on the binomial factor $\binom{10}{3}$. Moreover, under natural conditions the number of possible food types is usually much larger than 10, and the number of simultaneously available ones is also larger than 3, thus the survival range of NGT is probably even broader (extending to higher values of R_{food} and R_{migr}) than what our calculations

predict. We have attempted to survey the effect of more realistic, but computationally more demanding parameter sets as presented in Fig.2.8a and 2.8b, and found that the survival range of NGT indeed became broader. In Fig.2.8a we increased the number of possible food types from 10 to 12, while also increasing the number present at any on time from 3 to 4, consequently also setting $l_{\text{opt}} = 4$. In Fig.2.8.b we have decreased the mutation rate to one fifth of the value considered previously considered.

Characterizing transformation by a single rate coefficient r_t is also a rather crude simplification, since many bacteria become competent only under certain conditions. Fortunately, the limits for the survival of NGT calculated above depend very weakly on the value of r_t . Besides, the parameter range where a better regulated competence will persist is expected to be even larger than it is in our simplified model. Up to now we have only considered food sources that fluctuate with a single characteristic rate R_{food} . To check what happens in more complex situations, we have added an equal amount of constant food source (**XYZ**) to the fluctuating ones. These simulations have confirmed that the survival range of NGT remains virtually unaffected, indicating that the ability to lose and reload a few types of intermittently beneficial “operational” genes is advantageous to the bacteria and sufficient to maintain NGT.

Migration between bacterial populations in nature clearly depends on distance and, as a consequence, a wide variety of migration rates are usually present in the metapopulation. Populations close to each other (with strongly correlated environmental fluctuations or with intense inter-population migration), however, may be grouped together and considered as a large effective population. For NGT to persist it is sufficient that a reasonable number of such “effective” populations exist in the metapopulation, a requirement that does not seem unrealistic in face of the highly varied conditions under which bacteria prevail on Earth.

Finally one very important question remains which we have not addressed so far: is selection for shorter genome length necessary for the survival of NGT as a mechanism to reload genes? To answer this question we performed simulations where we gradually approached the limit $l_{\text{opt}} \rightarrow \infty$, where selection for genome length disappears. We found that if we took into consideration the length of the **t** gene l_t in terms of the mutation rate, the **t** gene survived (see Fig.2.8c). Because smaller mutation rates imply longer convergence times in our simulation, we implemented the value of $l_t = 0.1$, not by decreasing the mutation rate of the **t** gene, but by increasing those of the food genes by a factor of ten. From these results we can conclude that while selection for shorter genome length substantially increases the size of the parameter range where NGT is able to survive solely as a mechanism to reload genes, it is not indispensable. We may further argue that the relatively small

region in Fig.2.8c where the t gene persists, should become much larger for more realistic values of the mutation rate and the number of fluctuating food types (cf. Fig.2.8a and 2.8b).

2.8 Discussion

While it is, of course, clear that the reduction of any realistic population structure to a manageable construction is always an approximation, it has not been clearly established what the relevant degrees of freedom are in terms of evolutionary dynamics. Mean field approximations are a classic method of statistical and condensed matter physics, and are routinely used to circumvent intractable combinatorial problems which arise in many-body systems. Cluster mean field approximations of sufficient precision [45] have been developed that adequately describe the evolutionary dynamics of explicitly structured populations through systematically approximating the combinatorial complexity of the entire topology with that of small motif of appropriate symmetry. The effects of more minimal effective topologies have, however, not been investigated previously. In the above we have shown that straightforward hierarchical application of the mean field approximation (the assumption of a well-mixed system) surprisingly unveils a new level of complexity.

In the broader context of ecological and population genetics research on structured populations, our model can be described as a metapopulation model. The term “metapopulation” is, however, often used for any spatially structured population [49], and models thereof. More restrictive definitions of the term are often implied in the context of ecology and population genetics literature.

The foundations of the classic metapopulation concept were laid down by Levin’s vision of a “metapopulation” as a population of ephemeral local populations prone to extinction. A classic metapopulation persists, like an ordinary population of mortal individuals, in a balance between “deaths” (local extinctions) and “births” (establishment of new populations at unoccupied sites) [49]. This classic framework is most widespread in the ecology literature, a less often employed extension is the concept of a structured metapopulation, where the state of the individual populations is considered in more detail, this is more similar to our concept of hierarchical mixing, but differs in considering the possibility of local extinctions.

The effects of finite population size and migration, which our model considers, has been of more central concern in the population genetics literature. The analog of Levin’s classic metapopulation concept is often referred to as the “finite-island” model [50] the effective population genetic parameters describing which, have been explored in detail [94]. The

study of the population genetics of spatially subdivided populations in fact predates Levin, Wright having emphasised the capacity of drift in small populations to bring about genetic differentiation in the face of selection and/or migration several decades prior [50].

Our hierarchical mixing model treats the coevolutionary dynamics of evolutionary games on structured populations in a manner similar to the most simple population genetic models of spatially subdivided populations, focusing on the parallel effects of selection, drift and migration. It goes beyond these models both in considering the effects of frequency dependent selection (and the strategic aspects of the evolutionary dynamics this implies) and in using a self-consistent approach to describe the global state of the subdivided population. The effective population structure described by our hierarchical mixing model can be thought of as a population of individuals, interactions among which are specified by the edges of a hierarchically organized random graph. The fundamental difference in our picture is that the edges of this graph of interactions are not considered to be fixed, but are instead in a constant state of change, being present with a different probability between pairs of individuals who share the same local population and between pairs of individuals who do not (Fig.2.1.). We consider annealed randomness, which in contrast to the usual quenched picture of fixed edges is insensitive to the details of topology. Our approach we believe best facilitates the exploration of the effects of changing the relative strengths of drift and migration in the context of evolutionary games on structured populations.

Examining the effects of hierarchical mixing in the context of the evolution of cooperation we demonstrated that biased influx coupled with drift can result in cooperation being favored, provided the ratio of benefit to cost exceeds the local population size. This result bears striking resemblance to that of Ohtsuki et al. [47], who were able to calculate the fixation probability of a randomly placed mutant for any two-person, two-strategy game on a regular graph and found that cooperation is favored provided the ratio of benefit to cost exceeds the degree of the graph. Our results demonstrate that this rule extends to the minimal spatial structure induced by hierarchical levels of mixing.

Applying our model of spatial structure to the repeated prisoners dilemma revealed that a constant influx of defectors can help to stabilize cycles of cooperation, defection, and reciprocity through preventing the emergence of an intermittent period of ALLC domination in the population, which would present a situation that “leaves the door wide open” to domination by defectors. While previous work has been done on the effects of “forcing” cooperation [95], the idea that an influx of defectors can in fact stabilize the role of reciprocity in promoting cooperation has not been proposed previously. It seems highly unlikely that this mechanism can be explained in terms of kin or multilevel (group) selection,

the similarities between which in structured populations have recently been the subject of intensive debate (see e.g. Refs. [96] and [97] or Refs. [98] and [99]). Kin selection can operate whenever interactions occurring among individuals who share a more recent common ancestor than individuals sampled randomly from the whole population [99] are relevant. In our case it is the interaction between defectors, arriving from the global scale, and TFT players present at the local scale that is important, and not the interaction between individuals in the local population, who may be thought of as sharing a recent common ancestor due to local dispersal. Also, while the concept of multilevel selection presents a promising framework for the study of evolution of cooperation, it must nonetheless be possible to derive it from “first principles” – just as kin selection can be cast as an emergent effect of local dispersal.

While there has been considerable work on studying the evolutionary games on graphs and highly symmetric spatial structures, very little attention has been paid to the effects of more minimal effective population structures, despite their widespread application in ecology and population genetics, fields from which evolutionary game theory was born and must ultimately reconnect with. We believe that the minimal population structure that such a hierarchical mean field theory describes is potentially more relevant in a wide range of natural systems, than more subtle setups with a delicate dependence on the details and symmetries of the topology. We showed through several examples that such structure is sufficient for the emergence of some phenomena previously only observed for explicit spatial embedding, demonstrating the potential of our model to identify robust effects of population structure on the dynamics of evolutionary games that do not depend on the details of the underlying topology. The practical advantage of our approach, lies in its ability to readily determine whether or not some feature of a structured population depends on the topological details of local interactions.

Straying farther afield, we have applied our minimal model of population structure coupled with environmental fluctuations the problem of the maintenance of sex in bacteria. In light of simulations described above (and presented in our publication Ref. [51]), we believe that the existence of NGT is facilitated by its role as a vehicle to reload genes, as both the necessary minimal population structure and environmental fluctuations can reasonably be assumed to be present in Nature. We argue that the short-term advantage that sustains NGT long enough for its evolutionary effects to emerge, lies, at least in part, in providing mobility to variably selected genes. It allows individuals to reload genes lost from a population – due to long disuse – but still available in the metapopulation, bringing together genes from the collective gene pool of the species with locally adapted genomes.

This advantage prevails if spatio-temporal fluctuations [100] in the environment (imposing variable selection pressure on different populations of the same species) exist in parallel with weak migration between the populations (allowing genetic mixing). Whether or not natural bacterial populations actually experience the kind of population subdivision and inter-population migration necessary for our model to be applicable remains to be demonstrated experimentally. Some examples which may easily fulfill these conditions, however, readily come to mind, e.g. experience shows that any perishable substance that is a potential food source for bacteria is promptly colonized; one may also consider the intestinal flora of grazing animals, herds of which cover large distances while occasionally encountering each other at locations, such as water sources, where migration may occur between the microbial populations resident in their intestines.

Provided that the above conditions may be rather general, our results compel us to imply that the ability of active DNA uptake may easily have evolved through the gradual specialization of a more general transport mechanism [83], driven not by the need to serve the cell with nucleotides for food, but by an advantage conferred through providing homologous sequences for restoring genes eroded by formerly adaptive deletions. In other words, NGT is not an accidental byproduct of nutrient uptake, but has come into existence in order to counterbalance gene loss, which inevitably occurs in highly economized genomes under fluctuating selection pressure. It should be emphasized, however, that this does not exclude the advantage transformation confers through enabling access to DNA as a source of nutrients, which most probably helps sustain NGT.

A fundamental long-term effect of genetic mixing by NGT, with important implications for the understanding of the evolution of eukaryotic sex [80, 5, 101], is that it prevents bacterial species from falling apart into independently evolving clonal lineages [102], by facilitating genetic mixing between genomes of sufficient homology. The dynamic nature of bacterial genomes has been long been suggested by the rapid spread of antibiotic resistance and other pathogenic traits [91]. Recent efforts to assay microbial diversity by environmental sequencing have strikingly revealed³, however, that we have only so far glimpsed the tip of the iceberg in terms of prokaryotic diversity using laboratory cultures. It is even more striking that this observation extends to the level of genes. Predicted Open Reading Frames in the GOS data set⁴ have nearly doubled the number of known protein sequences, and on-

³*Metagenomics* studies sample the genome sequences in a given environment. In the last decade such studies [103] have lead to estimates, based on ribosomal RNA (rRNA) genes taken directly from the environment, that cultivation methods find less than 1% of the bacteria and archaea species.

⁴The Global Ocean Sampling expedition, initiated by J. Craig Venter, sampled a several-thousand km transect from the North Atlantic through the Panama Canal and ending in the South Pacific.

going sequencing continues to add new *gene families* at a linear rate [104] – indicating that we are only in the early stages of essaying the diversity of life on the planet. Refining our concept of a prokaryotic species, in light of the evidence for environment specific genomic diversity, frequent horizontal transfer and seemingly open pan genomes, looks to be central to this endeavour.

Chapter 3

The mapping between genotype and phenotype

"It is time for evolutionary theory to catch up with empirical paleontology, to confront the phenomena of evolutionary non-change, and to incorporate it into our theory.. "

S. J. Gould and N. Eldredge, Punctuated equilibria: the tempo and mode of evolution reconsidered. *Paleobiology* **3**, 115 (1977)

3.1 Introduction

Our concept of evolutionary change is based on two aspects of biological organization: phenotype and genotype. An organism's phenotype is the sum of its physical, organizational and behavioral manifestations during its lifetime. An organism's genotype is the heritable repository of information that instructs the production of molecules whose interactions, in conjunction with the environment, generate and maintain the phenotype [3]. Selection acts on phenotypes, causing the amplification of underlying genotypes in the population that give rise to phenotypes with higher fitness. Variation, however, is generated by mutations of the genotype. Consequently, evolutionary change through natural selection requires phenotypic variation among organisms that reflects genetic variation. It has been clear for some time that genetic variation is abundant in most species, how it translates into phenotypic variation, however, remains mostly unknown. The dynamics of selection (cf. section 1.2), by itself, has nothing to say about how evolutionary innovations come about. A model of the genotype-phenotype relationship is needed, to illuminate how genetic change maps into phenotypic change – to attain an understanding of the process of evolutionary change

through natural selection.

In general, the mapping between genotype and phenotype is accomplished through a range of complex internal processes and interactions with the environment, collectively called development [3]. Despite the complexity and diversity of developmental processes, simple and potentially universal features of this mapping have been identified that explain fundamental features of evolutionary phenomena; features that do not result naturally from an adaptations framework that considers phenotypes to be completely malleable by natural selection. Grounding patterns of phenotypic evolution in biophysical principles and mechanisms, we can hope to uncover constraints imposed on evolutionary trajectories of change by the process that translates genetic variation into phenotypic variation – by development [3].

The simple assumption that the genotype-phenotype map is degenerate, i.e., that there are many more genotypes than (distinguishable) phenotypes results in an emergent evolutionary dynamics that reconciles the concept that natural selection is the driving force of evolution with two of the most striking and fundamental observations regarding the evolution of life on Earth. With (i) *punctuated equilibrium*, the paleontological observation that most species, during their geological history, either do not change appreciably, only fluctuate mildly in morphology [7, 105, 106] and with (ii) *the neutral model of molecular evolution*, the observation at the molecular-level that the majority of evolutionary change is caused by random fixation of selectively neutral or nearly neutral mutants rather than positive Darwinian selection [12, 105, 107]. Further consideration of the topology of "accessibility", the distribution of evolutionary paths in genotype space connecting any two phenotypes, induced by model genotype-phenotype maps, has also been used to explain constraints to variation, the origin of novelties and directionality in evolution [3, 106].

Our pursuits here are more modest. If we accept that the road to evolutionary change is paved by long episodes of stasis, this implies that the potential for evolutionary innovation depends on the details of this micro-evolutionary stationary state determined by a highly degenerate genotype-phenotype map. Consequently, to understand the potential for innovation, we must understand the effects of genetic variation on a phenotype that selection acts to keep invariant, keeping in mind that it is neutral mutations that prepare the ground for later evolutionary adaptation during intermittent periods of phenotypic stasis. This endeavor to understand the invariance of phenotype in the face of variation in genotype (called genetic robustness) will be the focus of the remaining part of this chapter. We describe first van Nimwegen et al.'s seminal results on the neutral evolution of genetic robustness, followed by our own contributions. In section 3.3 we describe our results pub-

lished in Ref. [108] concerning the effects of recombination on the evolution of genetic robustness. In section 3.4 we reexamine Borenstein and Ruppin's results [109] and present evidence (published in Ref. [110]) that genetic robustness observed in miRNA sequences is the byproduct of selection for environmental robustness.

3.1.1 Genetic robustness

The RNA folding map is characterized by a number of remarkable statistical regularities with profound evolutionary consequences.

Walter Fontana, "The Topology of the Possible" in: "Understanding Change: Models, Methodologies and Metaphors." 2005, Macmillan, UK

Robustness is the invariance of phenotypes in the face of perturbation [111]. We must distinguish between robustness with respect to two types of perturbations: heritable and non-heritable. Invariance in the face heritable mutations, called genetic robustness, will be our primary concern here. As we will show, it can arise either as an emergent property of the evolutionary dynamics or as a byproduct of selection for non heritable perturbations, i.e., environmental robustness. Genetic robustness, the preservation of an optimal phenotype in the face of heritable perturbation (point mutations, deletions, insertions etc.), is critical to the understanding of evolution as phenotypically expressed genetic variation is the fuel of natural selection.

Empirical evidence has long been present that the magnitude of genetic effects on phenotype depends strongly on genetic background, the effects of the same mutation can be larger in one genetic background and smaller in another. The idea that wild-type genotypes are mutationally robust, i.e., show invariance in the face of mutations (more generally heritable perturbations), goes back to Waddington [112], who originally introduced the concept as canalization. Genetic robustness has been found across different levels of organization from individual genes, through simple genetic circuits [113, 114] to entire organisms (approximately 80% of yeast single knockouts have no obvious effect in rich medium [115]). The origin of the observed robustness has, however, remained a source of contention. The three main hypotheses regarding the potential origin of genetic robustness predate the concept itself, and fall along the lines of the famous debate between members of the modern synthesis (in particular Wright, Haldane and Fisher) surrounding the origin of dominance¹

¹Dominance can be understood as a simple case of robustness at the level of a single gene. A dominant phenotype is more robust against mutational (and perhaps environmental) perturbations than a recessive phenotype.

[111, 116]:

1. The most straightforward explanation, favoured by Wright, was that robustness evolves *directly*, through natural selection [117].
2. An alternative *congruent* hypotheses, put forward in the context of dominance by Haldane, proposes that the evolution of genetic robustness is a correlated byproduct of selection for environmental robustness, i.e., invariance in the face of nonheritable perturbations, e.g. temperature, salinity or internal factors such as fluctuations in the concentration of gene products during development [118].
3. A third view holds that genetic robustness is *intrinsic*, arising simply because the buffering of a character with respect to mutations is the necessary or likely consequence of character adaptation, in the context of dominance Wright [119] and later Kacser and Burns [120] argued that it arises as an inevitable, passive consequence of enzyme biochemistry and selection for increased metabolic flux².

An important stepping stone in solving this problem is establishing a theoretical understanding of under what conditions and to what extent natural selection can lead to the evolution of robustness. The question of which mechanism gave birth to and maintains the genetic robustness observed in biological systems must, however, ultimately be resolved by using empirical means – using actual biological data.

Several theoretical and simulation studies have addressed robustness in a wide range of contexts ranging from gene redundancy [121] to model regulatory networks [122, 123, 113, 124, 125]. Two pioneering studies, by Montville et al. [126] and Borenstein and Ruppin [109] have managed to step beyond computer simulations. Borenstein and Ruppin looked for evidence of excess mutational robustness present in RNA secondary structure [127], while Montville et al. relied on the expectation that high mutation rates present among RNA viruses should favour mutational robustness [128]. They used, respectively, microRNA sequences from diverse taxa and *in vitro* evolution experiments [126] to find evidence to support the hypotheses that genetic robustness can evolve directly. Further work on *in vitro* evolution experiments has provided additional evidence showing that if a population is highly polymorphic robustness can evolve directly [129, 130].

The theoretical underpinnings of these studies is provided by the results of van Nimwegen et al. (see Ref. [131] and the following section), who solved the quasispecies equations

²Adaptation to utilize transiently available substrates can be argued to result in selection for such high enzyme concentrations and/or catalytic efficiencies that render the effects of reducing enzyme concentrations by half – due to harboring the recessive allele – negligible under ambient conditions.

describing the evolution of a population on a network of phenotypically neutral sequences. They demonstrated (as described in detail in section 3.2), that provided a sufficiently polymorphic population, mutational robustness can evolve directly. The necessary mutation rates and/or population sizes were found to be very large in simulation studies using RNA secondary structure as a genotype-phenotype map [131, 132, 108], direct evolution of increased neutrality requiring the product of the effective population size N and the mutation rate per nucleotide u to be well in excess of one. Such high mutation rates can only readily be found among RNA viruses, are extraordinary even among unicellular organisms (*Prochlorococcus*³ $2N\mu \approx 2.$, *E. coli* $2N\mu \approx 0.2$, *S. cerevisiae* $4N\mu \approx 0.09$) and completely unheard of among multicellular eukaryotes possessing RNA silencing mechanisms and microRNA genes (*A. thaliana* $4N\mu \approx 0.012$, *D. melanogaster* $4N\mu \approx 0.015$, *C. elegans* $4N\mu \approx 0.013$, *C. intestinalis* $4N\mu \approx 0.012$, *M. musculus* $4N\mu \approx 0.001$, *H. sapiens* $4N\mu \approx 0.001$) [134].

3.2 The neutral evolution of genetic robustness

According to the classic results of Haldane [135] later extended by Kimura and Maruyama [136] the average fitness of an asexually reproducing population (in the limit of very large populations) depends only on the mutation rate and is independent of the details of the fitness landscape [136]. This result, however, as stated in Kimura and Maruyama's paper, only holds under the assumption that the fittest genotype does not have any neutral sites. If neutral mutations are taken into account the average fitness of the population will depend on both the mutation rate and the details of the fitness landscape [128]. A growing body of work has explored the effect of neutral mutations on mutation-selection balance in infinite populations (quasispecies) and the balance of mutation-selection and drift in finite populations. Quasispecies theory [128] and simulations of finite populations of genotypes evolving with phenotype defined by RNA secondary structure [127, 137] or simple lattice models of protein folding [138] have established that a selective pressure to evolve robustness against mutations exists. The net effect of this selection pressure is to concentrate the population in regions of genotype space where the density of neutral sequences is higher, selecting individual sequences with an increased robustness against mutations.

In order to demonstrate how this counterintuitive effect occurs, let us consider a space

³*Prochlorococcus* is a genus of very small ($0.6 \mu m$) marine cyanobacteria and is thought to be the most plentiful species on Earth [133]. A single millilitre of surface seawater may contain 100,000 cells or more. Estimates have put the number of individuals at 10^{29} worldwide, responsible for up to 20% of the oxygen production of the biosphere.

of discrete macromolecular sequences (genotypes) of L monomers, e.g. DNA or RNA molecules. In this space of genotypes, two sequences are “neighbors” if they only differ in their sequence by a single point mutation – the change of a single monomer. The phenotype of a sequence is some functionally relevant characteristic of the encoded sequence, e.g. the secondary structure of a microRNA stem-loop, the binding affinity of a transcription factor in the case of transcription binding site, or catalytic efficiency in the case of a amino acid sequence encoding some enzyme. In general, there are many more genotypes than (distinguishable) phenotypes. A network of genotypes that share a common phenotype we will refer to as a *neutral network*. Phenotypes have *neutral networks* of different size and topology [139, 140].

Let us follow van Nimwegen et al. [131] and focus on the case where selection acts to maintain some phenotype. In this case all mutations are either neutral (if they do not lead off the neutral network of the preferred phenotype) or significantly deleterious (if they lead off the neutral network of the preferred phenotype). The course of evolution under such circumstances will clearly be in keeping with the neutralist picture of evolution, wherein beneficial mutations are rare and the majority of genetic variation is either neutral or at worst weakly deleterious [12]. In fact, in our simplified picture all observed evolutionary change, all fixed mutations, will be neutral. We will not be concerned here with the next macro-evolutionary step leading away from this phenotype, but only the distribution of the population during an intermittent period of equilibrium. It seem natural to assume that on long time scales the population will be uniformly distributed over all available sequences – the entire neutral network. As we demonstrate below, this is in fact not the case, population dynamics will, given sufficient sustained variation is present, concentrate the population on highly connected regions of the neutral network, on genotypes that exhibit increased tolerance to mutations.

In some genotype space that consists of all sequences of length L over some finite alphabet \mathcal{A} of A symbols (e.g. for RNA $\mathcal{A} = \{A, U, C, G\}$ and $A = 4$). The neutral network belonging to the phenotype preferred by selection can be regarded as a graph G embedded in the space of genotypes. The vertex set of G defined by all genotypes in the neutral network and two vertices are connected by an edge if they are neighbors (i.e., differ at only a single position).

In their treatment van Nimwegen et al. assume as their evolutionary process a discrete generation selection-mutation dynamics with constant population size N , that is in effect identical to the Wright-Fisher process: each generation N new individuals are chosen with replacement and probabilities proportional to their fitness from the previous generation.

Each chosen individual suffers a mutation with probability μL , i.e., one of the L symbols is chosen with uniform probability for replacement by one of the $A - 1$ other symbols.

3.2.1 The infinite population limit

In the limit of infinite population size it is possible to explicitly solve for the asymptotic distribution of the population over the neutral network G of size M . We know that once the population has reached a stationary state there will be: (i) a stationary proportion of individuals X in each generation with genotypes that belong to G and (ii) a stationary average fitness $\langle s \rangle$. In order to dissect the balance between selection and mutation that defines this steady state let us consider separately the effects of these two forces separately. Under selection alone the proportion of individuals with genotypes that belong to G would increase from X to $sX/\langle s \rangle$, where s is the selective advantage of the preferred genotype. Mutation by itself would cause some fraction $\langle \nu \rangle$ of the N total individuals to: (i) to remain on the neutral network (i.e., retain the phenotype preferred by selection, with or without a change in genotype), while (ii) some fraction $1 - \langle \nu \rangle$ would “fall off” (i.e., mutate to some genotype with a phenotype distinct from that preferred by selection). Further, some fraction Y of the N individuals in the population, which are currently among the $(1 - X)M$ not on G , would mutate back onto the neutral network (i.e., mutate to some genotype with the phenotype preferred by selection). In the stationary state this implies:

$$X = \frac{s}{\langle s \rangle} \langle \nu \rangle X + Y \quad (3.1)$$

If the mutation rate is not too large and selective advantage of the preferred phenotype is larger than the inverse of the population size, i.e., $s > 1/N$, selection dominates drift and we may in general neglect the fraction of back mutations y simply due to the fact that $1 - X \approx 0$. This leads to

$$X = \frac{s}{\langle s \rangle} \langle \nu \rangle X \quad (3.2)$$

In order to connect this formula with the topology of the neutral network let us turn to some individual with a genotype $v \in G$ and vertex degree d_v (the number of neighbors that have the preferred phenotype). The probability of remaining on G depends only on the vertex degree d and is given by

$$\nu(d) = 1 - \mu L \left[1 - \frac{d}{(A - 1)L} \right], \quad (3.3)$$

If in the stationary state a fraction x_i of the population is located on vertex v then we may express the fraction of individuals that remain on the neutral network by averaging over the population fractions for all $i \in G$

$$\langle \nu \rangle = \sum_{i \in G} \nu(d_i) \frac{x_i}{X} = 1 - \mu L \left[1 - \frac{\sum_{i \in G} d_i x_i / X}{(A-1)L} \right] \quad (3.4)$$

from which by using Eq. (3.2) we get the expression

$$D = \sum_{i \in G} d_i \frac{x_i}{X} = L(A-1) \left[1 - \frac{s - \langle s \rangle}{\mu L s} \right]. \quad (3.5)$$

that relates the average vertex degree D and average fitness $\langle s \rangle$.

Using the above we may obtain a complete description of the stationary state in the form of $\{x_i\}$, the stationary population fractions on G in the stationary state by solving the M equations

$$x_i = (1 - \mu L) \frac{s}{\langle s \rangle} x_i + \frac{\mu L}{(A-1)L} \sum_{j \in [i]_G} \frac{s}{\langle s \rangle} x_j, \quad (3.6)$$

where $[i]_G$ denotes the neighbors of i in G . Introducing the adjacency matrix of G : for any $i \in G$

$$G_{ij} = \begin{cases} 1 & : j \in G \\ 0 & : j \notin G \end{cases} \quad (3.7)$$

and using Eqs. (3.5) and (3.2) we may rewrite Eq. (3.6) as the eigenvector equation

$$D x_i = G_{ij} x_j. \quad (3.8)$$

As G is none negative and the neutral network is connected according to the Peron-Frobenius theorem, the stationary population fractions $\{x_i\}$ are given by the principle eigenvector of the adjacency matrix of G , while the average neutrality D is equal to the spectral radius R of the adjacency matrix of G – this is the central result of van Nimwegen et al.

We can compare van Nimwegen et al.'s result to the classic results of Haldane on genetic load [135]. Haldane showed that the genetic load \mathcal{L} , i.e., the relative reduction in average fitness that results from mutation, is equal to the proportion of deleterious mutations per generation. The above result implies

$$\mathcal{L} = \frac{s - \langle s \rangle}{s} = \mu L \left[1 - \frac{R}{(A-1)L} \right], \quad (3.9)$$

that is the genetic load is reduced by a factor that depends on only the spectral radius of the adjacency matrix of G , and we recover the classic result in the absence of neutrality ($R = 0$).

3.2.2 Finite populations

It is clear that stationary distribution $\{x_i\}$ derived in the infinite population can not remain a good approximation for arbitrarily small population size. Let us consider the case where $\mu N \ll 1$. In this limit all individuals in the population will at any given time, with high probability, share the same genotype. Each generation with probability $\mu L N$ a mutant will enter the population, if this mutant is not on G it will rapidly go extinct, if, however, it is on G it will with probability $1/N$ take over the population (cf. section 1.2). In the event this takeover occurs the population will have completed a step of a random-walk on G , a random walk that attempts a step at a rate μL in one of the $L(A - 1)$ directions around v , but only accepts steps that keep the population on G . This is the random walk of a “blind ant” that, starting from some sequence that is part of the neutral network G , at each time step chooses one of the $L(A - 1)$ random neighbors and attempts to make a step onto this sequence. Our blind ant, however, only completes the step if this new node is also in G , if this not the case the blind ant remains in place. It can easily be shown that the limiting distribution of the blind ant random walk on G is uniform, i.e., asymptotically spends the same time on all $v \in G$. The asymptotic average neutrality in this case corresponds to the average degree of G :

$$D_0 = \frac{1}{M} \sum_{i \in G} d_i \quad (3.10)$$

To develop a qualitative understanding of the effects of finite population size in the crossover regime $\mu N \approx 1$, where finite population size effects are still substantial, but the population is already polymorphic a significant fraction of the time, let us, following van Nimwegen et al., consider a second type of random walk, that of a “myopic ant”. The myopic ant, starting from some sequence $v \in G$ that is part of the neutral network G , will before each step look around to see which neighbors are in G and choose one of these d_v sequences to step onto. The limiting distribution of the myopic ant random walk on G is such that the asymptotic fraction time spent on a given $v \in G$ is proportional to the degree of neutrality ($x_v \propto d_v$).

In the case of finite population sizes in the crossover regime will, as demonstrated by numerical simulations (cf. Ref. [131, 132, 108] and section 3.3.2 below), average neutrality will be reduced in comparison to the infinite population size limit (will be less than the

spectral radius, i.e., $D < R$). Similar to the myopic random walk introduced above the population will only “see” a local region of the neutral network.

3.3 The effects of recombination

In the following we examine the effects of recombination on the evolution of mutational robustness on networks of neutral genotypes and its dependence on the topology of the network presented in our publication [108]. We consider separately the effects of recombination in the limit of infinite populations and for finite populations.

Previous work on the effects of recombination on the population dynamics on networks of neutral genotypes is scarce, Xia *et al.* [141] have shown that in a simple lattice model of protein folding recombination leads to increased thermodynamic stability, but have not addressed the evolution of mutational robustness. In a wider context the effects of recombination on the evolution of robustness has been studied in highly simplified model developmental networks [123], where simulation results suggest that recombination between model gene networks imposes selection for mutational robustness, and that negative epistasis evolves as a byproduct of this selection. Even though the evolution of mutational robustness in networks of neutral genotypes under recombination has not been previously studied in detail, the observation that recombination has a contracting property (i.e., it always creates genotypes that are within the boundaries of the current mutational cloud) has led to the expectation that it should concentrate the population in those regions of genotype space in which the density of neutral genotypes is highest [128].

In order to gauge the most general effects of recombination, we can proceed by defining ensembles of random networks of neutral genotypes. As we show below, this allows us to demonstrate that recombination leads to the concentration of the population in highly neutral region of the genotype space, and hence the evolution of mutational robustness, in these generic neutral networks. Following the line of thought developed in the previous section we subsequently turn to finite populations effects. Similar to the case presented in the previous section, where only mutation was present (cf. also [131]) the evolution of mutational robustness is found to require a sufficiently polymorphic population (i.e., the product of the population size and the mutation rate, μN , must be greater than one). Setting the stage for the next section we use a scaled down analog of microRNA stem-loop hairpin structures. We demonstrate that only provided a sufficiently polymorphic population, will recombination lead to the evolution of mutational robustness on larger and more realistic neutral networks.

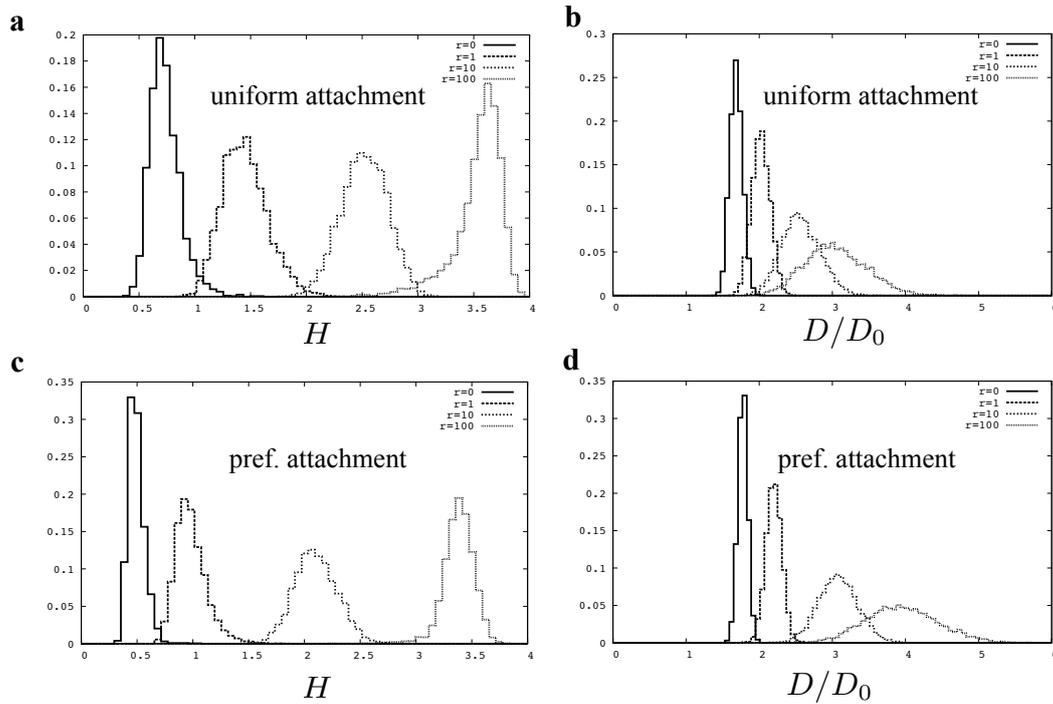


Figure 3.1: Histograms of the entropy H and mutational robustness enhancement D/D_0 for different values of r . Numerically calculating the stationary distribution of the population on 10^5 neutral networks $M = 200$ genotypes of length $L = 20$ randomly drawn from the uniform attachment ensemble (**a**, **b**) and preferential attachment ensembles (**b**, **c**) indicated (cf. Ref. [108]) that recombination leads to significant enhancement of mutational robustness under very general conditions. Comparison of the results for the two ensembles suggests that preferential attachment networks, where genotypes of higher centrality are more neutral, evolve higher levels of mutational robustness.

Considering, similar to section 3.2 the genotype space of sequences of length L over some alphabet \mathcal{A} we assume that genotype space contains a neutral network of high, but equal fitness genotypes – those sharing a common preferred phenotype. We assume, further, that the majority of the population is concentrated on this neutral network, i.e., that the remaining part of genotype space consist of genotypes with markedly lower fitness, such that to good approximation all such phenotypes maybe considered lethal. While making this assumption from the outset may seem more strict than the approach in the previous section, it is in fact in one to one correspondence with the assumption of $Y \approx 0$ in section 3.2.1. For our evolutionary process we assume a selection-mutation-recombination dynamics with constant population size N . We consider all mutations or recombination events leading off the neutral subset of genotypes G with high fitness to be fatal. Each individual of the population suffers mutations at a rate μL and undergoes recombination with a random member of the population at a rate ρL . Individuals who acquire a fatal genotype (one not part of the neutral network) are replaced through reproduction of a random individual of the population.

In the limit of large populations the stationary distribution of the population on a given neutral network G only depends on the ratio $r = \rho/\mu$ and the neutral network G . Considering one-point recombination crossover events and denoting the number of neutral genotypes in G by M and the frequency of genotype $i \in \{0, M\}$ in the population by x_i the time evolution of our system in the $N \rightarrow \infty$ limit is given by:

$$\dot{x}_i = \sigma x_i - (\mu L + \rho L)x_i + \mu \sum_{j=1}^M G_{ij}x_j + \rho \sum_{k=1}^M \sum_{l=1}^M R_{ikl}x_kx_l, \quad (3.11)$$

where the mutation operator is defined as $M_{ij} = 1$ if genotype i can be derived from genotype j by replacing a single letter and zero otherwise, while the recombination operator R_{ikl} equals the number of recombination (one-point crossover) events through which genotypes k and l yield genotype i , and

$$\sigma = (\mu L + \rho L) - \sum_{i=1}^M \left(\mu \sum_{j=1}^M M_{ij}x_j + \rho \sum_{k=1}^M \sum_{l=1}^M R_{ikl}x_kx_l \right). \quad (3.12)$$

is the uniform growth rate of the population, that compensates for the disappearance of lethal genotypes generated by mutation and recombination, ensuring that the population size remains constant, i.e., $\sum_i \dot{x}_i = 0$ at all times. We compute the limit distribution of the population ($x_i^{\text{st.}}$) by numerically computing the stationary solution of equations (3.11).

For finite N stochasticity resulting from the discrete nature of the reproduction process must also be taken into consideration. Considering discrete generations (see section 1.2), we may proceed by solving equation (3.11) while introducing sampling noise (i.e., drift) by sampling the population at unit time intervals, the mutation and recombination rates will then be in units of events per generation (see Methods section 4.2). The resulting stochastic population dynamics depends on the product of the mutation rate and population size μN , the ratio of the mutation and recombination rates r , and the structure of G .

Computing R_{ikl} requires enumerating all $M \times M \times L$ possible recombination events. This is only feasible for M not larger than a few thousand genotypes. For larger neutral networks calculation of the infinite population limit stationary distribution is not currently tractable, it is still possible, however, to simulate the finite population dynamics as the complexity of these computations scales with $\mathcal{O}(N^2L)$ and not $\mathcal{O}(M^2L)$. In these simulations we average over several runs starting from random initial conditions.

3.3.1 The infinite population limit

In section 3.2.1 we have shown, following van Nimwegen et al. [131], that in the limit of large populations the average number of neutral single mutant neighbors in a population in selection-mutation balance is larger than the average number of neutral neighbors in the network. The population tends to concentrate in parts of the network with enhanced neutrality. If we introduce recombination and numerically compute the stationary distribution of the population in selection-mutation-recombination balance using equation (3.11) we observe that recombination concentrates the population even further. This observation can be quantified by comparing the entropy of the stationary population distributions $x_i^{\text{st.}}$ defined by:

$$H = \ln \frac{1}{M} - \sum_{i=1}^M x_i^{\text{st.}} \ln x_i^{\text{st.}}, \quad (3.13)$$

for different values of r . To assess the mutational robustness of the population we have to compute the average number of neutral neighbors of a random individual in the population:

$$D = \sum_{i=1}^M x_i^{\text{st.}} d_i, \quad (3.14)$$

where d_i is the number genotypes in the neutral network that can be obtained from genotype i by replacing a single letter. To obtain a measure of the extent to which excess mutational robustness emerges solely as an effect of the population dynamics we have to compare this

value to the average neutrality of the network:

$$D_0 = \frac{1}{M} \sum_{i=1}^M d_i. \quad (3.15)$$

To compute the above averages in Ref. [108] we generated 10^5 random neutral networks with $M = 200$ and $L = 20$ of both types and averaged H and D/D_0 over them. We also generated networks with larger M (and L) values and found qualitatively similar results. Due to the finite nature of both the mutation rate and population size in natural populations (together quantified by μN) it is networks of relatively small M that are most relevant biologically, in the sense that natural populations with finite μN are in effect restricted to some relatively small region of the neutral network for time scales long enough for local selection-mutation-recombination balance to be achieved. This implies that the extent to which robustness is evolved is determined by local topology. The exploration of the entire neutral network occurs on a much longer time scale. This long time scale exploration is of secondary importance from the perspective of the evolution of mutational robustness as the topology of the neutral network beyond the mutational cloud is in effect invisible to the population dynamics.

As shown in Fig.3.1.a. and Fig.3.1.c. recombination leads to a similar increase in entropy in both uniform attachment and preferential attachment random networks. The entropy increase as a function of r is slightly less significant for preferential attachment networks $H = 0.726 \pm 0.128, 0.982 \pm 0.14, 2.102 \pm 0.19, 3.379 \pm 0.143$ for $r = 0, 1, 10, 100$, respectively, then for uniform attachment networks, $H = 0.492 \pm 0.76, 1.449 \pm 0.202, 2.518 \pm 0.207, 3.541 \pm 0.219$ for $r = 0, 1, 10, 100$, respectively.

If we look, however, at the increase in average mutational robustness (Fig.3.1.b. and Fig. 3.1.d.) we find that as a function of r it is significantly larger for preferential attachment networks $D/D_0 = 1.775 \pm 0.067, 2.200 \pm 0.108, 3.079 \pm 0.268, 3.951 \pm 0.494$ for $r = 0, 1, 10, 100$, respectively, then for uniform attachment networks, $D/D_0 = 1.684 \pm 0.087, 2.024 \pm 0.125, 2.569 \pm 0.260, 3.045 \pm 0.422$ for $r = 0, 1, 10, 100$, respectively. This suggest that in preferential attachment networks, were the expected neutrality of genotypes with high centrality reaches higher levels, evolve higher levels of mutational robustness.

3.3.2 Finite populations

Similar to section 3.2.2, where only mutation was considered, the evolution of mutational robustness in the presence of recombination requires that the population be sufficiently

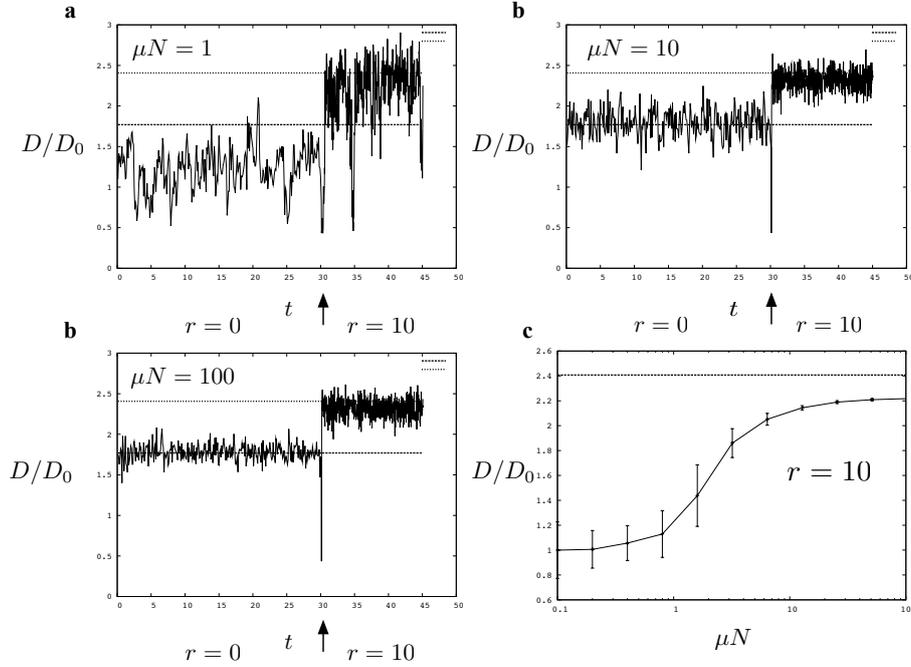


Figure 3.2: Simulations on a random uniform attachment network with $M = 200$, $L = 20$ with different values of μN (cf. Ref. [108]) show that the extent to which mutational robustness is evolved increases as μN becomes larger. **(a-c)** snapshots of the time evolution of mutational robustness enhancement D/D_0 in a population evolving on the same random uniform attachment network with different values of μN , time is indicated in units of $2N$ generations. Recombination was turned on at $t = 30$ (indicated by the arrow). **(d)** mutational robustness enhancement D/D_0 as a function of μN for the same network with $r = 10$ as calculated from 100 simulations with random initial conditions, where the population was allowed to evolve for $100 \times 2N$ generations, the error bars indicate the variance in the time averages over runs with random initial conditions. The dashed line indicates the value of the mutational robustness enhancement D/D_0 in the $N \rightarrow \infty$ limit, throughout.

polymorphic, i.e., that the product of the mutation rate and the population be greater than unity. If, however, this condition is satisfied the presence of recombination leads to the evolution of increased mutational robustness under rather general circumstances. To quantify the extent to which mutational robustness increases in Ref. [108] we numerically calculated the time average over the stochastic population dynamics. Due to the separation of the time scales at which the population attains local selection-mutation-recombination balance and the much longer time-scale over which it explores the entire neutral network, averaging over a set of random initial conditions was used to attain a computationally tractable approximation of the long time average.

Performing simulations (presented in Ref. [108]) for different values of μN showed that D/D_0 approaches its infinite population as μN is increased (see Fig.3.2.). This indicates that recombination concentrates the population on local regions of higher neutrality for smaller μN as well. In order to examine the effects of recombination on a more realistic neutral network in Ref. [108] we also performed simulations on a scaled down version of a microRNA stem-loop (Fig.3.3.). We found that for all values of μN the extent to which the population evolves mutational robustness is higher in the presence of recombination than in its absence.

3.4 Congruent evolution of robustness in microRNA

One of few empirical results on the evolutionary origins of genetic robustness was recently presented by Borenstein and Ruppin [109]. In their study Borenstein and Ruppin examined microRNA (miRNA) precursor sequences from several eukaryotic species, and found significantly increased mutational robustness in comparison with a sample of random RNA sequences with the same stem-loop structure. The observed robustness was found to be uncorrelated with traditional measures of environmental robustness – implying that miRNA sequences show evidence of the direct evolution of genetic robustness. These findings are surprising as theoretical results indicate that the direct evolution of robustness requires high mutation rates and/or large effective population sizes only found among RNA viruses, not multicellular eukaryotes.

3.4.1 MicroRNA stem-loops as a molecular phenotype

MicroRNAs are small non-coding RNAs found in both multi-cellular animals and land plants [142, 143, 144]. In both kingdoms they act as negative regulators of translation that

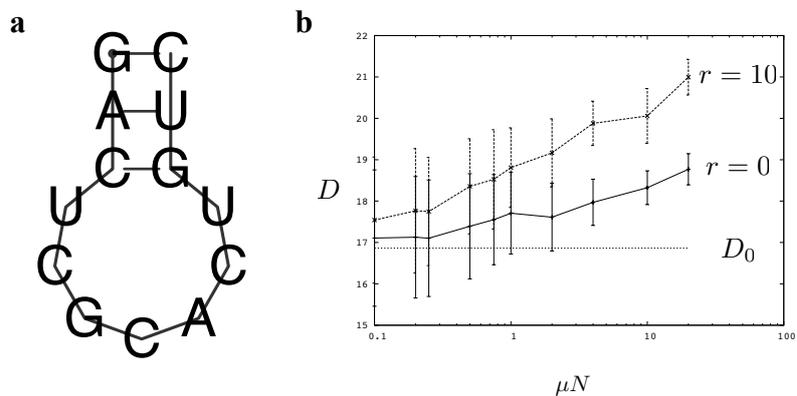


Figure 3.3: To investigate a more realistic neutral network in Ref. [108] we performed simulations using a scaled down analog of microRNA stem-loop hairpin structures (a) consisting of a 7 nucleotide long loop and a 3 nucleotide long stem region. (b) The extent of mutational robustness was found to be higher in the presence of recombination ($r = 10$) then without it ($r = 0$) for all values of μN . Simulations for different values of μN were performed for a set of 20 random initial conditions starting from which the population was allowed to evolve for $100 \times 2N$ generations, the error bars indicate the variance in the time averages over runs with random initial conditions. The dashed line indicates the value of D_0 throughout.

silence genes via complementary interactions with mRNAs [145, 146]. With thousands of miRNA genes identified and genome sequences of diverse eukaryotes available for comparison, the opportunity has emerged for insights into the origin, evolution and molecular mechanism of RNA interference (RNAi) [147]. In the eukaryotic cell long double-stranded RNAs (dsRNA) can arise from several sources including: the replication of RNA viruses, self-annealing cellular transcripts, and miRNA genes. These long dsRNA are cleaved by Dicer to form a RNA duplex of 21-25 nucleotides, with miRNA gene transcript subject to preprocessing involving the excision of a central hairpin forming precursor sequence by Drosha [148]. In the case of dsRNA of non-miRNA origin the endonuclease Dicer functions as a molecular ruler to cleave the dsRNA at 21-25 nucleotide intervals. By contrast, in the case of miRNA precursor sequences a specific 21-25 nucleotide subsequence is cleaved. In both cases, these small RNAs are subsequently incorporated into RNA-induced silencing complexes that use the small RNA as guides for the sequence-specific silencing of complementary messenger RNAs.

The hairpin like secondary structure of precursor stem-loops plays a crucial role in the maturation process and is under evolutionary constraint to conserve its structure. Borenstein and Ruppin used the novel and ingenious method of generating for each miRNA sequence a random sample of sequences with identical minimum free-energy (MFE) structure to uncover traces of adaptation. To compare the mutational robustness of miRNA precursor sequences to random sample sequences with identical MFE structure, they compared the single mutant neighborhood of a given miRNA precursor sequence to the single mutant neighborhood of the sample sequences. Calculating the average distance of the MFE structure of each single mutant sequence to the MFE structure of the original sequence, for both stem-loop and sample sequences, they demonstrated that miRNA precursor sequences have single mutant neighborhoods with sequences that fold into more similar MFE structures compared to sequences in the single mutant neighborhoods of sample sequences with identical MFE structure. While a similar comparison of the folding minimum free-energy showed a comparable, but lower bias, the finding that the two were only weakly correlated allowed the authors to conclude that the observed bias is a result of direct selection for mutational robustness. Their results were reexamined by Shu et al. [149], who argued that mutational robustness among miRNA precursors may be the correlated byproduct of selection for environmental robustness, but found only a moderately higher correlation using a different measure of mutational robustness.

3.4.2 Robustness of microRNA sequences

In Ref. [110] we assessed the environmental and mutational robustness of 3641 unique miRNA precursor sequences available at the time, and for each sequence compared them to a random sample of sequences with the same MFE structure. The idea of looking for signs of adaptation for increased robustness among miRNA precursor sequences by comparing the robustness of naturally occurring sequences to that of random sequences with the same secondary structure is conceptually similar to the approach used to support the argument that the genetic code has evolved to minimize mutational load [150, 151, 101]. In the case of the genetic code the authors took the common genetic code, and, for each codon, calculated the change in polarity of the encoded amino acid caused by replacing each of the three nucleotides, one after the other. In order to determine whether the genetic code is adapted to minimize mutational load they proceeded by comparing the mean squared change caused by the replacement of a single nucleotide in the common genetic code to 10000 randomly generated codes with the same redundancies. They found that only two of the random codes were more conservative than the common code with respect to polarity distances between neighboring amino acids.

We undertook a similar program in the case of miRNA precursor sequences. Each miRNA gene encodes a short ≈ 22 nucleotide sequence that is partially complementary to the mRNA of proteins regulated by the particular miRNA gene. For the proper short sequence to be excised by the proteins of the RNA interference pathway, and hence for the miRNA gene to be functional, a larger part of the miRNA sequence, called the miRNA precursor sequence, must fold into the proper secondary structure. In order to determine whether a miRNA precursor sequence is adapted to minimize the effects of mutational and/or environmental perturbations, i.e., to maximize mutational and/or environmental robustness, we compared the mutational and environmental robustness of each miRNA precursor sequence (robustness measures we used are defined below) to the mutational and environmental robustness of a random sample of sequences with identical structural phenotype (i.e., identical MFE structure).

To generate a random sample of sequences with given MFE structure we first used, starting from a random sequence, stochastic minimization of the free-energy of the target structure to find a sequence with the desired MFE structure. This method by itself, however, yields a *biased* sample of sequences (see Fig.4.1a and b) and must be supplemented by an additional randomization step (see Methods section 4.4). To measure the mutational robustness of a given sequence we used the measures introduced by Borenstein and Ruppin [109]: (i) the structural distance based mutational robustness measure η_s of an RNA

sequence of length L is defined by $\eta_s = 1/(3L) \sum_{i=0}^{3L} (L - d_i)/L$, where d_i is the base-pair distance between the secondary structure of mutant i and the native sequence (given by the number of base pairs present in one structure, but not the other), and the sum goes over all $3L$ single mutant neighbors and (ii) the more stringent measure η_n is simply defined as the fraction of neutral single mutant neighbors, i.e., those that have identical MFE structure to the original sequence. In order to quantify the level of excess mutational robustness among miRNA precursor sequences we counted, for each miRNA precursor sequence, the number of sample sequences that have higher mutational robustness according to a given measure (see Methods section 4.4) and used this to calculate the rank scores r_s and r_n , defined as the fraction of sample sequences with identical or higher robustness according to η_s and η_n , respectively. To facilitate an overview of the extent of excess mutational robustness we also calculated the average of the rank scores over all miRNA precursor sequences \bar{r}_s and \bar{r}_n as well as the fraction of miRNA precursor sequences with higher than average robustness (i.e., rank-scores < 0.5) R_s and R_n and the fraction of sequences with statistically significant increased robustness (i.e., rank-scores < 0.05 , see Methods section 4.4) S_s and S_n , respectively, according to a give measure. The statistical significance of both rank scores for individual miRNA precursor sequences as well as that of the finding a given fraction of robust sequences among a group of sequences was determined as detailed in Methods section 4.4.

Reexamining the mutational robustness of miRNA precursor sequences in comparison to an *unbiased* sample of sequences with identical MFE structure we found that miRNA precursor sequences – in contrast to the results of Borenstein and Ruppín – do not have significantly more neutral single mutant neighbors than sample sequences, but do show a statistically significant increase in robustness measured according to η_s (see Fig.4.1b and Table 3.4.2). In other words, native miRNA precursor sequences have on average the same number of single mutant neighbors with MFE structures identical to their own, as random sample sequences with the same structure. The MFE structure of those single mutant neighbors that are not identical to their own are, on the other hand, significantly more similar than in the case of sample sequences.

The presence of excess mutational robustness is, by itself, insufficient to determine whether mutational robustness has evolved as a result of direct selection or in congruence with selection for environmental robustness. As established previously [109], there is evidence for excess thermodynamic robustness, robustness to thermal fluctuations, as evidenced by a significantly lower than chance minimum folding energy among miRNA precursor sequences. Defining the environmental robustness measure η_E simply as minus

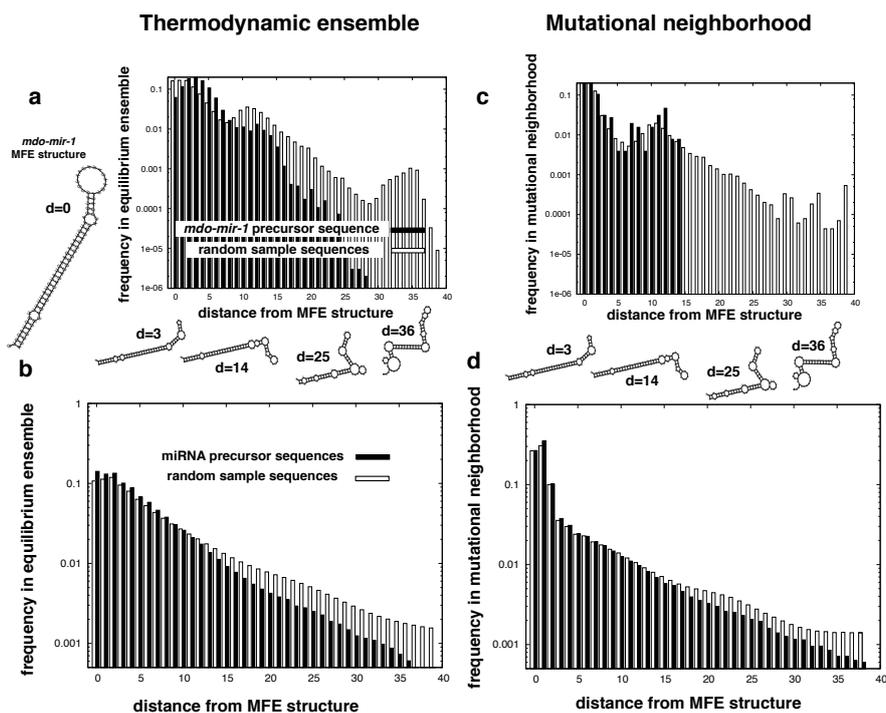


Figure 3.4: In order to examine the robustness of miRNA precursor sequences to thermal fluctuations in work presented in Ref. [110] we sampled the equilibrium thermodynamic ensemble of structures. Sampling 10^6 structures for each miRNA precursor sequence and each member of the random sequence sample, we binned structures according to their distance from the MFE structure. (a) For, e.g. the *Monodelphis domestica* miRNA precursor sequence *mdo-mir-1* examining the distribution of structures as a function of the base-pair distance shows that the averaged random sequence sample distribution (white bars) has a much larger fraction of structures that are drastically different from the MFE structure, compared to the distribution of structures for the original miRNA precursor sequence (black bars). (b) Examining the averaged distribution of stem-loop (black bars) and random corresponding random sequence sample distributions (white bars) shows that there is a general tendency among miRNA precursor sequences for increased thermodynamic robustness, i.e., of avoiding structures that are highly dissimilar to the MFE structure. A strikingly similar effect can be observed if we examine the distribution of structures in the mutational neighborhood. Analogous to (a), in (c) we binned, according to their distance from the MFE structure of the wild type, the MFE structures of all (3L) single point mutants for the *Monodelphis domestica* miRNA precursor sequence *mdo-mir-1* (black bars) as well as the MFE structures for each sequence in the single mutant neighborhood of sample sequences (white bars). The distribution of structures in both the thermodynamic ensemble (a) and the mutational neighborhood (c) of the *mdo-mir-1* miRNA precursor sequence have a significantly smaller fraction of structures that are highly dissimilar than sample sequences with identical MFE structure. Comparing the averaged distribution of stem-loop (black bars) and random corresponding random sequence sample distributions (white bars) in the mutational neighborhood (d) to similar averaged distributions in the thermodynamic ensembles of the same sequences (b) shows that the tendency among miRNA precursor sequences for increased robustness is present both in the mutational neighborhood and the thermodynamic ensemble, i.e., miRNA precursor sequences show excess robustness in the face of both thermal and mutational perturbation.

the minimum folding energy we also find $\bar{r}_E = 0.278$, $R_E = 0.796$, $S_E = 0.220$ using unbiased sampling. The correlation between r_s and r_E across miRNA precursor sequences is, however, rather low with a Pearson's correlation coefficient of $c(r_s, r_E) = 0.217$ and $c(r_n, r_E) = 0.071$. The minimum folding energy is a somewhat crude measure of thermodynamic robustness and does not even reflect the excess mutational robustness according to the measure η_s . There is no good reason to assume that a low MFE in itself confers environmental robustness, as even for relatively high free energies a given sequence may none the less with high probability fold into structures sufficiently similar to the MFE structure to remain functional. The large number of miRNA precursor sequences that exhibit excess mutational robustness as measured by the structural similarity based measure η_s suggests that a strict adherence to the MFE structure is not necessary to retain functionality – a sufficiently similar, but not necessarily identical, secondary structure is enough to guarantee the excision of the proper subsequence. This is further supported by the fact that folding free energy alone is not sufficient to discriminate miRNA precursors, as well as recent evidence that a diverse set of structural features are needed for successful cleavage of a miRNA precursor sequence [152].

To construct an appropriate measure of thermodynamic robustness that also reflects this observation we would need to know the extent of similarity that is required to retain functionality – indeed we would require detailed knowledge of the interaction between the RNA substrate and the enzymes of the RNA interference pathway to establish an appropriate measure of structure similarity. As such information is not at present available, we chose to use the most simple and widely employed structure similarity measure, the base-pair distance used above. In order to determine the extent of similarity required to retain functionality we defined the threshold thermodynamic robustness measure $\eta_t(d_{\text{th.}})$ as a function of the threshold distance $d_{\text{th.}}$, by equating it with the probability in the equilibrium thermodynamic ensemble of structures that have base-pair distances equal to or less than a threshold $d_{\text{th.}}$ with respect to the MFE structure, i.e.,

$$\eta_t(d_{\text{th.}}) = \sum_{i \in \Omega} H(d_{\text{th.}} - d_i) \frac{e^{-G_i/kT}}{Z}, \quad (3.16)$$

where the sum goes over the set of all possible structures Ω , d_i denotes the base-pair distance of structure i to the MFE structure, $Z = \sum_{i \in \Omega} e^{-G_i/kT}$ is the partition sum and $H(x)$ is the unit step function, i.e., $H(x) = 0$ if $x < 0$ and $H(x) = 1$ if $x \geq 0$.

Examining the thermodynamic robustness of miRNA precursor sequences in comparison to an unbiased sample of sequences with identical MFE structure we found that miRNA

Table 3.1: Phylogenetic breakdown of different measures of robustness

group / species	\bar{r}_n	R_n	S_n	\bar{r}_s	R_s	S_s	$\bar{r}_t(25)$	$R_t(25)$	$S_t(25)$	$c(r_s, r_t(25))$	# of seqs.
all	0.44	0.59	0.08	0.29	0.78	0.17	0.31	0.74	0.28	0.73	3641
vertebrate	0.46	0.55	0.06	0.31	0.78	0.13	0.29	0.75	0.30	0.76	2215
invertebrate	0.37	0.68	0.10	0.21	0.88	0.27	0.22	0.84	0.36	0.73	488
landplant	0.41	0.63	0.11	0.31	0.75	0.21	0.40	0.63	0.19	0.68	848
virus	0.38	0.66	0.09	0.23	0.84	0.18	0.21	0.85	0.32	0.65	82
<i>Homo sapiens</i>	0.48	0.53	0.04	0.32	0.75	0.11	0.28	0.76	0.33	0.74	471
<i>Mus musculus</i>	0.46	0.56	0.06	0.33	0.76	0.12	0.31	0.74	0.27	0.79	373
<i>Drosophila melanogaster</i>	0.40	0.64	0.08	0.22	0.88	0.24	0.23	0.82	0.35	0.74	78
<i>Caenorhabditis elegans</i>	0.30	0.78	0.18	0.20	0.89	0.37	0.23	0.82	0.34	0.75	114
<i>Arabidopsis thaliana</i>	0.39	0.62	0.11	0.29	0.78	0.19	0.43	0.60	0.15	0.75	131
Epstein-Barr virus	0.31	0.78	0.00	0.16	0.96	0.22	0.16	0.87	0.48	0.81	23

Average rank-scores that indicate significantly increased according to both measures discussed in Methods section 4.4 (p -value $< 10^{-3}$) are given in bold.

precursor sequences have significantly more structures in their equilibrium thermodynamic ensemble that are similar to the MFE structure than sample sequences (see Fig.3.4a,b and Table 3.4.2). In other words, miRNA precursor sequences tend to adapt more similar structures as a result of thermal fluctuations than random sample sequences with the same structure. Calculating the average rank score $\bar{r}_t(d_{th.})$ and the fraction of robust $R_t(d_{th.})$ and significantly robust $S_t(d_{th.})$ miRNA precursor sequences, with respect to the measure $\eta_t(d_{th.})$ (Fig.3.5a) and examining the distribution of structures as a function of the base-pair distance for individual miRNA precursor sequences (see e.g. Fig.3.4a) indicates that above a threshold distance $d_{th.} \approx 20$ the measures start to saturate, yielding an estimate of the required similarity to retain function. The correlation between the rank-score of miRNA precursor sequences according to the distance similarity based mutational robustness measure and the threshold thermodynamics measure is high for all threshold values. This is the direct result of the high degree of similarity between the distribution of structures in the thermodynamic ensemble and the mutational neighborhood (Fig.3.4). The average rank score $\bar{r}_t(d_{th.})$ and the fraction of robust $R_t(d_{th.})$ and significantly robust $S_t(d_{th.})$ miRNA precursor sequences with respect to the threshold thermodynamic robustness measure indicate a markedly larger extent of excess robustness than their counterparts for mutation robustness, i.e., \bar{r}_s , R_s and S_s (see Fig.3.5a,b and Table 3.4.2) above $d_{th.} > 20$.

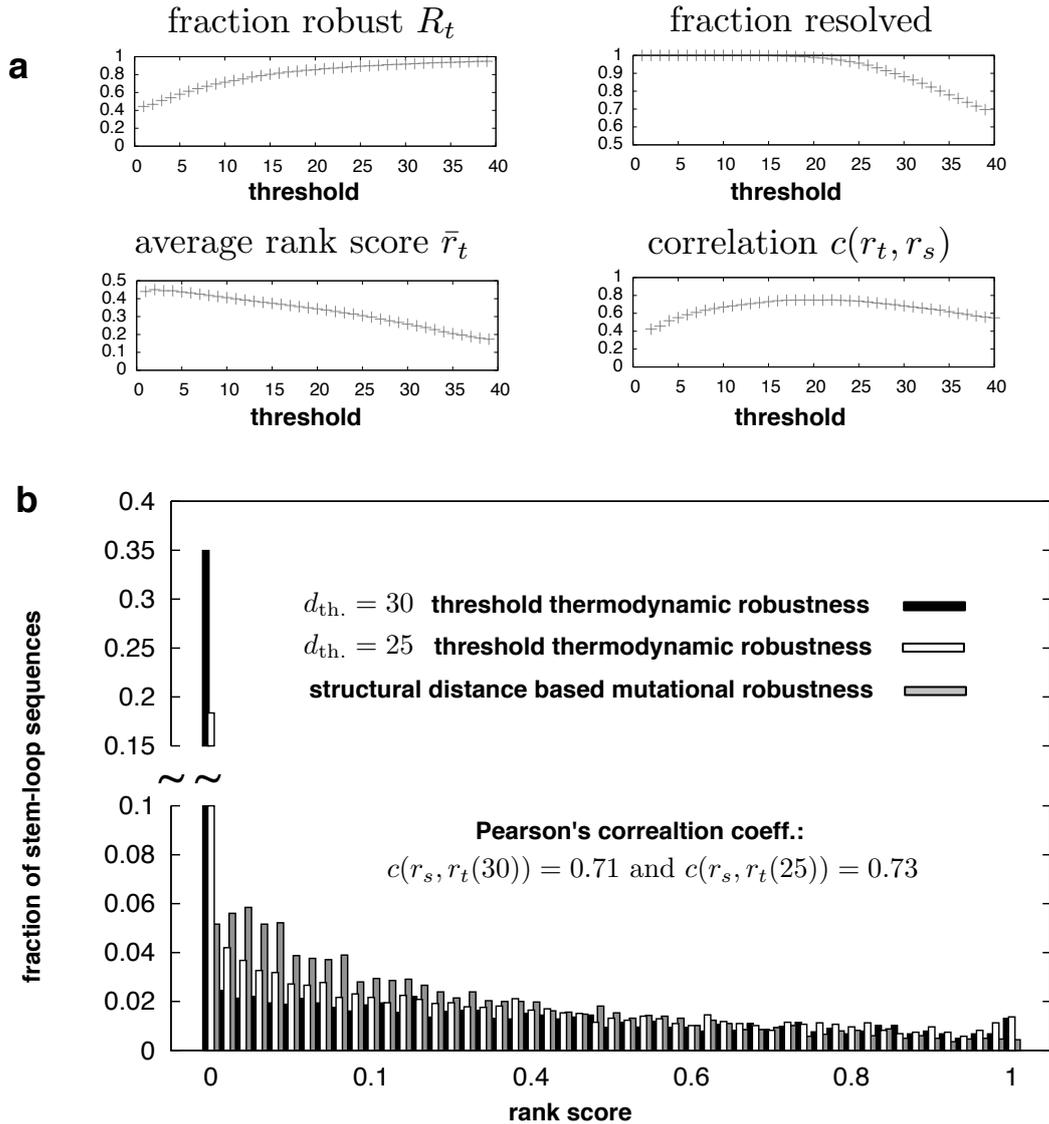


Figure 3.5: To quantify the extent of thermodynamic robustness present in miRNA precursor sequences in Ref. [110], we examined the rank statistics of the threshold thermodynamic robustness $\eta_t(d_{th.})$ which we defined as the cumulative frequency of structures in the thermodynamic ensemble that are equal to or less than a distance threshold $d_{th.}$. (a) For threshold values larger than $d_{th.} \approx 20$ the average rank of miRNA precursor sequences with respect to $\eta_t(d_{th.})$, denoted by $\bar{r}_t(d_{th.})$, becomes larger than \bar{r}_s , while remaining highly correlated with it. For thresholds $d_{th.} > 25$ the $\eta_t(d_{th.})$ values for an increasingly larger fraction of random sample sequences becomes indistinguishable from the $\eta_t(d_{th.})$ value of the original miRNA precursor sequence due to a lack of structures with $d > d_{th.}$ in our finite resolution sample of the equilibrium ensemble. The fraction of sequences which remain distinguishable are labeled *fraction resolved*. (b) Among the majority of the random sample sequences that are resolved, however, the threshold robustness of large number of the miRNA precursor sequences becomes highly significant (black and white bars) compared to the mutation robustness as measured by η_s (grey bars), whilst retaining a high correlation among r_s and $r_t(d_{th.})$.

3.4.3 The case for congruent evolution of robustness

The results presented above demonstrate the correlated presence of excess environmental (thermodynamic) and genetic (mutational) robustness among miRNA precursor sequences as measured according to, respectively, η_s and $\eta_t(d_{\text{th.}})$. A rather general causality between environmental and genetic robustness in the context of RNA secondary structure has been suggested by Ancel and Fontana [153], who studied the dynamics of an *in silico* population of RNA sequences evolving towards a predefined target shape. They found that a correlation exists between the set of shapes in the plastic repertoire of a sequence and the set of dominant (minimum free energy) shapes in its genetic neighborhood. They argue that this statistical property of the RNA genotype-phenotype map, which they call plastogenetic congruence, traps populations in regions where most genetic variation is phenotypically neutral. In other words, RNA sequences explore a similar repertoire of suboptimal structures as a result of perturbations due to mutations and perturbations resulting from thermal fluctuations, and selection for a given target structure favours sequences with higher robustness to perturbations of both type.

Since, in contrast to genetic robustness, environmental robustness does not require high values of μN , as it is a property of the sequence and not its mutational neighborhood, we contend that the observed bias in mutational robustness is in fact the result of the *congruent* evolution of environmental and genetic robustness.

3.4.4 The temperature of mutations

The observation that mutational and thermal perturbations cause RNA sequences to explore a similar repertoire of suboptimal structures can be quantified. In fact as we show below the effects of point mutations can with striking accuracy cast as increased effective temperature.

The response to thermodynamic fluctuations of a sequence s is reflected in the distribution

$$p_T(d|s, T) = \sum_{i \in \Omega(d, s, T)} \frac{\exp(-G(i, s, T)/kT)}{Z(s, T)} \quad (3.17)$$

that gives the probability structures at temperature T with structural distance d from the MFE structure of s in the thermodynamic ensemble ($\Omega(d, s, T)$ is the subset of all possible structures Ω that are at distance d from the MFE structure at temperature T of s , $G(i, t, T)$ is the free energy of structure i on sequence t and $Z(v, T) = \sum_{i \in \Omega} \exp(-G(i, v, T)/kT)$).

The combined effects of mutational and thermodynamic perturbations, on the other

hand, are reflected in the distribution of structures over the single mutational neighborhood $[s]_1$ of s . More precisely the probability of structures over $[s]_1$ with structural distance d from the MFE structure of s at the same thermodynamic temperature T . We may calculate this distribution by averaging over the set of $3L$ single mutant neighbor sequences $[s]_1$ of sequence s :

$$p_{\text{TM}}(d|s, T) = \sum_{i \in \Omega(d, s, T)} \frac{1}{3L} \sum_{t \in [s]_1} \frac{\exp(-G(i, t, T)/kT)}{Z(t, T)}. \quad (3.18)$$

Provided that mutational and thermal perturbations have sufficiently similar effects we may hope to obtain the distribution $p_{\text{TM}}(d|s, T)$ from $p_{\text{T}}(d|s, T)$ by some simple transformation. As base-pairing energies are of the order of a few kT we may try to recover the seemingly more complex mutational distribution $p_{\text{TM}}(d|s, T)$ over $[s]_1$ from the more simple $p_{\text{T}}(d|s, T)$ by introducing an *effective mutational temperature*. This effective temperature T_{eff} , however, will not be completely analogous to the thermodynamics temperature T as it will not effect the entropic contributions to the free energy, but merely the apparent equilibrium energy scale. That is a T_{eff} must be included in the denominator in the exponent of the Boltzmann factor, but not in the calculation of the free energy:

$$p_{\text{TE}}(d|s, T, T_{\text{eff}}) = \sum_{i \in \Omega(d, T)} \frac{\exp(-G(i, s, T)/kT_{\text{eff}})}{Z_{\text{eff}}(s, T, T_{\text{eff}})}, \quad (3.19)$$

where $Z_{\text{eff}}(v, T, T_{\text{eff}}) = \sum_{i \in \Omega} \exp(-G(i, v, T)/kT_{\text{eff}})$.

As show in Fig.3.6 the combined distribution $p_{\text{TM}}(d|s, T)$ may be recovered using a single effective temperature T_{eff} with striking accuracy in the case of the miRNA precursor sequence data set considered above (for details see Methods section 4.4.2). The existence of an effective mutational temperature implies a tight and rather general connection between robustness to mutational and environmental perturbations. To demonstrate this connection let us denote the probability of some structure i on sequence s at thermodynamic temperature T in the thermodynamic ensemble as

$$P_{\text{T}}(i, s, T) = \frac{\exp(-G(i, s, T)/kT)}{Z(s, T)}, \quad (3.20)$$

and mutational ensemble as

$$P_{\text{TM}}(i, s, T) = \frac{1}{3L} \sum_{t \in [s]_1} \frac{\exp(-G(i, t, T)/kT)}{Z(t, T)}. \quad (3.21)$$

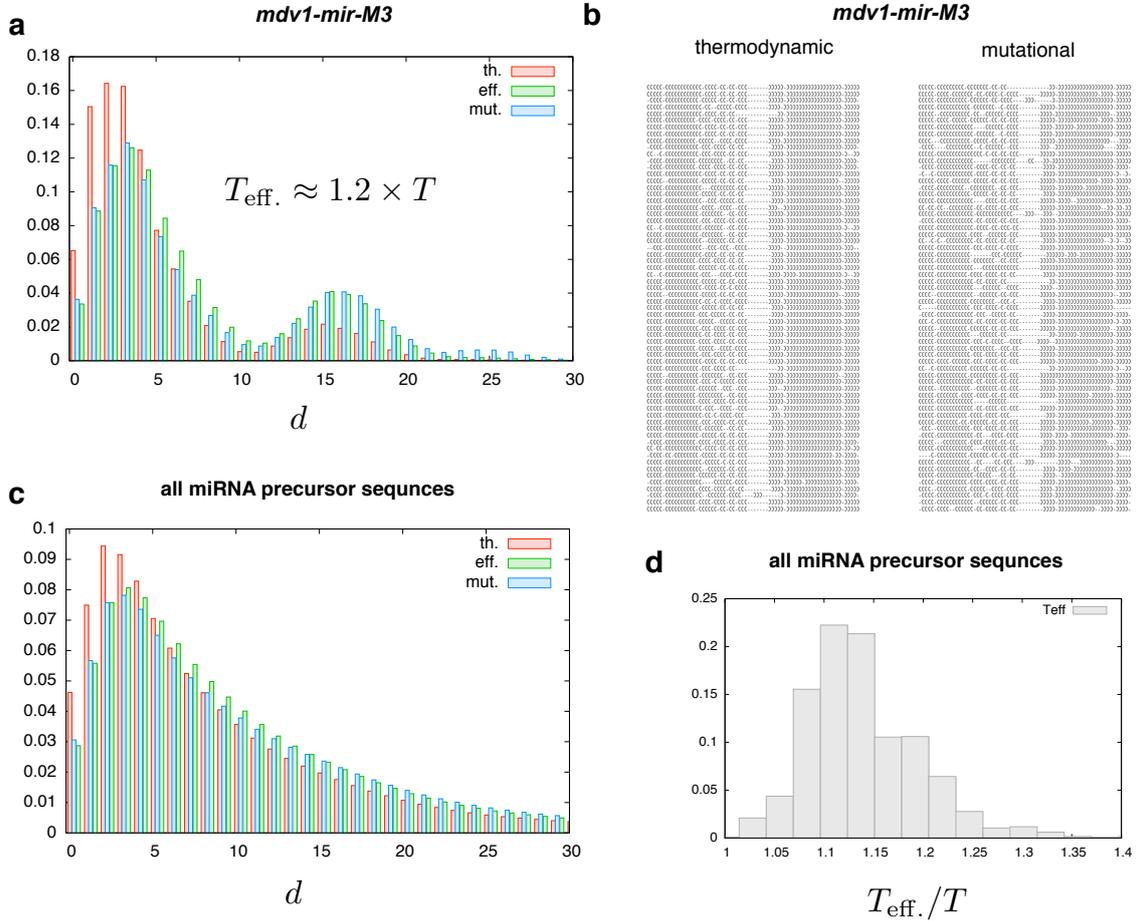


Figure 3.6: We estimated the *effective mutational temperature* T_{eff} that best reproduces $p_{\text{mut.}}(d|s, T)$ for all precursor sequences in the miRNA data set (cf. Methods section 4.4). **(a)** shows $p_{\text{T}}(d|s, T)$ (red bars), $p_{\text{TM}}(d|s, T)$ (blue bars) and $p_{\text{TE}}(d|s, T, T_{\text{eff}})$ (green bars) corresponding to the best fit $T_{\text{eff}} = 1.98T$ for the miRNA sequence *mdv1-mir-M3* from Mareks disease virus. **(b)** shows a random sample of the thermodynamic and combined mutational and thermodynamic ensembles. **(c)** shows the average over all sequences in the miRNA precursor data set of $p_{\text{T}}(d|s, T)$ (red bars), $p_{\text{TM}}(d|s, T)$ (blue bars) and $p_{\text{TE}}(d|s, T, T_{\text{eff}})$ (green bars) with the distribution of best fit temperatures shown in **(d)**. See text and Methods section 4.4.2 for details.

The existence of some T_{eff} implies that $P_{\text{TM}}(i, s, T)$ can be approximated as

$$P_{\text{TM}}(i, s, T) \approx P_{\text{TE}}(i, s, T, T_{\text{eff}}) = \frac{\exp(-G(i, s, T)/kT_{\text{eff}})}{Z_{\text{eff.}}(s, T, T_{\text{eff}})}. \quad (3.22)$$

Provided selection favours some subset of functional structures $\mathcal{S} \subset \Omega$ we may measure the robustness with respect to thermodynamics perturbations as

$$\eta_{\text{T}}(\mathcal{S}|s, T) = \sum_{i \in \mathcal{S}} P_{\text{T}}(i, s, T) = \sum_{i \in \mathcal{S}} e^{-\beta G(i, s, T)} / \sum_{i \in \Omega} e^{-\beta G(i, s, T)}, \quad (3.23)$$

where we have introduced the notation $\beta = kT$. This expression is equivalent to $\eta_t(d_{\text{th.}})$ defined in Eq. (3.16) when \mathcal{S} is the set of structures that have base-pair distances equal to or less than a threshold $d_{\text{th.}}$ with respect to the MFE structure of s at temperature T . The joint robustness to thermodynamic and mutational perturbations is measured by the analogous expression

$$\eta_{\text{TM}}(\mathcal{S}|s, T) = \sum_{i \in \mathcal{S}} P_{\text{T}}(i, s, T). \quad (3.24)$$

The existence of an effective temperature implies

$$\eta_{\text{TM}}(\mathcal{S}|s, T) \approx \sum_{i \in \mathcal{S}} P_{\text{TE}}(i, s, T, T_{\text{eff}}) = \sum_{i \in \mathcal{S}} e^{-\beta_{\text{eff.}} G(i, s, T)} / \sum_{i \in \Omega} e^{-\beta_{\text{eff.}} G(i, s, T)}, \quad (3.25)$$

where $\beta_{\text{eff.}} = kT_{\text{eff.}}$.

To separate out mutational robustness we introduce the measure of mutational robustness

$$\eta_{\text{M}}(\mathcal{S}|s, T) = \frac{\eta_{\text{TM}}(\mathcal{S}|s, T)}{\eta_{\text{T}}(\mathcal{S}|s, T)} \approx \frac{\eta_{\text{TE}}(\mathcal{S}|s, T, T_{\text{eff.}})}{\eta_{\text{T}}(\mathcal{S}|s, T)}. \quad (3.26)$$

In agreement with the results in previous section, looking at Fig.3.6 we can see that provided $\mathcal{S} = \Omega(d \leq d_{\text{th.}}, s, T)$: (i) a single effective mutational temperature T_{eff} faithfully reproduces the effects of point mutations, and (ii) $T_{\text{eff.}}/T > 1$ values among miRNA precursor sequences imply that selection for robustness with respect to thermodynamics perturbations causes mutational robustness to emerge (and vice versa).

Selection will favour sequences that have thermodynamic ensembles where the cumulative probability of functional sequences $\eta_{\text{T}}(\Omega(d \leq d_{\text{th.}}, s, T)|s, T)$ is larger. The cumulative probability of functional sequences in the combined ensemble $\eta_{\text{TM}}(\Omega(d \leq d_{\text{th.}}, s, T)|s, T)$ is determined by the the cumulative probability in the thermodynamic ensemble and the effects of an effective temperature, i.e., is given by $\eta_{\text{TE}}(\Omega(d \leq d_{\text{th.}}, s, T)|s, T, T_{\text{eff.}})$. As a consequence of the fact that $T_{\text{eff.}}/T > 1$, selection for sequences with increased

$\eta_T(\Omega(d \leq d_{th.}, s, T)|s, T)$ implies $\eta_M(\Omega(d \leq d_{th.}, s, T)|s, T) > 1$. That is the sequences for which the effects of thermodynamic perturbations are reduced over the set of functional structures $\Omega(d \leq d_{th.}, s, T)$ are also the sequences for which the effects of mutational perturbations are reduced over $\Omega(d \leq d_{th.}, s, T)$. In general, selection for thermodynamic robustness favours sequences for which the functional subset of structures \mathcal{S} reside in narrower and deeper free energy valleys of configuration space; sequences for which an effective temperature increase – corresponding to mutational perturbation – will push less probability out of the narrower and deeper free energy valleys.

Our work on the effective mutational temperature is still in its preliminary stages. Further investigation of different data sets, e.g. diverse tRNA sequences such as tRNA from prokaryotes adapted to different temperatures, and also of different genotype-phenotype relationships, e.g. tractable protein folding models, is necessary to establish how universal this our concept of an effective mutational temperature is. None the less the possibility of such a simple relationship between the variability resulting from thermal perturbations and mutational ones has the possibility to offer important *quantitative* insight into the origins of genetic robustness and evolutionary novelty.

3.5 Discussion

We examined the effects of recombination on the extent to which populations evolving on neutral networks exhibit mutational robustness. Our results show that recombination leads to enhanced mutational robustness under very general circumstances. We calculated the stationary limit distribution of populations evolving on neutral networks drawn from two different random ensembles in the infinite population limit. Comparison of the results showed that populations in which recombination is present are more sensitive to details of the topology than populations where only mutation is present. In particular, neutral networks where genotypes of high centrality exhibit larger neutrality evolve greater mutational robustness.

While our results indicate that significant mutational robustness readily evolves in the presence recombination, this result must be considered with the caveat that evolution of mutational robustness through neutral dynamics requires sufficient polymorphism to be present in the population. Our results shown, however, that provided this condition is met, recombination leads to increased values of mutational robustness in comparison to populations where only mutation is present.

Examining microRNA genes of several eukaryotic species, we demonstrated that the

sampling method used by Borenstein and Ruppin introduced significant bias that lead to an overestimation of robustness. Introducing a novel measure of environmental robustness based on the equilibrium thermodynamic ensemble of secondary structures of the miRNA precursor sequences, we demonstrated that the biophysics of RNA folding induces a high level of correlation between genetic (mutational) and environmental (thermodynamic) robustness. In light of theoretical considerations, we believe that this correlation strongly suggests that genetic robustness observed in miRNA sequences is the byproduct of selection for environmental robustness.

The correlation between the response to heritable (mutational) and nonheritable (thermodynamic) perturbation, and hence the congruent evolution of genetic and environmental robustness may extend to other systems with genotype-phenotype maps different from RNA secondary structure. In particular, Xia and Levitt [141] have found compelling evidence of the correlated evolution of increased thermodynamic stability and the number of neutral neighbors in lattice protein models. Understanding the relationship between sequence, structure, and function is, and will remain to be in the foreseeable future, a central theme in both molecular and evolutionary biology. A comprehensive view of how the relationship between sequence, structure, and function is shaped during the course of evolution must take into consideration both the potential correlations that arise from the physics of the structure-sequence relationship as well as the relevant population genetic conditions in the context of which it takes on the role of a genotype-phenotype map.

The relevance of the biophysics of folding to evolutionary problems –ranging from relationship between genetic robustness and evolvability to the tempo and mode of evolutionary change – has only recently been considered in any detail by evolutionary theory. Conversely, population genetic effects – such as the importance of genetic drift and the potentially universal effects of thermodynamic stability on organismal fitness – have been largely absent from the biochemical approach to protein and RNA evolution [154].

Chapter 4

Methods

4.1 Genome organization dynamics

Population Dynamics

The population dynamics simulations the results of which are discussed in section 2.7 (and published in Ref. [51]) were carried out in a manner that allowed the separate treatment of every bacteria possessing any of the possible model genotypes. The frequency of individual bacteria of each genotype was calculated by solving the 2^{10+1} differential equations describing the number of bacteria $n(G)$ in each of the 2^{10+1} genome states G :

$$\begin{aligned} \frac{dn(G)}{dt} = & r_r(G, \{F\}, \{n\})n(G) - r_w n(G) \\ & - r_m(l_f(G) + \delta_{t,G})n(G) + [\text{mut.in}] \\ & - r_t(10 + 1)\delta_{t,G}n(G) + [\text{trf.in}] + [\text{migr.in}], \end{aligned}$$

where $l_f(G)$ is the number of functional model food genes in genotype G , $\delta_{x,G}$ is equal to unity if G contains an intact model gene x and zero otherwise. $n(G)$ is treated as a continuous variable, but with a lower cutoff at $n(G) = 1$ to mimic the discrete nature of bacterial populations. This lower cutoff is implemented for each subpopulation with $n(G) < 1$, by resetting $n(G)$ to 1 with probability $n(G)$, or to 0 with probability $1 - n(G)$, at each time step of the simulation. The first term on the right hand side describes the reproduction of bacteria as detailed below. The second term corresponds to the washout of bacteria from the population. The third and fourth terms deal with mutation, the former one considers bacteria with genotype G that undergo mutation, while the later one $[\text{mut.in}]$ is the sum of contributions from all bacteria that mutate into state G . The fifth and sixth

terms describe transformation, the former one corresponding to those bacteria in state G that undergo transformation, while the later one [trf.in] represents the complicated sum of all transformations that lead to state G . The last term [migr.in] corresponding to migration is also described below.

Reproduction Rate

The fitness of individual bacteria with a given genotype G is influenced both by available food $\{F\} \equiv F_1, F_2, F_3$, as well as genome length $l(G) = l_f(G) + l_t \delta_{t,G}$, where l_t denotes the length of the **t** gene relative to the model food genes (and is considered to be less than unity, as explained below). Available food types F_i were each considered to be constantly replenished in the environment, each type being characterized by a strength S_i corresponding to the number of bacteria the given food type could sustain at the maximum division rate r_r^{\max} . Each bacterium with a functional copy of model gene f_i (corresponding to a given food type F_i) present in its genome ($\delta_{f_i,G} = 1$) received an “equal share” $S_i / \sum_{G'} \delta_{f_i,G'} n(G')$ of this food, while those bacteria that did not have an intact copy of the model gene ($\delta_{f_i,G} = 0$) in their genome received none. In general the speed at which a given bacterium can divide $r_r(G, \{F\}, \{n\})$ is proportional to the total amount of food it can utilize $\sum_{i=1}^3 \delta_{f_i,G} S_i / \sum_{G'} \delta_{f_i,G'} n(G')$, but also decreases by a factor of $2/(1 + l(G)/l_{\text{opt}})$ as the total number of intact model food genes present in its genome increases, where $l_{\text{opt}} = 3$ is the optimal genome size. The denominator of the latter factor takes into account that the time necessary for a bacterium to divide has a component proportional to genome length. Bacteria that possess only those three model genes that are necessary for the utilization of the three food types available at a given moment are the ones that divide the fastest (i.e., have the highest fitness), their genomes being the most highly economized. As outlined in section 2.7 real bacteria may possess several dozen gene groups necessary for the utilization of fluctuating resources, but as we are not able to treat numerically more than a few model gene–food pairs, we have chosen the parameters of the fitness function such that it decreases rapidly as the number of model genes grow (e.g. by 33% if twice the optimal number of model genes is present), that is we, in this sense, consider each model gene to correspond to several gene groups. Further, as the gene group responsible for DNA uptake (the **t** model gene) is only one of several dozens and in reality consist of a relatively small number of genes, we have included it with a smaller relative length ($l_t = 0.1$) in the total genome length. Due to practical constraints, however, the mutation rate of the **t** gene was considered (except for the data presented in Fig.2.8c.) to be the same as that of the food genes (which corresponds to a tenfold increase in the mutation rate for the **t** gene).

Therefore, in our simulations the reproduction rate of a bacterium with genotype G and available foods $\{F\}$ was given by:

$$r_r(G, \{F\}, \{n\}) = \frac{2 r_r^{\max}}{1 + l(G)/l_{\text{opt}}} \min \left[1, \sum_{i=1}^3 \frac{\delta_{f_i, G} S_i}{\sum_{G'} \delta_{f_i, G'} n(G')} \right],$$

and we chose the food strength to be $S_i = 10^5$ for each available food type.

Migration

Since all food types are equivalent and consequently all combinations of these are equally likely, the number of mean field variables needed to describe the global genotype distribution can be reduced from 2^{10+1} to 10×2 corresponding to all genotypes G with 1, 2, \dots , 10 intact metabolic genes, with and without the t gene. To take into account the influx of bacteria from external populations we took the averages of these 10×2 types in the population with a sliding-growing time window (always encompassing the last quarter of the simulation) and subsequently calculated the migrational term [migr.in] for each genotype G by multiplying the corresponding mean field variable with the appropriate combinatorial factor and the migration rate R_{migr} .

4.2 Random neutral networks

In Ref. [108], as described in section 3.3, we investigated the effects of recombination on two types of random neutral network ensembles. In both cases we choose an alphabet of size 4 and generated random neutral networks consisting of M genotypes of length L . The first type of random networks (hereafter referred to as uniform attachment networks) were generated by choosing a random seed genotype and adding a random neighboring genotype (one that can be obtained from it by a single mutation). This process was continued by selecting a random genotype from those already added to the network and adding one its neighbors not already in the network at random until the network contains M genotypes. The second type of random networks (hereafter referred to as preferential attachment networks) were grown at each step by choosing genotypes from those already added to the network with a probability proportional to the number of neighbors the genotype already has in the network, and adding a neighbor not contained in the network at random until the network contained M genotypes.

This choice of random ensembles was motivated by the expected differences in the

topology of networks belonging to the two ensembles. Preferential networks have a topology where the most connected genotypes are also more centrally located on average.

A numerical scheme for the Wright-Fisher process

In the case of the simulations performed in Ref. [108], the results of which are described in section 3.3, sampling consisted of choosing N individuals and assigning each a genotype $i \in G$ with probability x_i and subsequently updating x_i accordingly and integrating equation 3.11 for one generation time before performing sampling again. For sufficiently large sample size the stochastic dynamics of our model only depends on the combined parameters μN and r . Appropriately large sample size N was chosen by increasing N (while keeping μN and r fixed) until the effects of sample size became negligible ($N \gtrsim 1000$).

4.3 RNA secondary structure

RNA is a very appealing system to study the relationship between sequence (genotype) and structure (phenotype). It has a relatively simple structure being composed by only four kind of different monomers. The physical interactions determining its secondary structure (hydrogen bonding between G-C, A-U and G-U, stacking energies and entropic contributions from loops) are approximately one order of magnitude stronger than those determining the tertiary structure. The relative simplicity of RNA has allowed the development of very effective coarse grained models based on empirical parameters. These models allow us reliably and quickly calculate the minimum free energy secondary structure into which an RNA molecule is most likely to fold, and can be used to fully analyse the energy landscape of moderately large (up to few hundred bp. long) RNA molecules.

4.3.1 The Vienna package

The Vienna RNA Package [155] consists of a C code library and several stand-alone programs for the prediction and comparison of RNA secondary structures. The package implements three different dynamic programming algorithms for the prediction of RNA secondary:

1. *The Zuker and Stiegler algorithm* that yields the single minimum free energy (MFE) structure [156].

2. *McCaskill's algorithm*, which calculates the partition function along with the base pairing probabilities in the thermodynamic ensemble [157].
3. *Wuchty et al.'s algorithm*, which generates all suboptimal structures within a given energy range of the optimal energy [158].

An algorithm to design sequences with a predefined structure (inverse folding) is also provided, the limitations of which for sampling random sequences with a given MFE structure is discussed below.

4.3.2 Scaled down microRNA neutral network

As discussed in detail in section 3.4.1 the structure of microRNA precursor stem-loops, have been shown to exhibit a significantly high level of robustness in comparison with random RNA sequences with similar stem-loop structures [109]. This observation makes the neutral networks with identical stem-loop like secondary structure a natural testing ground for the evolution of mutational robustness. To construct a neutral network on which the above population dynamics is computationally tractable, in Ref. [108] we constructed a scaled down analog of such stem-loop structures. We downloaded all currently available miRNA stem-loop precursor sequences from miRBase [159]. Analysis of the sequences and the corresponding structures indicated that hairpin-like stem-loops consisted of, on average, a 7.241-base-long loop and a 50.429-base-pair-long stem region. To first approximation we may consider the neutral network of a hairpin like structure to consist of large quasi-independent regions corresponding to mutations in the central loop region that are connected by more rare mutations in the stem region. Aiming to approximate the neutral networks of stem-loop structures by that of a single such region we used the Vienna RNA secondary structure prediction package [155] to find connected neutral networks (a set of genotypes that can be reached through single mutations) of the hairpin like structure with 3-base-pair-long stem and a 7-base-long loop region. In section 3.3 we present results for the $M = 37972$ connected neutral network that contains the sequence GACUCGCACUGUC.

4.4 MicroRNA sequences

Sequences

MicroRNA (miRNA) precursor sequences were downloaded from miRBase version 9.0 [159]. All 4361 miRNA genes were used, yielding 3641 unique miRNA precursor se-

quences.

Measuring thermodynamic robustness of structural ensembles

In order to calculate the thermodynamic robustness measure η_t , defined in section 3.4.2, we sampled the equilibrium thermodynamic ensemble of stem-loop and sample sequences using the stochastic backtracking routine from the Vienna RNA package producing 10^6 suboptimal structures per sequence, using the default temperature of 310 K . The average distance from the MFE structure in the thermodynamic ensemble can be calculated exactly with the help of base-pairing probabilities, which are available as a byproduct of partition function folding in the Vienna package, and were used to validate the sampling.

Statistics

Given a rank score r and sample size N a good estimate for the probability of observing an equal or lower rank score by chance is given by $(rN)/(N + 1) \approx r$. Following Ref. [109] rank scores of $r < 0.05$ are considered significantly robust. To determine if the robustness of miRNA precursor sequences according to some measure η has the same distribution as the robustness of sample sequences η' for a group of sequences, following Ref. [109] we test against the null hypothesis that they are drawn from identical distributions using the nonparametric Wilcoxon signed rank test. In contrast to Ref. [109], however, we do not consider as paired values η and the average of η' over all N sample sequences, $\langle \eta' \rangle$, as we found this to result in spuriously low p -values, but instead calculate the p -values for a given group of sequences by averaging over 1000 different sets of $\{\eta, \eta'\}$ pairs where in each set the η' values belong to a random sample sequence. As a complementary approach, we also tested the hypothesis that the distribution of rank scores of a group of sequences for a given robustness measure is uniform – as we would expect if miRNA precursor sequences were randomly sampled from the set of sequences with identical MFE structure – using a standard Kolmogorov-Smirnov goodness of fit test. We found the two significance analyses to be in good agreement indicating highly significant bias for higher values of $\eta_t(d_{th.})$ and η_s , but mostly no or only nonsignificant bias for higher η_m . The supplemental information accompanying Ref. [110] contains species level statistics and significance analyses, this is also available as part of the electronic supplement of the thesis at <http://angel.elte.hu/~ssolo/thesis.html>.

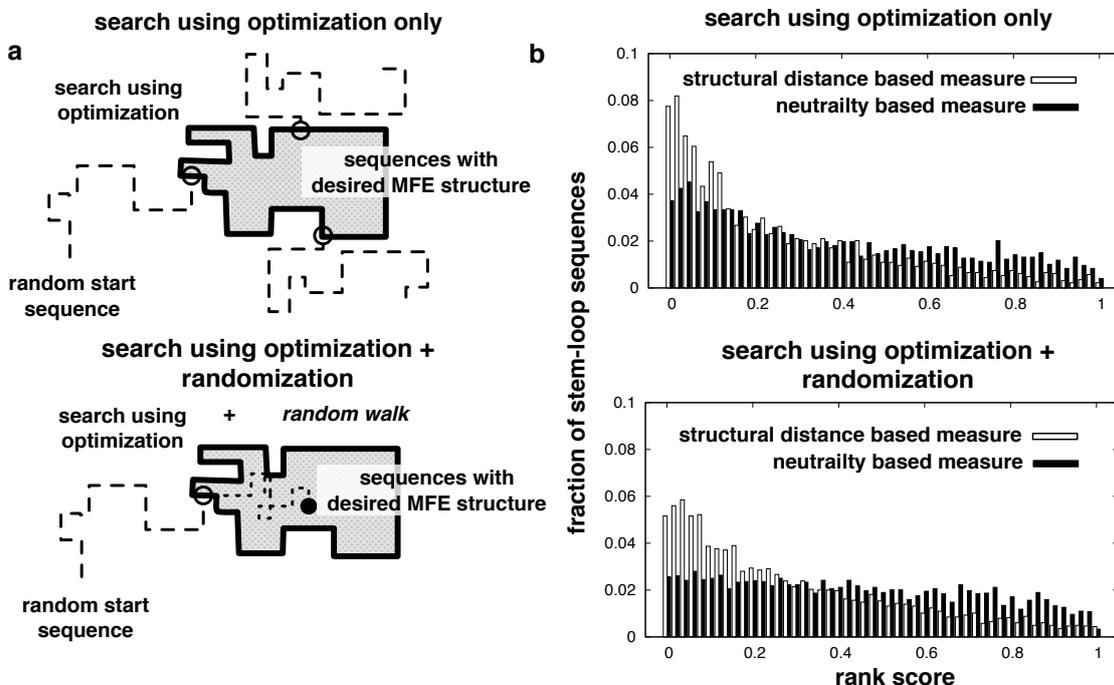


Figure 4.1: (a) Generating a random sample of sequences with a desired MFE structure by stochastic minimization of the free-energy of the desired fold (the method employed the RNAinverse program used by Borenstein and Ruppin [109] as well as Shu et al. [149]) results in a biased sample in which sequences with lower than average neutrality (higher than average number single mutant neighbors) are overrepresented. This can be avoided if after finding a sequence with the desired MFE structure a random walk is performed among sequences with the desired MFE structure. This random walk on the neutral network associated with the MFE structure mimics the sequence drift of a sequence evolving under the constraint to fold in to the desired MFE structure. (b) Rank score distributions for two measures of mutational robustness (η_s and η_n see text). Comparing the distributions derived from sampling using only stochastic optimization (top, $\bar{r}_s^{\text{biased}} = 0.25$, $R_s^{\text{biased}} = 0.83$, $\bar{r}_n^{\text{biased}} = 0.37$, $R_n^{\text{biased}} = 0.66$) to that derived from sampling with subsequent randomization (bottom, $\bar{r}_s = 0.29$, $R_s = 0.78$, $\bar{r}_n = 0.44$, $R_n = 0.59$) shows that increased neutrality is predominately an artifact of biased sampling, while the lower than average distance of MFE structures in the mutational neighborhood to the wild-type MFE structure becomes somewhat less pronounced, but is still significant.

4.4.1 Sampling of sequences with a given MFE structure

For each miRNA precursor sequence we produced a sample of random sequences by (i) using the stochastic optimization routine from the Vienna RNA package [155] to produce a sequence with MFE structure identical to that of the native sequence that is stored (ii) and subsequently randomizing this sequence by attempting $4L$ random nucleotide substitutions in a miRNA precursor sequence of length L , accepting a substitution if the resulting sequence's MFE structure remains unchanged. For each miRNA precursor sequence on average > 800 sample sequences with identical MFE structure was generated. Supplemental information accompanying Ref. [110] contains the robustness values for all 4361 genes associated with 3641 unique sequences we considered, this is also available as part of the electronic supplement of the thesis at <http://angel.elte.hu/~ssolo/thesis.html>.

4.4.2 Estimating the effective temperature of mutations

The distributions $p_T(d|s, T)$ and $p_{TM}(d|s, T)$ was estimated by sampling 10^6 structures per sequence from the Boltzmann ensemble of secondary structures over a sequence s and its single mutant neighbourhood $[s]_1$ at temperature T . To examine the similarity between these two distributions we sampled the Boltzmann ensemble of secondary structures over a sequence s in order to estimate the joint probability distribution $p_T(G, d|s, T)$ of structures with free energy G over s and with structural distance d from the MFE structure of s at temperature T . We used the expression

$$p_{TE}(d|s, T, T_{\text{eff.}}) = \frac{Z(s, T)}{Z(s, T, T_{\text{eff.}})} \int_0^\infty dG \frac{\exp(-G(i, s, T)/kT_{\text{eff.}})}{\exp(-G(i, s, T)/kT)} p_T(G, d|s, T) \quad (4.1)$$

to recover $p_{TE}(d|s, T, T_{\text{eff.}})$ from the joint distribution $p_T(d|s, T)$. To estimate $T_{\text{eff.}}$ we minimized

$$\sum_d (p_{TM}(d|s, T) - p_{TE}(d|s, T, T_{\text{eff.}}))^2 \quad (4.2)$$

for each sequence s in the miRNA data set.

Conclusion

*We shall not cease from exploration
And the end of all our exploring
Will be to arrive where we started
And know the place for the first time.*

T.S. Eliot – "Little Gidding"

The idea that the evolution of life on Earth has been (and continues to be) driven by natural selection ("survival of the fittest") is, to quote Karl Popper, "an immensely impressive and powerful theory" [160]. Popper's contention that the claim that natural selection completely explains evolution is a bold claim, a claim that is far from being established, is just as true today as it was three decades ago. The population genetic underpinnings of modern evolutionary theory have been well tested, and so has the phylogenetic theory of evolution which says that all life on Earth has evolved from a few (most probably a single) primitive unicellular organisms. Much remains to be understood, however, of how natural selection coupled with the nonadaptive forces of mutation, recombination and random genetic drift has led to the emergent complexity and diversity of life observed today. The challenge of integrating the vast information in the genetic heritage of extant life with the progression of forms in the fossil record lies before us, and so does the problem of deciphering the processes by which life has climbed the coevolutionary ladder[161]; reconstructing the steps of the recursive dance of life and environment in which new levels of organization emerge to exploit new opportunities in an environment transformed by the previous revolution.

It has been argued that the science of molecular evolution is in its golden age. To fulfill its promise, however, the relentless boom in cataloging patterns of diversity in organization and function at various levels of biological organization ranging from gene regulatory networks to cell signaling pathways and beyond must progress toward a quantitative synthesis of this ever growing body of knowledge. The road to such a quantitative synthesis will, I believe, be paved by a greater understanding of biological function and complexity as an emergent result of fundamental evolutionary processes. And will lead ultimately,

we may hope, to a theory of evolution that is able to explain long time scale evolutionary patterns and emergent levels of complexity, as the results of an evolutionary dynamic driven by natural selection. An evolutionary dynamic driven by natural selection that is: fueled by genetic variation, played out in an intricate coevolving dance of populations and environment, and subject to the physical constraints of function and interaction.

Bibliography

- [1] Spencer, H (1864) *The Principles of Biology*. (Williams and Norgate, London, UK).
- [2] Lynch, M (2007) *The Origins of Genomic Architecture* (Sinauer, Sunderland, USA).
- [3] Fontana, W (2005) *The Topology of the Possible in: "Understanding Change: Models, Methodologies and Metaphors."* (Palgrave Macmillan).
- [4] Darwin, C (1859) *The Origin of Species* (John Murray, London UK).
- [5] Smith, JM (1978) *The evolution of sex*. (Cambridge Univ. Press, Cambridge, UK).
- [6] Axelrod, R (1984) *The Evolution of Cooperation* (Basic Books, New York, USA).
- [7] Gould, SJ, N, E (1977) Punctuated equilibria: the tempo and mode of evolution reconsidered. *Paleobiology* 3:115.
- [8] Hutchinson, GE (1961) The paradox of the plankton. *American Naturalist* 95:137–145.
- [9] Malthus, TR (1798) *An Essay on the Principle of Population* (J. Johnson, London, UK).
- [10] Verhulst, PF (1838) Notice sur la loi que la population poursuit dans son accroissement. *Correspondance mathématique et physique* 10:113–121.
- [11] Gillespie, JH (2004) *Population Genetics: A Concise Guide* (John Hopkins University Press, Baltimore, USA).
- [12] Kimura, M (1983) *The Neutral Theory of Molecular Evolution* (Cambridge Univ. Press, Cambridge, UK).
- [13] Sella, G, Hirsh, AE (2005) The application of statistical physics to evolutionary biology. *Proc Natl Acad Sci U S A* 102:9541–9546.
- [14] Zuckerkandl, E, Pauling, L (1965) *Evolutionary divergence and convergence in proteins. in "Evolving Genes and Proteins"* (Academic Press, New York).
- [15] Kimura, M (1968) Evolutionary rate at the molecular level. *Nature* 217:624–626.
- [16] King, JL, Jukes, TL (1969) Non-darwinian evolution. *Science* 164:788–798.
- [17] Mustonen, V, Lassig, M (2007) Adaptations to fluctuating selection in drosophila. *Proc Natl Acad Sci U S A* 104:2277–2282.
- [18] Mustonen, V, Lassig, M (2009) From fitness landscapes to seascales: non-equilibrium dynamics of selection and adaptation. *Trends Genet* 25:111–119.

- [19] Liti, G et al. (2009) Population genomics of domestic and wild yeasts. *Nature* 458:337–41.
- [20] Fay, JC, Wyckoff, GJ, Wu, CI (2002) Testing the neutral theory of molecular evolution with genomic data from drosophila. *Nature* 415:1024–6.
- [21] Eyre-Walker, A (2006) The genomic rate of adaptive evolution. *Trends Ecol Evol (Amst)* 21:569–75.
- [22] Andolfatto, P (2005) Adaptive evolution of non-coding dna in drosophila. *Nature* 437:1149–52.
- [23] Bachtrog, D (2005) Sex chromosome evolution: molecular aspects of y-chromosome degeneration in drosophila. *Genome Res* 15:1393–1401.
- [24] Fisher, RA (1930) *The Genetical Theory of Natural Selection* (Oxford University Press, Oxford, UK).
- [25] Smith, JM, Price, G (1973) The logic of animal conflict. *Nature* 246:15–18.
- [26] Hofbauer, J, Sigmund, K (1991) *Evolutionary Game Theory and Population Dynamics* (Cambridge University Press, Cambridge, UK).
- [27] Nowak, MA, May, RM (1992) Evolutionary games and spatial chaos. *Nature* 359:826–829.
- [28] Lieberman, E, Hauert, C, Nowak, MA (2005) Evolutionary dynamics on graphs. *Nature* 433:312–316.
- [29] Axelrod, R, Hamilton, WD (1981) The evolution of cooperation. *Science* 211:1390–1396.
- [30] Nowak, MA, Bonhoeffer, S, May, RM (1994) Spatial games and the maintenance of cooperation. *Proc Natl Acad Sci U S A* 91:4877–4881.
- [31] Szabo, G, Hauert, C (2002) Phase transitions and volunteering in spatial public goods games. *Phys Rev Lett* 89:118101.
- [32] Kerr, B, Riley, MA, Feldman, MW, Bohannan, BJM (2002) Local dispersal promotes biodiversity in a real-life game of rock-paper-scissors. *Nature* 418:171–174.
- [33] Nowak, MA, Sigmund, K (2002) Biodiversity: Bacterial game dynamics. *Nature* 418:138–139.
- [34] Kirkup, BC, Riley, MA (2004) Antibiotic-mediated antagonism leads to a bacterial game of rock-paper-scissors in vivo. *Nature* 428:412–414.
- [35] Czarán, TL, Hoekstra, RF, Pagie, L (2002) Chemical warfare between microbes promotes biodiversity. *Proc Natl Acad Sci U S A* 99:786–790.
- [36] Lenski, RE, Riley, MA (2002) Chemical warfare from an ecological perspective. *Proc Natl Acad Sci U S A* 99:556–558.
- [37] Von Neumann, J, Morgenstern, O (1944) *Theory of games and economic behavior* (Princeton University Press, Princeton).
- [38] Nash, JF (1950) Equilibrium points in n-person games. *Proc Natl Acad Sci U S A* 36:48–49.

- [39] Traulsen, A, Claussen, JC, Hauert, C (2005) Coevolutionary dynamics: from finite to infinite populations. *Phys Rev Lett* 95:238701.
- [40] Traulsen, A, Claussen, JC, Hauert, C (2006) Coevolutionary dynamics in large, but finite populations. *Phys Rev E Stat Nonlin Soft Matter Phys* 74:011901.
- [41] Moran, PAP (1962) *The Statistical Processes of Evolutionary Theory* (Clarendon, New York, USA).
- [42] Van Kampen, NG (1981) *Stochastic processes in physics and chemistry*. (Elsevier, Amsterdam, The Neatherlands).
- [43] Szabó, G, Tóke, C (1998) Evolutionary prisoner's dilemma game on a square lattice. *Phys. Rev. E* 58:69.
- [44] Szabo, G, Vukov, J, Szolnoki, A (2005) Phase diagrams for an evolutionary prisoner's dilemma game on two-dimensional lattices. *Phys Rev E Stat Nonlin Soft Matter Phys* 72:047107.
- [45] Hauert, C, Gy, S (2005) Game theory and physics. *Am. J. Phys.* 73:405–414.
- [46] Santos, FC, Pacheco, JM (2005) Scale-free networks provide a unifying framework for the emergence of cooperation. *Phys Rev Lett* 95:098104.
- [47] Ohtsuki, H, Hauert, C, Lieberman, E, Nowak, MA (2006) A simple rule for the evolution of cooperation on graphs and social networks. *Nature* 441:502–505.
- [48] Hanski, IA, Gilpin, ME (1997) *Metapopultaion Biology* (Academic Press, New York, USA).
- [49] Hanski, IA (1998) Metapopulation dynamics. *Nature* 396:41.
- [50] Pannell, JR, Charlesworth, B (2000) Effects of metapopulation processes on measures of genetic diversity. *Philos Trans R Soc Lond B Biol Sci* 355:1851–1864.
- [51] Szollosi, GJ, Derenyi, I, Vellai, T (2006) The maintenance of sex in bacteria is ensured by its potential to reload genes. *Genetics* 174:2173–2180.
- [52] Szollosi, GJ, Derenyi, I (2008) Evolutionary games on minimally structured populations. *Phys Rev E Stat Nonlin Soft Matter Phys* 78:031919.
- [53] Wang, S, Szalay, MS, Zhang, C, Csermely, P (2008) Learning and innovative elements of strategy adoption rules expand cooperative network topologies. *PLoS ONE* 3:e1917.
- [54] Nowak, MA (2006) Five rules for the evolution of cooperation. *Science* 314:1560–1563.
- [55] Sinervo, B, Lively, CM (1996) The rock-paper-scissors game and the evolution of alternative male strategies. *Nature* 380.
- [56] Reichenbach, T, Mobilia, M, Frey, E (2006) Coexistence versus extinction in the stochastic cyclic lotka-volterra model. *Phys Rev E Stat Nonlin Soft Matter Phys* 74:051907.
- [57] Szolnoki, A, Szabo, G (2004) Phase transitions for rock-scissors-paper game on different networks. *Phys Rev E Stat Nonlin Soft Matter Phys* 70:037102.
- [58] Szabó, G, Szolnoki, A, Izsák, R (2004) Rock-scissors-paper game on regular small-world networks. *J. Phys. A* 37:2599–2606.

- [59] Nowak, MA, Sasaki, A, Taylor, C, Fudenberg, D (2004) Emergence of cooperation and evolutionary stability in finite populations. *Nature* 428:646–650.
- [60] Imhof, LA, Fudenberg, D, Nowak, MA (2005) Evolutionary cycles of cooperation and defection. *Proc Natl Acad Sci U S A* 102:10797–10800.
- [61] Molander, P (1985) The optimal level of generosity in a selfish, uncertain environment. *J. Conflict Resolut.* 29:611–618.
- [62] Taylor, PD, Day, T, Wild, G (2007) Evolution of cooperation in a finite homogeneous graph. *Nature* 447:469–472.
- [63] Vellai, T, Kovacs, AL, Kovacs, G, Ortutay, C (1999) Genome economisation and a new approach to the species concept in bacteria. *Proc. R. Soc. Lond. B* 266:1953–1958.
- [64] Mira, A, Ochman, H, Moran, NA (2001) Deletional bias and the evolution of bacterial genomes. *Trends Genet* 17:589–596.
- [65] Vellai, T, Takacs, K, Vida, G (1998) A new aspect to the origin and evolution of eukaryotes. *J Mol Evol* 46:499–507.
- [66] Avery, OT, MacLeod, CM, McCarty, M (1944) Studies on the chemical nature of the substance inducing transformation of pneumococcal types. *J. Exp. Med.* 79:137–158.
- [67] Solomon, JM, Grossman, AD (1996) Who's competent and when: regulation of natural genetic competence in bacteria. *Trends Genet* 12:150–155.
- [68] Thomas, CM, Nielsen, KM (2005) Mechanisms of, and barriers to, horizontal gene transfer between bacteria. *Nat Rev Microbiol* 3:711–721.
- [69] Levin, BR, Bergstrom, CT (2000) Bacteria are different: observations, interpretations, speculations, and opinions about the mechanisms of adaptive evolution in prokaryotes. *Proc Natl Acad Sci U S A* 97:6981–6985.
- [70] Feil, EJ et al. (2001) Recombination within natural populations of pathogenic bacteria: short-term empirical estimates and long-term phylogenetic consequences. *Proc Natl Acad Sci U S A* 98:182–187.
- [71] Feil, EJ (2004) Small change: keeping pace with microevolution. *Nat Rev Microbiol* 2:483–495.
- [72] Smith, JM, Dowson, CG, Spratt, BG (1991) Localised sex in bacteria. *Nature* 349:29–31.
- [73] Smith, JM (1993) The role of sex in bacterial evolution. *J Hered* 84:326–327.
- [74] Bernstein, H, Byerly, HC, Hopf, FA, Michod, RE (1985) Genetic damage, mutation, and the evolution of sex. *Science* 229:1277–1281.
- [75] Elena, SF, Lenski, RE (1997) Test of synergistic interactions among deleterious mutations in bacteria. *Nature* 390:395–398.
- [76] Butlin, R (2002) Evolution of sex: The costs and benefits of sex: new insights from old asexual lineages. *Nat Rev Genet* 3:311–317.

- [77] Otto, SP, Lenormand, T (2002) Resolving the paradox of sex and recombination. *Nat Rev Genet* 3:252–261.
- [78] Redfield, RJ, Schrag, MR, Dean, AM (1997) The evolution of bacterial transformation: sex with poor relations. *Genetics* 146:27–38.
- [79] Cohen, E, Kessler, DA, Levine, H (2005) Recombination dramatically speeds up evolution of finite populations. *Phys Rev Lett* 94:098102.
- [80] Holmes, EC, Urwin, R, Maiden, MC (1999) The influence of recombination on the population structure and evolution of the human pathogen neisseria meningitidis. *Mol Biol Evol* 16:741–749.
- [81] Vetsigian, K, Goldenfeld, N (2005) Global divergence of microbial genome sequences mediated by propagating fronts. *Proc Natl Acad Sci U S A* 102:7332–7337.
- [82] Redfield, RJ (2001) Do bacteria have sex? *Nat. Rev. Genet.* 2:634–639.
- [83] Woese, C (1998) The universal ancestor. *Proc Natl Acad Sci U S A* 95:6854–6859.
- [84] Ochman, H, Lawrence, JG, Groisman, EA (2000) Lateral gene transfer and the nature of bacterial innovation. *Nature* 405:299–304.
- [85] Koonin, EV (2003) Comparative genomics, minimal gene-sets and the last universal common ancestor. *Nat Rev Microbiol* 1:127–136.
- [86] Abby, S, Daubin, V (2007) Comparative genomics and the evolution of prokaryotes. *Trends Microbiol* 15:135–141.
- [87] Lerat, E, Daubin, V, Moran, NA (2003) From gene trees to organismal phylogeny in prokaryotes: the case of the gamma-proteobacteria. *PLoS Biol* 1:E19.
- [88] Welch, RA et al. (2002) Extensive mosaic structure revealed by the complete genome sequence of uropathogenic escherichia coli. *Proc Natl Acad Sci U S A* 99:17020–17024.
- [89] Chen, SL et al. (2006) Identification of genes subject to positive selection in uropathogenic strains of escherichia coli: a comparative genomics approach. *Proc Natl Acad Sci U S A* 103:5977–5982.
- [90] Tettelin, H et al. (2005) Genome analysis of multiple pathogenic isolates of streptococcus agalactiae: implications for the microbial "pan-genome". *Proc Natl Acad Sci U S A* 102:13950–13955.
- [91] Claverys, JP, Prudhomme, M, Mortier-Barriere, I, Martin, B (2000) Adaptation to the environment: Streptococcus pneumoniae, a paradigm for recombination-mediated genetic plasticity? *Mol Microbiol* 35:251–259.
- [92] Nakamura, Y, Itoh, T, Matsuda, H, Gojobori, T (2004) Biased biological function of horizontally transferred genes in prokaryotic genomes. *Nat. Genet.* 36:760–766.
- [93] Alm, RA et al. (1999) Genomic-sequence comparison of two unrelated isolates of the human gastric pathogen helicobacter pylori. *Nature* 397:176–180.
- [94] Whitlock, MC, Barton, NH (1997) The effective size of a subdivided population. *Genetics* 146:427–441.

- [95] Szabo, G, Antal, T, Szabo, P, Droz, M (2000) Spatial evolutionary prisoner's dilemma game with three strategies and external constraints. *Phys Rev E Stat Phys Plasmas Fluids Relat Interdiscip Topics* 62:1095–1103.
- [96] Killingback, T, Bieri, J, Flatt, T (2006) Evolution in group-structured populations can resolve the tragedy of the commons. *Proc Biol Sci* 273:1477–1481.
- [97] Grafen, A (2007) Detecting kin selection at work using inclusive fitness. *Proc Biol Sci* 274:713–719.
- [98] Traulsen, A, Nowak, MA (2006) Evolution of cooperation by multilevel selection. *Proc Natl Acad Sci U S A* 103:10952–10955.
- [99] Lehmann, L, Keller, L, West, S, Roze, D (2007) Group selection and kin selection: two concepts but one process. *Proc Natl Acad Sci U S A* 104:6736–6739.
- [100] Meyers, AM, Bull, JJ (2002) Fighting change with change: adaptive variation in an uncertain world. *Trends Ecol. Evol.* 17:551–557.
- [101] Szathmáry, E, Smith, JM (1997) *The major transitions in evolution* (Oxford University Press, Oxford, UK).
- [102] Dykhuizen, DE, Green, L (1991) Recombination in *Escherichia coli* and the definition of biological species. *J Bacteriol* 173:7257–7268.
- [103] Hugenholtz, P, Tyson, GW (2008) Microbiology: metagenomics. *Nature* 455:481–483.
- [104] Yooseph, S et al. (2007) The sorcerer ii global ocean sampling expedition: expanding the universe of protein families. *PLoS Biol* 5:e16.
- [105] Fontana, W, Schuster, P (1998) Continuity in evolution: on the nature of transitions. *Science* 280:1451–1455.
- [106] Stadler, BM, Stadler, PF, Wagner, GP, Fontana, W (2001) The topology of the possible: formal spaces underlying patterns of evolutionary change. *J Theor Biol* 213:241–274.
- [107] Wagner, A (2008) Neutralism and selectionism: a network-based reconciliation. *Nat Rev Genet* 9:965–974.
- [108] Szollosi, GJ, Derenyi, I (2008) The effect of recombination on the neutral evolution of genetic robustness. *Math Biosci* 214:58–62.
- [109] Borenstein, E, Ruppin, E (2006) Direct evolution of genetic robustness in microRNA. *Proc Natl Acad Sci U S A* 103:6593–6598.
- [110] Szollosi, GJ, Derenyi, I (2009) Congruent evolution of genetic and environmental robustness in micro-rna. *Mol Biol Evol* 26:867–874.
- [111] de Visser, JAGM et al. (2003) Perspective: Evolution and detection of genetic robustness. *Evolution* 57:1959–1972.
- [112] Waddington, CH (1957) *The Strategy of the Genes: A Discussion of Some Aspects of Theoretical Biology* (MacMillan, New York USA).

- [113] Ciliberti, S, Martin, OC, Wagner, A (2007) Robustness can evolve gradually in complex regulatory gene networks with varying topology. *PLoS Comput Biol* 3:e15.
- [114] Zauner, H, Sommer, RJ (2007) Evolution of robustness in the signaling network of pristinonchus vulva development. *Proc Natl Acad Sci U S A* 104:10086–10091.
- [115] Hillenmeyer, ME et al. (2008) The chemical genomic portrait of yeast: uncovering a phenotype for all genes. *Science* 320:362–365.
- [116] Mayo, O, Bürger, R (1997) Evolution of dominance: a theory whose time has passed? *Biol. Rev.* 72:97–110.
- [117] Fisher, RA (1928) The possible modifications of the response of the wild type to recurrent mutations. *American Naturalist* 62:115–126.
- [118] Haldane, J (1930) A note on fisher’s theory of dominance. *American Naturalist* 64.
- [119] Wright, S (1934) Physiological and evolutionary theories of dominance. *American Naturalist* 68:25–53.
- [120] Kacser, H, Burns, JA (1981) The molecular basis of dominance. *Genetics* 97:6639–6666.
- [121] Krakauer, DC, Plotkin, JB (2002) Redundancy, antiredundancy, and the robustness of genomes. *Proc Natl Acad Sci U S A* 99:1405–1409.
- [122] Siegal, ML, Bergman, A (2002) Waddington’s canalization revisited: developmental stability and evolution. *Proc Natl Acad Sci U S A* 99:10528–10532.
- [123] Azevedo, RBR, Lohaus, R, Srinivasan, S, Dang, KK, Burch, CL (2006) Sexual reproduction selects for robustness and negative epistasis in artificial gene networks. *Nature* 440:87–90.
- [124] Ciliberti, S, Martin, OC, Wagner, A (2007) Innovation and robustness in complex regulatory gene networks. *Proc Natl Acad Sci U S A* 104:13591–13596.
- [125] Crombach, A, Hogeweg, P (2008) Evolution of evolvability in gene regulatory networks. *PLoS Comput Biol* 4:e1000112.
- [126] Montville, R, Froissart, R, Remold, SK, Tenailon, O, Turner, PE (2005) Evolution of mutational robustness in an rna virus. *PLoS Biol* 3:e381.
- [127] Wagner, A, Stadler, PF (1999) Viral rna and evolved mutational robustness. *J Exp Zool* 285:119–127.
- [128] Wilke, CO, Adami, C (2003) Evolution of mutational robustness. *Mutat Res* 522:3–11.
- [129] Sanjuan, R, Cuevas, JM, Furio, V, Holmes, EC, Moya, A (2007) Selection for robustness in mutagenized rna viruses. *PLoS Genet* 3:e93.
- [130] Bloom, JD et al. (2007) Evolution favors protein mutational robustness in sufficiently large populations. *BMC Biol* 5:29.
- [131] van Nimwegen, E, Crutchfield, JP, Huynen, M (1999) Neutral evolution of mutational robustness. *Proc Natl Acad Sci U S A* 96:9716–9720.
- [132] Forster, R, Adami, C, Wilke, CO (2006) Selection for mutational robustness in finite populations. *J Theor Biol* 243:181–190.

- [133] Partensky, F, Hess, WR, Vaulot, D (1999) Prochlorococcus, a marine photosynthetic prokaryote of global significance. *Microbiol Mol Biol Rev* 63:106–127.
- [134] Lynch, M, Conery, JS (2003) The origins of genome complexity. *Science* 302:1401–1404.
- [135] Haldane, JBS (1937) The effect of variation of fitness. *American Naturalist* 71:337–350.
- [136] Kimura, M, Maruyama, T (1966) The mutational load with epistatic gene interactions in fitness. *Genetics* 54:1337–1351.
- [137] Huynen, MA, Hogeweg, P (1994) Pattern generation in molecular evolution: exploitation of the variation in rna landscapes. *J Mol Evol* 39:71–79.
- [138] Bornberg-Bauer, E, Chan, HS (1999) Modeling evolutionary landscapes: mutational stability, topology, and superfunnels in sequence space. *Proc Natl Acad Sci U S A* 96:10689–10694.
- [139] Reidys, C, Stadler, PF, Schuster, P (1997) Generic properties of combinatorial maps: neutral networks of rna secondary structures. *Bull Math Biol* 59:339–397.
- [140] Jorg, T, Martin, OC, Wagner, A (2008) Neutral network sizes of biological rna molecules can be computed and are not atypically small. *BMC Bioinformatics* 9:464.
- [141] Xia, Y, Levitt, M (2002) Roles of mutation and recombination in the evolution of protein thermodynamics. *Proc Natl Acad Sci U S A* 99:10382–10387.
- [142] Lagos-Quintana, M, Rauhut, R, Lendeckel, W, Tuschl, T (2001) Identification of novel genes coding for small expressed rnas. *Science* 294:853–858.
- [143] Lau, NC, Lim, LP, Weinstein, EG, Bartel, DP (2001) An abundant class of tiny rnas with probable regulatory roles in caenorhabditis elegans. *Science* 294:858–862.
- [144] Lee, RC, Ambros, V (2001) An extensive class of small rnas in caenorhabditis elegans. *Science* 294:862–864.
- [145] Bartel, DP (2004) Micrnas: genomics, biogenesis, mechanism, and function. *Cell* 116:281–297.
- [146] Zamore, PD, Haley, B (2005) Ribo-gnome: the big world of small rnas. *Science* 309:1519–1524.
- [147] Shabalina, SA, Koonin, EV (2008) Origins and evolution of eukaryotic rna interference. *Trends Ecol Evol* 23:578–587.
- [148] Jinek, M, Doudna, JA (2009) A three-dimensional view of the molecular machinery of rna interference. *Nature* 457:405–412.
- [149] Shu, W, Bo, X, Ni, M, Zheng, Z, Wang, S (2007) In silico genetic robustness analysis of microrna secondary structures: potential evidence of congruent evolution in microrna. *BMC Evol Biol* 7:223.
- [150] Di Giulio, M (1989) The extension reached by the minimization of the polarity distances during the evolution of the genetic code. *J Mol Evol* 29:288–293.

- [151] Haig, D, Hurst, L (1999) A quantitative measure of error minimization in the genetic code. *J Mol Evol* 49:708.
- [152] Ritchie W, Legendre M, GD (2007) Rna stem-loops: to be or not to be cleaved by rnase iii. *RNA* 13:457–462.
- [153] Ancel, LW, Fontana, W (2000) Plasticity, evolvability, and modularity in rna. *J Exp Zool* 288:242–283.
- [154] DePristo, MA, Weinreich, DM, Hartl, DL (2005) Missense meanderings in sequence space: a biophysical view of protein evolution. *Nat Rev Genet* 6:678–687.
- [155] Hofacker, IL et al. (1994) Fast folding and comparison of rna secondary structures. *Monatshfte für Chemie* 125:167–188.
- [156] Zuker, M, Stiegler, P (1981) Optimal computer folding of large rna sequences using thermodynamics and auxiliary information. *Nucleic Acids Res* 9:133–148.
- [157] McCaskill, JS (1990) The equilibrium partition function and base pair binding probabilities for rna secondary structure. *Biopolymers* 29:1105–1119.
- [158] Wuchty, S, Fontana, W, Hofacker, IL, Schuster, P (1999) Complete suboptimal folding of rna and the stability of secondary structures. *Biopolymers* 49:145–165.
- [159] Griffiths-Jones, S, Grocock, RJ, van Dongen, S, Bateman, A, Enright, AJ (2006) mirbase: microRNA sequences, targets and gene nomenclature. *Nucleic Acids Res* 34:D140–4.
- [160] Popper, K (1978) Natural selection and the emergence of mind. *Dialectica* 32:339–355.
- [161] Lenton, TM, Schellnhuber, HJ, Szathmary, E (2004) Climbing the co-evolution ladder. *Nature* 431:913.

Summary

Darwin laid the foundations of evolutionary biology on two central ideas: (i) all species are related to one another through a history of common descent, and (ii) the exquisite match between a species and its environment is explained by natural selection. These two ideas suggest a dichotomy of long and short time scales in evolution. We explore this dichotomy from a statistical physics perspective; undertaking a bottom up program of trying to discern statistical laws, and understand emergent complexity, from the behaviour of large numbers of simple individuals subject to shared constraints.

The first half of the thesis is concerned with the effects of population structure. We derive the dynamics governing the evolution of a minimally structured population from fundamental individual level stochastic processes. This allows us to demonstrate that a continuous transition leading to the dominance of cooperation exists in populations with hierarchical levels of interaction. Applying our model of spatial structure to the repeated prisoner's dilemma we uncovered a novel and counterintuitive mechanism by which the constant influx of defectors sustains cooperation. Further exploring the phase space of the repeated prisoner's dilemma and also of the "rock-paper-scissor" game we find indications of rich structure and are able to reproduce several effects observed in other models with explicit spatial embedding.

Turning to the problem of natural genetic transformation (NGT), a sexual process by which bacteria actively take up exogenous DNA and use it to replace homologous chromosomal sequences. We develop a novel simulation approach for the long-term dynamics of genome organization (involving the loss and acquisition of genes) in a bacterial species consisting of a large number of spatially distinct populations subject to independently fluctuating ecological conditions. We show that in the presence of weak inter-population migration NGT is able to subsist as a mechanism to reload locally lost, intermittently selected genes from the collective gene pool of the species through DNA uptake from migrants. The machinery of transformation survives under a wide range of model parameters readily encompassing real-world biological conditions. These findings imply that the primary role of NGT is not to serve the cell with food, but to provide homologous sequences for restoring genes that have disappeared from or become degraded in the local population.

In the second half of the thesis we examine the effects of the degeneracy the genotype-phenotype mapping. It has previously been demonstrated that populations of genotypes evolving on the neutral networks corresponding to all genotypes with the same secondary structure only through neutral mutations can evolve mutational robustness, by concentrating the population on regions of high neutrality. We introduce recombination, and demonstrate, through numerically calculating the stationary distribution of an infinite population on ensembles of random neutral networks that mutational robustness is significantly enhanced. We simulated finite populations of genotypes evolving on random neutral networks chosen from random network ensembles with different topology, and also a scaled down microRNA neutral network. We show that even in finite populations recombination will still act to focus the population on regions of locally high neutrality.

Finally, we turn to the problem of the origin of empirically observed genetic robustness. It unclear whether genetic robustness has evolved directly by natural selection or is a correlated byproduct of selection for environmental robustness. We demonstrated that sampling method used previously introduced a significant bias that lead to an overestimation of genetic robustness among miRNA. We develop a novel measure of environmental robustness based on the equilibrium thermodynamic ensemble of secondary structures of the miRNA precursor sequences. Using this measure we demonstrate that the biophysics of RNA folding induces a high level of correlation between genetic (mutational) and environmental (thermodynamic) robustness. Based on the small effective population sizes among multicellular eukaryotes we argue that the correlation between mutational and thermodynamic robustness strongly suggests that genetic robustness observed in miRNA sequences is the byproduct of selection for environmental robustness.

Összefoglaló

Darwin két gondolattal alapozta meg az evolúcióbíológia tudományát: (i) a fajok eredete közös, azok rokonok a közös leszármazás folytán, (ii) a fajok és környezetük közötti szoros illeszkedés a természetes szelekció eredménye. Ez a két gondolat az evolúció folyamatában a hosszú és rövid időskálák kettősségét sugallja. Ezt a kettősséget vizsgáljuk statisztikus fizikai szemzőgből: statisztikus statisztikus törvényszerűségek kikövetkeztetését és a kibontakozó komplexitást megértését, kísérelve meg nagyszámú egyszerű, közös kényszereknek kitett egyed viselkedésének vizsgálatán keresztül.

Az értekezés első fele a populáció struktúrájának hatásaival foglalkozik. Alapvető, egyedszintű folyamatoktól kiindulva levezetjük a minimálisan strukturált populációt leíró dinamikát. Ezt felhasználva megmutatjuk, hogy az együttműködés dominanciájához vezető, folytonos fázisátalakulás játszódik le minimálisan strukturált populációkban. Populációstruktúra-modellünket az az ismételt fogolydilemma-játékokra alkalmazva egy váratlan mechanizmusra derül fény, amely során az "árulók" folyamatos beáramlása tartja fent az együttműködést. Az ismételt fogolydilemma és a "kő-papír-olló" játékok fázisátalakulását továbbvizsgálva gazdag szerkezetre utaló eredményeket kapunk, és képesek vagyunk több hatást reprodukálni, melyet eddig csak explicit térbeli beágyazás során figyeltek meg.

Ezt követően a természetes genetikai transzformáció (natural genetic transformation: NGT) fennmaradásának problémájával foglalkozunk. Az NGT egy szexuális folyamat, mely során a baktériumok külső eredetű DNS-t vesznek fel és cserélnek ki kromoszomális szekvenciákkal. Újszerű szimulációs megközelítést alkalmazunk a hosszútávú genomorganizációs dinamika (gének felvétele és elvesztése) leírására egy bakteriális fajban, amely faj nagyszámú térben elkülönült, független környezeti fluktuációknak kitett populációból áll. Megmutatjuk, hogy gyenge populáció közti migráció jelenlétében az NGT képes fennmaradni, lokálisan már elveszett, váltakozóan szelektált géneket "visszatölteni" képes mechanizmusként. A transzformáció mechanizmusa a modellparaméterek széles tartományában képes fennmaradni, amely tartomány könnyedén magában foglalja a tényleges biológiai körülményeket.

Az értekezés második felében a degenerált genotípus-fenotípus leképezés hatásait vizsgáljuk. Mint azt már mások megmutatták, egy adott fenotípushoz tartozó genotípusokból álló neutrális hálózaton, kizárólag neutrális mutációkon keresztül evolválódva, a populáció a neutrális hálózat fokozott neutralitással rendelkező régióin koncentrálnak, tehát a mutációs robusztusság közvetlen evolúcióját tapasztaljuk. A végtelen populációs határeset stacionárius állapotát numerikusan kiszámítva megmutatjuk, hogy a rekombináció jelentősen emeli a mutációs robusztusság mértékét. Ezt követően véges populációk dinamikáját vizsgáljuk különböző topológiájú, véletlen neutrális hálózatokon és egy lekicsinyített mikroRNS neutrális hálózaton. Megmutatjuk, hogy a neutrális hálózat ugyan csökkent mértékben, de véges populációkban is képes lokálisan neutrálisabb régiókban koncentrálni a populációt.

Befejezésül az empirikusan megfigyelhető genetikai robusztusság problémájával foglalkozunk. Nem egyértelmű, hogy ez a genetikai robusztusság közvetlenül evolválódott-e, vagy pedig a környezeti robusztusság korrelált mellékterméke. Megmutatjuk, hogy az eddig alkalmazott mintavételezési eljárás a genetikai robusztusság túlbecsüléséhez vezetett. A környezeti robusztusság egy új mértékét bevezetve megmutatjuk, hogy az RNS folding erős korrelációt indukál a genetikai (mutációs) és környezeti (termodinamikai) robusztusság között. A többsejtű eukariótákra jellemző, kis effektív populációméreteket figyelembe véve úgy véljük, a mutációs és termodinamikai robusztusság közötti korreláció erős bizonyíték arra, hogy a mikroRNS-ek körében tapasztalható genetikai robusztusság a termodinamikai robusztusságra való szelekció mellékterméke.