



# **From simple molecules to complex structures**

## **A constructive role for neutrality in evolution?**

Peter Schuster

Institut für Theoretische Chemie, Universität Wien, Austria

and

The Santa Fe Institute, Santa Fe, New Mexico, USA



Alife XI Conference

Winchester, 04.– 08.08.2008

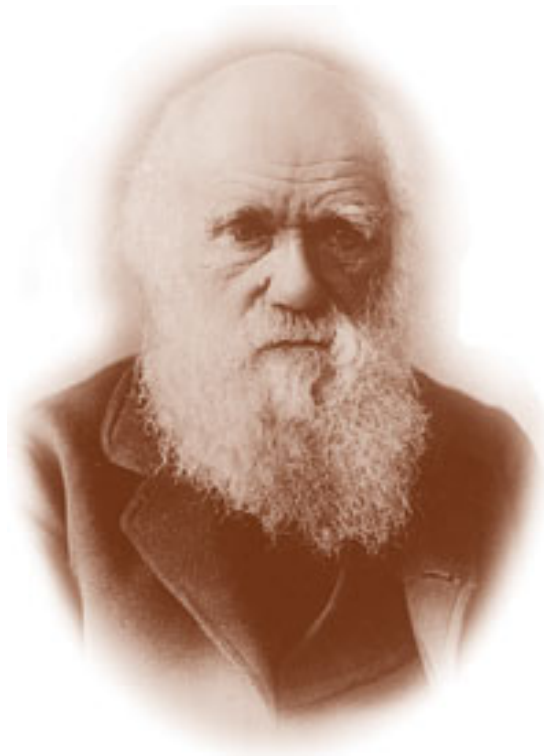
Web-Page for further information:

<http://www.tbi.univie.ac.at/~pks>

# What is neutrality ?

Selective neutrality =  
= several genotypes having the **same fitness**.

Structural neutrality =  
= several genotypes forming molecules with  
the **same structure**.



ON  
THE ORIGIN OF SPECIES

BY MEANS OF NATURAL SELECTION,

OR THE

PRESERVATION OF FAVOURED RACES IN THE STRUGGLE  
FOR LIFE.

By CHARLES DARWIN, M.A.,

FELLOW OF THE ROYAL, GEOLOGICAL, LINNEAN, ETC., SOCIETIES;  
AUTHOR OF 'JOURNAL OF RESEARCHES DURING H. M. S. BEAGLE'S VOYAGE  
ROUND THE WORLD.'

LONDON:  
JOHN MURRAY, ALBEMARLE STREET.

1859.

*The right of Translation is reserved.*

This preservation of favourable individual differences and variations, and the destruction of those which are injurious, I have called Natural Selection, or the Survival of the Fittest. Variations neither useful nor injurious would not be affected by natural selection, and would be left either a fluctuating element, as perhaps we see in certain polymorphic species, or would ultimately become fixed, owing to the nature of the organism and the nature of the conditions.

Charles Darwin. *The Origin of Species*. Sixth edition. John Murray. London: 1872



Motoo Kimuras Populationsgenetik der neutralen Evolution.

Evolutionary rate at the molecular level.  
*Nature* **217**: 624-626, 1955.

*The Neutral Theory of Molecular Evolution.*  
Cambridge University Press. Cambridge,  
UK, 1983.

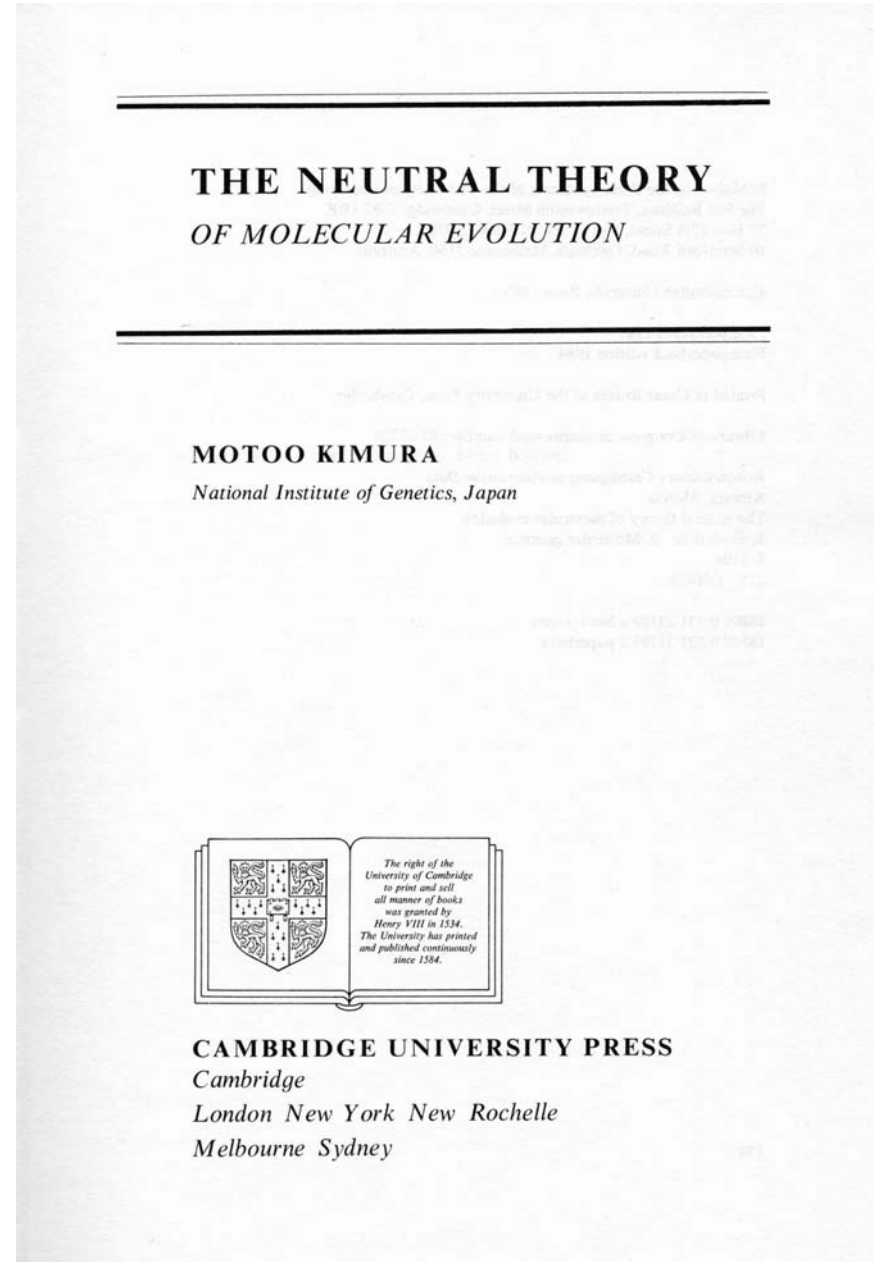
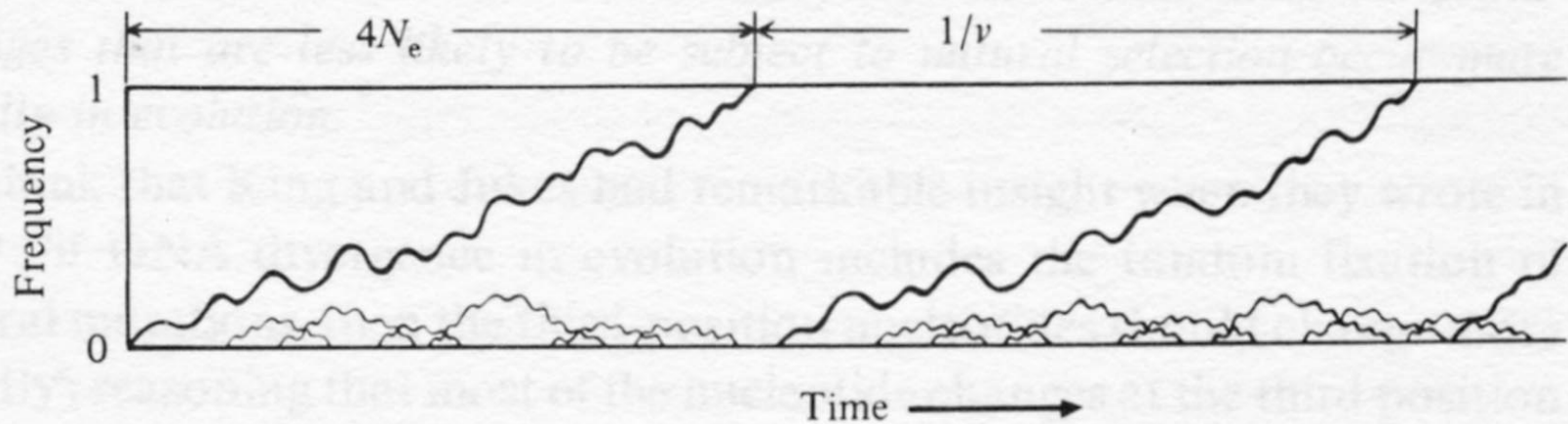


Fig. 3.1. Behavior of mutant genes following their appearance in a finite population. Courses of change in the frequencies of mutants destined to fixation are depicted by thick paths.  $N_e$  stands for the effective population size and  $v$  is the mutation rate.



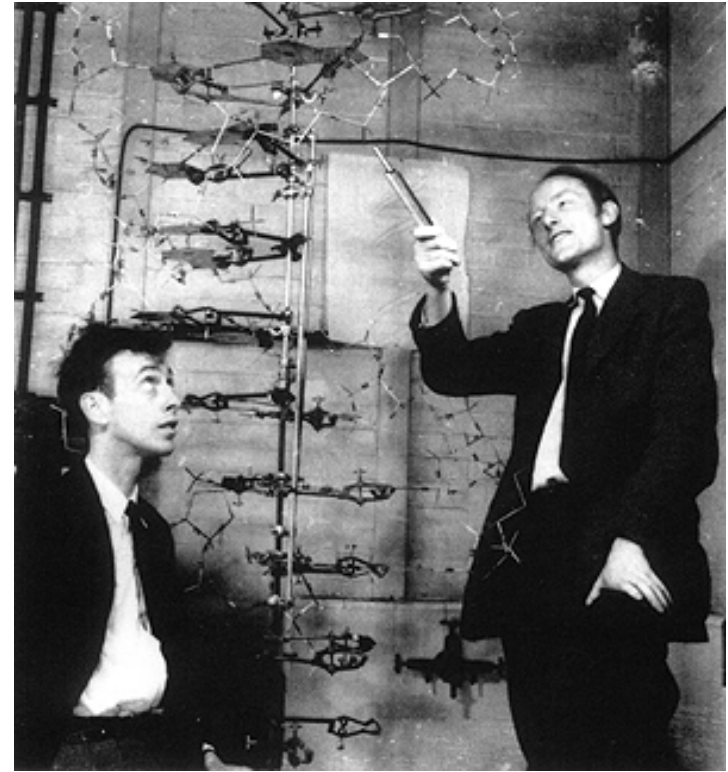
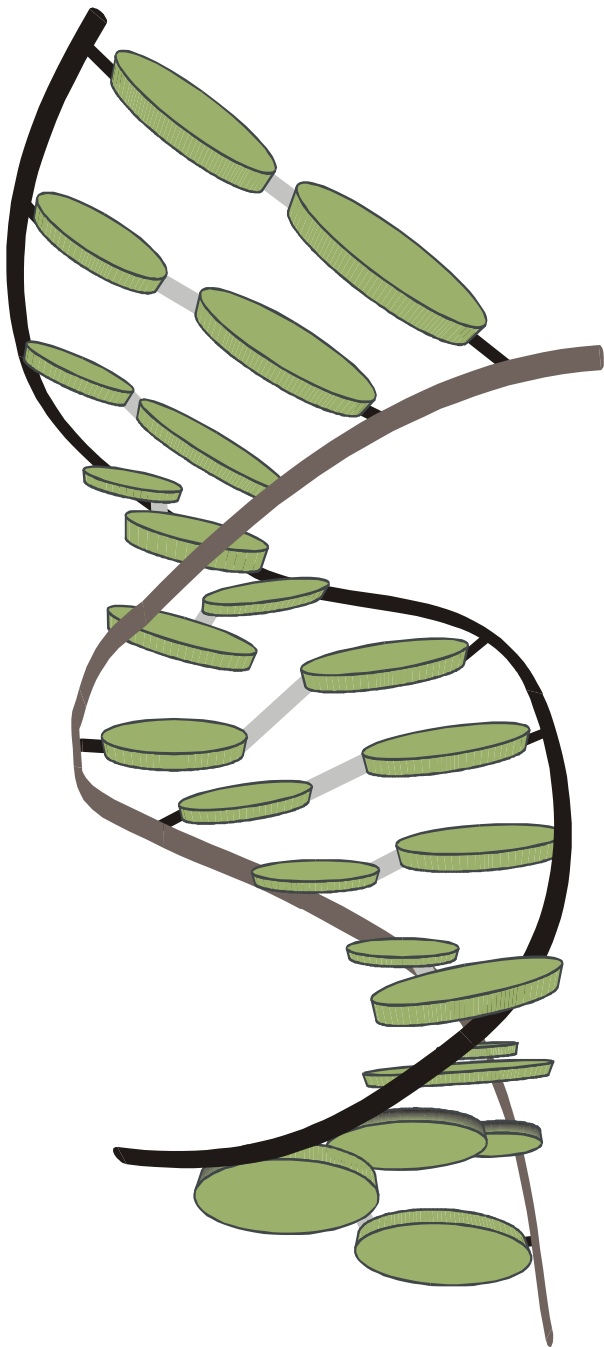
The average time of replacement of a dominant genotype in a population is the reciprocal mutation rate,  $1/v$ , and therefore independent of population size.

Fixation of mutants in neutral evolution (Motoo Kimura, 1955)



1. The chemistry of Darwinian evolution
2. RNA sequences and structures
3. Consequences of neutrality
4. Evolutionary optimization of RNA structure
5. Complexity in biology

1. **The chemistry of Darwinian evolution**
2. RNA sequences and structures
3. Consequences of neutrality
4. Evolutionary optimization of RNA structure
5. Complexity in biology

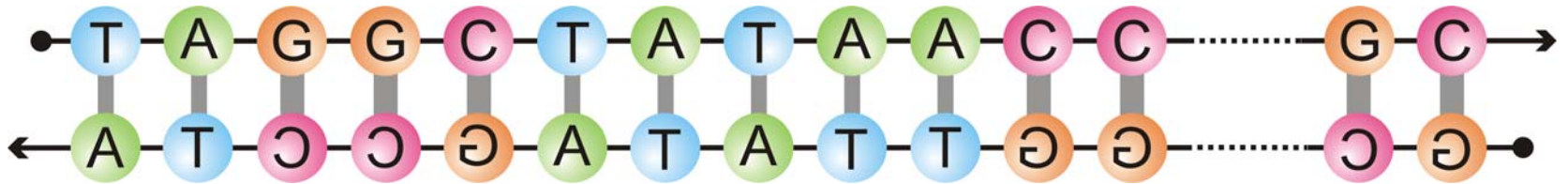


James D. Watson, 1928-, and Francis H.C. Crick, 1916-2004

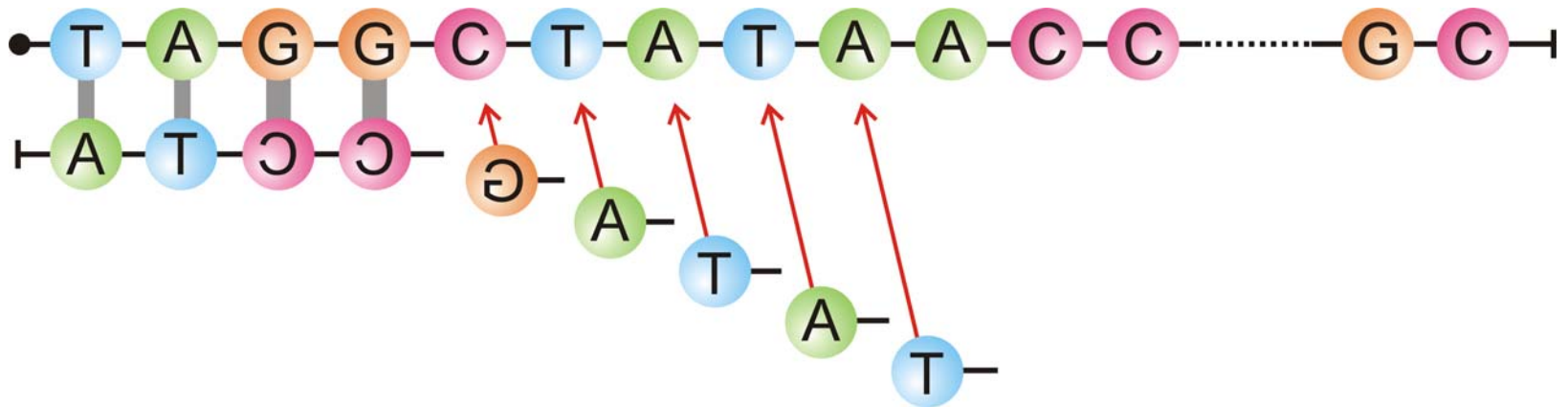
Nobel prize 1962

**1953 – 2003 fifty years double helix**

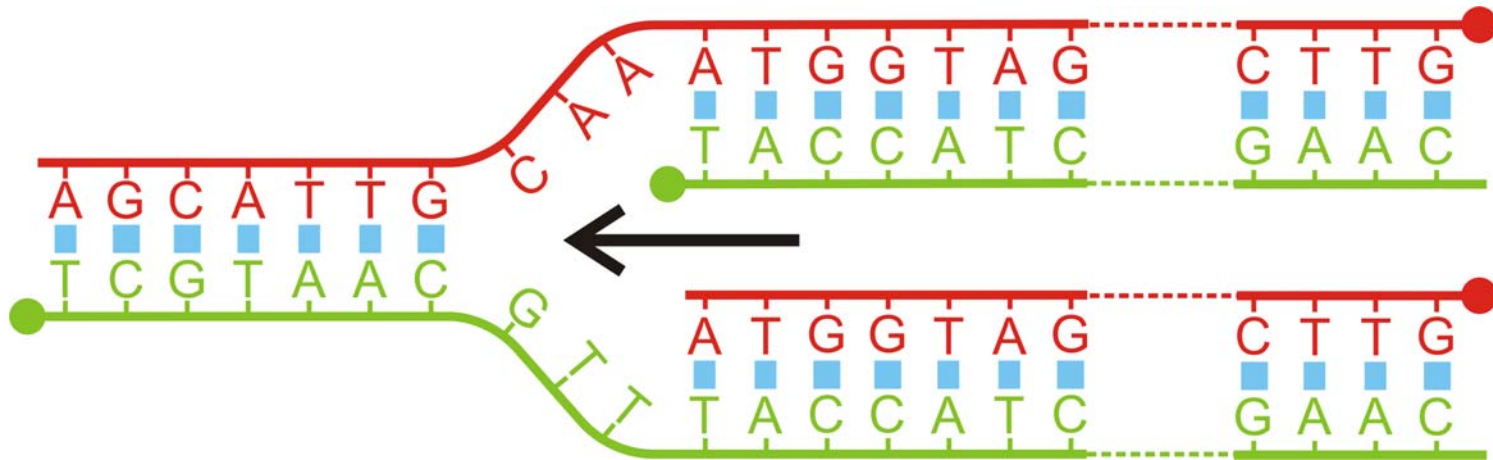
The three-dimensional structure of a short double helical stack of B-DNA



|  
o  
|

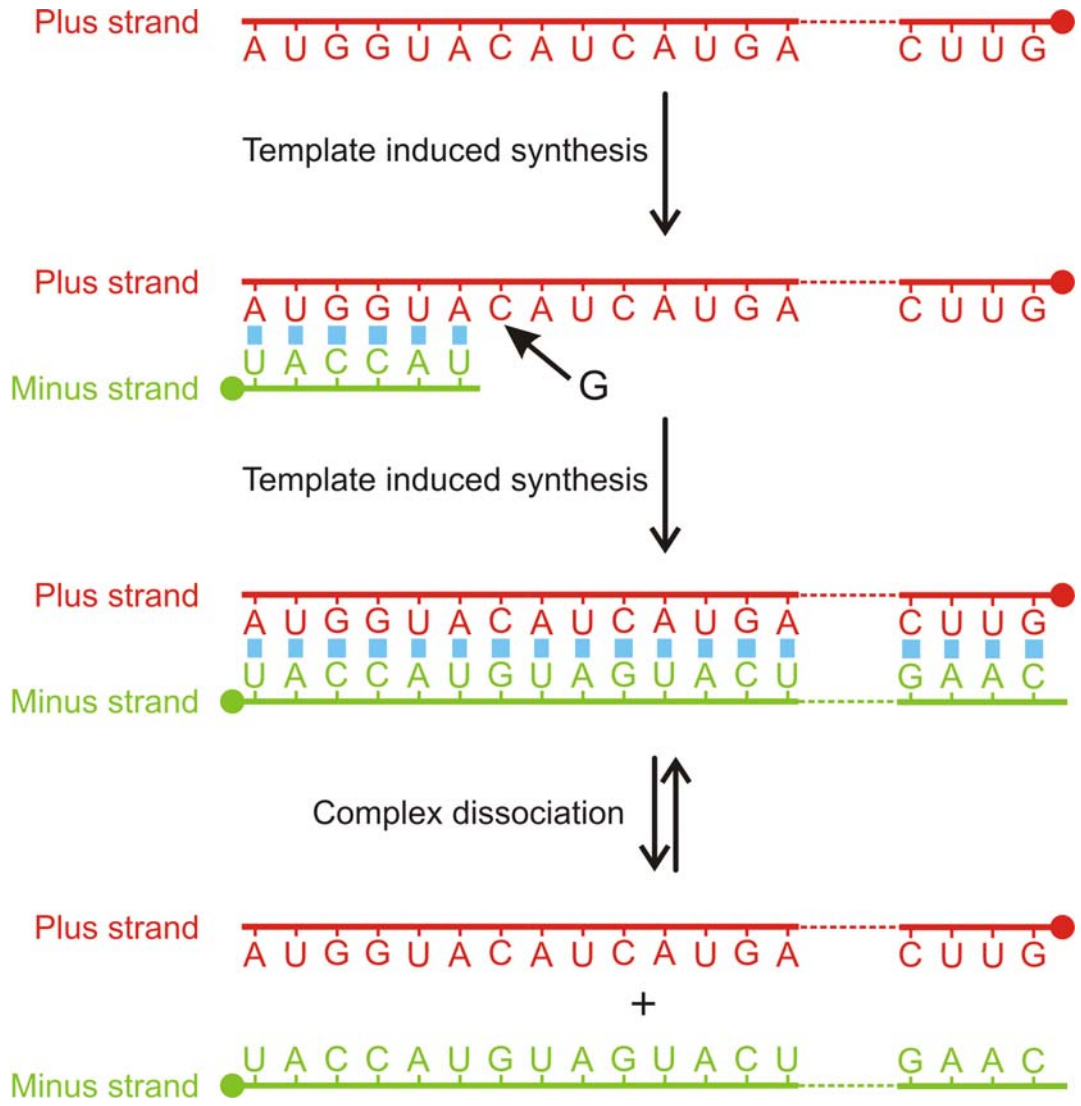


DNA structure and DNA replication



,'Replication fork' in DNA replication

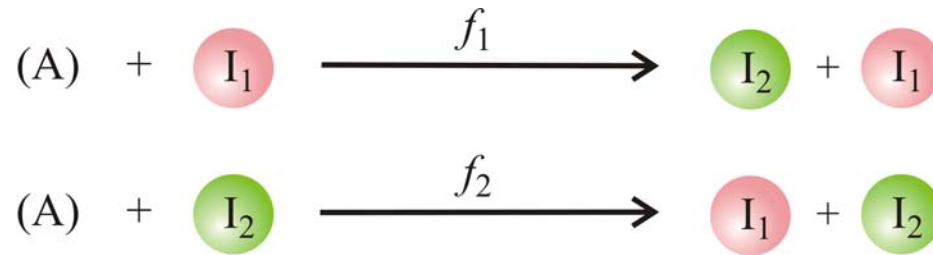
The mechanism of DNA replication is ,semi-conservative'



Complementary replication is the simplest copying mechanism of RNA.

Complementarity is determined by Watson-Crick base pairs:

**G≡C** and **A=U**



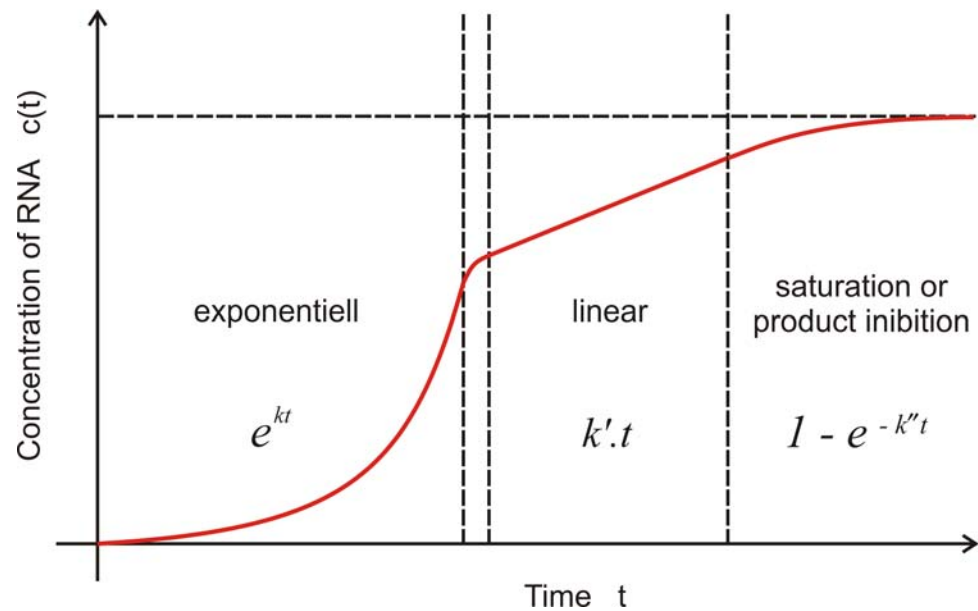
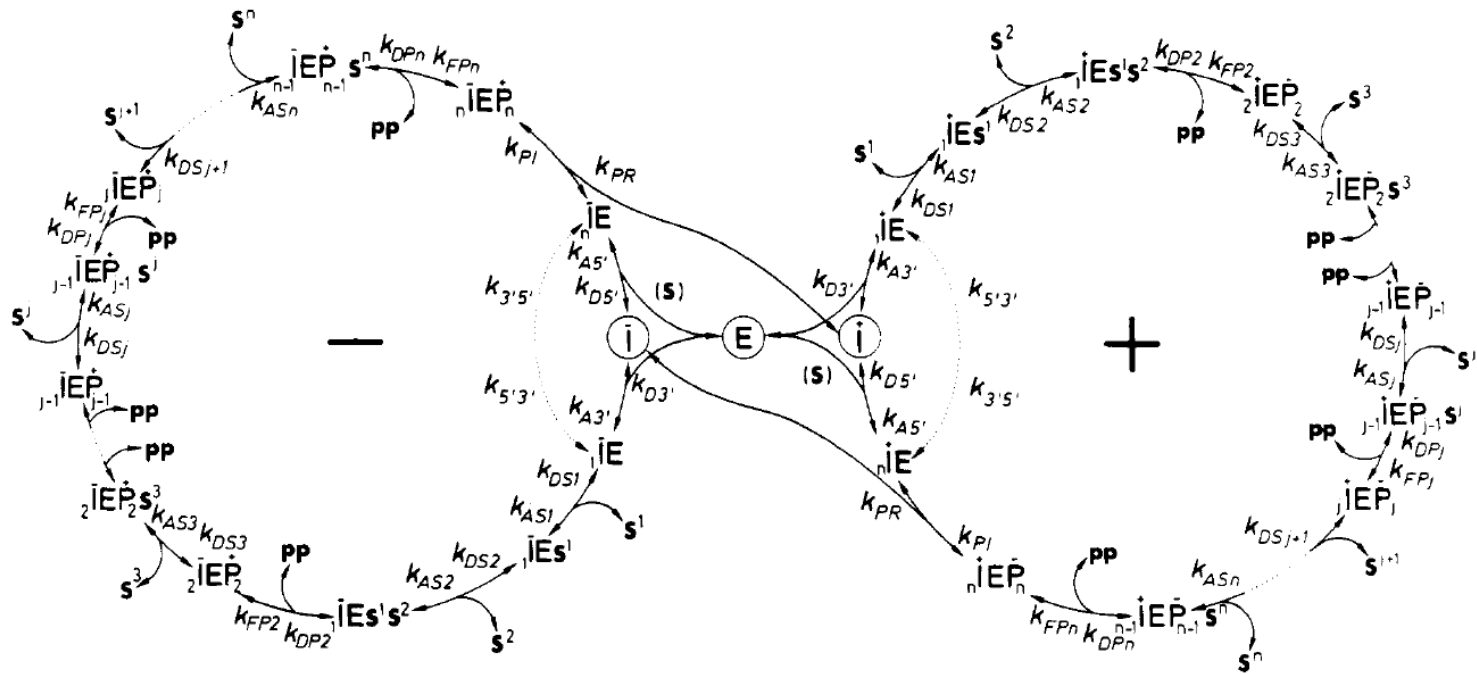
$$\frac{dx_1}{dt} = f_2 x_2 \quad \text{and} \quad \frac{dx_2}{dt} = f_1 x_1$$

$$x_1 = \sqrt{f_2} \xi_1, \quad x_2 = \sqrt{f_1} \xi_2, \quad \zeta = \xi_1 + \xi_2, \quad \eta = \xi_1 - \xi_2, \quad f = \sqrt{f_1 f_2}$$

$$\eta(t) = \eta(0) e^{-ft}$$

$$\zeta(t) = \zeta(0) e^{ft}$$

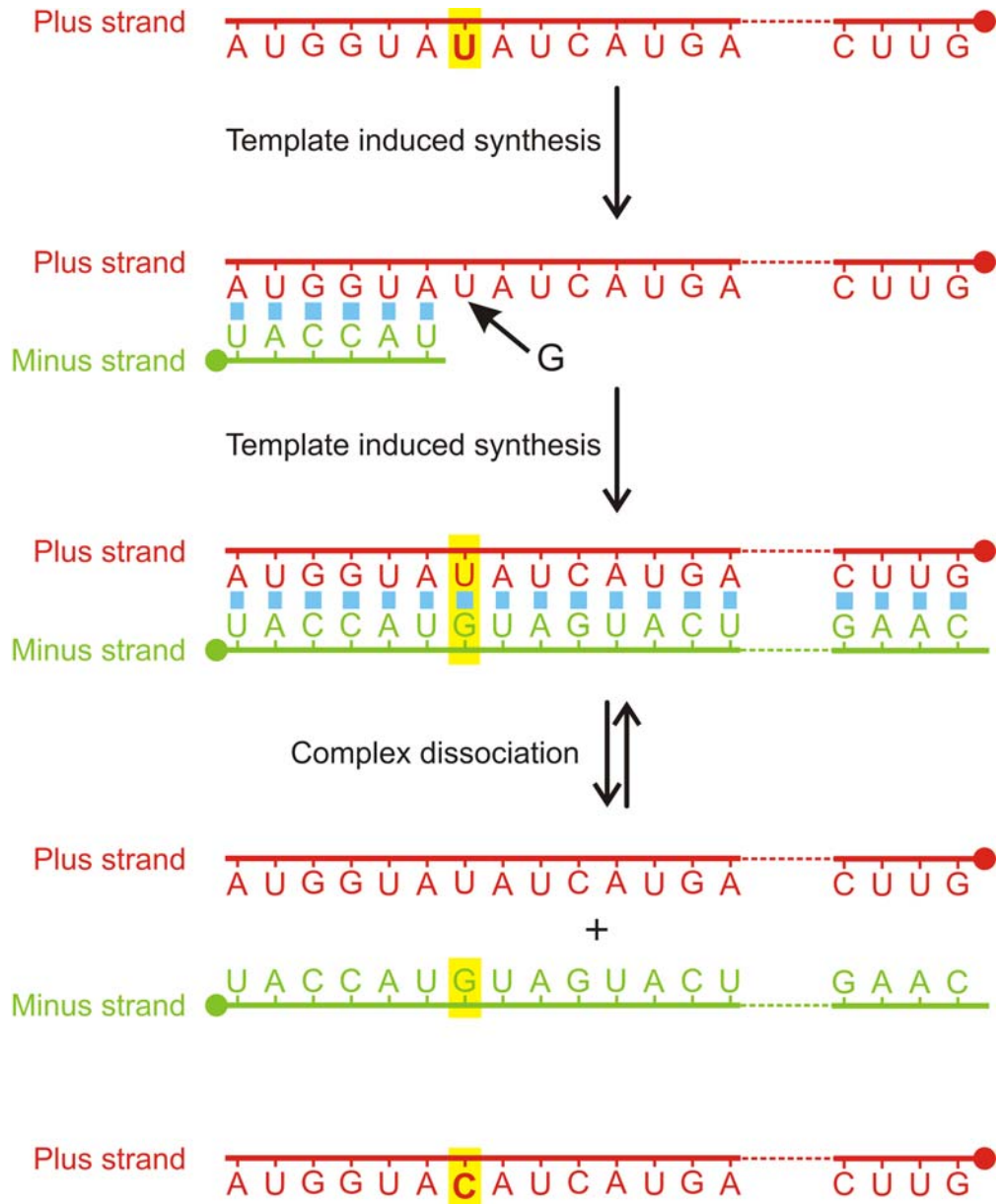
Complementary replication as the simplest molecular mechanism of reproduction



## Kinetics of RNA replication

C.K. Biebricher, M. Eigen, W.C. Gardiner, Jr.  
*Biochemistry* **22**:2544-2559, 1983





## Evolution of RNA molecules based on Q $\beta$ phage

D.R.Mills, R.L.Peterson, S.Spiegelman, *An extracellular Darwinian experiment with a self-duplicating nucleic acid molecule*. Proc.Natl.Acad.Sci.USA **58** (1967), 217-224

S.Spiegelman, *An approach to the experimental analysis of precellular evolution*. Quart.Rev.Biophys. **4** (1971), 213-253

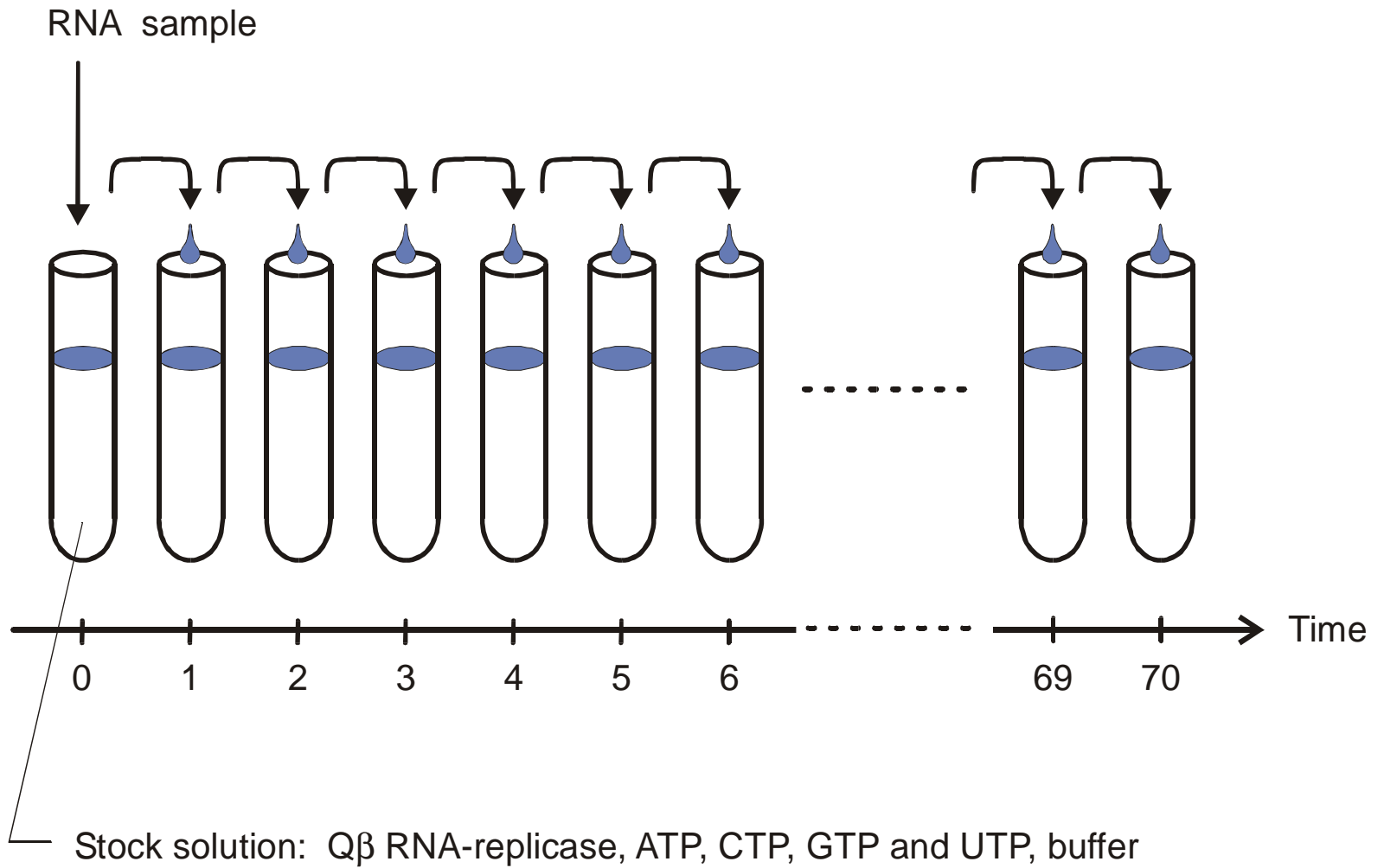
C.K.Biebricher, *Darwinian selection of self-replicating RNA molecules*. Evolutionary Biology **16** (1983), 1-52

G.Bauer, H.Otten, J.S.McCaskill, *Travelling waves of in vitro evolving RNA*. Proc.Natl.Acad.Sci.USA **86** (1989), 7937-7941

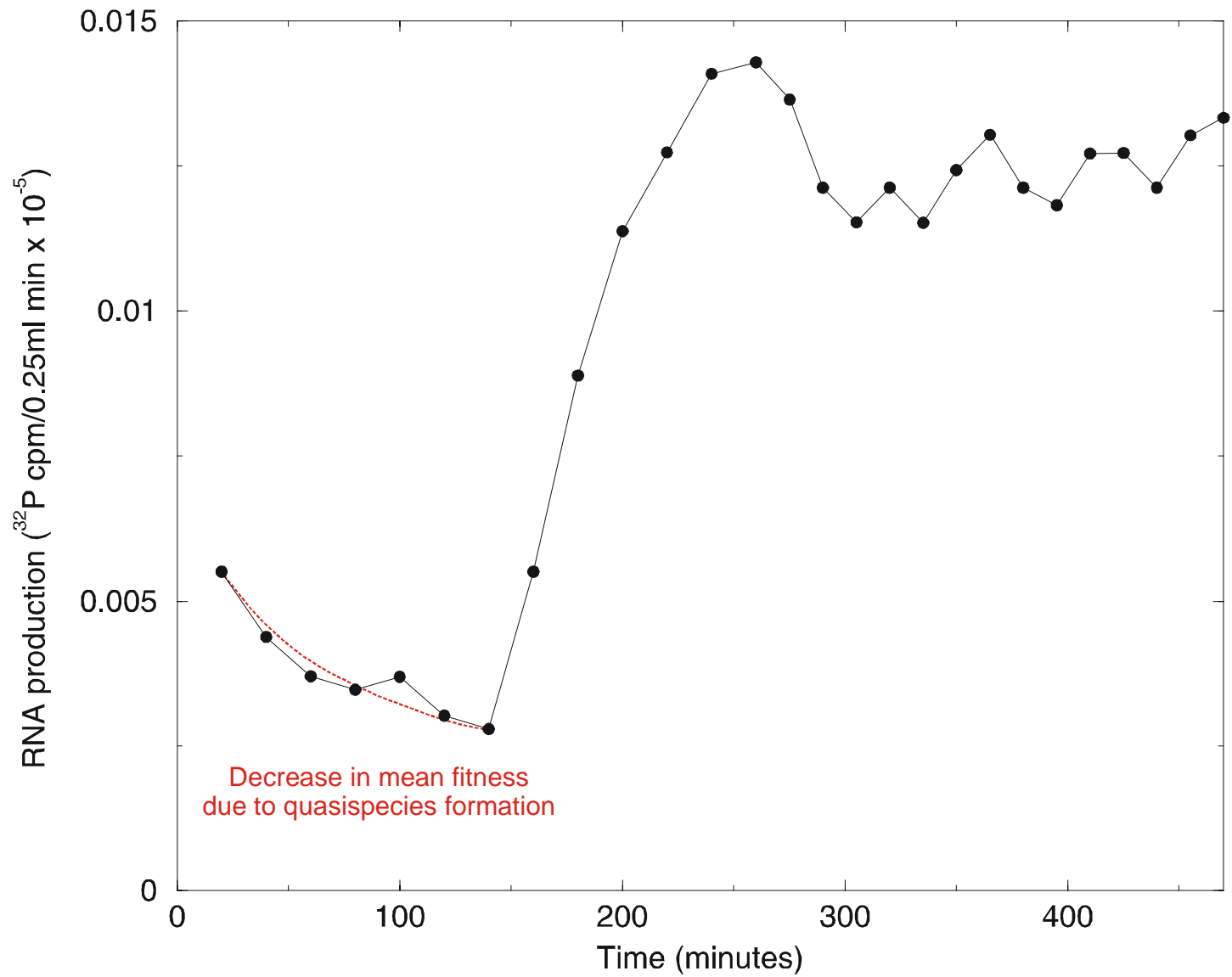
C.K.Biebricher, W.C.Gardiner, *Molecular evolution of RNA in vitro*. Biophysical Chemistry **66** (1997), 179-192

G.Strunk, T.Ederhof, *Machines for automated evolution experiments in vitro based on the serial transfer concept*. Biophysical Chemistry **66** (1997), 193-202

F.Öhlenschläger, M.Eigen, *30 years later – A new approach to Sol Spiegelman's and Leslie Orgel's in vitro evolutionary studies*. Orig.Life Evol.Biosph. **27** (1997), 437-457



Anwendung der seriellen Überimpfungstechnik auf RNA-Evolution in Reagenzglas

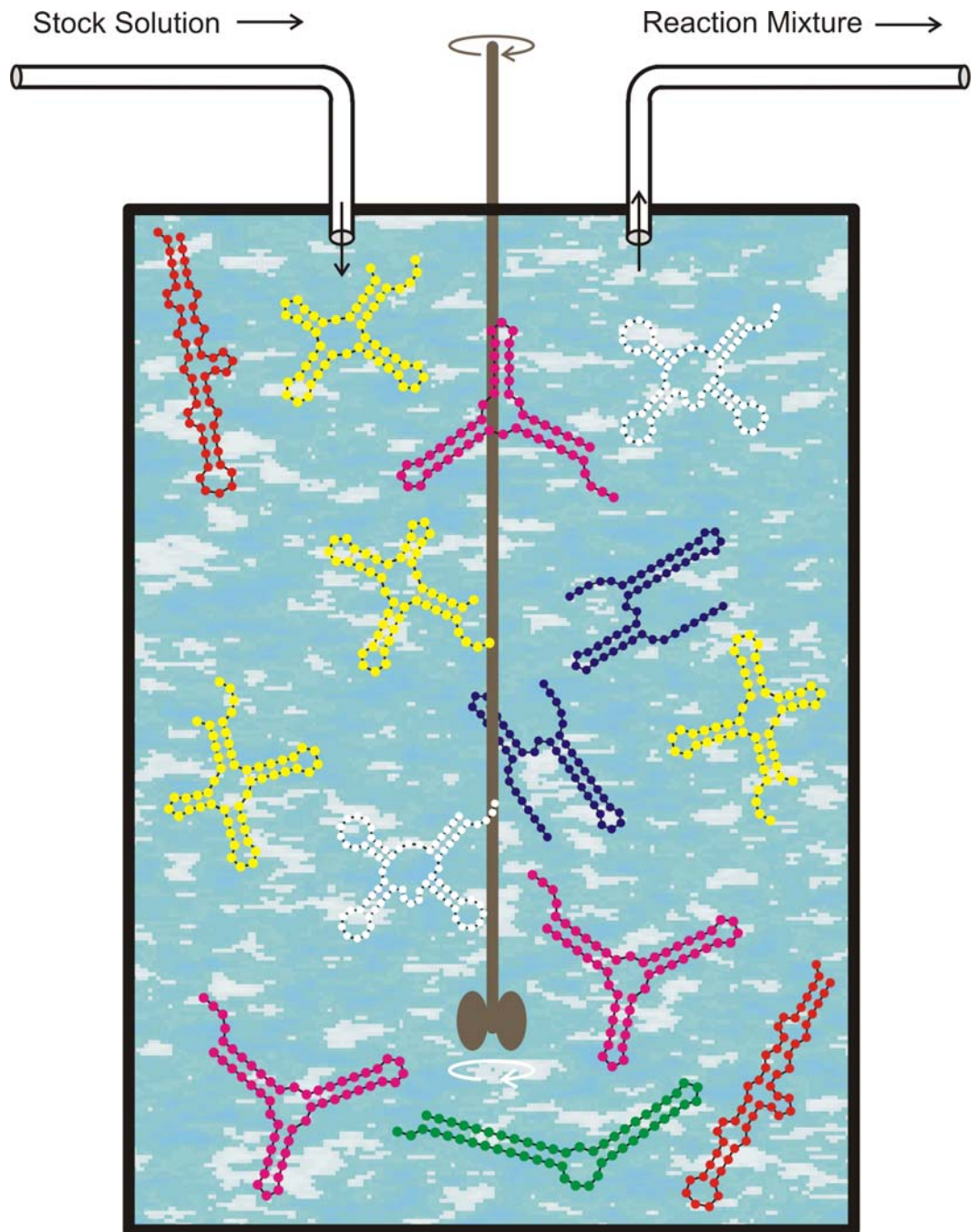


The increase in RNA production rate during a serial transfer experiment

**Stock solution:**

activated monomers, **ATP, CTP, GTP, UTP (TTP)**;  
a replicase, an enzyme that performs complementary replication;  
buffer solution

The flowreactor is a device for **studies** of evolution *in vitro* and *in silico*.



## Evolutionary design of RNA molecules

A.D. Ellington, J.W. Szostak, *In vitro selection of RNA molecules that bind specific ligands*. Nature **346** (1990), 818-822

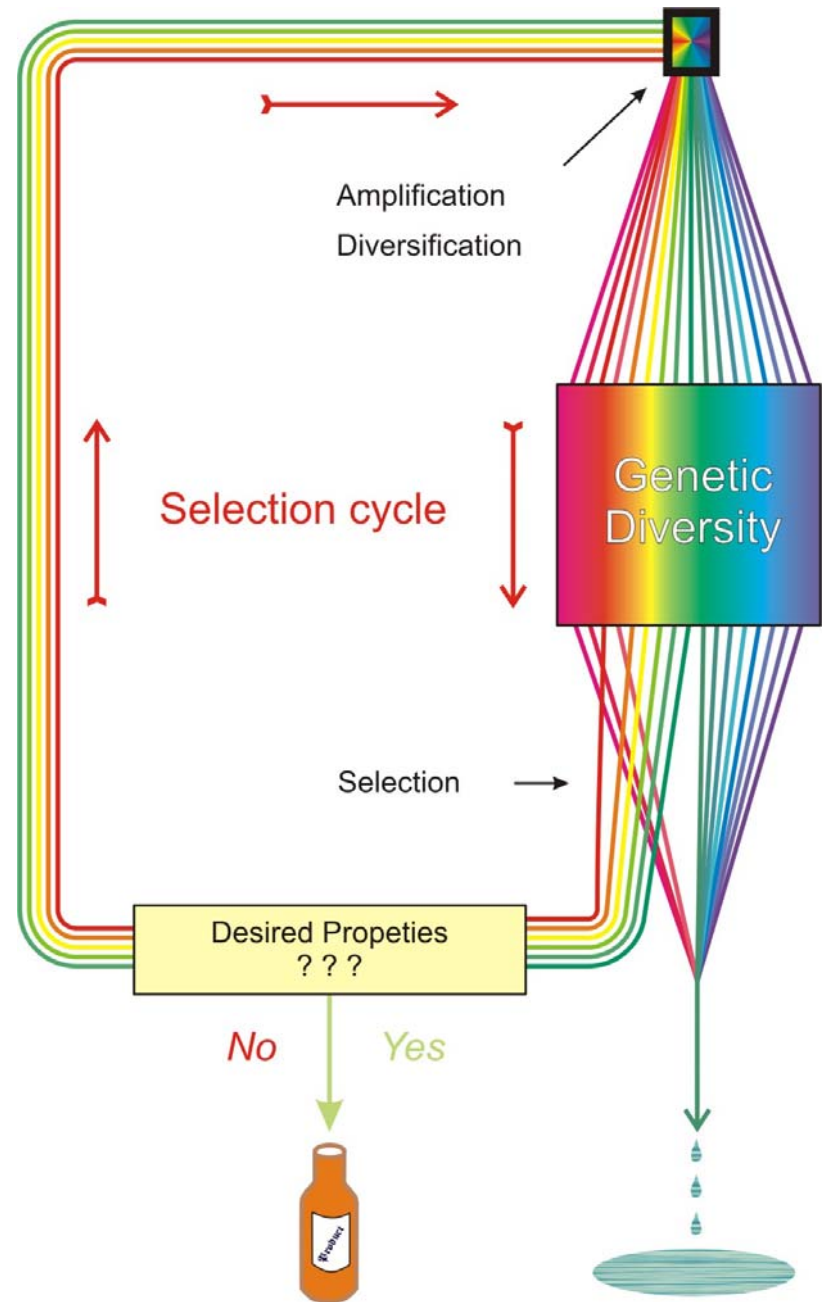
C. Tuerk, L. Gold, *SELEX - Systematic evolution of ligands by exponential enrichment: RNA ligands to bacteriophage T4 DNA polymerase*. Science **249** (1990), 505-510

D.P. Bartel, J.W. Szostak, *Isolation of new ribozymes from a large pool of random sequences*. Science **261** (1993), 1411-1418

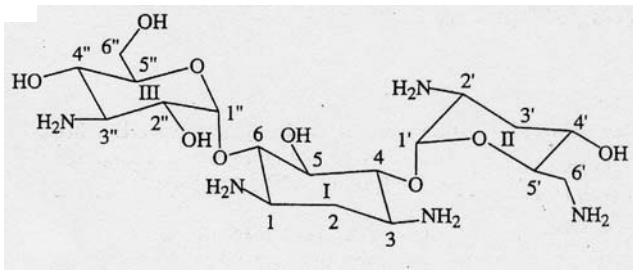
R.D. Jenison, S.C. Gill, A. Pardi, B. Poliski, *High-resolution molecular discrimination by RNA*. Science **263** (1994), 1425-1429

Y. Wang, R.R. Rando, *Specific binding of aminoglycoside antibiotics to RNA*. Chemistry & Biology **2** (1995), 281-290

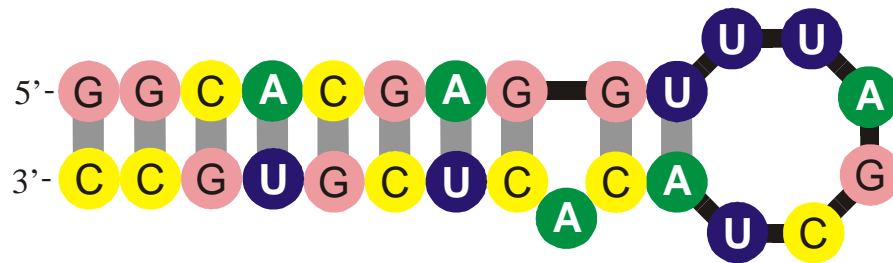
L. Jiang, A. K. Suri, R. Fiala, D. J. Patel, *Saccharide-RNA recognition in an aminoglycoside antibiotic-RNA aptamer complex*. Chemistry & Biology **4** (1997), 35-50



An example of 'artificial selection' with RNA molecules or 'breeding' of biomolecules



tobramycin

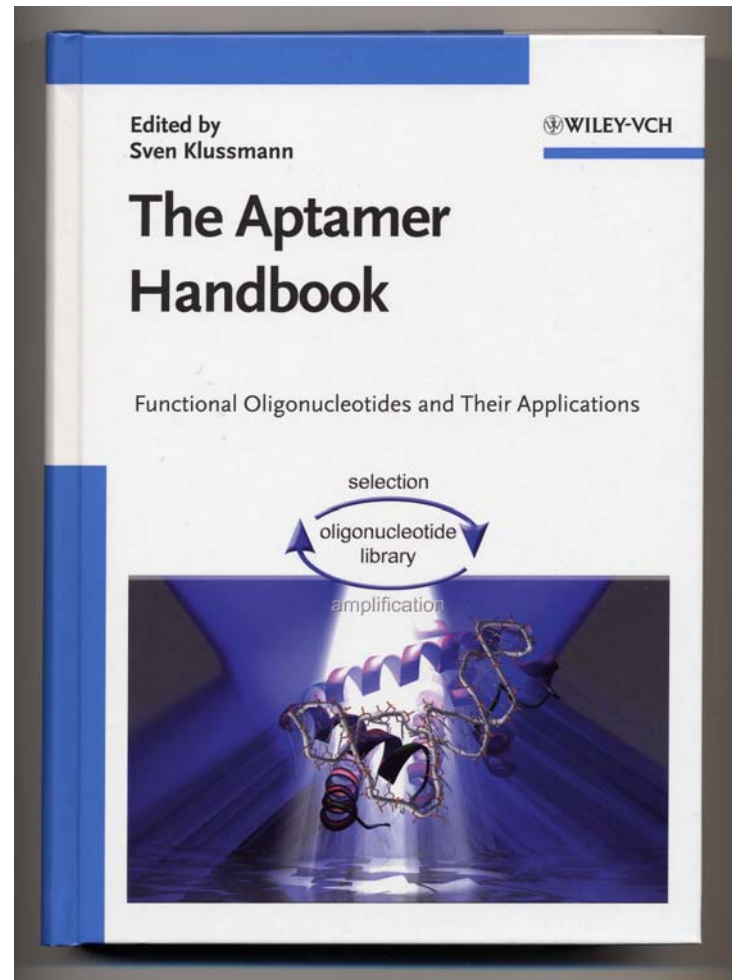
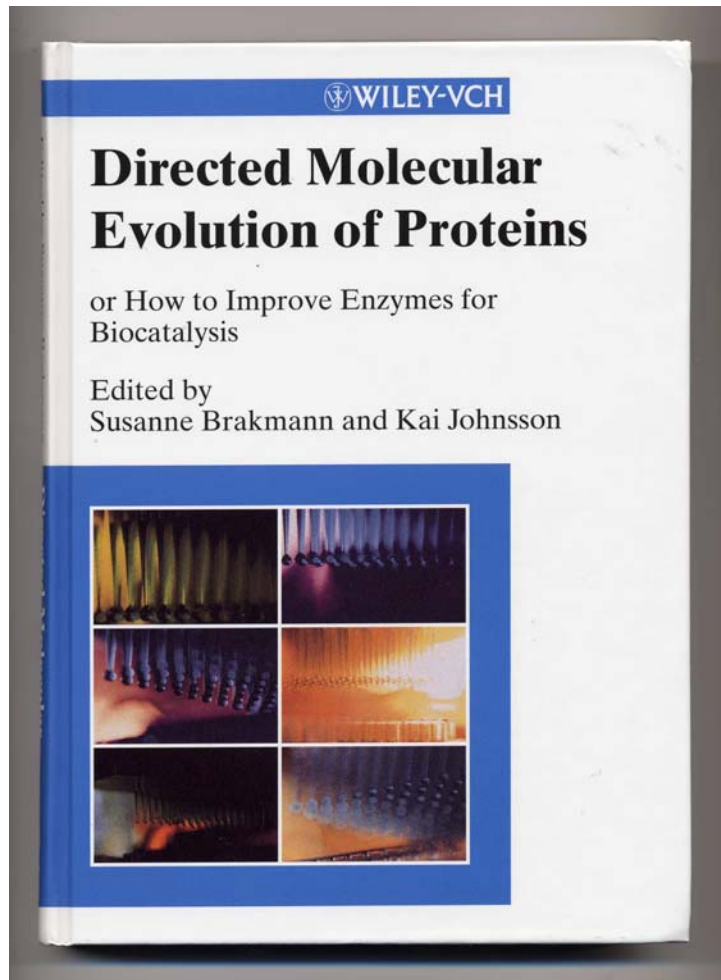


RNA aptamer

Formation of secondary structure of the tobramycin binding RNA aptamer with  $K_D = 9 \text{ nM}$

L. Jiang, A. K. Suri, R. Fiala, D. J. Patel, *Saccharide-RNA recognition in an aminoglycoside antibiotic-RNA aptamer complex*. *Chemistry & Biology* 4:35-50 (1997)



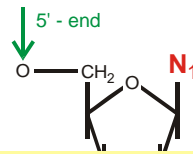


Application of molecular evolution to problems in biotechnology

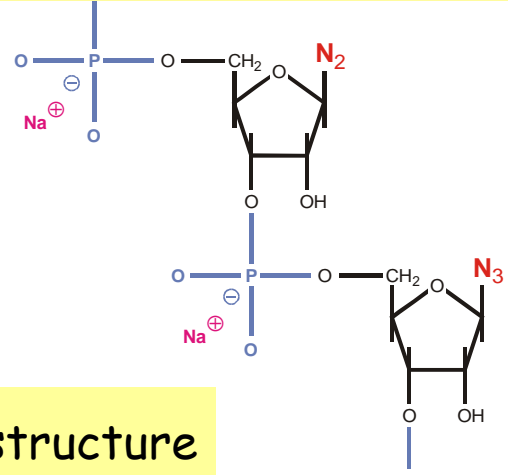
## Results from molecular evolution in laboratory experiments:

- Evolutionary optimization does not require cells and occurs in molecular systems too.
- *In vitro* evolution allows for production of molecules for predefined purposes and gave rise to a branch of biotechnology.

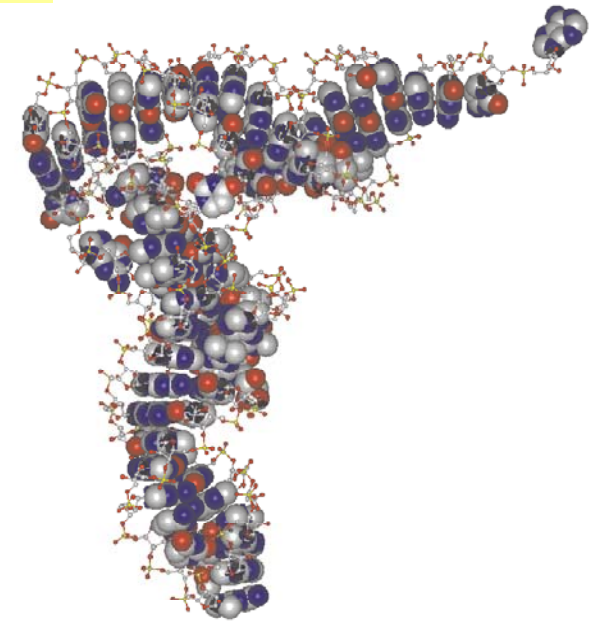
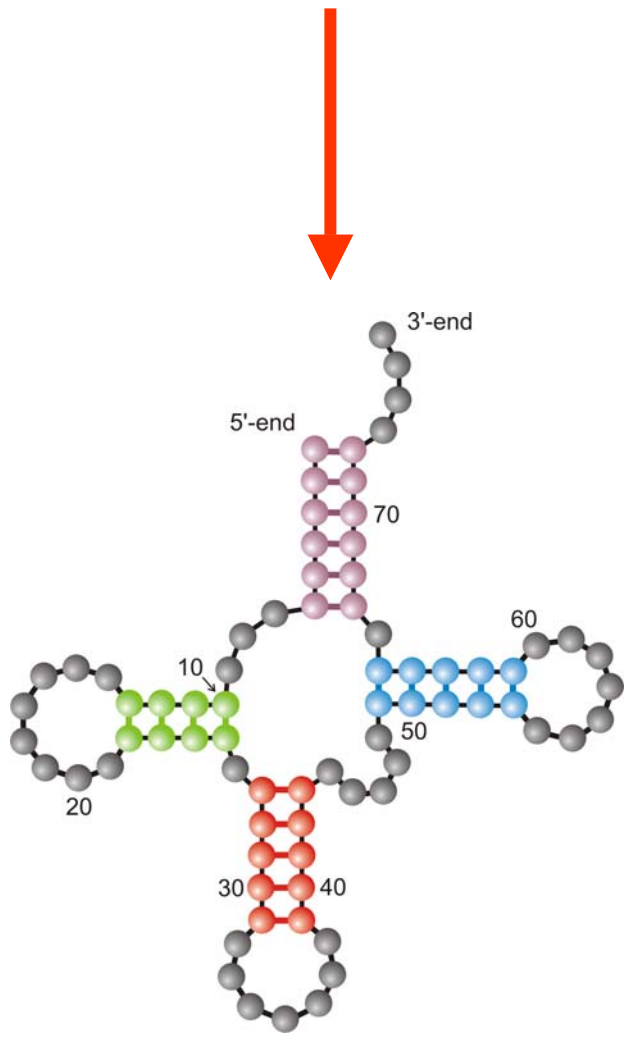
1. The chemistry of Darwinian evolution
2. **RNA sequences and structures**
3. Consequences of neutrality
4. Evolutionary optimization of RNA structure
5. Complexity in biology

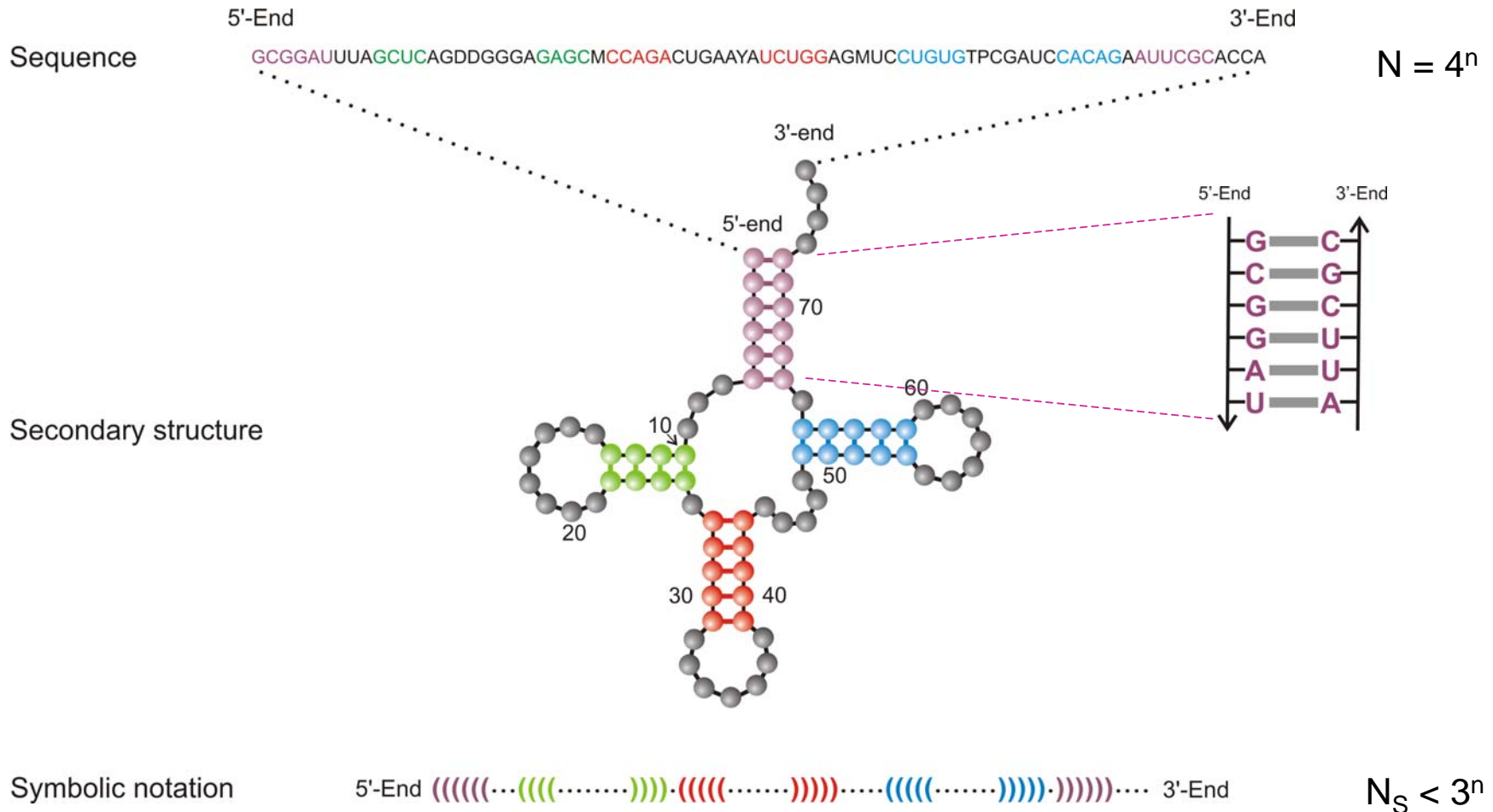


5'-end **GCGGAUUUAGCUC**AGUUGGGAGAG**CGCCAGACUGAAGAUCUGG**AGGUC**CUGUGUUCGAUCCACAGAAUUCGCACCA** 3'-end



Definition of RNA structure

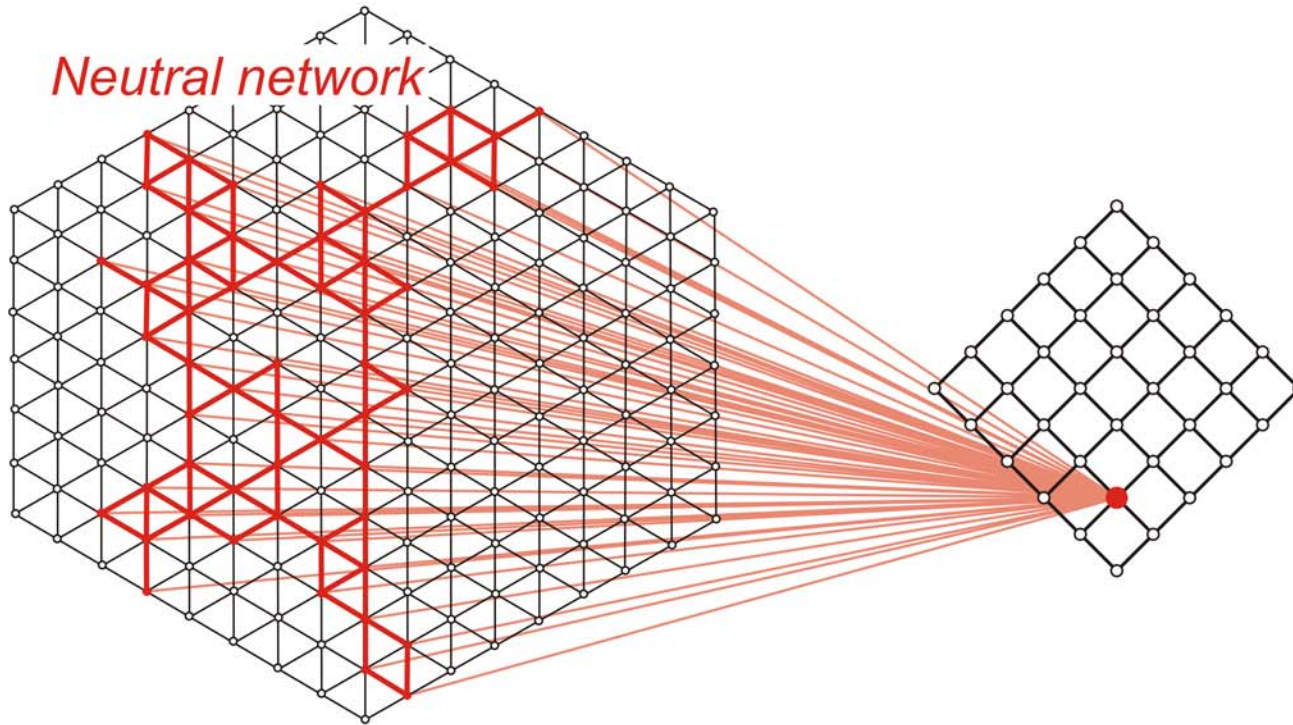




Criterion: Minimum free energy (mfe)

Rules:  $\_ (\_ ) \_ \in \{AU, CG, GC, GU, UA, UG\}$

A symbolic notation of RNA secondary structure that is equivalent to the conventional graphs



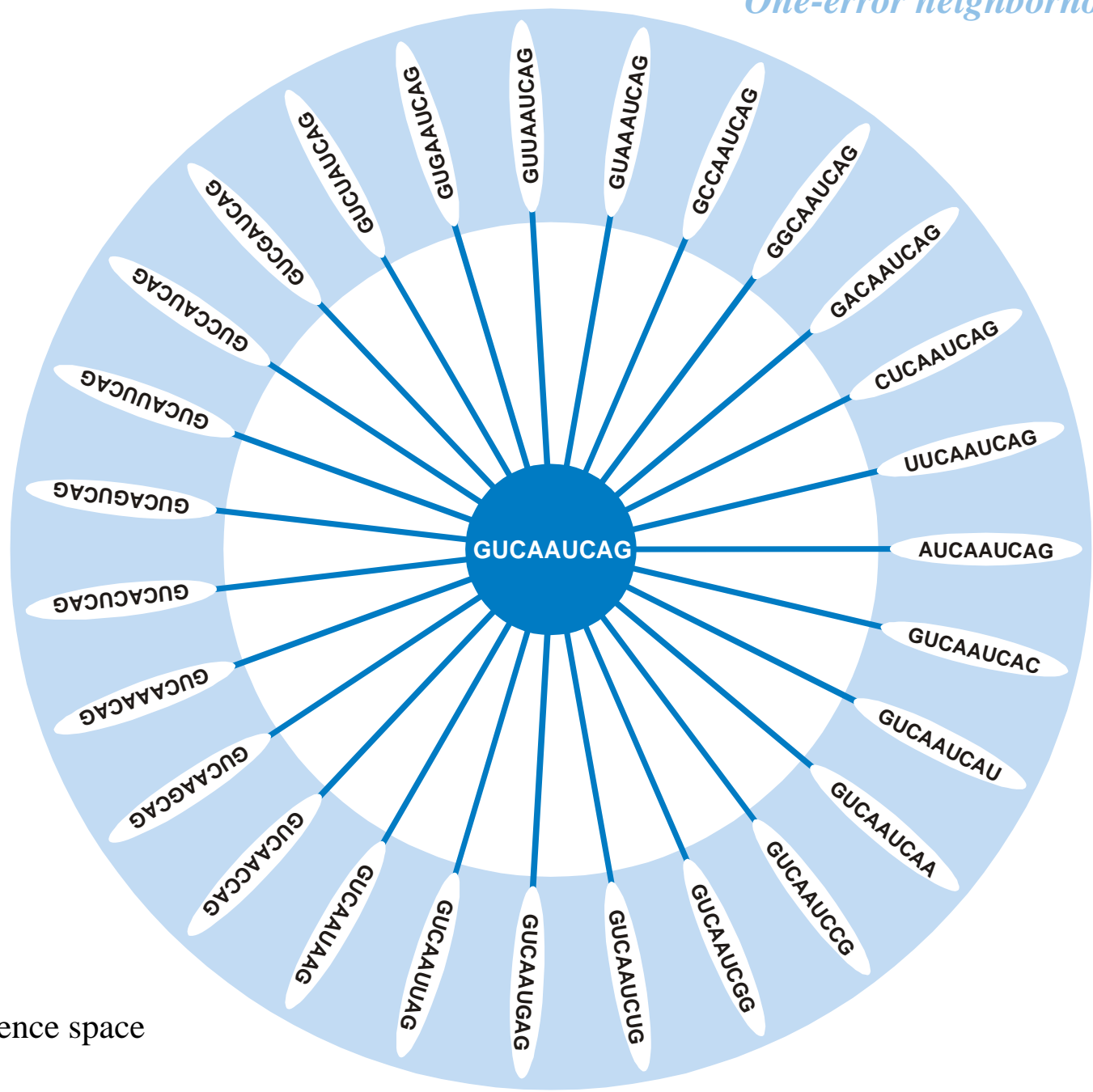
Sequence space

Structure space

many genotypes

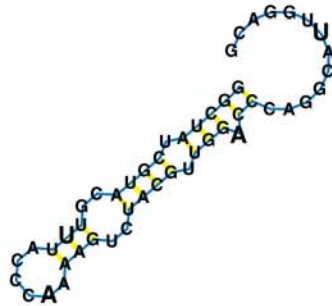
⇒

one phenotype



The surrounding of **GUCAAUCAG** in sequence space

GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG

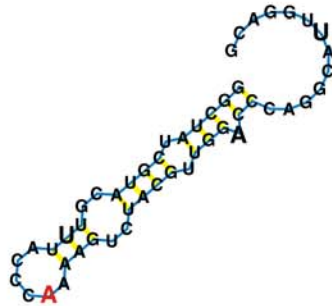


One error neighborhood – Surrounding of an RNA molecule of chain length  $n=50$  in sequence and shape space

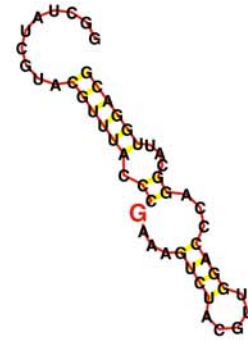


GGCUAUCGUACGUUUACCCGAAAGUCUACGUUGGACCCAGGCAUUGGACG

GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG

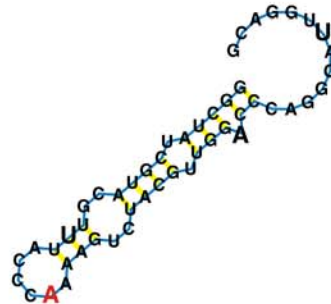


One error neighborhood – Surrounding of an RNA molecule of chain length  $n=50$  in sequence and shape space

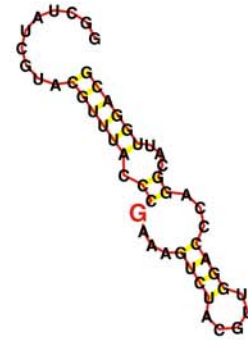


GGCUAUCGUACGUUUACCCGAAAGUCUACGUUGGACCCAGGCAUUGGACG

GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG



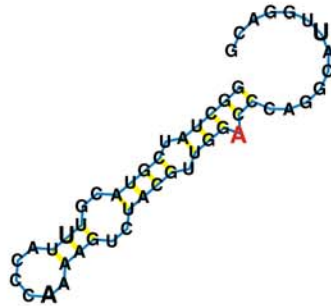
One error neighborhood – Surrounding of an RNA molecule of chain length  $n=50$  in sequence and shape space



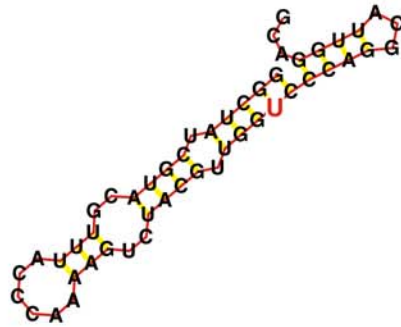
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGG**U**CCAGGCAUUGGACG

GGCUAUCGUACGUUUACCC**G**AAAGUCUACGUUGGACCCAGGCAUUGGACG

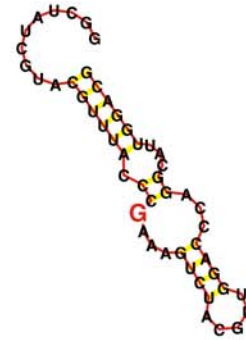
GGCUAUCGUACGU**U**UACCCAAAAGUCUACGUUGG**A**CCAGGCAUUGGACG



One error neighborhood – Surrounding of an RNA molecule of chain length  $n=50$  in sequence and shape space

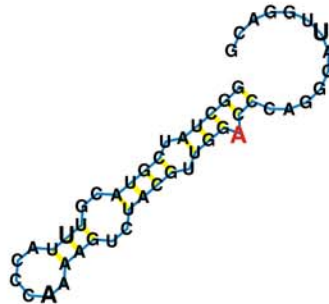


GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGG**U**CCAGGCAUUGGACG

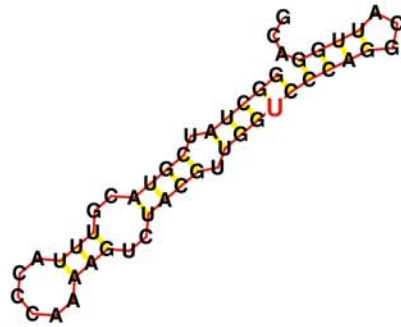


GGCUAUCGUACGUUUACCC**G**AAAGUCUACGUUGGACCCAGGCAUUGGACG

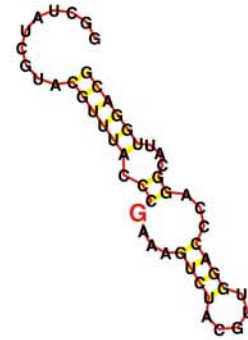
GGCUAUCGUACGU**U**UACCCAAAAGUCUACGUUGG**A**CCAGGCAUUGGACG



One error neighborhood – Surrounding of an RNA molecule of chain length  $n=50$  in sequence and shape space



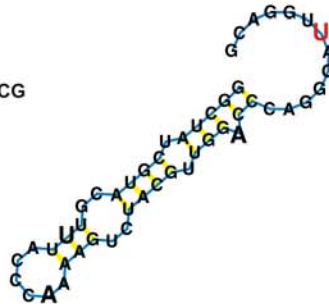
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGG**U**CCAGGCAUUGGACG



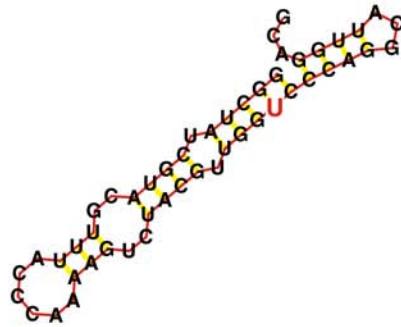
GGCUAUCGUACGUUUACCC**G**AAAGUCUACGUUGGACCCAGGCAUUGGACG

GGCUAUCGUACGU**U**UACCCAAAAGUCUACGUUGGACCCAGGCA**U**UGGACG

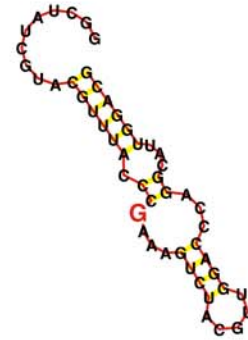
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCA**C**UGGACG



One error neighborhood – Surrounding of an RNA molecule of chain length  $n=50$  in sequence and shape space



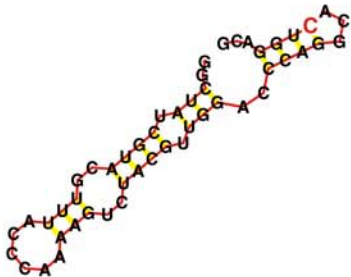
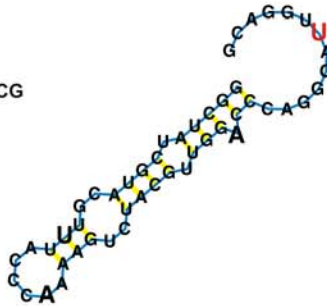
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGG**U**CCAGGCAUUGGACG



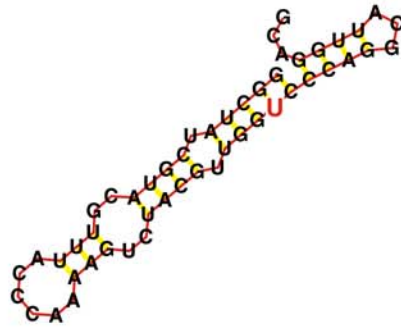
GGCUAUCGUACGUUUACCC**G**AAAGUCUACGUUGGACCCAGGCAUUGGACG

GGCUAUCGUACGU**U**UACCCAAAAGUCUACGUUGGACCCAGGCA**U**UGGACG

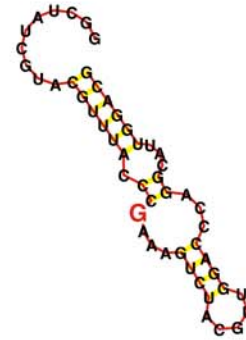
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCA**C**UGGACG



One error neighborhood – Surrounding of an RNA molecule of chain length  $n=50$  in sequence and shape space



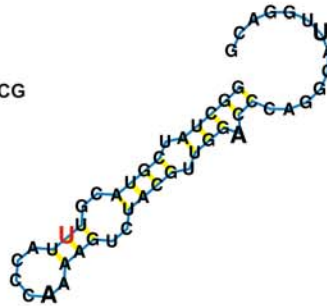
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGG**U**CCAGGCAUUGGACG



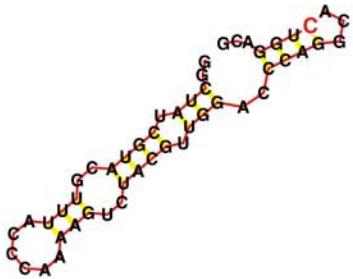
GGCUAUCGUACGUUUACCC**G**AAAGUCUACGUUGGACCCAGGCAUUGGACG

GGCUAUCGUACGU**U**UACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG

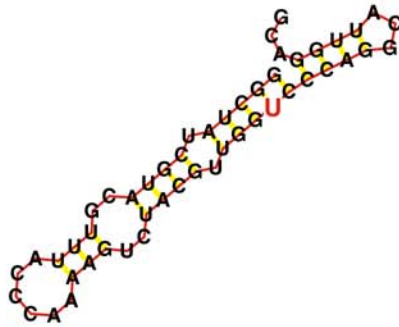
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCA**C**UGGACG



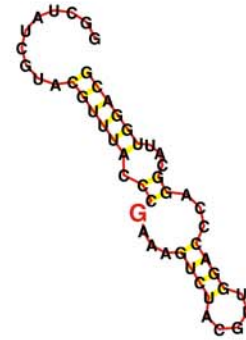
GGCUAUCGUACGU**G**UACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG



One error neighborhood – Surrounding of an RNA molecule of chain length  $n=50$  in sequence and shape space



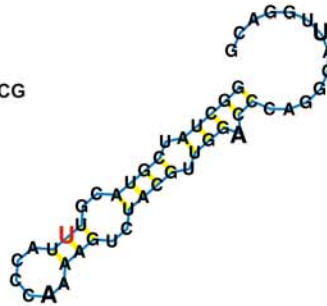
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGG**U**CCAGGCAUUGGACG



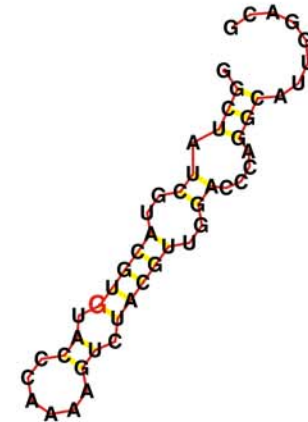
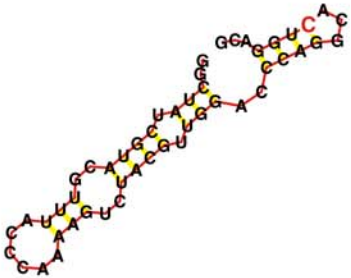
GGCUAUCGUACGUUUACCC**G**AAAGUCUACGUUGGACCCAGGCAUUGGACG

GGCUAUCGUACGU**U**UACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG

GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCA**C**UGGACG



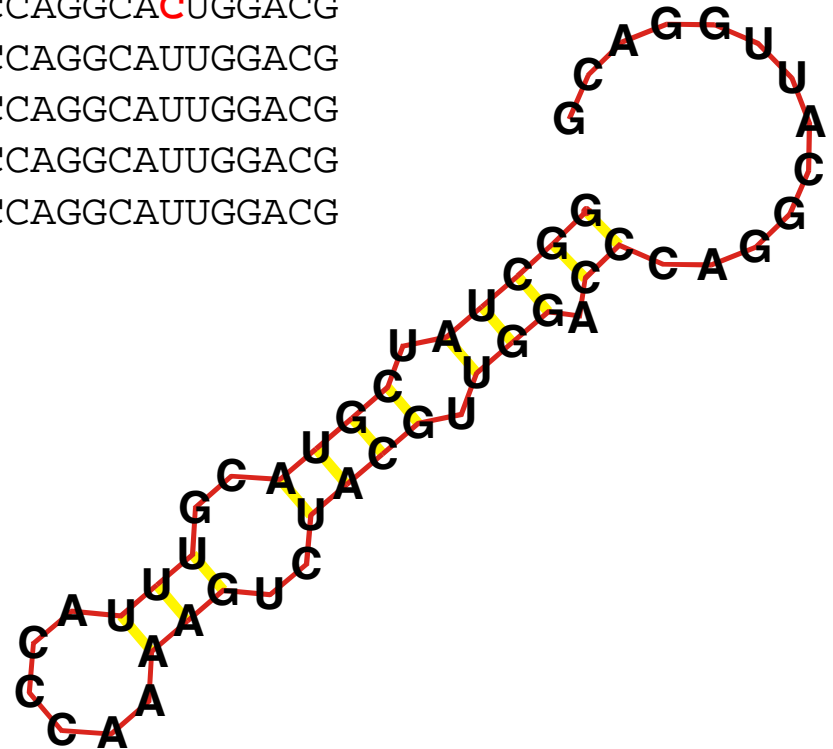
GGCUAUCGUACGU**G**UACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG



One error neighborhood – Surrounding of an RNA molecule of chain length  $n=50$  in sequence and shape space



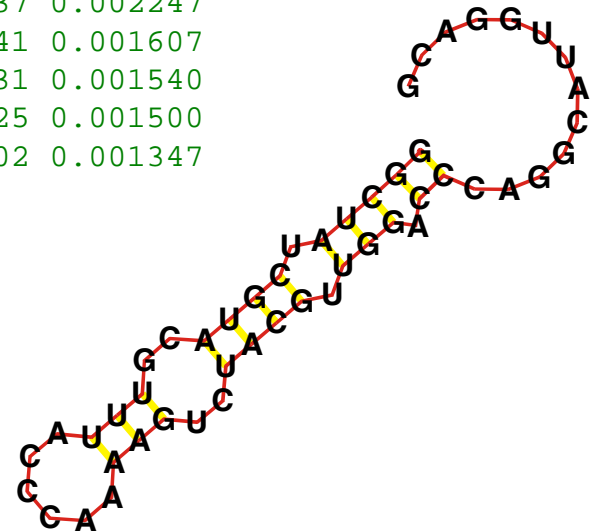
GGCUAUCGUAU**U**GUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG  
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUA**A**GACG  
GGCUAUCGUACGUUUAC**U**CAAAGUCUACGUUGGACCCAGGCAUUGGACG  
GGCUAUCGUACG**C**UUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG  
GGC**C**AUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG  
**GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG**  
GGCUAUCGUACGU**G**UACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG  
GGCUA**A**CGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG  
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCC**U**GGCAUUGGACG  
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCA**C**UGGACG  
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGG**U**CCCAGGCAUUGGACG  
GGCUA**G**CGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG  
GGCUAUCGUACGUUUACCC**G**AAAGUCUACGUUGGACCCAGGCAUUGGACG  
GGCUAUCGUACGUUUACCCAAAAG**C**CUACGUUGGACCCAGGCAUUGGACG



One error neighborhood – Surrounding of an RNA molecule of chain length  $n=50$  in sequence and shape space

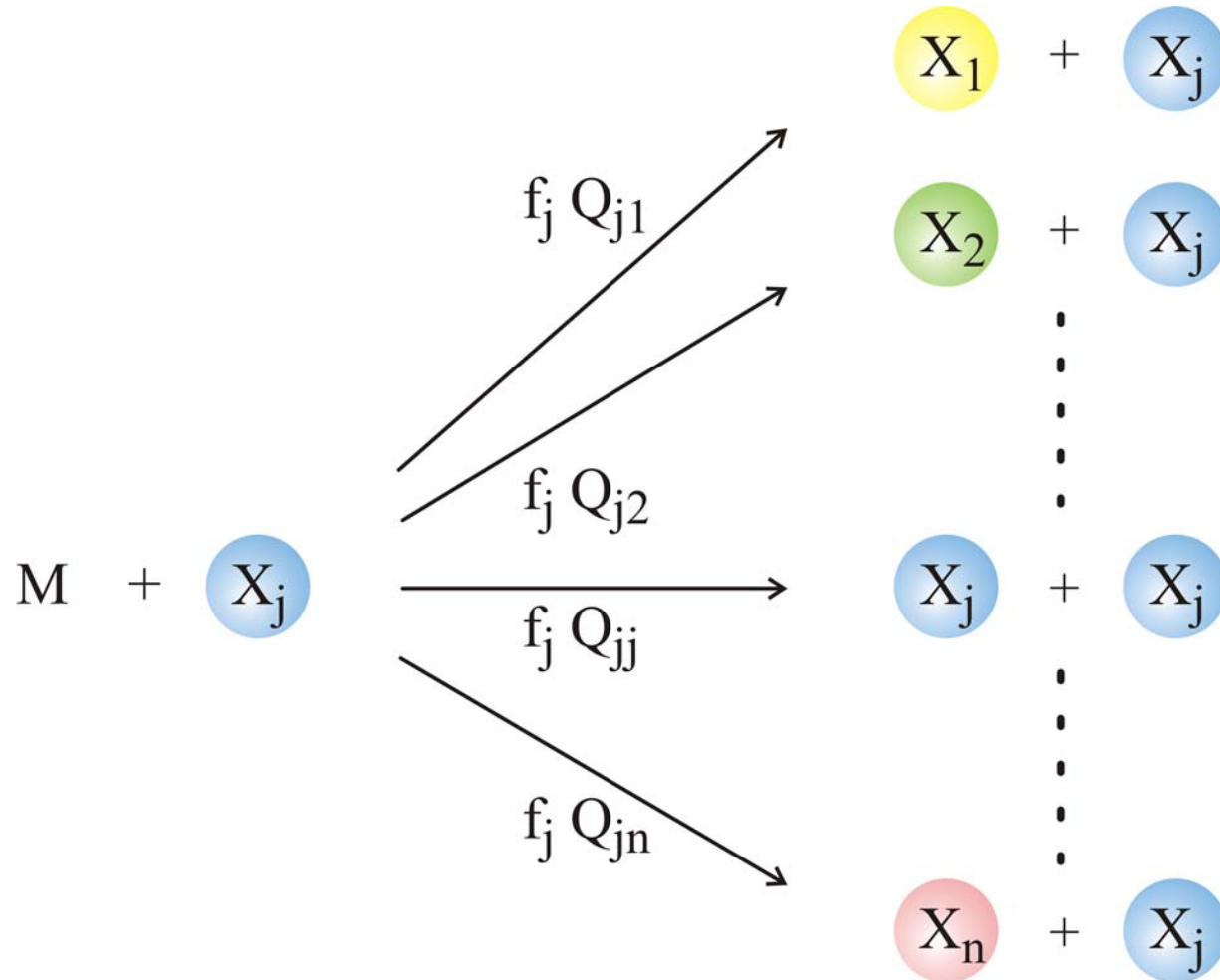
	Number	Mean Value	Variance	Std.Dev.
Total Hamming Distance:	150000	11.647973	23.140715	4.810480
Nonzero Hamming Distance:	99875	16.949991	30.757651	5.545958
Degree of Neutrality:	50125	<b>0.334167</b>	0.006961	<b>0.083434</b>
Number of Structures:	<b>1000</b>	<b>52.31</b>	85.30	<b>9.24</b>

1	(((((.((((..(((.....))))..))))..)))..)).....	50125	0.334167
2	..(((.((((..(((.....))))..))))..))).....	2856	0.019040
3	(((((.((((..(((.....))))..))))..))).....	2799	0.018660
4	(((((.((((..(((.....))))..))))..))).....	2417	0.016113
5	(((((.((((..(((.....))))..))))..))).....	2265	0.015100
6	(((((.((((..(((.....))))..))))..))).....	2233	0.014887
7	(((((..(((..(((.....))))..))))..))).....	1442	0.009613
8	(((((.((((..(((.....))))..))))..))).....	1081	0.007207
9	(((((..(((..(((.....))))..))))..))).....	1025	0.006833
10	(((((.((((..(((.....))))..))))..))).....	1003	0.006687
11	..(((.((((..(((.....))))..))))..))).....	963	0.006420
12	(((((.((((..(((.....))))..))))..))).....	860	0.005733
13	(((((.((((..(((.....))))..))))..))).....	800	0.005333
14	(((((.((((..(((.....))))..))))..))).....	548	0.003653
15	(((((.((((.....))))..))))..))).....	362	0.002413
16	(((((.((((..(((.....))))..))))..))).....	337	0.002247
17	(((((.((((..(((.....))))..))))..))).....	241	0.001607
18	(((((.((((..(((.....))))..))))..))).....	231	0.001540
19	(((((..(((..(((.....))))..))))..))).....	225	0.001500
20	(((((..(((..(((.....))))..))))..))).....	202	0.001347

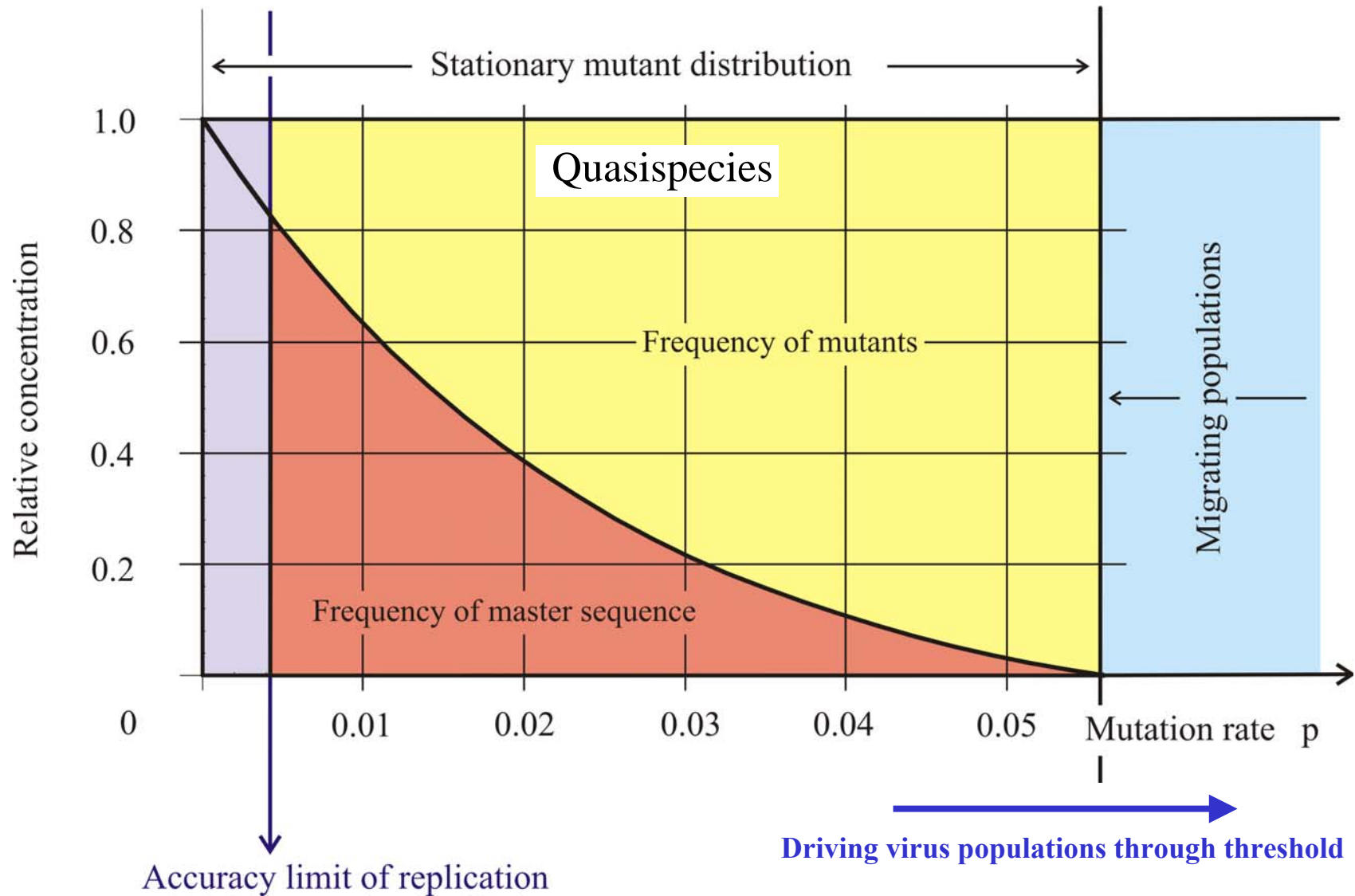


Shadow – Surrounding of an RNA structure in shape space:  
**AUGC** alphabet, chain length n=50

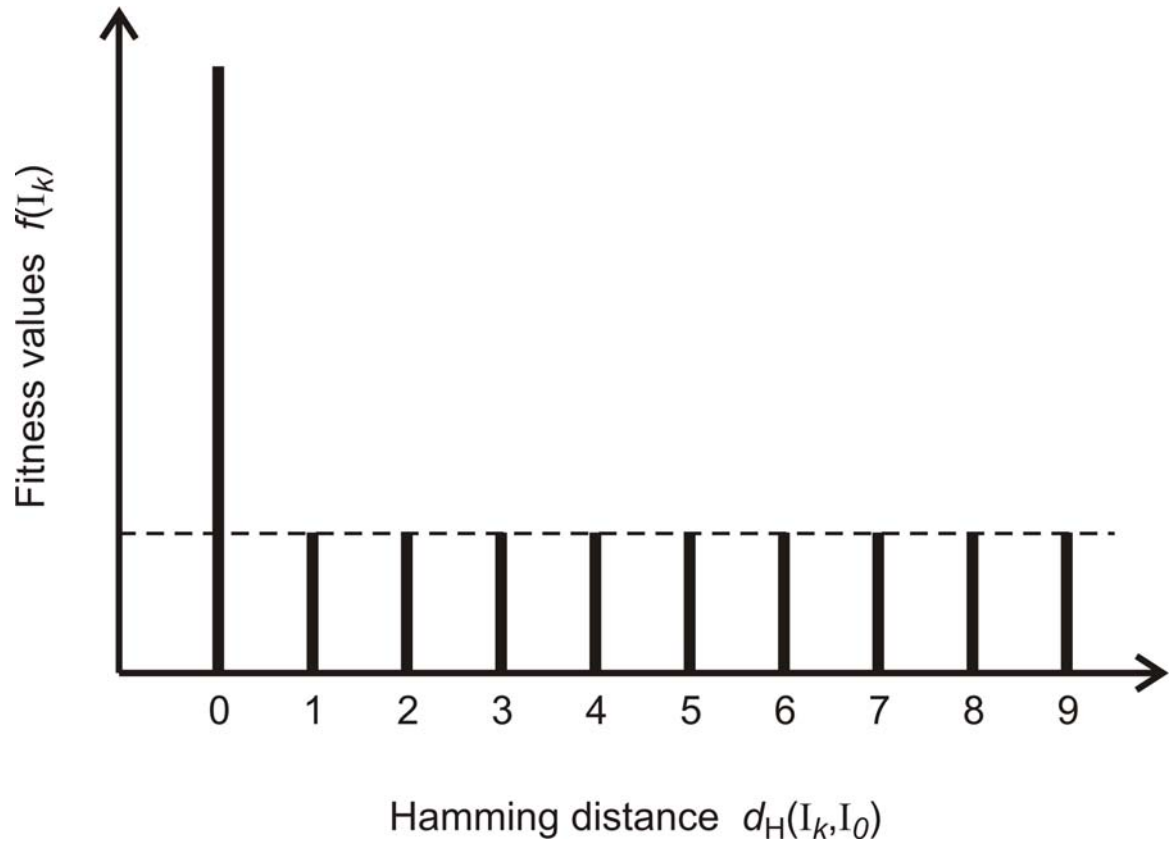
1. The chemistry of Darwinian evolution
2. RNA sequences and structures
- 3. Consequences of neutrality**
4. Evolutionary optimization of RNA structure
5. Complexity in biology



Chemical kinetics of replication and mutation as parallel reactions



The error threshold in replication



A fitness landscape showing an error threshold

SELF-REPLICATION WITH ERRORS

A MODEL FOR POLYNUCLEOTIDE REPLICATION\*\*

Jörg SWETINA and Peter SCHUSTER\*

Institut für Theoretische Chemie und Strahlenchemie der Universität, Währingerstraße 17, A-1090 Wien, Austria

Received 4th June 1982  
 Revised manuscript received 23rd August 1982  
 Accepted 30th August 1982

Key words: Polynucleotide replication; Quasi-species; Point mutation; Mutant class; Stochastic replication

A model for polynucleotide replication is presented and analyzed by means of perturbation theory. Two basic assumptions allow handling of sequences up to a chain length of  $n = 80$  explicitly: point mutations are restricted to a two-digit model and individual sequences are subsumed into mutant classes. Perturbation theory is in excellent agreement with the exact results for long enough sequences ( $n > 20$ ).

1. Introduction

Eigen [8] proposed a formal kinetic equation (eq. 1) which describes self-replication under the constraint of constant total population size:

$$\frac{dx_i}{dt} = x_i \sum_j w_{ij} x_j - \frac{x_i}{c} \phi; i = 1, \dots, n \quad (1)$$

By  $x_i$  we denote the population number or concentration of the self-replicating element  $I_i$ , i.e.,  $x_i = [I_i]$ . The total population size or total concentration  $c = \sum_i x_i$  is kept constant by proper adjustment of the constraint  $\phi = \sum_i \sum_j w_{ij} x_j x_i$ . Characteristically, this constraint has been called 'constant organization'. The relative values of diagonal

( $w_{ii}$ ) and off-diagonal ( $w_{ij}, i \neq j$ ) rates, as we shall see in detail in section 2, are related to the accuracy of the replication process. The specific properties of eq. 1 are essentially based on the fact that it leads to exponential growth in the absence of constraints ( $\phi = 0$ ) and competitors ( $n = 1$ ).

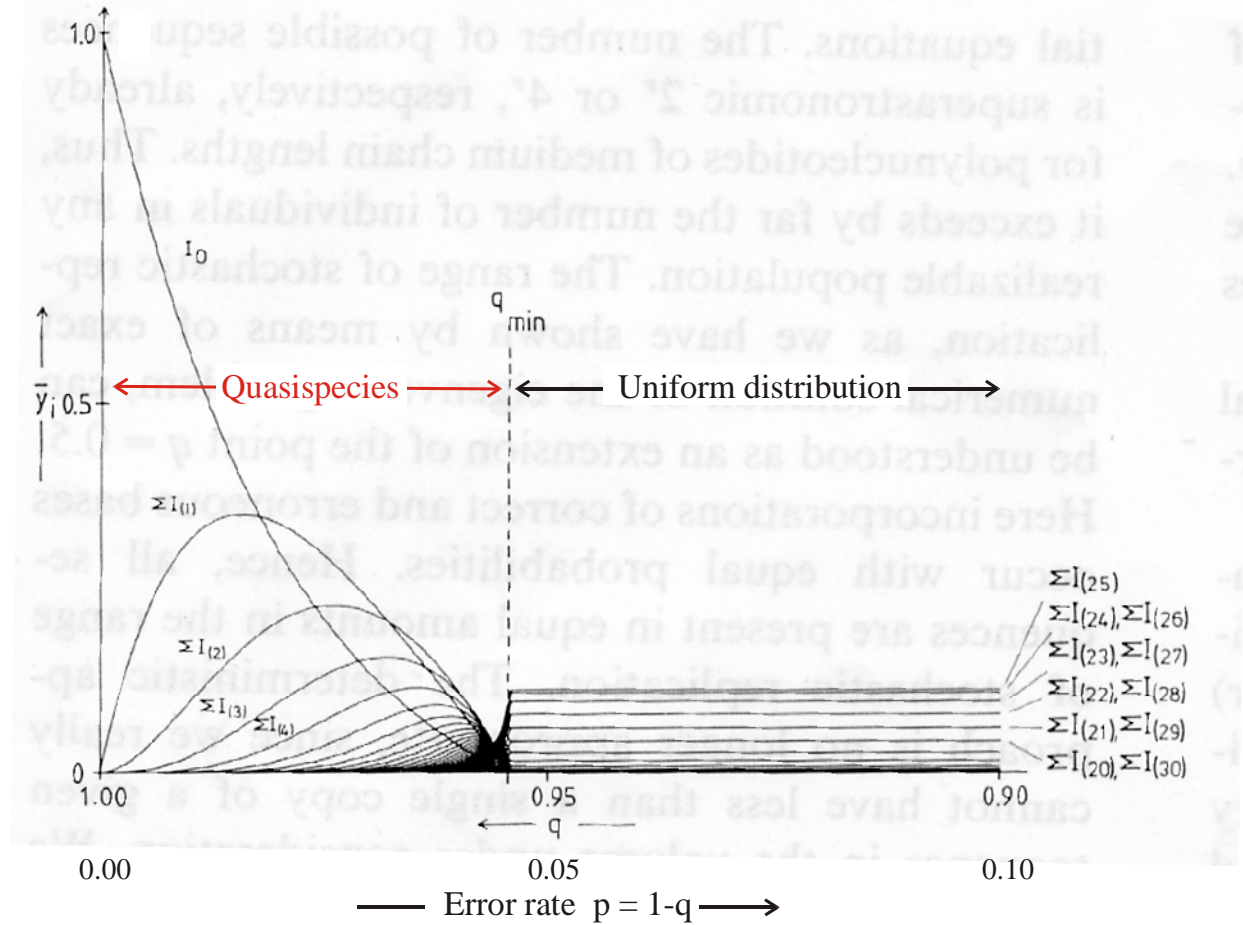
The non-linear differential equation, eq. 1 - the non-linearity is introduced by the definition of  $\phi$  at constant organization - shows a remarkable feature: it leads to selection of a defined ensemble of self-replicating elements above a certain accuracy threshold. This ensemble of a master and its most frequent mutants is a so-called 'quasi-species' [9]. Below this threshold, however, no selection takes place and the frequencies of the individual elements are determined exclusively by their statistical weights.

Rigorous mathematical analysis has been performed on eq. 1 [7,15,24,26]. In particular, it was shown that the non-linearity of eq. 1 can be removed by an appropriate transformation. The eigenvalue problem of the linear differential equation obtained thereby may be solved approximately by the conventional perturbation technique

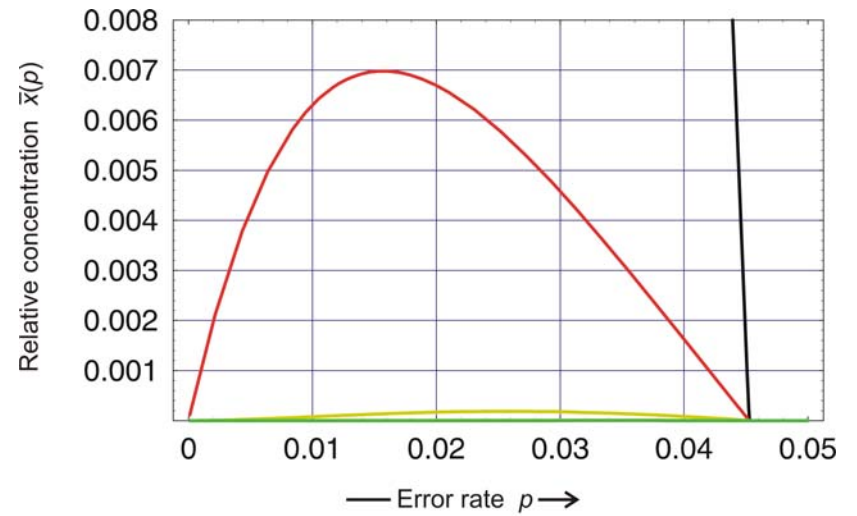
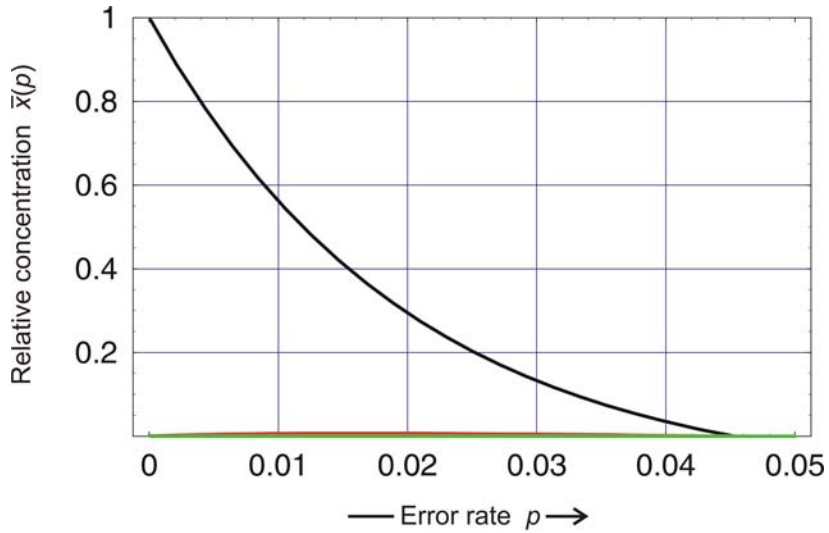
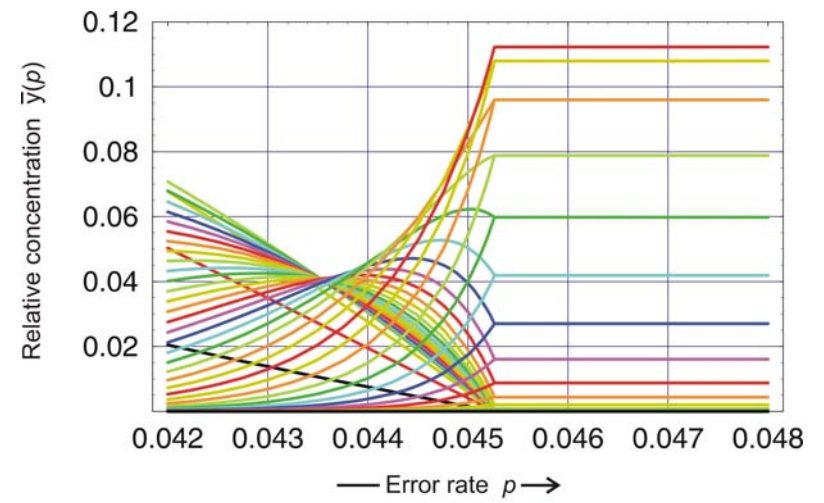
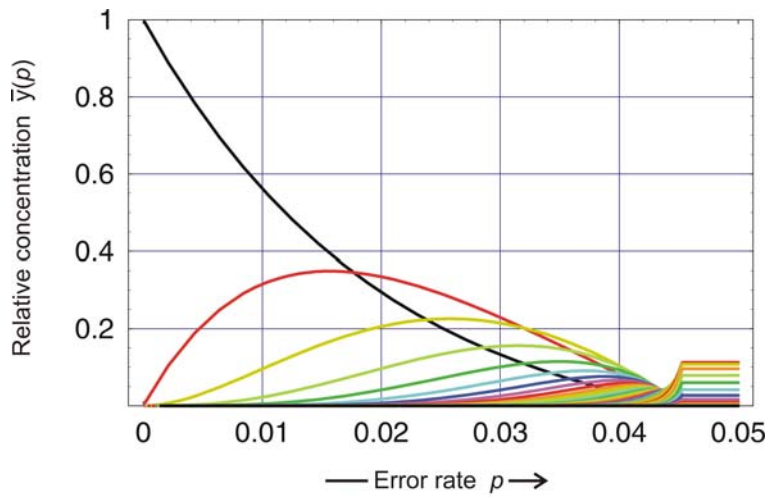
\* Dedicated to the late Professor B.L. Jones who was among the first to do rigorous mathematical analysis on the problems described here.

\*\* This paper is considered as part II of Model Studies on RNA replication. Part I is by Gassner and Schuster [14].

† All summations throughout this paper run from 1 to  $n$  unless specified differently:  $\Sigma_i = \Sigma_{i=1}^n$  and  $\Sigma_{i,j} = \Sigma_{i=1}^n + \Sigma_{j=1}^n$ , respectively.

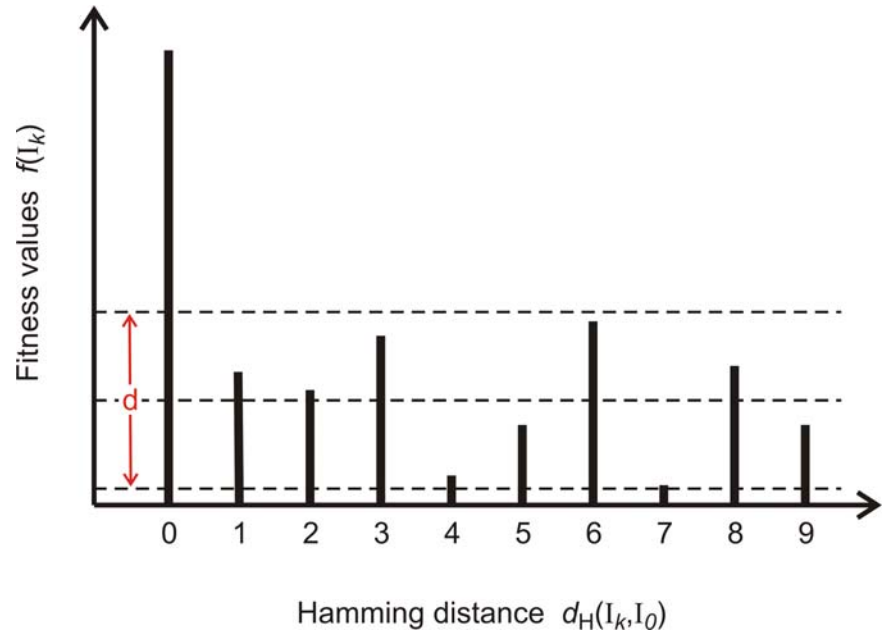
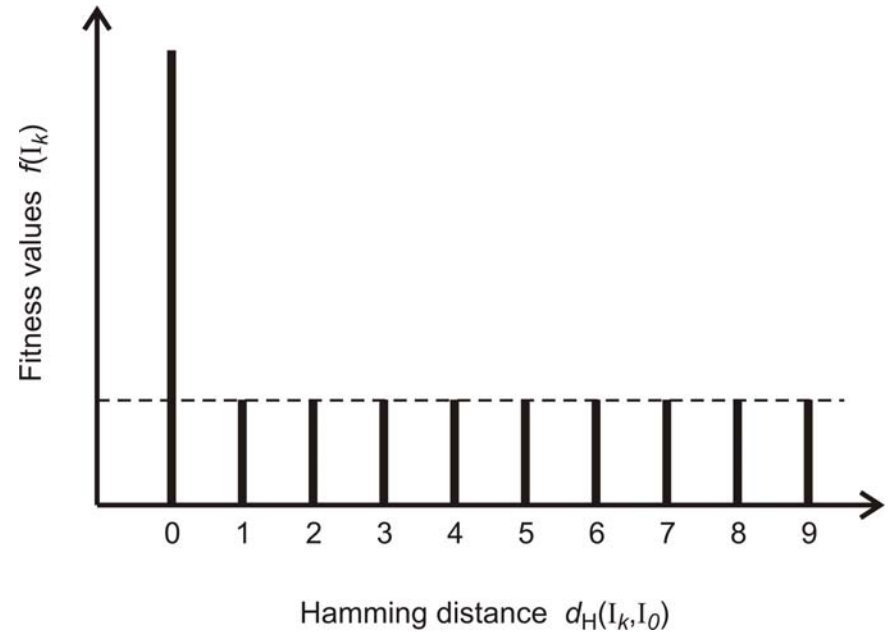


Stationary population or **quasispecies** as a function of the mutation or error rate  $p$

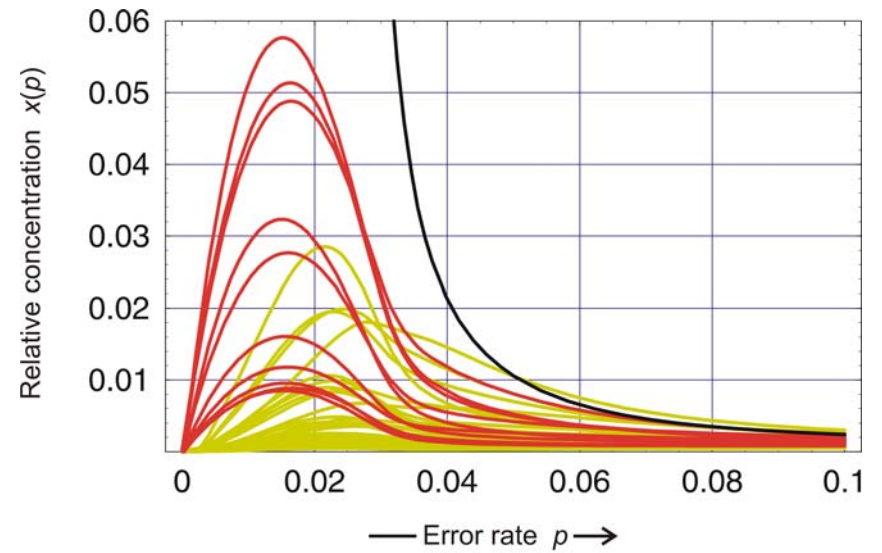
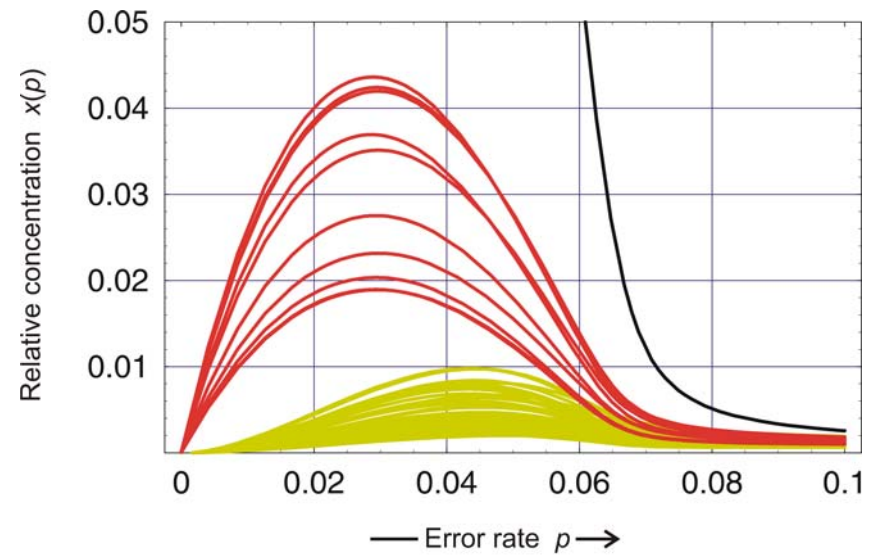
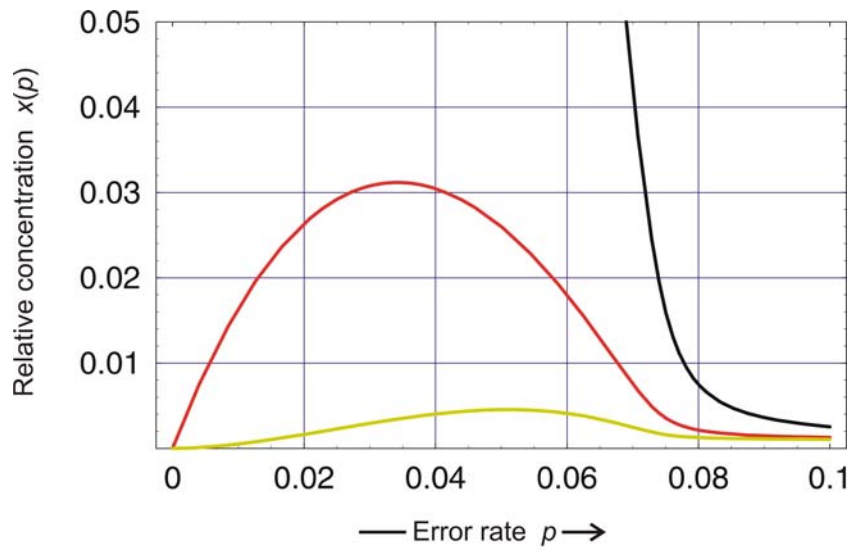


Error threshold on a single peak fitness landscape with  $n = 50$  and  $\sigma = 10$





Fitness landscapes showing error thresholds



Error threshold: Individual sequences

$n = 10$ ,  $\sigma = 2$  and  $d = 0, 1.0, 1.85$

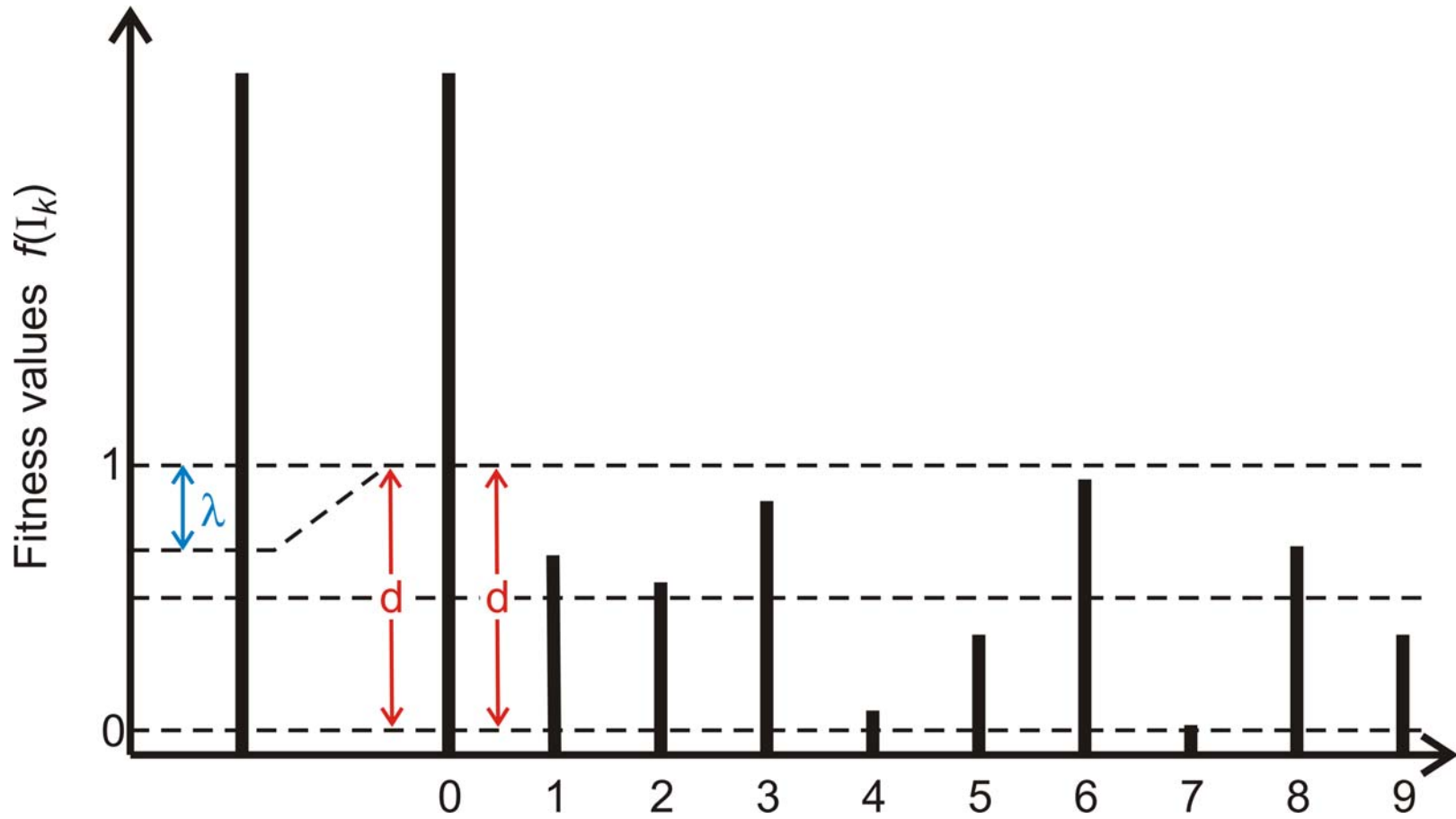
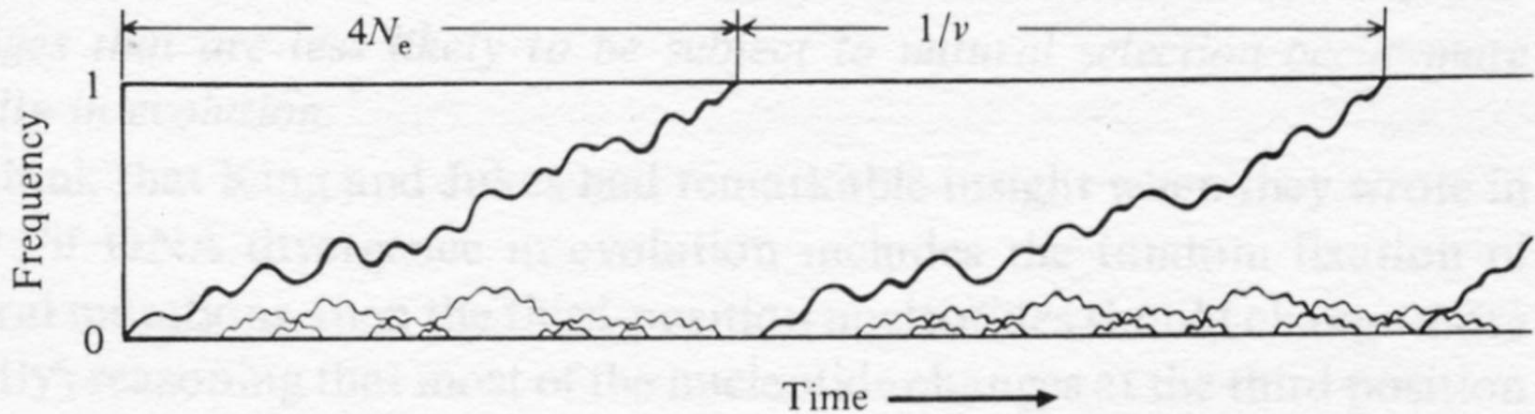


Fig. 3.1. Behavior of mutant genes following their appearance in a finite population. Courses of change in the frequencies of mutants destined to fixation are depicted by thick paths.  $N_e$  stands for the effective population size and  $v$  is the mutation rate.



Motoo Kimura

Is the Kimura scenario correct for frequent mutations?

## STATIONARY MUTANT DISTRIBUTIONS AND EVOLUTIONARY OPTIMIZATION

■ PETER SCHUSTER and JÖRG SWETINA  
Institut für theoretische Chemie  
und Strahlenchemie der Universität Wien,  
Währingerstraße 17,  
A 1090 Wien,  
Austria

Molecular evolution is modelled by erroneous replication of binary sequences. We show how the selection of two species of equal or almost equal selective value is influenced by its nearest neighbours in sequence space. In the case of perfect neutrality and sufficiently small error rates we find that the Hamming distance between the species determines selection. As the error rate increases the fitness parameters of neighbouring species become more and more important. In the case of almost neutral sequences we observe a critical replication accuracy at which a drastic change in the "quasispecies", in the stationary mutant distribution occurs. Thus, in frequently mutating populations fitness turns out to be an ensemble property rather than an attribute of the individual.

In addition we investigate the time dependence of the mean excess production as a function of initial conditions. Although it is optimized under most conditions, cases can be found which are characterized by decrease or non-monotonous change in mean excess productions.

*1. Introduction.* Recent data from populations of RNA viruses provided direct evidence for vast sequence heterogeneity (Domingo *et al.*, 1987). The origin of this diversity is not yet completely known. It may be caused by the low replication accuracy of the polymerizing enzyme, commonly a virus specific, RNA dependent RNA synthetase, or it may be the result of a high degree of selective neutrality of polynucleotide sequences. Eventually, both factors contribute to the heterogeneity observed. Indeed, mutations occur much more frequently than previously assumed in microbiology. They are by no means rare events and hence, neither the methods of conventional population genetics (Ewens, 1979) nor the neutral theory (Kimura, 1983) can be applied to these virus populations. Selectively neutral variants may be close with respect to Hamming distance and then the commonly made assumption that the mutation backflow from the mutants to the wilde type is negligible does not apply.

A kinetic theory of polynucleotide evolution which was developed during the past 15 years (Eigen, 1971; 1985; Eigen and Schuster, 1979; Eigen *et al.*, 1987; Schuster, 1986); Schuster and Sigmund, 1985) treats correct replication and mutation as parallel reactions within one and the same reaction network

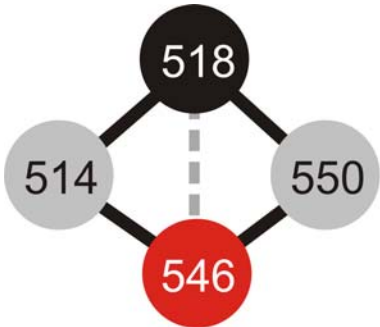


Neutral network

$\lambda = 0.01, s = 367$

$$d_H = 1$$

$$\lim_{p \rightarrow 0} x_1(p) = x_2(p) = 0.5$$



Neutral network

$\lambda = 0.01, s = 877$

$$d_H = 2$$

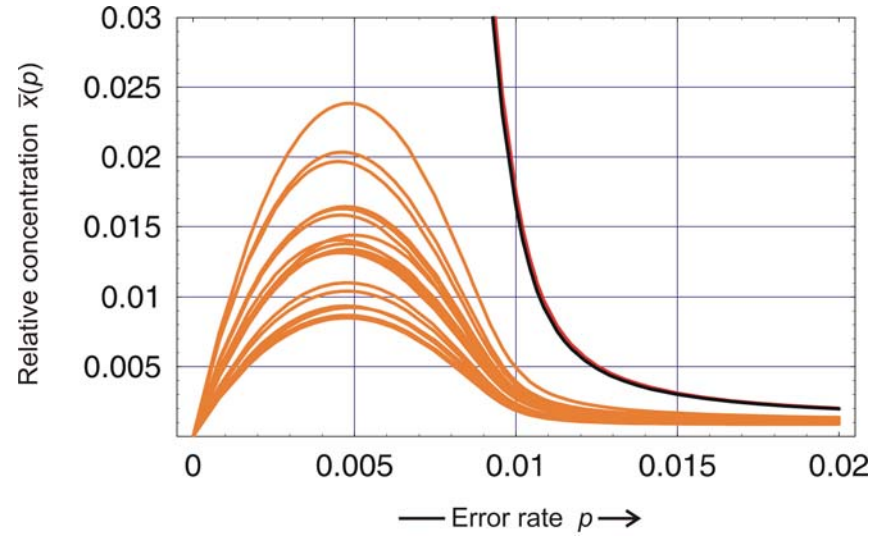
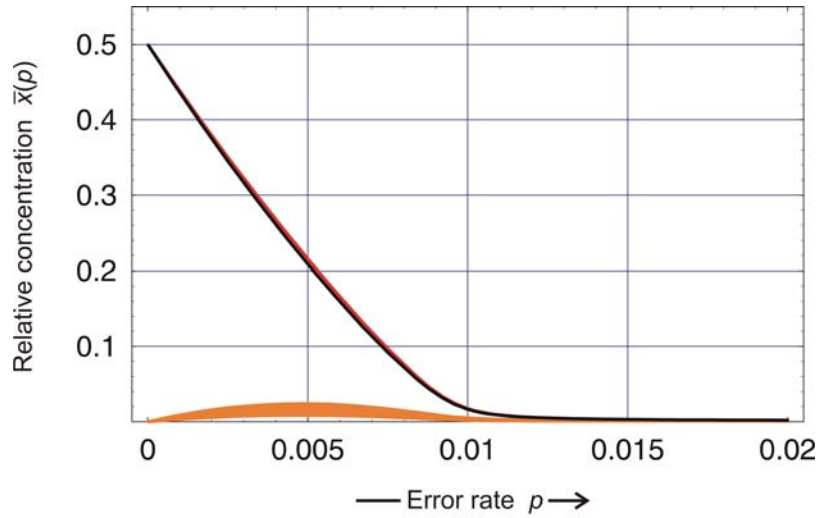
$$\lim_{p \rightarrow 0} x_1(p) = a$$

$$\lim_{p \rightarrow 0} x_2(p) = 1 - a$$

$$d_H = 3$$

random fixation in the sense of  
Motoo Kimura

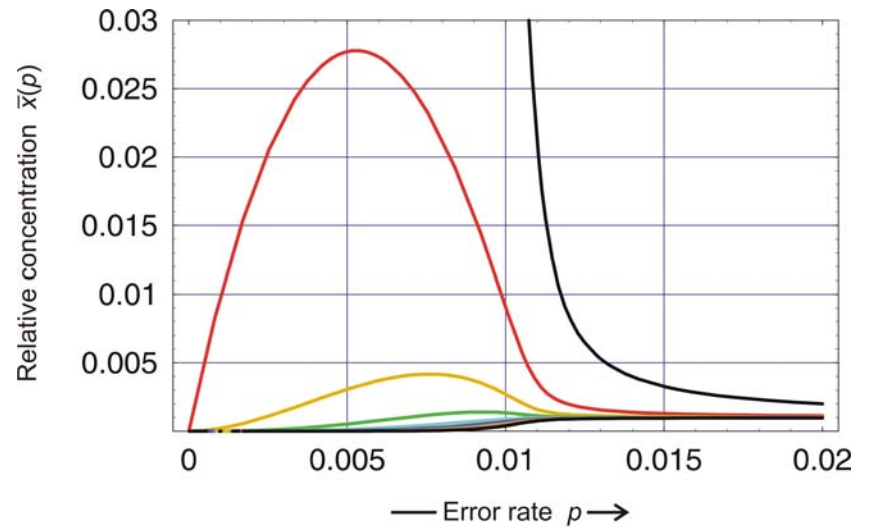
Pairs of genotypes in neutral replication networks

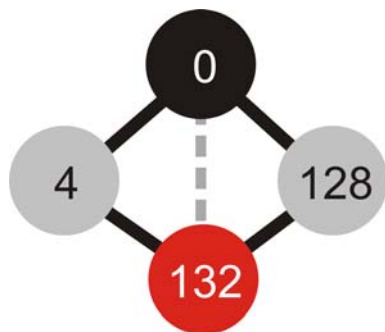


Neutral network  
 $\lambda = 0.01, s = 367$

Neutral network: Individual sequences

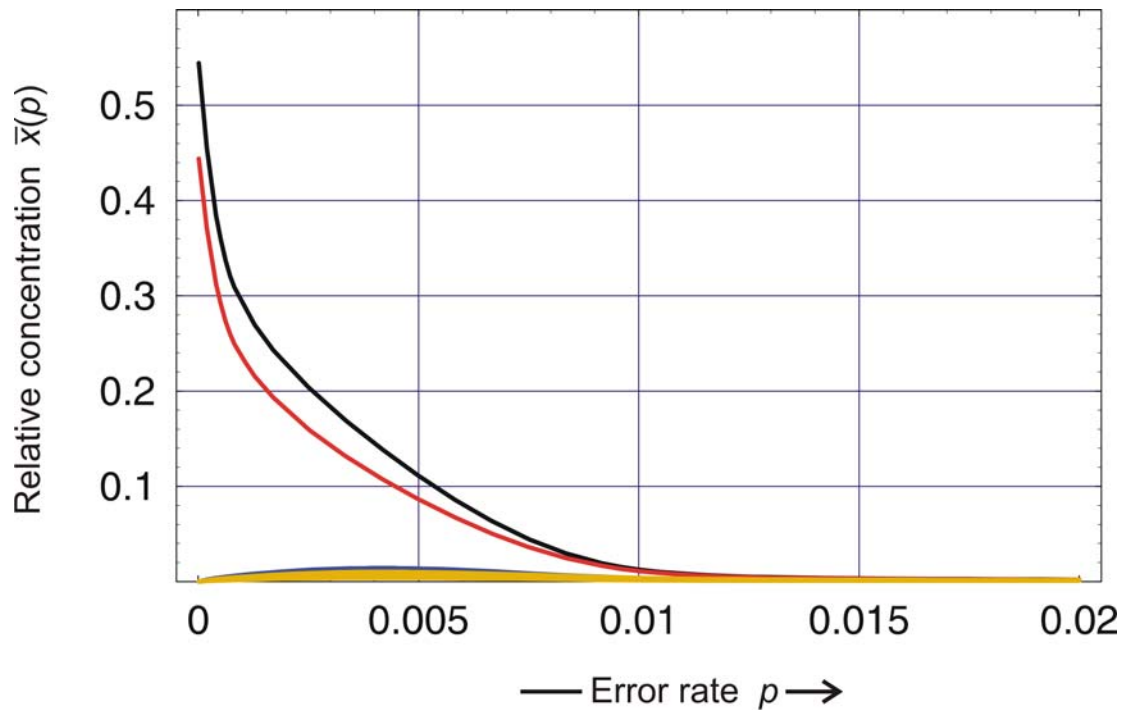
$n = 10, \sigma = 1.1, d = 1.0$





Neutral network

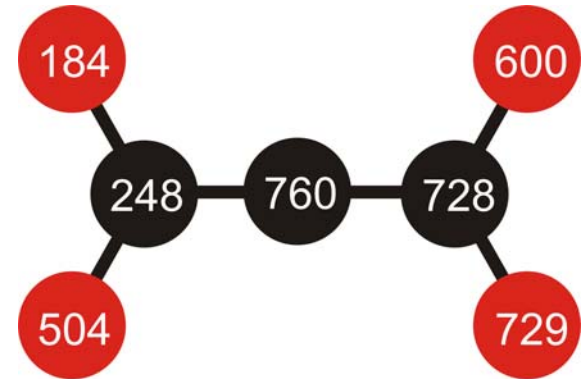
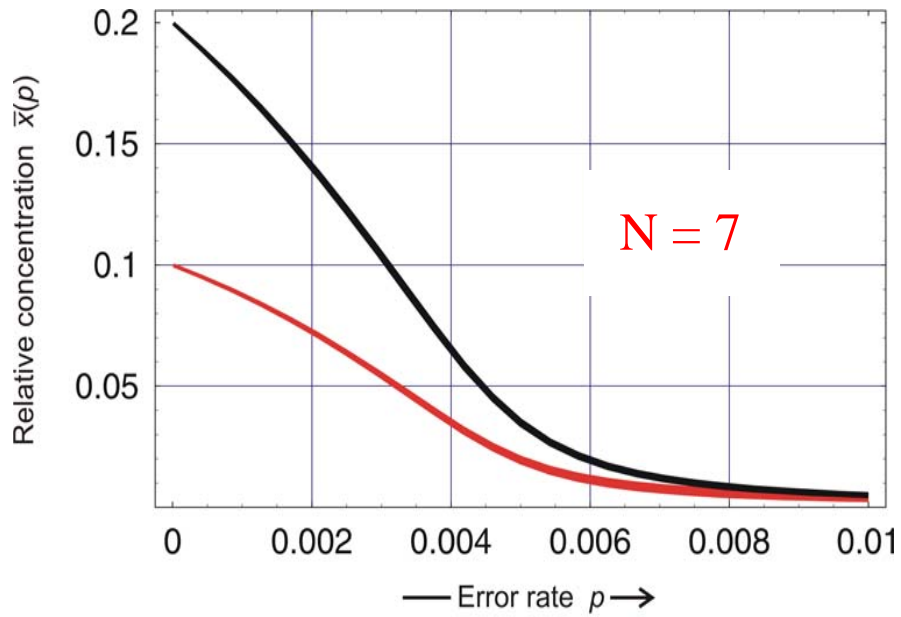
$\lambda = 0.01$ ,  $s = 877$



Neutral network: Individual sequences

$n = 10$ ,  $\sigma = 1.1$ ,  $d = 1.0$

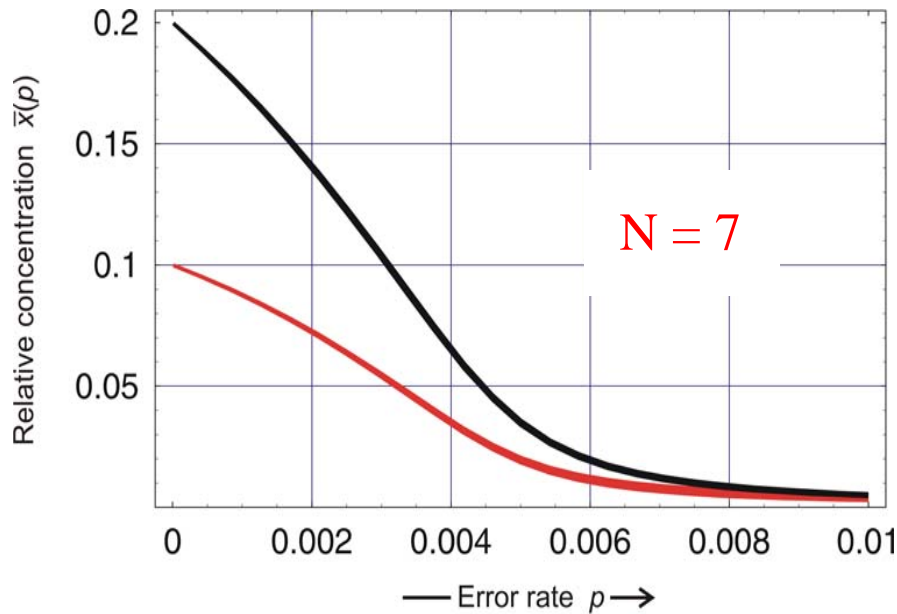




Neutral network

$\lambda = 0.10, s = 229$

Neutral networks with increasing  $\lambda$ :  $\lambda = 0.10, s = 229$



Perturbation matrix  $W$

$$W = \begin{pmatrix} f & 0 & \varepsilon & 0 & 0 & 0 & 0 \\ 0 & f & \varepsilon & 0 & 0 & 0 & 0 \\ \varepsilon & \varepsilon & f & \varepsilon & 0 & 0 & 0 \\ 0 & 0 & \varepsilon & f & \varepsilon & 0 & 0 \\ 0 & 0 & 0 & \varepsilon & f & \varepsilon & \varepsilon \\ 0 & 0 & 0 & 0 & \varepsilon & f & 0 \\ 0 & 0 & 0 & 0 & \varepsilon & 0 & f \end{pmatrix}$$

Eigenvalues of  $W$

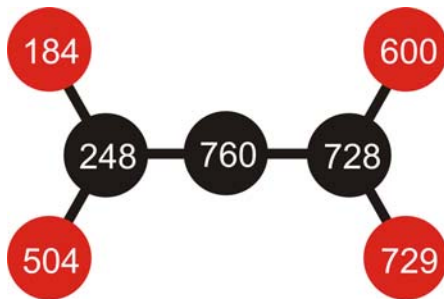
$$\lambda_0 = f + 2\varepsilon,$$

$$\lambda_1 = f + \sqrt{2}\varepsilon,$$

$$\lambda_{2,3,4} = f,$$

$$\lambda_5 = f - \sqrt{2}\varepsilon,$$

$$\lambda_6 = f - 2\varepsilon.$$



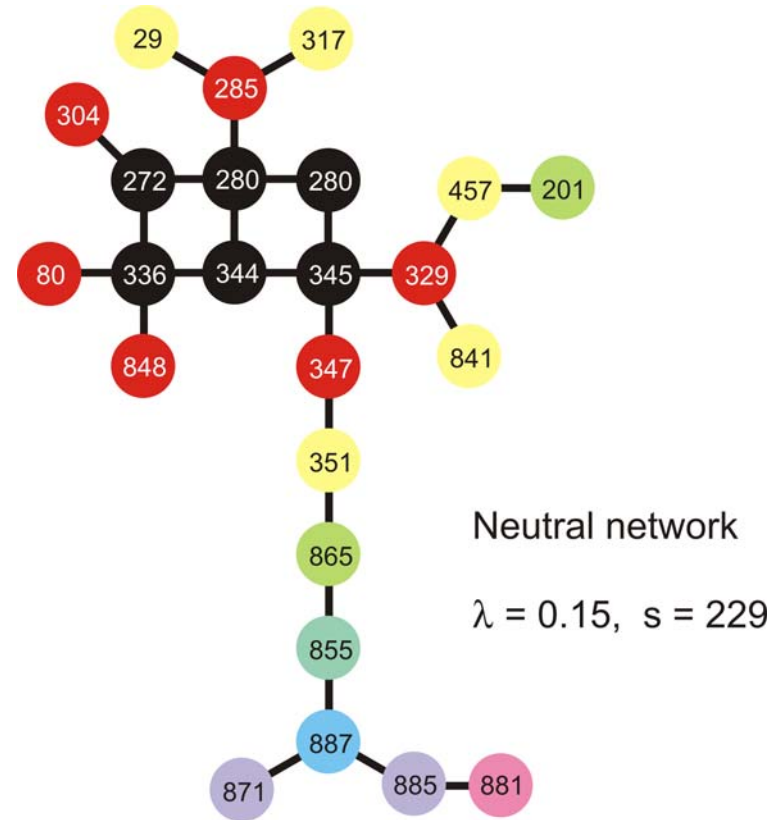
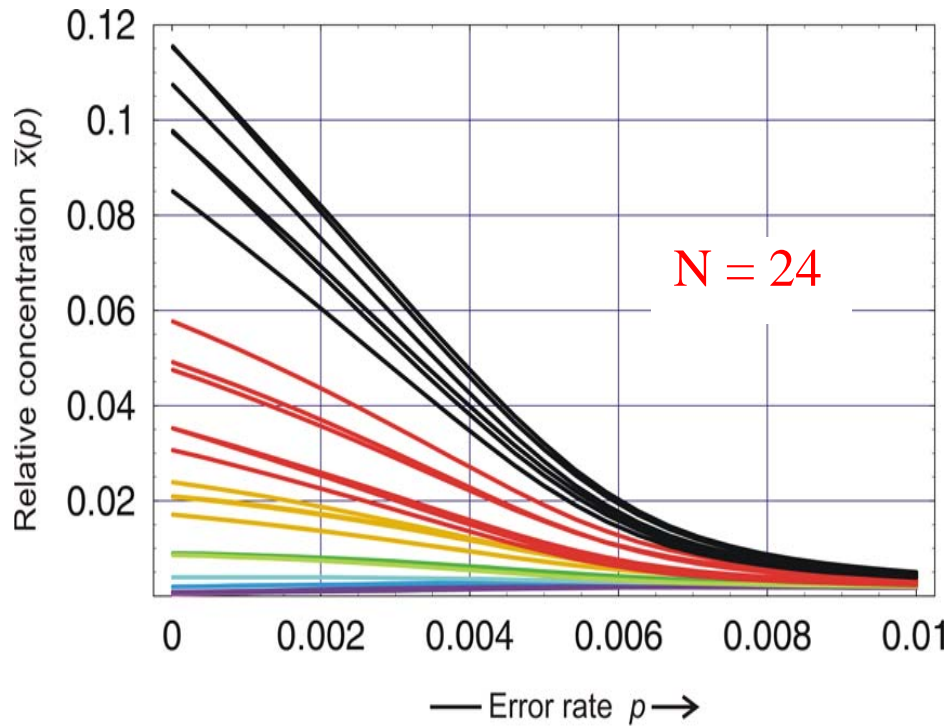
Neutral network

$$\lambda = 0.10, s = 229$$

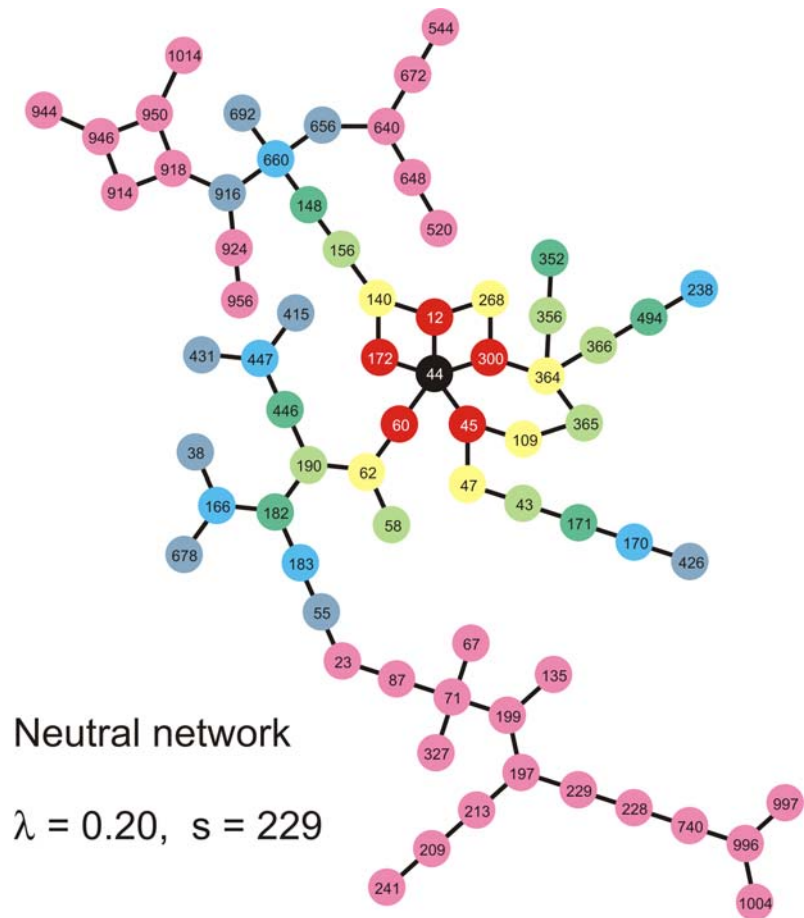
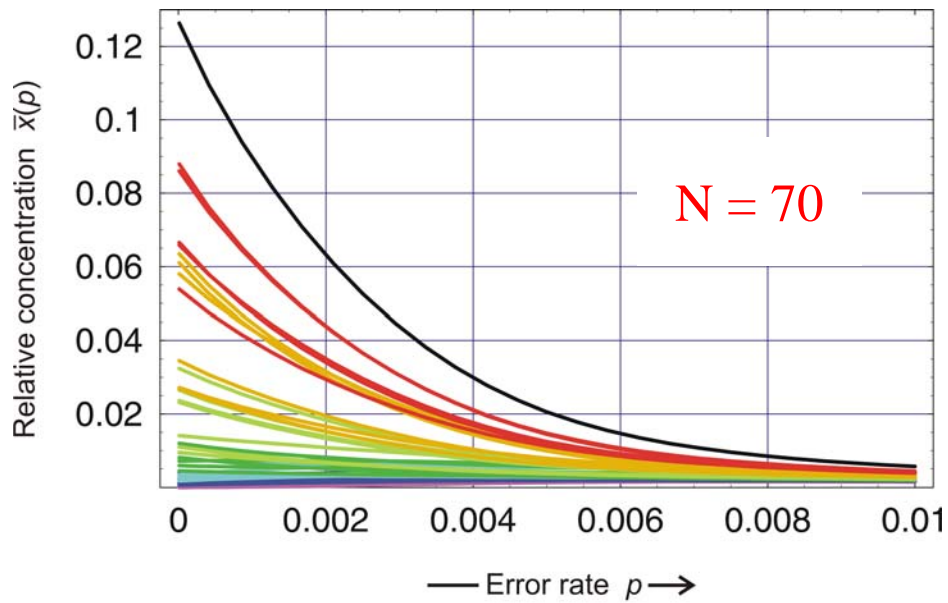
Largest eigenvector of  $W$

$$\xi_0 = (0.1, 0.1, 0.2, 0.2, 0.2, 0.1, 0.1).$$

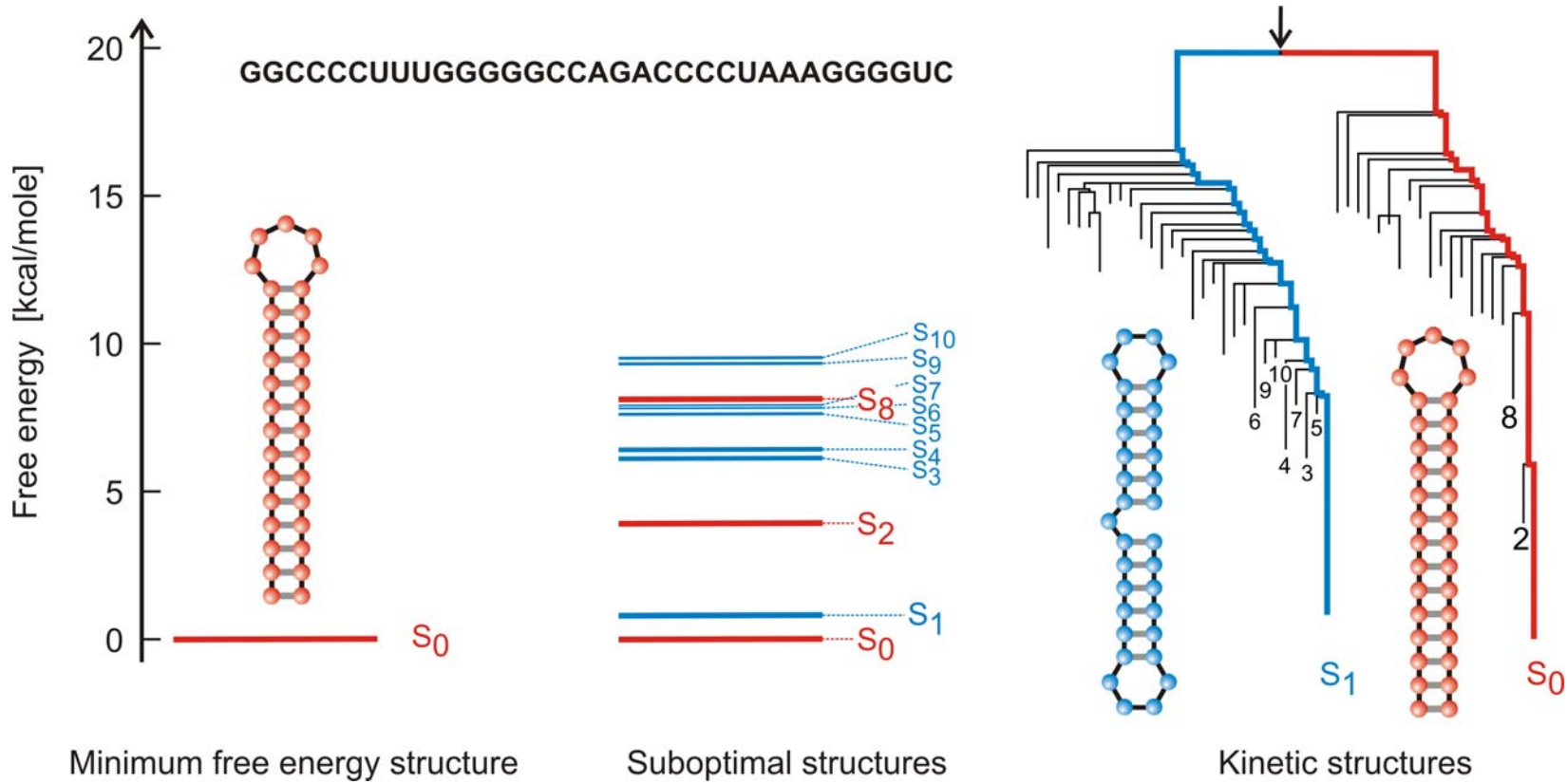
Neutral networks with increasing  $\lambda$ :  $\lambda = 0.10, s = 229$



Neutral networks with increasing  $\lambda$ :  $\lambda = 0.15, s = 229$



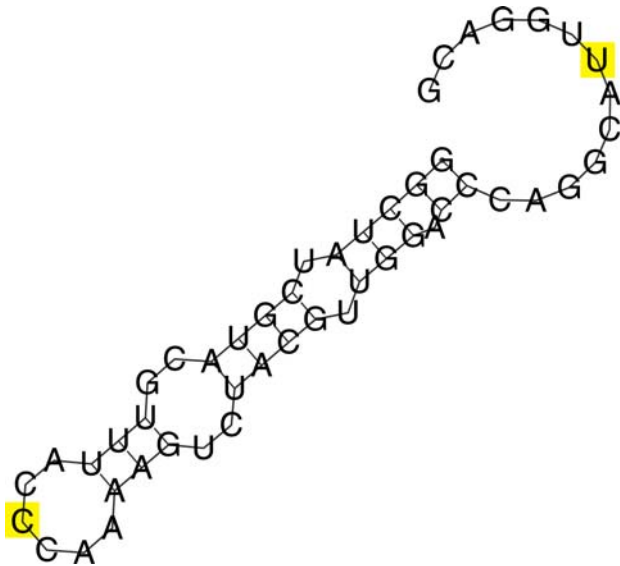
Neutral networks with increasing  $\lambda$ :  $\lambda = 0.20, s = 229$



Extension of the notion of structure

GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG

((((((.((((..(((.....))))..))))..)))..))..... -7.30  
 .....(((((((.(.....((((.....))))))..))..)))) -6.70  
 .....(((((((.(.....((((.....))))))..))..)))) -6.60  
 ..(((.((((..(((.....))))..))))..))..(((.....)))... -6.10  
 ((((((.((((..(((.....))))..))))..)))..))..(.....). -6.00  
 ((((((.((((..(((.....))))..))))..)))..))..(.....). -6.00  
 .(((.(.....((((.....))))..))..))..(.....)..... -6.00

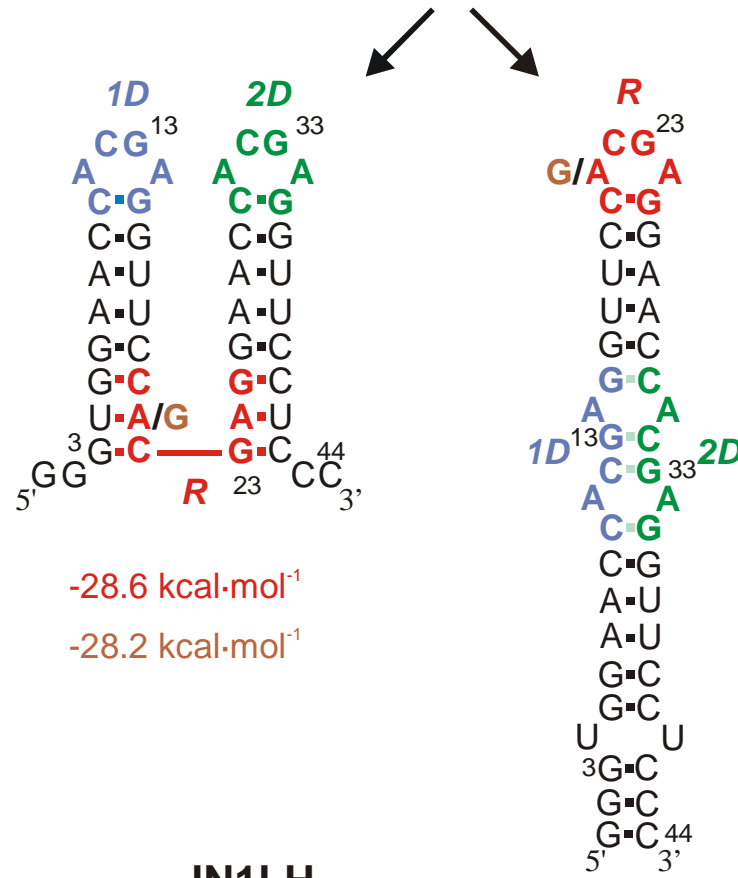


GGCUAUCGUACGUUUACCAAAAAGUCUACGUUGGACCCAGGCAUUGGACG

((((((.((((..(((.....))))..))))..)))..))..... -7.30  
 .(((.(.....((((.....))))..))..))..... -6.50  
 .(((.....((((.....))))..))..(.....)))..... -6.30  
 ..(((.((((..(((.....))))..))))..))..(((.....)))... -6.10  
 ((((((.((((..(((.....))))..))))..)))..))..(.....). -6.00  
 ((((((.((((..(((.....))))..))))..)))..))..(.....). -6.00  
 .(((.....((((.....))))..))..))..(.....)..... -6.00

GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG

((((((.((((..(((.....))))..))))..)))..))..... -7.30  
 ..(((.((((..(((.....))))..))))..))..(((.....)))... -7.20  
 .....(((((((.(.....((((.....))))))..))..)))) -6.70  
 .....(((((((.(.....((((.....))))))..))..)))) -6.60  
 ((((((.((((..(((.....))))..))))..)))..))..(.....)..... -6.50  
 (.(((.((((..(((.....))))..))))..)))..(.....)))... -6.30  
 .(((.(.....((((.....))))..))..))..(.....)))... -6.30  
 .....(((.(.....((((.....))))..))..))..(.....)))... -6.30  
 (.(((.(.....((((.....))))..))..))..(.....)))... -6.10  
 .....((.....((((.....))))..))..(.....)))... -6.10  
 .....(((.(.....((((.....))))..))..))..(.....)))... -6.10  
 ((((((.((((..(((.....))))..))))..)))..))..(.....). -6.00  
 ((((((.((((..(((.....))))..))))..)))..))..(.....). -6.00  
 .(((.(.....((((.....))))..))..))..(.....)..... -6.00  
 .....(((.(.....((((.....))))..))..))..(.....)..... -6.00



-28.6 kcal·mol<sup>-1</sup>  
 -28.2 kcal·mol<sup>-1</sup>

-28.6 kcal·mol<sup>-1</sup>  
 -31.8 kcal·mol<sup>-1</sup>

## An RNA switch

JN1LH

J.H.A. Nagel, C. Flamm, I.L. Hofacker, K. Franke,  
 M.H. de Smit, P. Schuster, and C.W.A. Pleij.

Structural parameters affecting the kinetic competition of  
 RNA hairpin formation. *Nucleic Acids Res.* **34**:3568-3576,  
 2006.

- minus the background levels observed in the HSP in the control (Sar1-GDP-containing) incubation that prevents COPII vesicle formation. In the microsome control, the level of p115-SNARE associations was less than 0.1%.
46. C. M. Carr, E. Grote, M. Munson, F. M. Hughson, P. J. Novick, *J. Cell Biol.* **146**, 333 (1999).
  47. C. Ungermann, B. J. Nichols, H. R. Pelham, W. Wickner, *J. Cell Biol.* **140**, 61 (1998).
  48. E. Grote and P. J. Novick, *Mol. Biol. Cell* **10**, 4149 (1999).
  49. P. Uetz et al., *Nature* **403**, 623 (2000).
  50. GST-SNARE proteins were expressed in bacteria and purified on glutathione-Sepharose beads using standard methods. Immobilized GST-SNARE protein (0.5  $\mu$ M) was incubated with rat liver cytosol (20 mg) or purified recombinant p115 (0.5  $\mu$ M) in 1 ml of NS buffer containing 1% BSA for 2 hours at 4°C with rotation. Beads were briefly spun (3000 rpm for 10 s) and sequentially washed three times with NS buffer and three times with NS buffer supplemented with 150 mM NaCl. Bound proteins were eluted three times in 50  $\mu$ l of 50 mM tris-HCl (pH 8.5), 50 mM reduced glutathione, 150 mM NaCl, and 0.1% Triton X-100 for 15 min at 4°C with intermittent mixing, and elutes were pooled. Proteins were precipitated by MeOH/CH<sub>2</sub>Cl<sub>2</sub> and separated by SDS-polyacrylamide gel electrophoresis (PAGE) followed by immunoblotting using p115 mAb 13F12.
  51. V. Rybin et al., *Nature* **383**, 266 (1996).
  52. K. G. Hardwick and H. R. Pelham, *J. Cell Biol.* **119**, 513 (1992).
  53. A. P. Newman, M. E. Groesch, S. Ferro-Novick, *EMBO J.* **11**, 3609 (1992).
  54. A. Spang and R. Schekman, *J. Cell Biol.* **143**, 589 (1998).
  55. M. F. Rexach, M. Latterich, R. W. Schekman, *J. Cell Biol.* **126**, 1133 (1994).
  56. A. Mayer and W. Wickner, *J. Cell Biol.* **136**, 307 (1997).
  57. M. D. Turner, H. Plutner, W. E. Balch, *J. Biol. Chem.* **272**, 13479 (1997).
  58. A. Price, D. Seals, W. Wickner, C. Ungermann, *J. Cell Biol.* **148**, 1231 (2000).
  59. X. Cao and C. Barlowe, *J. Cell Biol.* **149**, 55 (2000).
  60. G. G. Tall, H. Hama, D. B. DeWald, B. F. Horadzovsky, *Mol. Biol. Cell* **10**, 1873 (1999).
  61. C. G. Burd, M. Peterson, C. R. Cowles, S. D. Emr, *Mol. Biol. Cell* **8**, 1089 (1997).
  62. M. R. Peterson, C. G. Burd, S. D. Emr, *Curr. Biol.* **9**, 159 (1999).
  63. M. G. Waters, D. O. Clary, J. E. Rothman, *J. Cell Biol.* **118**, 1015 (1992).
  64. D. M. Walter, K. S. Paul, M. G. Waters, *J. Biol. Chem.* **273**, 29565 (1998).
  65. N. Hui et al., *Mol. Biol. Cell* **8**, 1777 (1997).
  66. T. E. Kreis, *EMBO J.* **5**, 931 (1986).
  67. H. Plutner, H. W. Davidson, J. Saraste, W. E. Balch, *J. Cell Biol.* **119**, 1097 (1992).
  68. D. S. Nelson et al., *J. Cell Biol.* **143**, 319 (1998).
  69. We thank G. Waters for p115 cDNA and p115 mAbs; G. Warren for p97 and p47 antibodies; R. Scheller for rbt1, membrin, and sec22 cDNAs; H. Plutner for excellent technical assistance; and P. Tan for help during the initial phase of this work. Supported by NIH grants GM 33301 and GM42336 and National Cancer Institute grant CA58689 (W.E.B.), a NIH National Research Service Award (B.D.M.), and a Wellcome Trust International Traveling Fellowship (B.B.A.).

20 March 2000; accepted 22 May 2000

## One Sequence, Two Ribozymes: Implications for the Emergence of New Ribozyme Folds

Erik A. Schultes and David P. Bartel\*

We describe a single RNA sequence that can assume either of two ribozyme folds and catalyze the two respective reactions. The two ribozyme folds share no evolutionary history and are completely different, with no base pairs (and probably no hydrogen bonds) in common. Minor variants of this sequence are highly active for one or the other reaction, and can be accessed from prototype ribozymes through a series of neutral mutations. Thus, in the course of evolution, new RNA folds could arise from preexisting folds, without the need to carry inactive intermediate sequences. This raises the possibility that biological RNAs having no structural or functional similarity might share a common ancestry. Furthermore, functional and structural divergence might, in some cases, precede rather than follow gene duplication.

Related protein or RNA sequences with the same folded conformation can often perform very different biochemical functions, indicating that new biochemical functions can arise from preexisting folds. But what evolutionary mechanisms give rise to sequences with new macromolecular folds? When considering the origin of new folds, it is useful to picture, among all sequence possibilities, the distribution of sequences with a particular fold and function. This distribution can range very far in sequence space (1). For example, only seven nucleotides are strictly conserved among the group I self-splicing introns, yet secondary (and presumably tertiary) structure within the core of the ribozyme is preserved (2). Because these dis-

parate isolates have the same fold and function, it is thought that they descended from a common ancestor through a series of mutational variants that were each functional. Hence, sequence heterogeneity among divergent isolates implies the existence of paths through sequence space that have allowed neutral drift from the ancestral sequence to each isolate. The set of all possible neutral paths composes a "neutral network," connecting in sequence space those widely dispersed sequences sharing a particular fold and activity, such that any sequence on the network can potentially access very distant sequences by neutral mutations (3-5).

Theoretical analyses using algorithms for predicting RNA secondary structure have suggested that different neutral networks are interwoven and can approach each other very closely (3, 5-8). Of particular interest is whether ribozyme neutral networks approach each other so closely that they intersect. If so, a single sequence would be capable of folding into two different conformations, would

have two different catalytic activities, and could access by neutral drift every sequence on both networks. With intersecting networks, RNAs with novel structures and activities could arise from previously existing ribozymes, without the need to carry non-functional sequences as evolutionary intermediates. Here, we explore the proximity of neutral networks experimentally, at the level of RNA function. We describe a close apposition of the neutral networks for the hepatitis delta virus (HDV) self-cleaving ribozyme and the class III self-ligating ribozyme.

In choosing the two ribozymes for this investigation, an important criterion was that they share no evolutionary history that might confound the evolutionary interpretations of our results. Choosing at least one artificial ribozyme ensured independent evolutionary histories. The class III ligase is a synthetic ribozyme isolated previously from a pool of random RNA sequences (9). It joins an oligonucleotide substrate to its 5' terminus. The prototype ligase sequence (Fig. 1A) is a shortened version of the most active class III variant isolated after 10 cycles of *in vitro* selection and evolution. This minimal construct retains the activity of the full-length isolate (10). The HDV ribozyme carries out the site-specific self-cleavage reactions needed during the life cycle of HDV, a satellite virus of hepatitis B with a circular, single-stranded RNA genome (11). The prototype HDV construct for our study (Fig. 1B) is a shortened version of the antigenomic HDV ribozyme (12), which undergoes self-cleavage at a rate similar to that reported for other antigenomic constructs (13, 14).

The prototype class III and HDV ribozymes have no more than the 25% sequence identity expected by chance and no fortuitous structural similarities that might favor an intersection of their two neutral networks. Nevertheless, sequences can be designed that simultaneously satisfy the base-pairing requirements

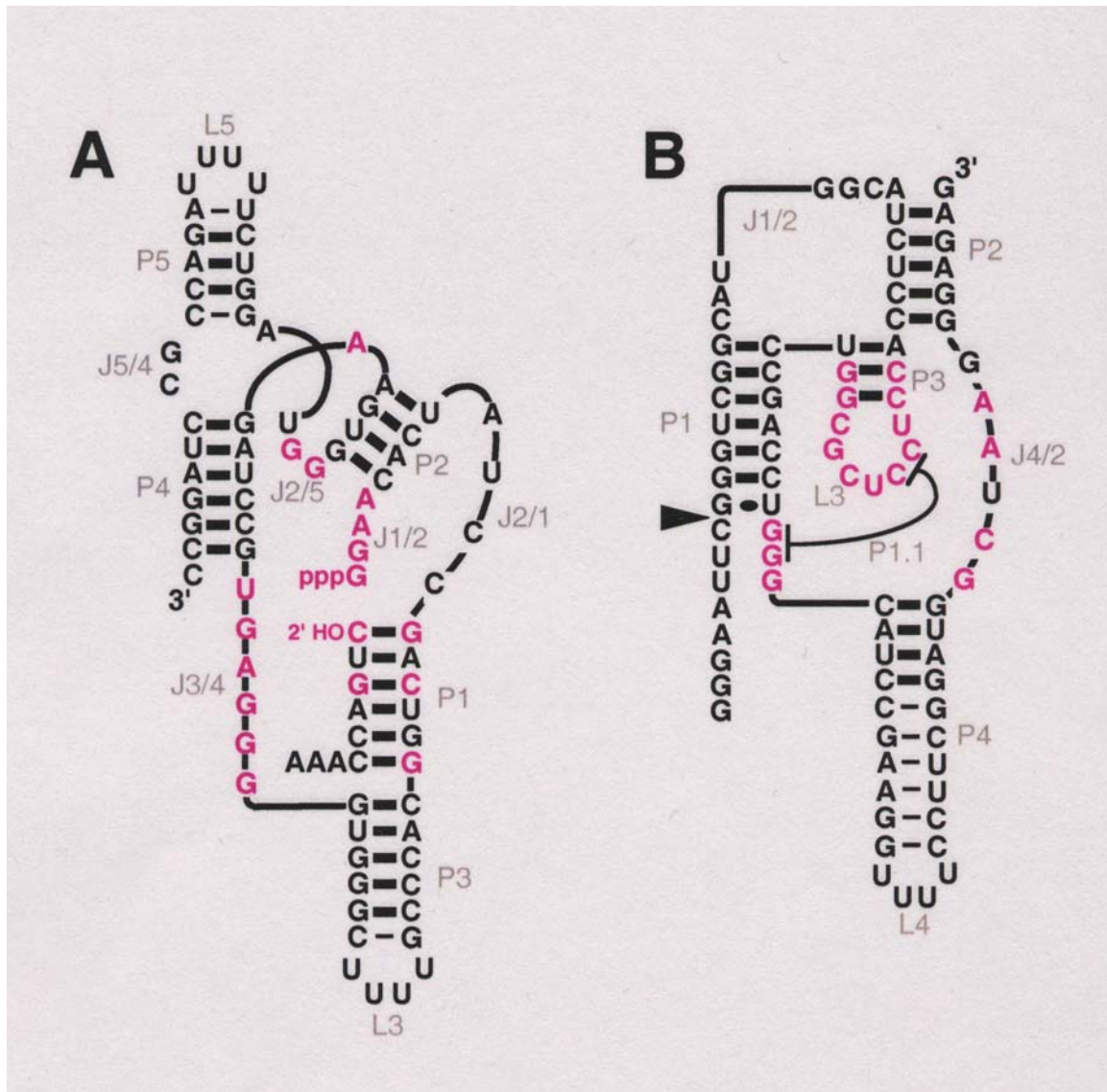
## A ribozyme switch

E.A.Schultes, D.B.Bartel, *Science*  
**289** (2000), 448-452

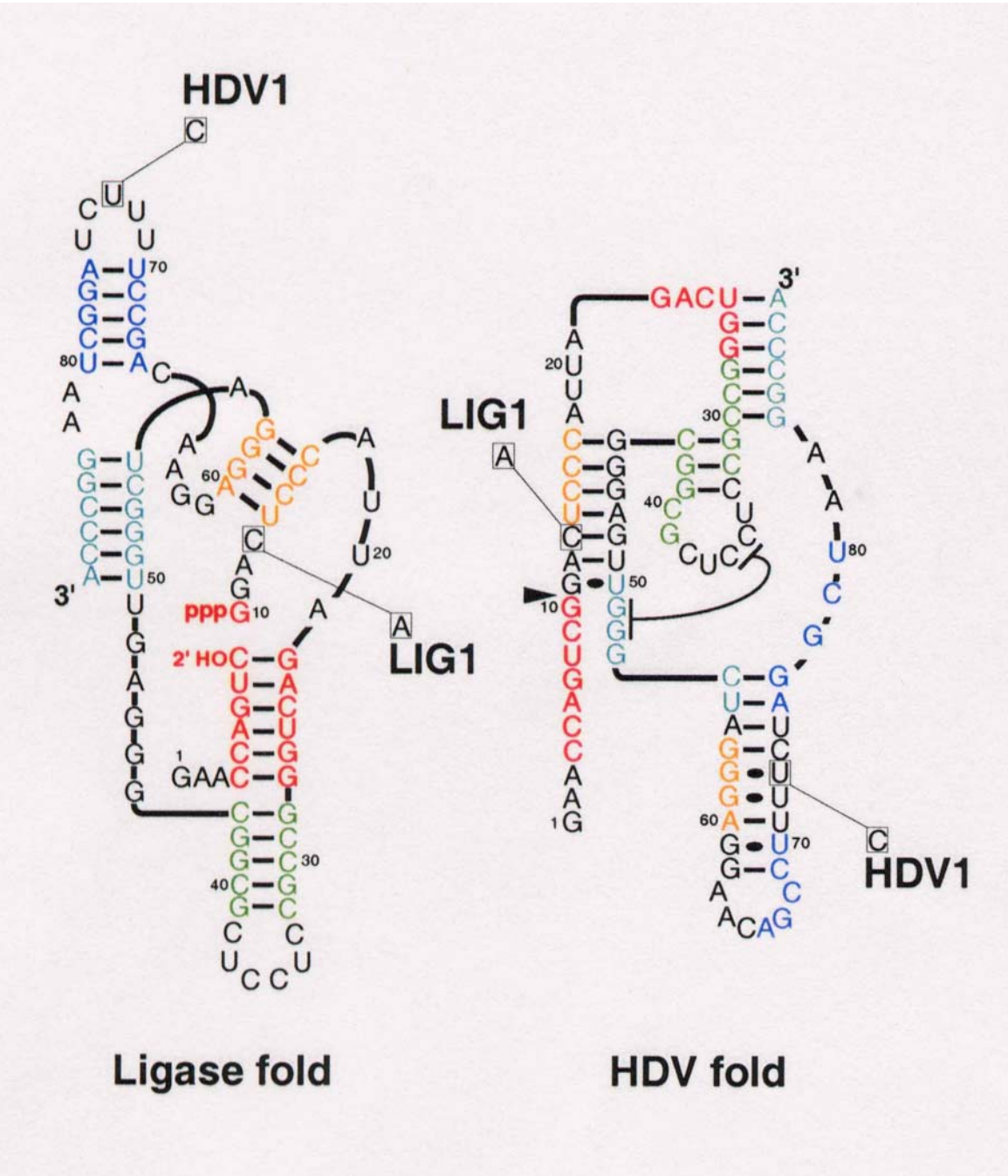
Whitehead Institute for Biomedical Research and Department of Biology, Massachusetts Institute of Technology, 9 Cambridge Center, Cambridge, MA 02142, USA.

\*To whom correspondence should be addressed. E-mail: dbartel@wi.mit.edu



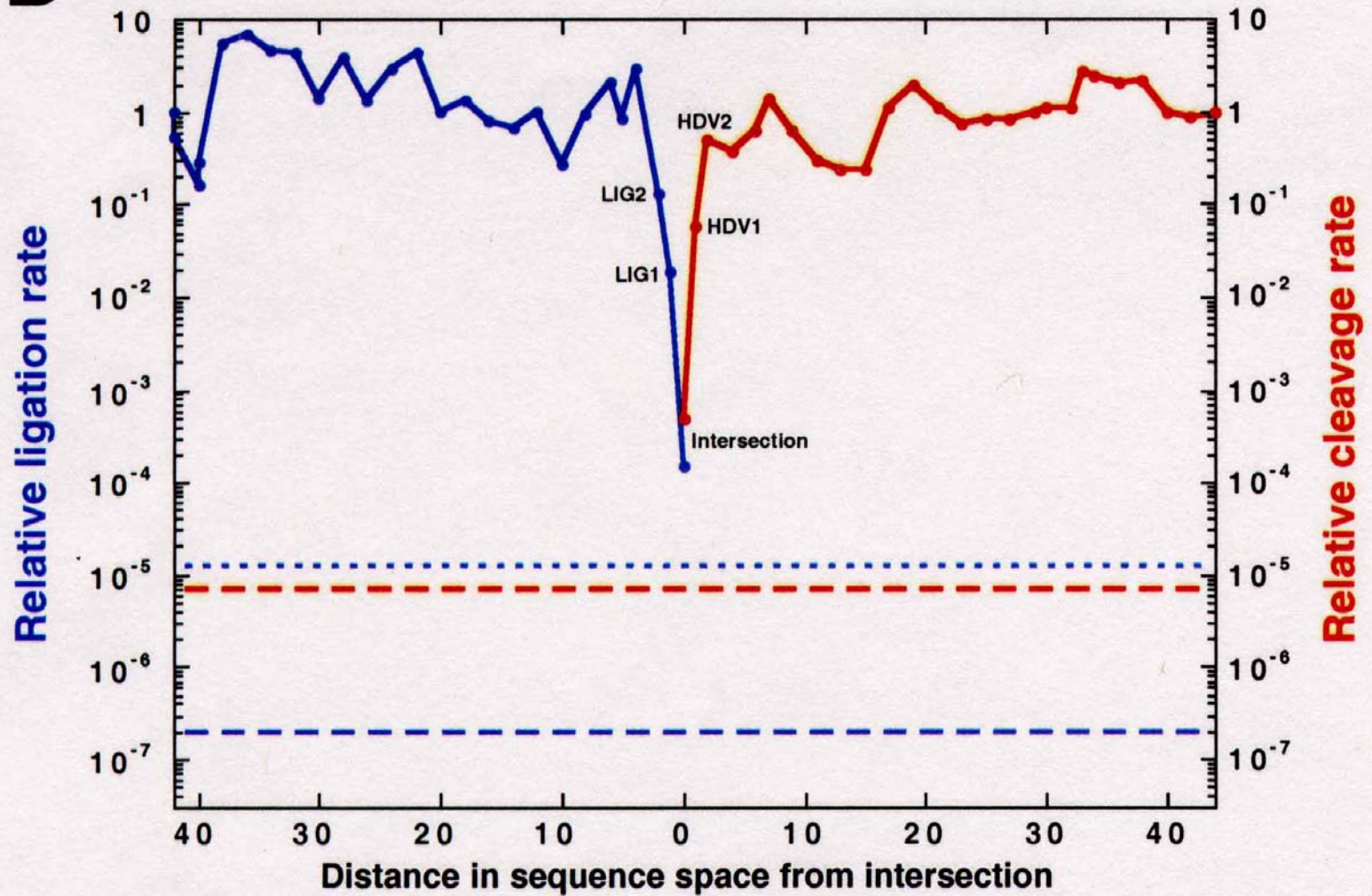


Two ribozymes of chain lengths  $n = 88$  nucleotides: An artificial ligase (**A**) and a natural cleavage ribozyme of hepatitis- $\delta$ -virus (**B**)



The sequence at the *intersection*:

An RNA molecules which is 88 nucleotides long and can form both structures

**B**

Two neutral walks through sequence space with conservation of structure and catalytic activity

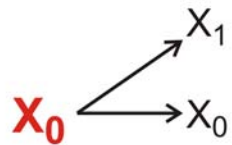
## Results from replication kinetics and *in silico* RNA evolution:

- Novel antiviral strategies were developed from known molecular mechanisms of virus evolution.
- Direct evidence that neutrality is increasing the repertoire of structures and properties in populations.

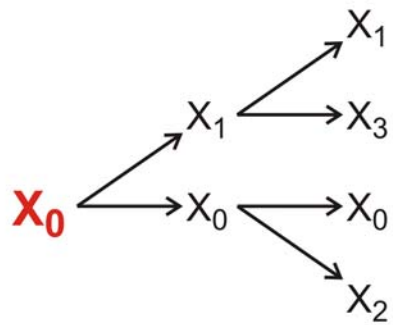
1. The chemistry of Darwinian evolution
2. RNA sequences and structures
3. Consequences of neutrality
4. **Evolutionary optimization of RNA structure**
5. Complexity in biology

$X_0$

Evolution of RNA molecules as a Markow process and its analysis by means of the relay series

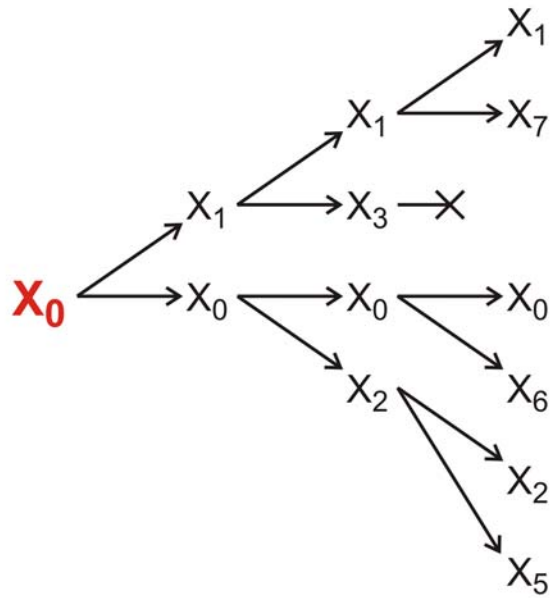


Evolution of RNA molecules as a Markov process and its analysis by means of the relay series

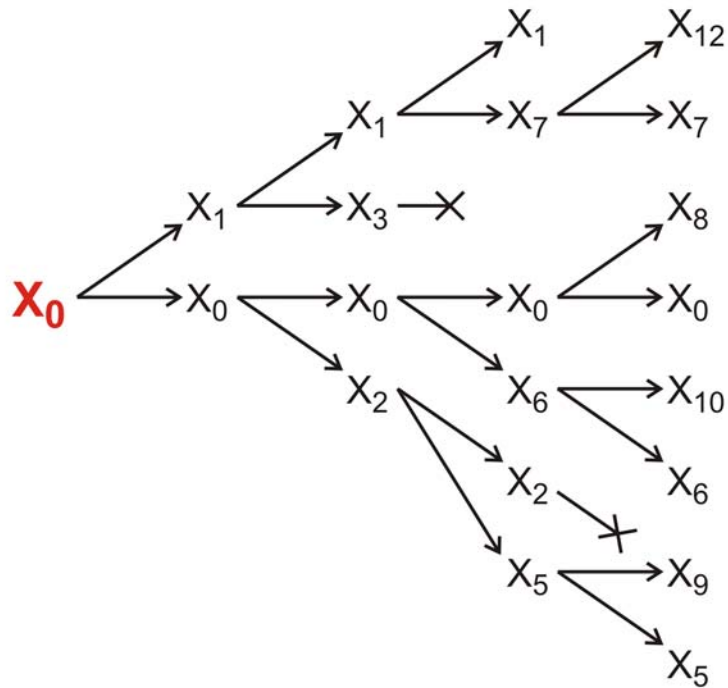


Evolution of RNA molecules as a Markov process and its analysis by means of the relay series

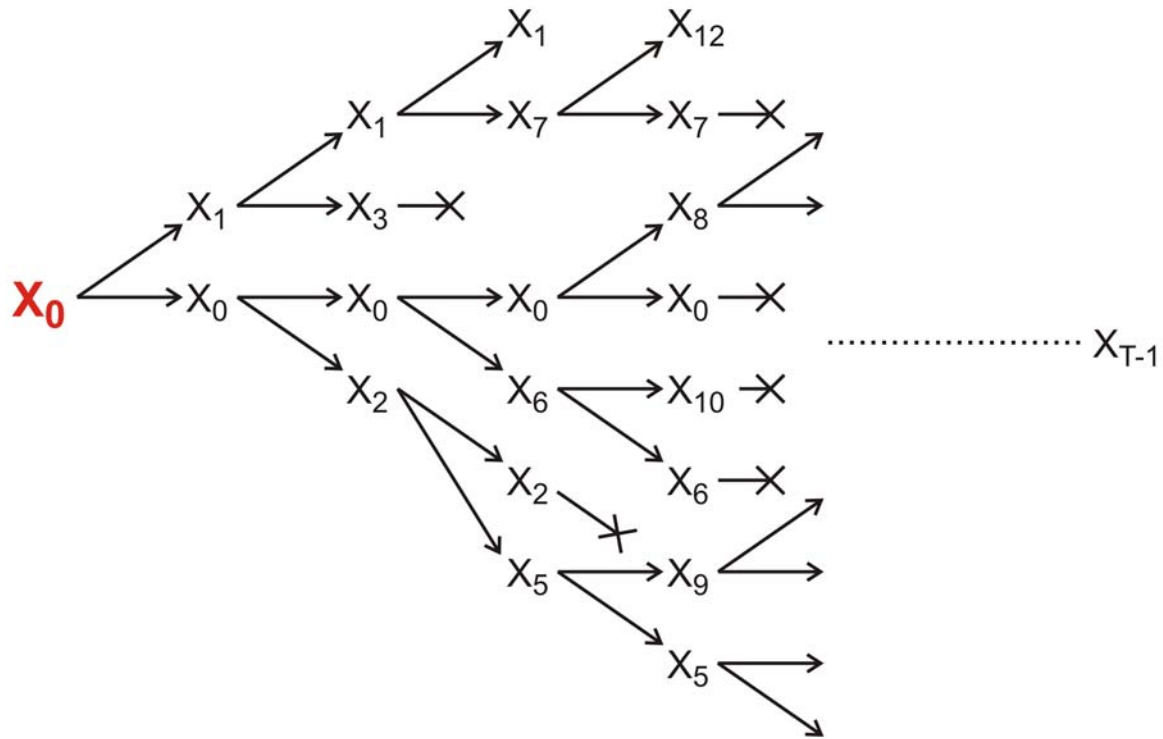




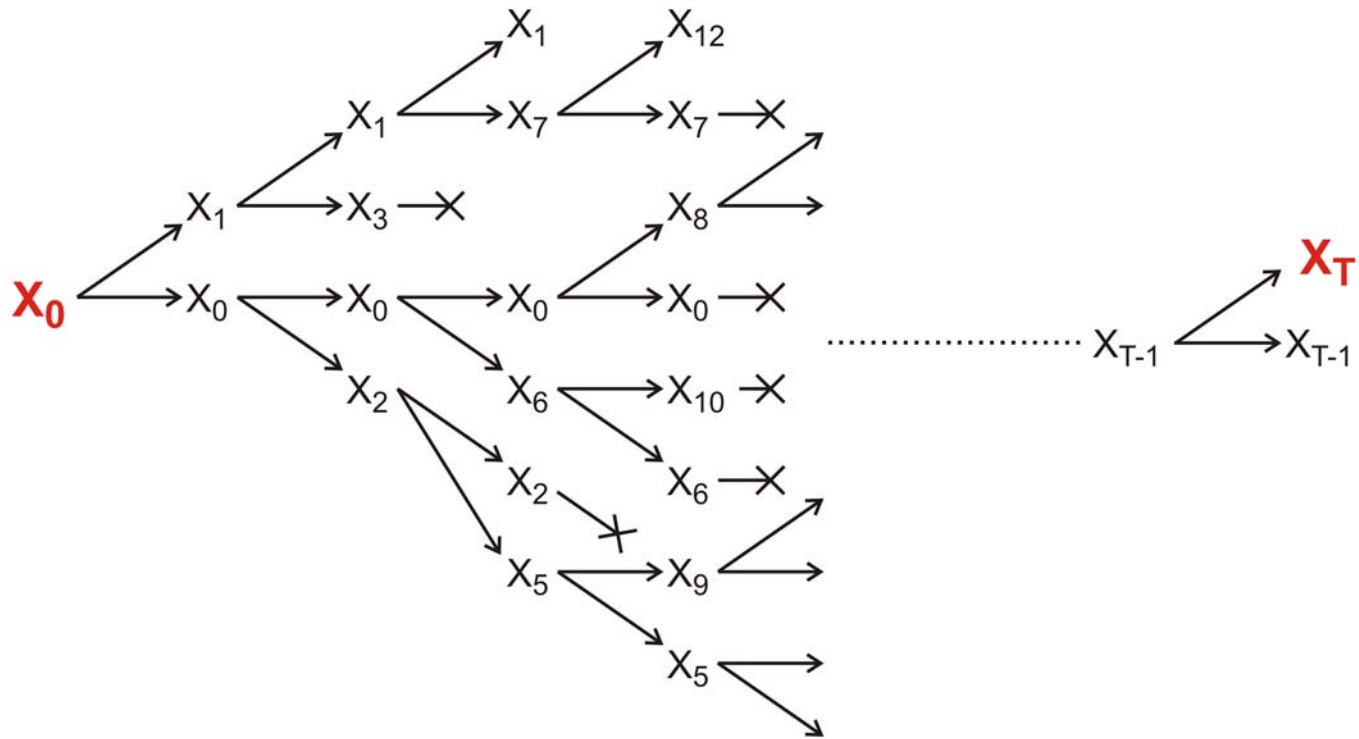
Evolution of RNA molecules as a Markov process and its analysis by means of the relay series



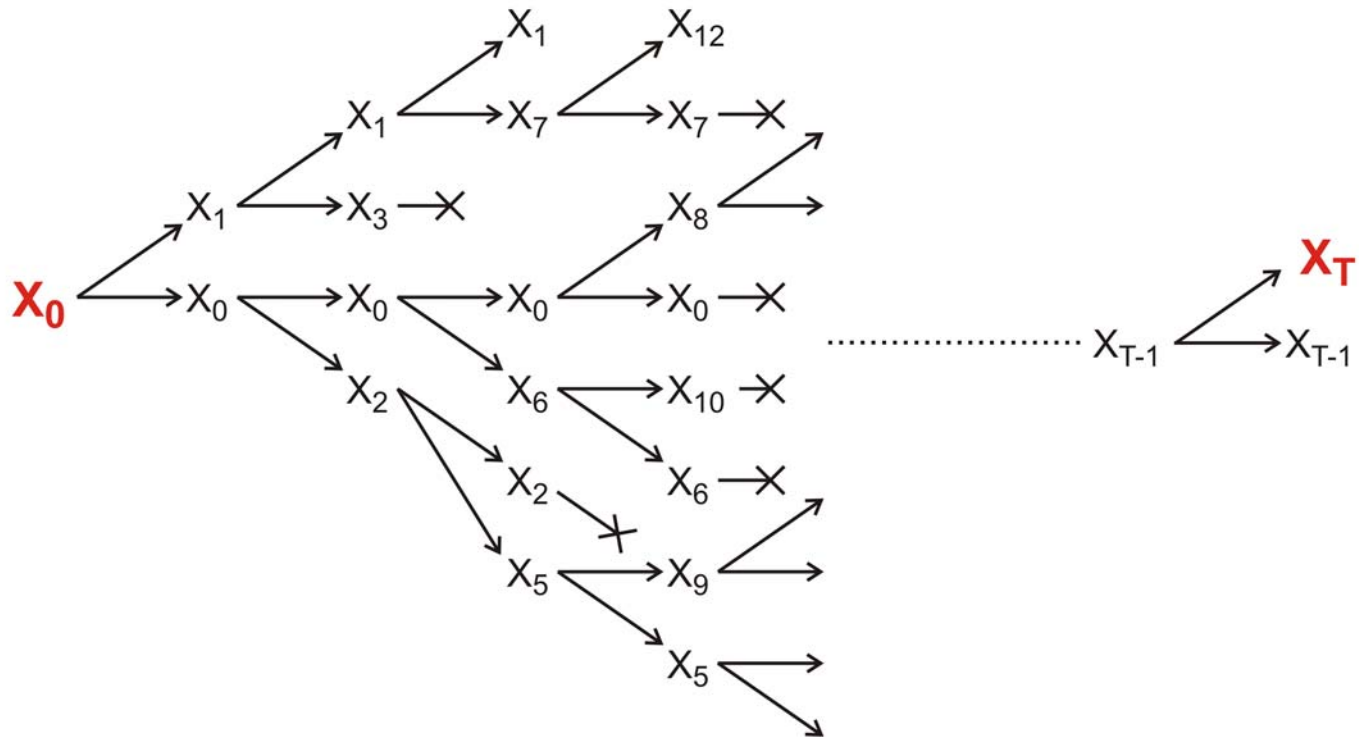
Evolution of RNA molecules as a Markov process and its analysis by means of the relay series



Evolution of RNA molecules as a Markow process and its analysis by means of the relay series

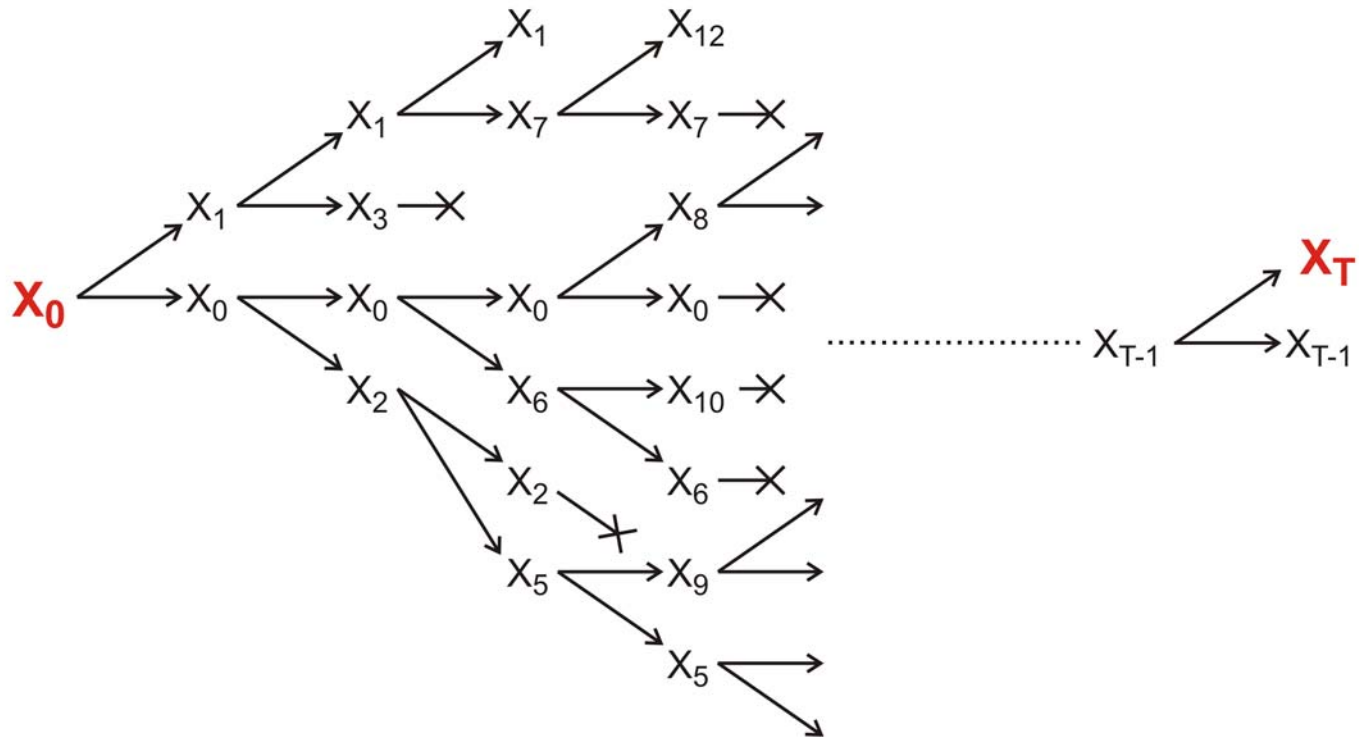


Evolution of RNA molecules as a Markow process and its analysis by means of the relay series



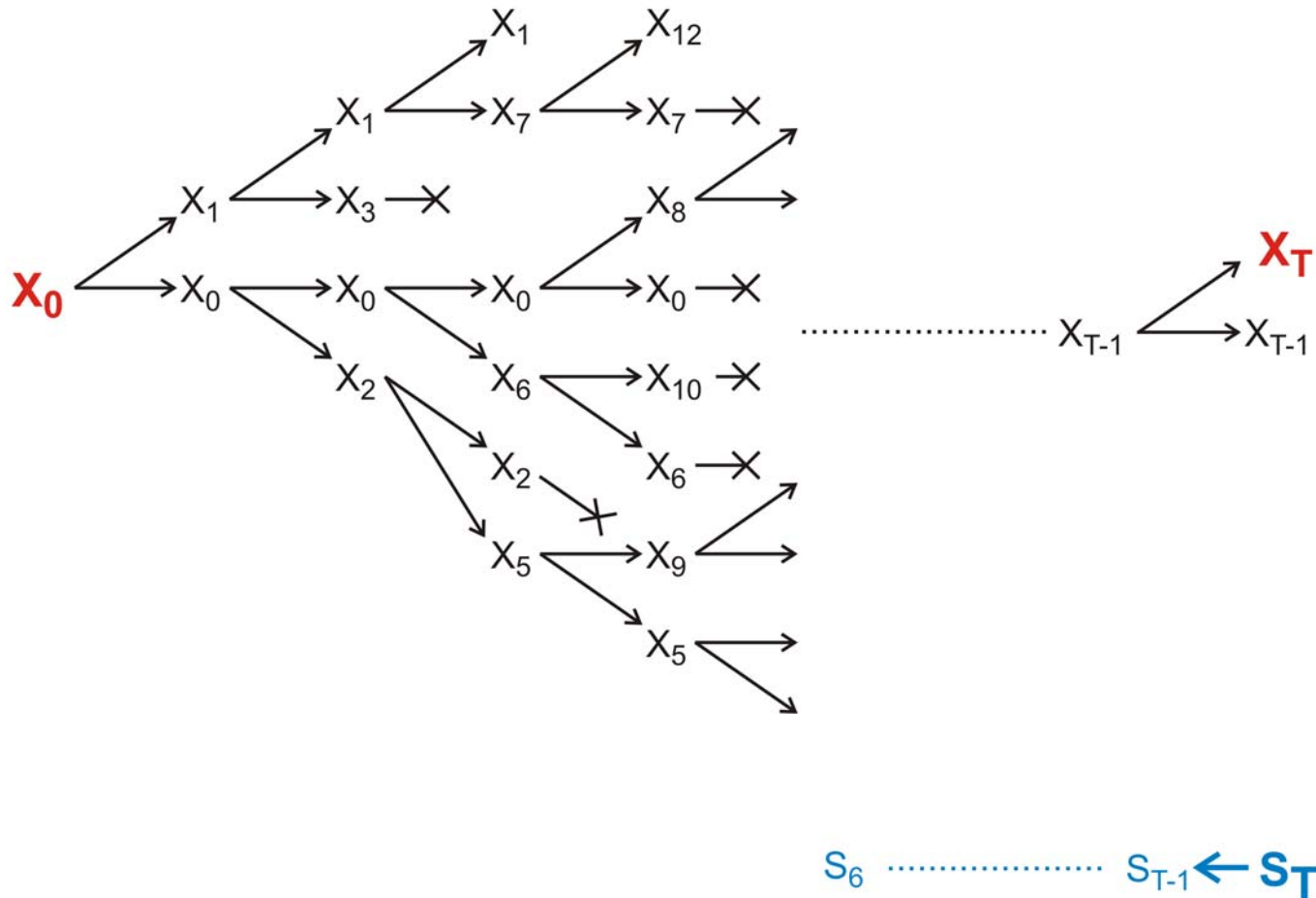
$S_T$

Evolution of RNA molecules as a Markow process and its analysis by means of the relay series

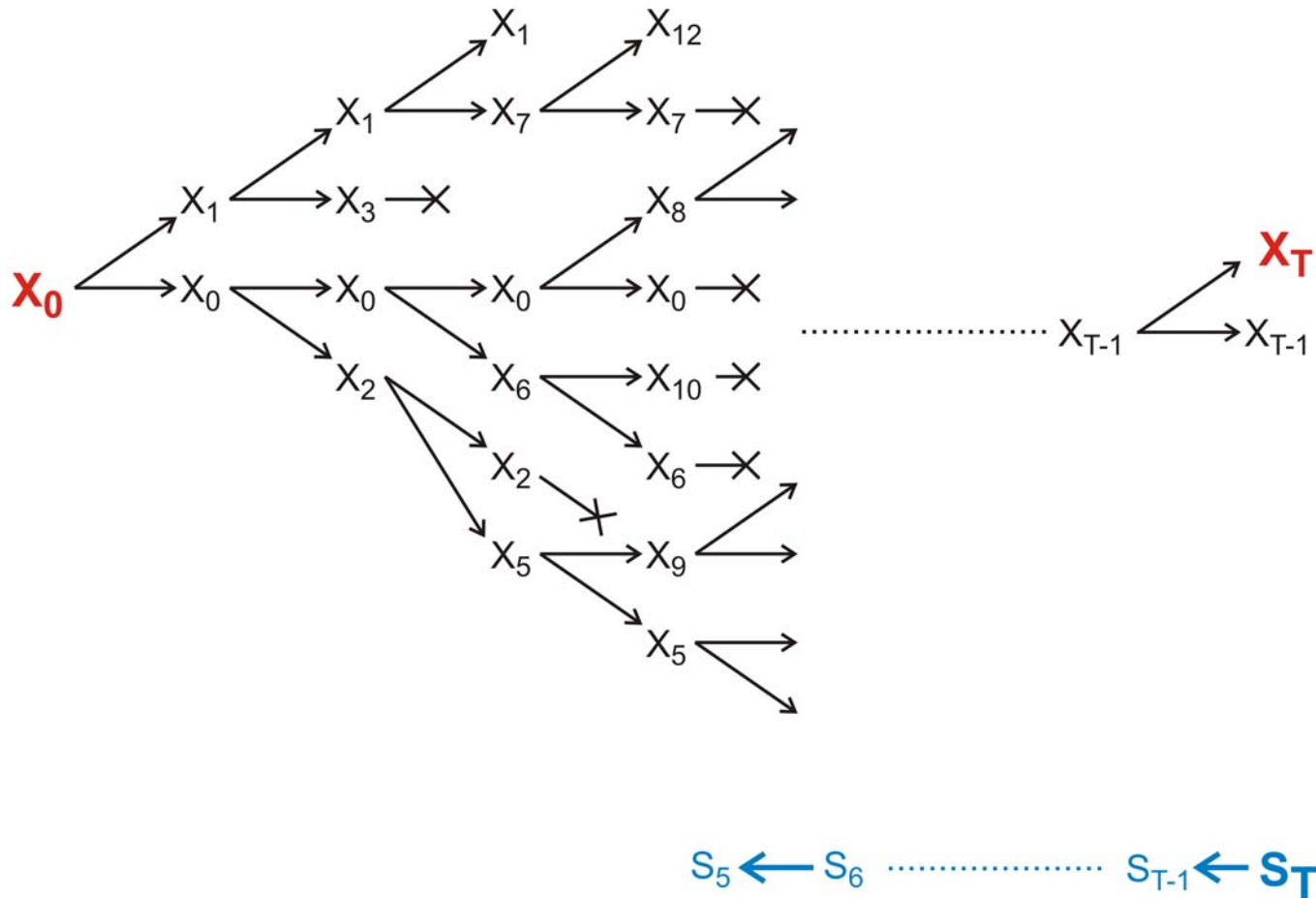


$S_{T-1} \leftarrow S_T$

Evolution of RNA molecules as a Markow process and its analysis by means of the relay series

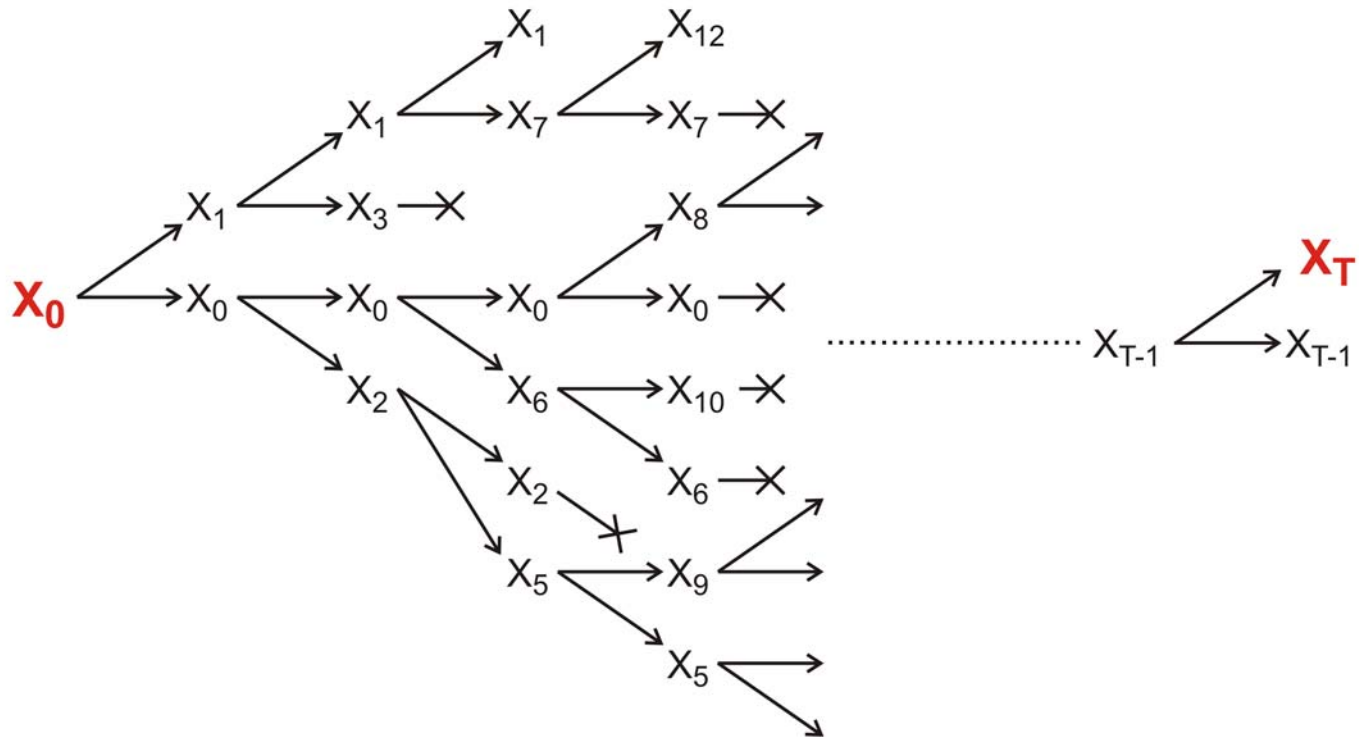


Evolution of RNA molecules as a Markow process and its analysis by means of the relay series

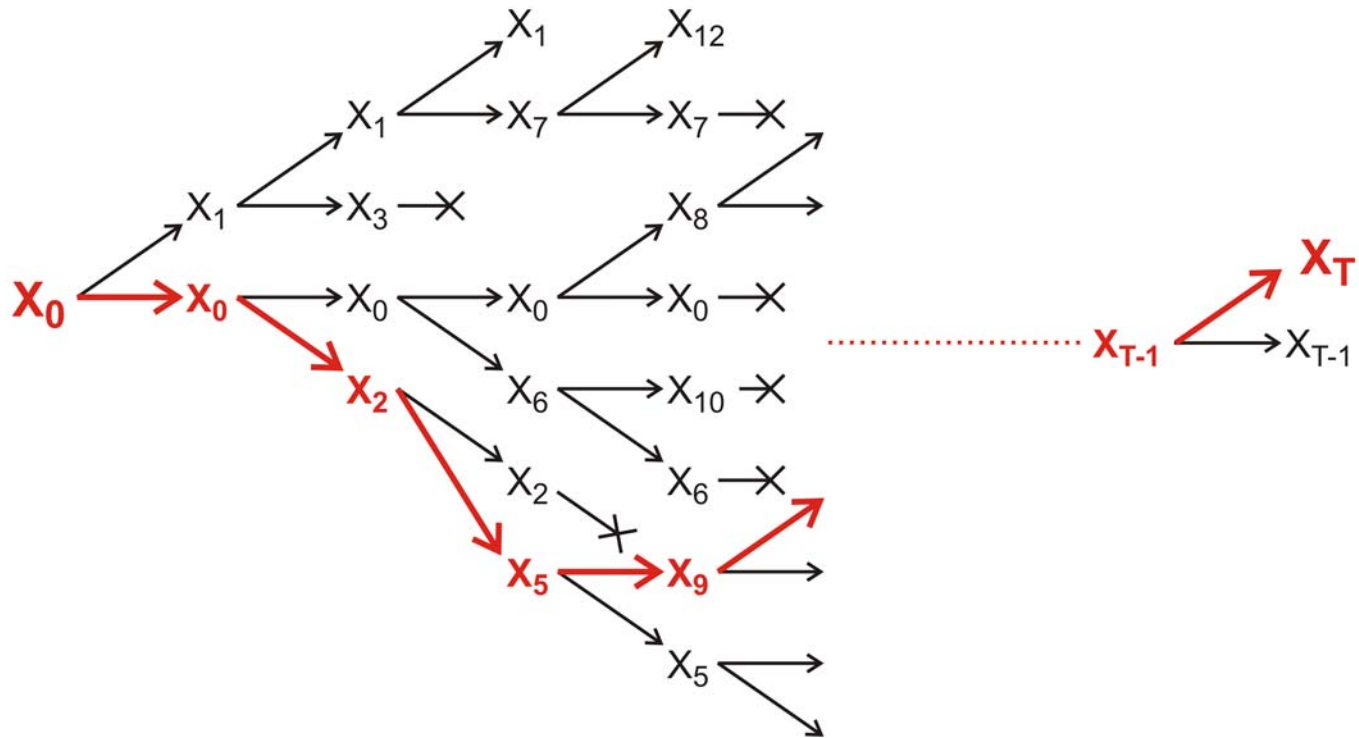


Evolution of RNA molecules as a Markow process and its analysis by means of the relay series





Evolution of RNA molecules as a Markow process and its analysis by means of the relay series



Evolution of RNA molecules as a Markow process and its analysis by means of the relay series

random individuals. The primer pair used for genomic DNA amplification is 5'-TCTCCCTGGATTCT-CATTTA-3' (forward) and 5'-TCTTTGTCTTCTGT-TGCACC-3' (reverse). Reactions were performed in 25  $\mu$ l using 1 unit of Taq DNA polymerase with each primer at 0.4  $\mu$ M, 200  $\mu$ M each dATP, dTTP, dCTP, and dGTP, and PCR buffer [10 mM Tris-HCl (pH 8.3), 50 mM KCl, 1.5 mM MgCl<sub>2</sub>] in a cycle condition of 94°C for 1 min and then 35 cycles of 94°C for 30 s, 55°C for 30 s, and 72°C for 30 s followed by 72°C for 6 min. PCR products were purified (Qiagen), digested with Xmn I, and separated in a 2% agarose gel.

32. A nonsense mutation may affect mRNA stability and result in degradation of the transcript [L. Maquat, *Am. J. Hum. Genet.* **59**, 279 (1996)].

33. Data not shown; a dot blot with poly (A)<sup>+</sup> RNA from 50 human tissues (The Human RNA Master Blot, 7770-1, Clontech Laboratories) was hybridized with a probe from exons 29 to 47 of *MYO15* using the same condition as Northern blot analysis [13].

34. Smith-Magenis syndrome (SMS) is due to deletions of 17p11.2 of various sizes, the smallest of which includes *MYO15* and perhaps 20 other genes [6]; K-S Chen, L. Potocki, J. R. Lupski, *MROD Res. Rev.* **2**, 122 (1996). *MYO15* expression is easily detected in the pituitary gland (data not shown). Haploinsufficiency for *MYO15* may explain a portion of the SMS

phenotype such as short stature. Moreover, a few SMS patients have sensorineural hearing loss, possibly because of a point mutation in *MYO15* in trans to the SMS 17p11.2 deletion.

35. R. A. Fiedel, data not shown.

36. K. B. Avraham *et al.*, *Nature Genet.* **11**, 369 (1995); X-Z. Liu *et al.*, *ibid.* **17**, 268 (1997); F. Gibson *et al.*, *Nature* **374**, 62 (1995); D. Weil *et al.*, *ibid.*, p. 60.

37. RNA was extracted from cochlea (membranous labyrinth) obtained from human fetuses at 18 to 22 weeks of development in accordance with guidelines established by the Human Research Committee at the Brigham and Women's Hospital. Only samples without evidence of degradation were pooled for poly (A)<sup>+</sup> selection over oligo(dT) columns. First-strand cDNA was prepared using an Advantage RT-for-PCR kit (Clontech Laboratories). A portion of the first-strand cDNA (4%) was amplified by PCR with Advantage cDNA polymerase mix (Clontech Laboratories) using human *MYO15*-specific oligonucleotide primers (forward, 5'-GCATGACCTGCGGGTAAT-GCG-3'; reverse, 5'-CTCAAGGCTTCTGGCATGGT-GCTCGCTGCG-3'). Cycling conditions were 40 s at 94°C, 40 s at 66°C (3 cycles), 60°C (5 cycles), and 55°C (29 cycles); and 45 s at 68°C. PCR products were visualized by ethidium bromide staining after fractionation in a 1% agarose gel. A 688-bp PCR

product is expected from amplification of the human *MYO15* cDNA. Amplification of human genomic DNA with this primer pair would result in a 2903-bp fragment.

38. We are grateful to the people of Bengkala, Bali, and the two families from India. We thank J. R. Lupski and K.-S. Chen for providing the human chromosome 17 cosmid library. For technical and computational assistance, we thank N. Dietrich, M. Ferguson, A. Gupta, E. Sorbello, R. Torzkadash, C. Varner, M. Walker, G. Bouffard, and S. Beckstrom-Stenberg (National Institutes of Health Intramural Sequencing Center). We thank J. T. Hinnant, I. N. Arhya, and S. Winata for assistance in Bali, and J. Barber, S. Sullivan, E. Green, D. Drayna, and T. Battey for helpful comments on this manuscript. Supported by the National Institute on Deafness and Other Communication Disorders (NIDCD) (Z01 DC 00335-01 and Z01 DC 00338-01 to T.B.F. and E.R.W. and R01 DC 03402 to C.G.M.), the National Institute of Child Health and Human Development (R01 HD30428 to S.A.C.) and a National Science Foundation Graduate Research Fellowship to F.J.P. This paper is dedicated to J. B. Snow Jr. on his retirement as the Director of the NIDCD.

9 March 1998; accepted 17 April 1998

## Continuity in Evolution: On the Nature of Transitions

Walter Fontana and Peter Schuster

To distinguish continuous from discontinuous evolutionary change, a relation of nearness between phenotypes is needed. Such a relation is based on the probability of one phenotype being accessible from another through changes in the genotype. This nearness relation is exemplified by calculating the shape neighborhood of a transfer RNA secondary structure and provides a characterization of discontinuous shape transformations in RNA. The simulation of replicating and mutating RNA populations under selection shows that sudden adaptive progress coincides mostly, but not always, with discontinuous shape transformations. The nature of these transformations illuminates the key role of neutral genetic drift in their realization.

A much-debated issue in evolutionary biology concerns the extent to which the history of life has proceeded gradually or has been punctuated by discontinuous transitions at the level of phenotypes (1). Our goal is to make the notion of a discontinuous transition more precise and to understand how it arises in a model of evolutionary adaptation.

We focus on the narrow domain of RNA secondary structure, which is currently the simplest computationally tractable, yet realistic phenotype (2). This choice enables the definition and exploration of concepts that may prove useful in a wider context. RNA secondary structures represent a coarse level of analysis compared with the three-dimensional structure at atomic resolution. Yet, secondary structures are empir-

ically well defined and obtain their biophysical and biochemical importance from being a scaffold for the tertiary structure. For the sake of brevity, we shall refer to secondary structures as "shapes." RNA combines in a single molecule both genotype (replicable sequence) and phenotype (selectable shape), making it ideally suited for *in vitro* evolution experiments (3, 4).

To generate evolutionary histories, we used a stochastic continuous time model of an RNA population replicating and mutating in a capacity-constrained flow reactor under selection (5, 6). In the laboratory, a goal might be to find an RNA aptamer binding specifically to a molecule (4). Although in the experiment the evolutionary end product was unknown, we thought of its shape as being specified implicitly by the imposed selection criterion. Because our intent is to study evolutionary histories rather than end products, we defined a target shape in advance and assumed the replication rate of a sequence to be a function of

the similarity between its shape and the target. An actual situation may involve more than one best shape, but this does not affect our conclusions.

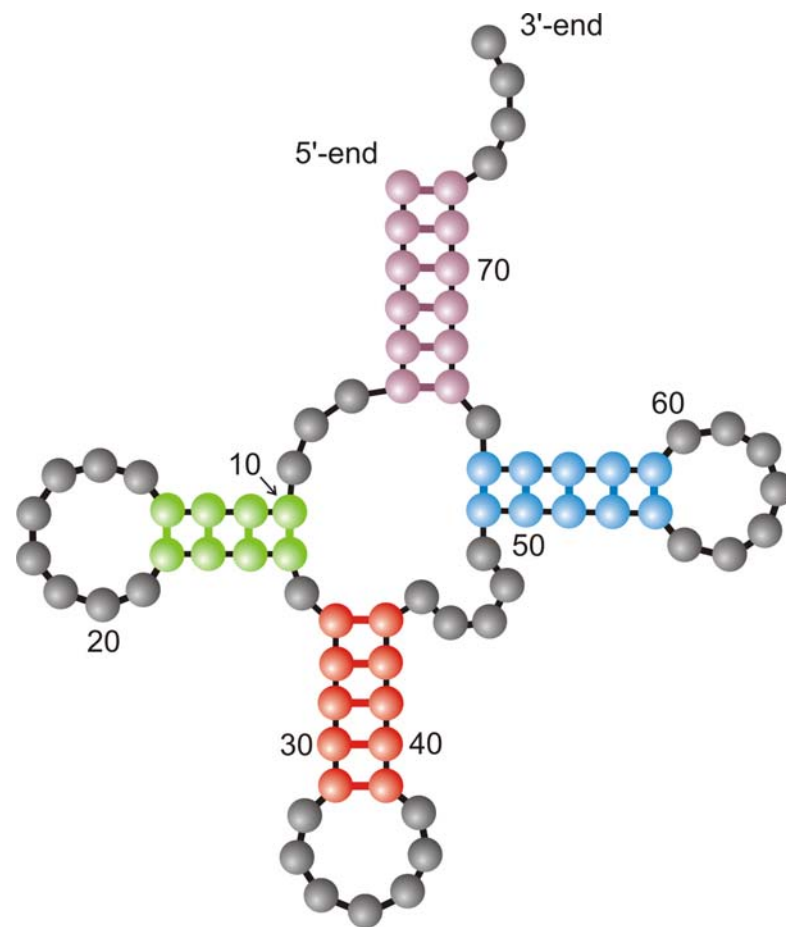
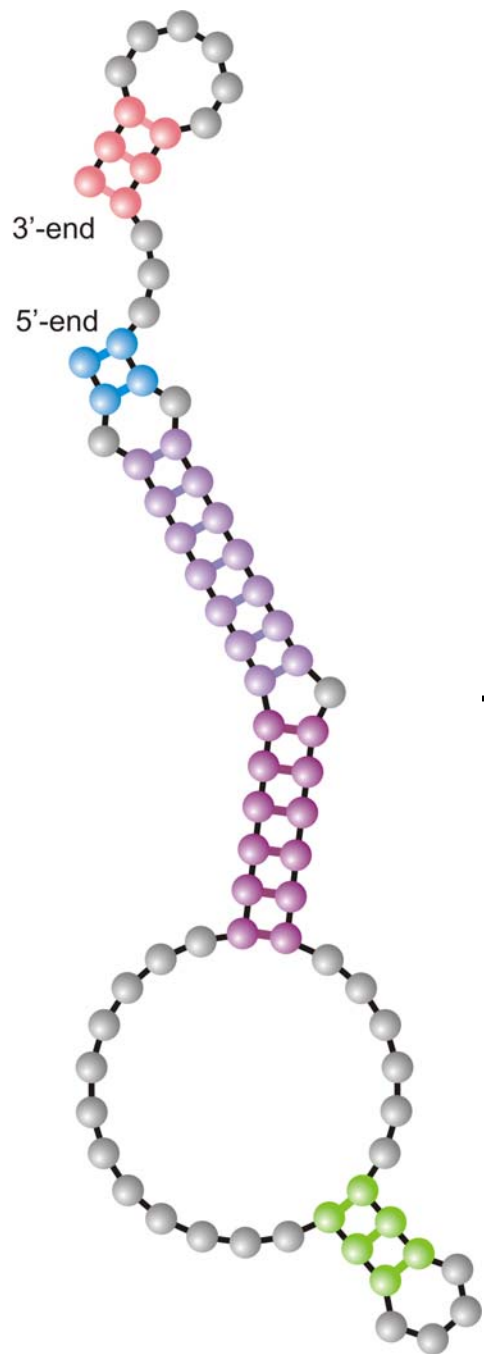
An instance representing in its qualitative features all the simulations we performed is shown in Fig. 1A. Starting with identical sequences folding into a random shape, the simulation was stopped when the population became dominated by the target, here a canonical tRNA shape. The black curve traces the average distance to the target (inversely related to fitness) in the population against time. Aside from a short initial phase, the entire history is dominated by steps, that is, flat periods of no apparent adaptive progress, interrupted by sudden approaches toward the target structure (7). However, the dominant shapes in the population not only change at these marked events but undergo several fitness-neutral transformations during the periods of no apparent progress. Although discontinuities in the fitness trace are evident, it is entirely unclear when and on the basis of what the series of successive phenotypes itself can be called continuous or discontinuous.

A set of entities is organized into a (topological) space by assigning to each entity a system of neighborhoods. In the present case, there are two kinds of entities: sequences and shapes, which are related by a thermodynamic folding procedure. The set of possible sequences (of fixed length) is naturally organized into a space because point mutations induce a canonical neighborhood. The neighborhood of a sequence consists of all its one-error mutants. The problem is how to organize the set of possible shapes into a space. The issue arises because, in contrast to sequences, there are

## Evolution *in silico*

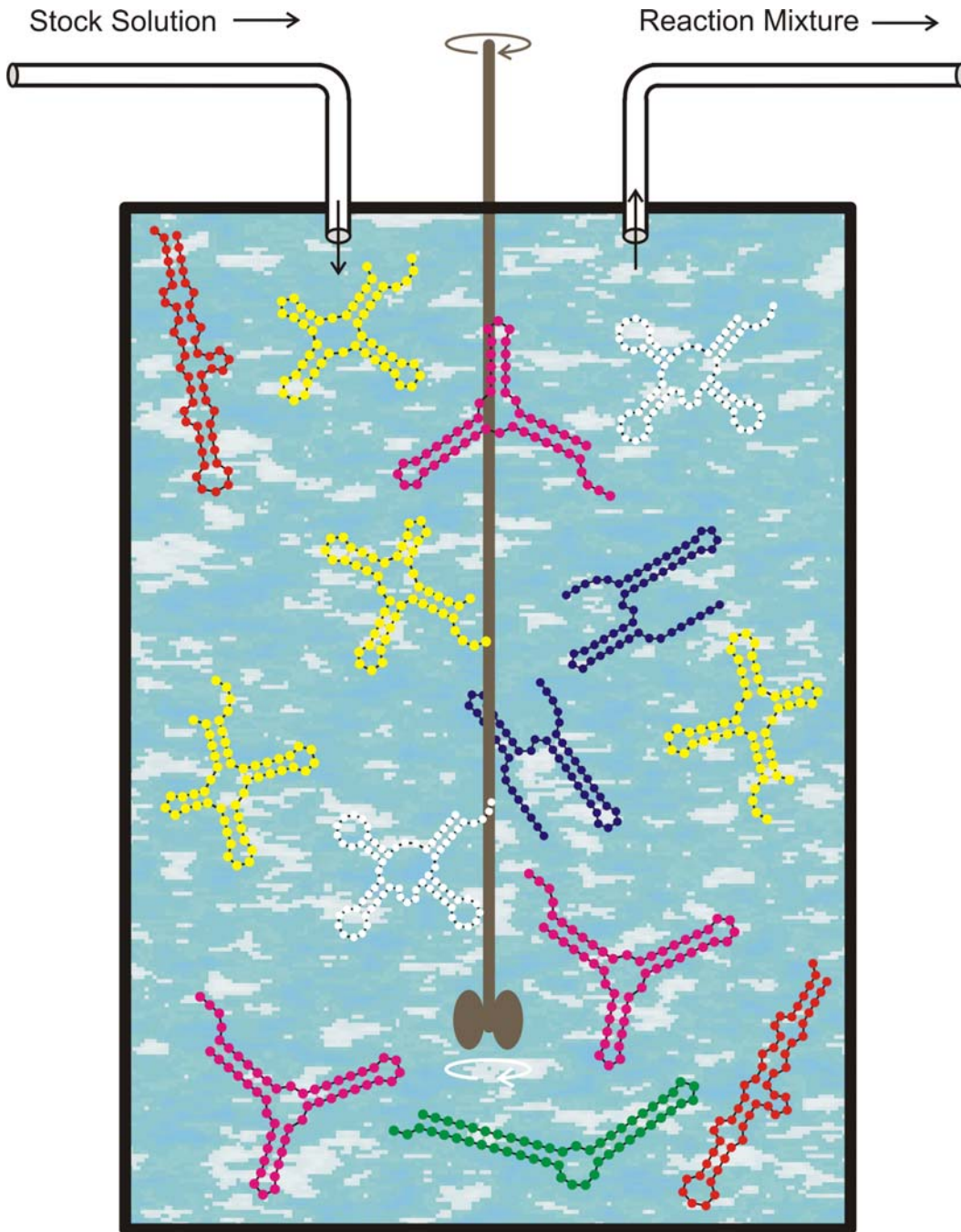
W. Fontana, P. Schuster,  
*Science* **280** (1998), 1451-1455

Institut für Theoretische Chemie, Universität Wien, Währingerstrasse 17, A-1090 Wien, Austria, Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA, and International Institute for Applied Systems Analysis (IIASA), A-2361 Laxenburg, Austria.



Structure of  
randomly chosen  
initial sequence

Phenylalanyl-tRNA as  
target structure



## Replication rate constant

(Fitness):

$$f_k = \gamma / [\alpha + \Delta d_S^{(k)}]$$

$$\Delta d_S^{(k)} = d_H(S_k, S_\tau)$$

**Selection pressure:**

The population size,

$N = \#$  RNA molecules,

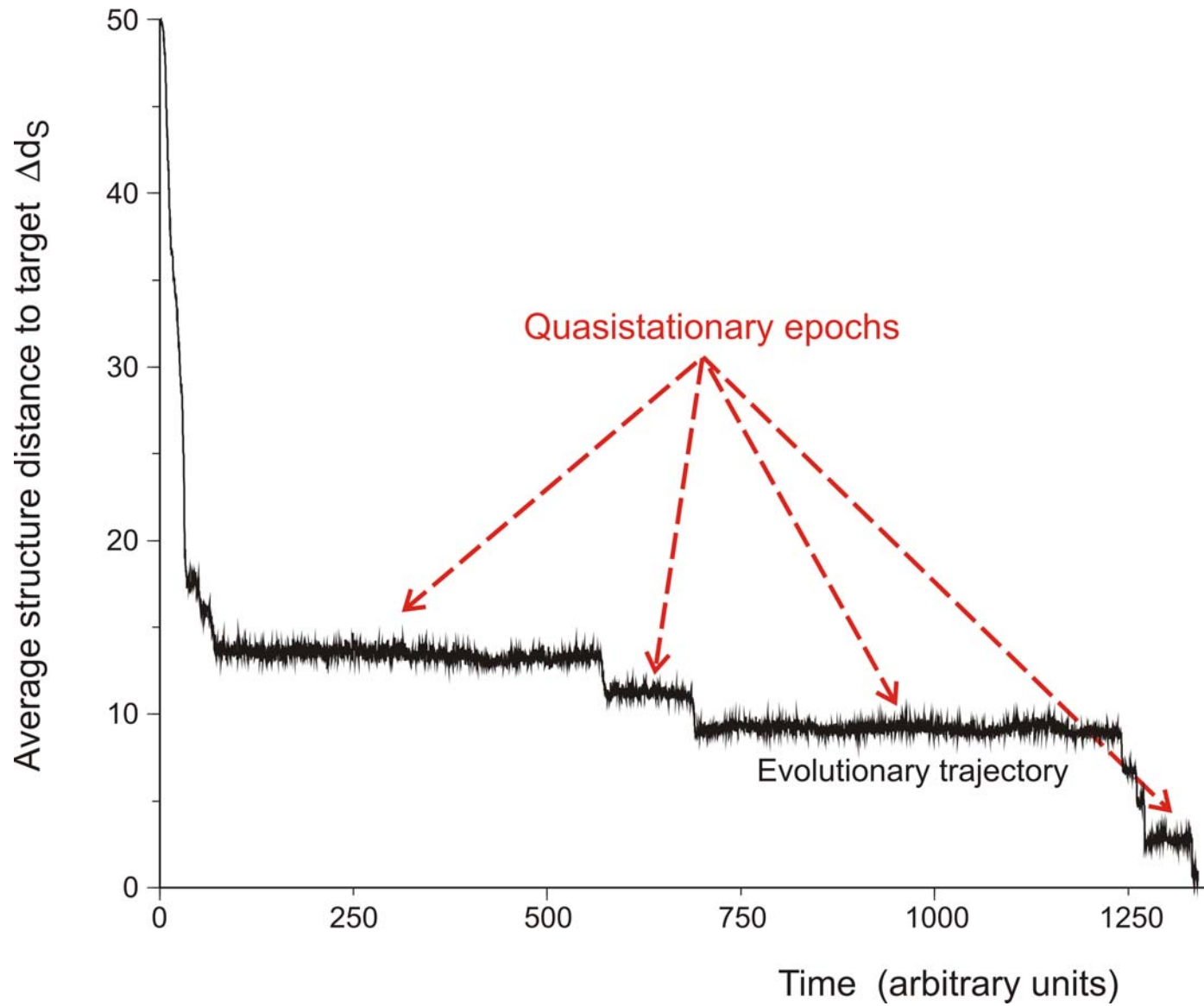
is determined by the flux:

$$N(t) \approx \bar{N} \pm \sqrt{\bar{N}}$$

**Mutation rate:**

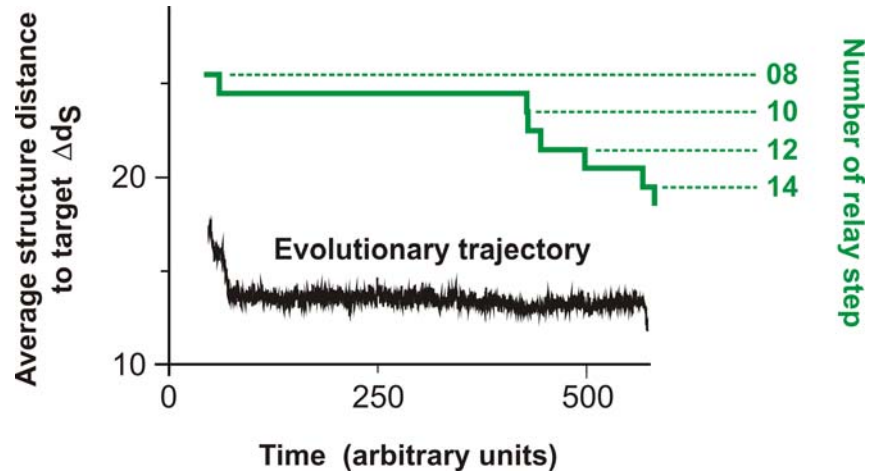
$$p = 0.001 / \text{Nucleotide} \times \text{Replication}$$

The flow reactor as a device for studying the evolution of molecules *in vitro* and *in silico*.



*In silico* optimization in the flow reactor: Evolutionary Trajectory

**28 neutral point mutations** during a long quasi-stationary epoch



```

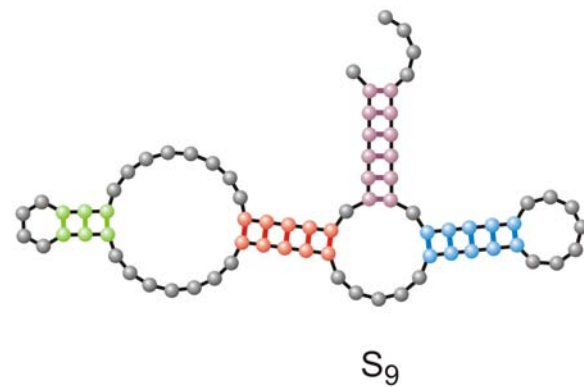
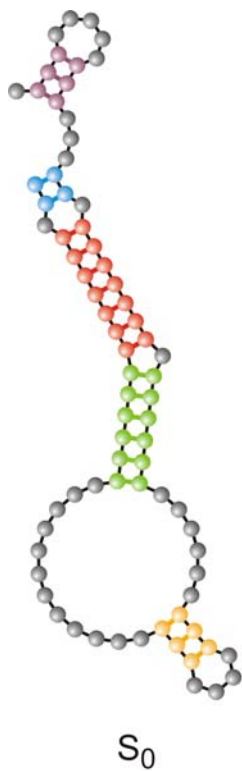
entry  GGUAUGGGCGUUGAAUAGUAGGGUUUAAACCAAUCGGCAACGAUCUCGUGUGCGCAUUUCAUAUCCCGUACAGAA
8      .(((((((((((((. . . . . (((. . . . .)))) . . . . .)))))) . . . . .(((((. . . . .))))))))) . . . .
exit   GGUAUGGGCGUUGAAUAUAGGGUUUAAACCAAUCGGCCAACGAUCUCGUGUGCGCAUUUCAUAUCCAUAACAGAA
entry  GGUAUGGGCGUUGAAUAAUAGGGUUUAAACCAAUCGGCCAACGAUCUCGUGUGCGCAUUUCAUAUACCAUACAGAA
9      .((((((. ((((. . . . . (((. . . . .)))) . . . . .)))) . . . . .(((((. . . . .))))).)))) . . . .
exit   UGGAUGGACGUUGAAUAAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACACCGUCCCAAG
entry  UGGAUGGACGUUGAAUAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACACCGUCCCAAG
10     .(((((. ((((. . . . . (((. . . . .)))) . . . . .)))) . . . . .(((((. . . . .))))).)))) . . . .
exit   UGGAUGGACGUUGAAUAAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACACCGUCCCAAG
  
```

**Transition inducing point mutations**  
change the molecular structure

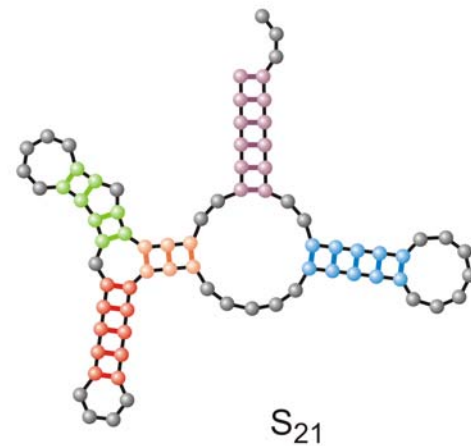
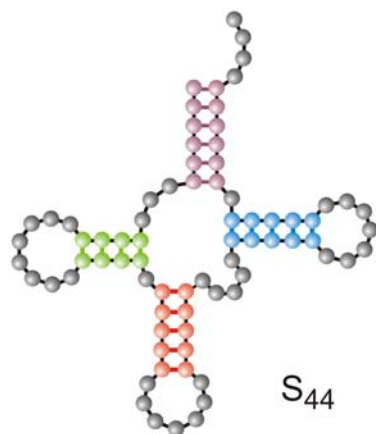
**Neutral point mutations** leave the  
molecular structure unchanged

Neutral genotype evolution during phenotypic stasis

Randomly chosen  
initial structure



Phenylalanyl-tRNA  
as target structure

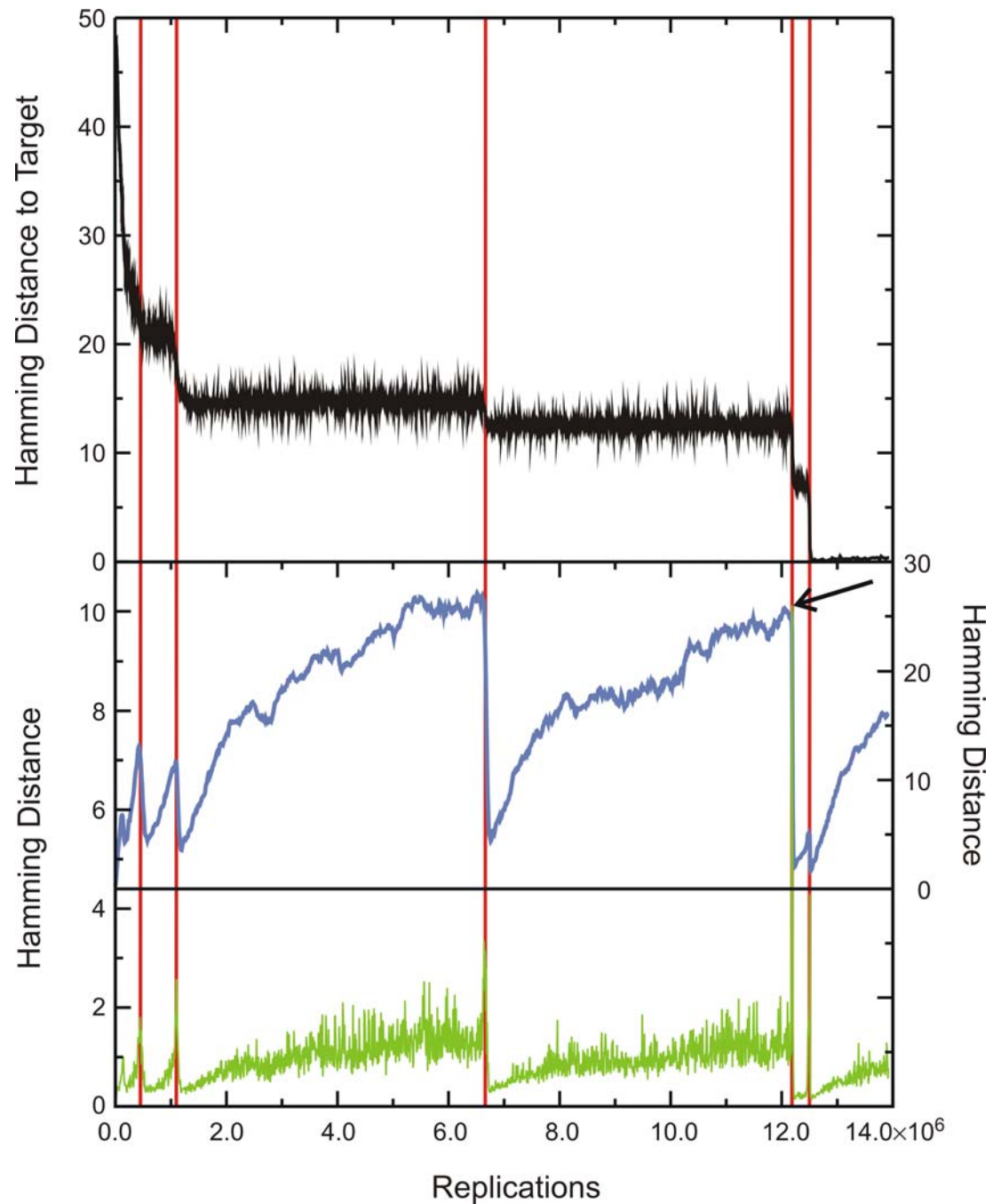


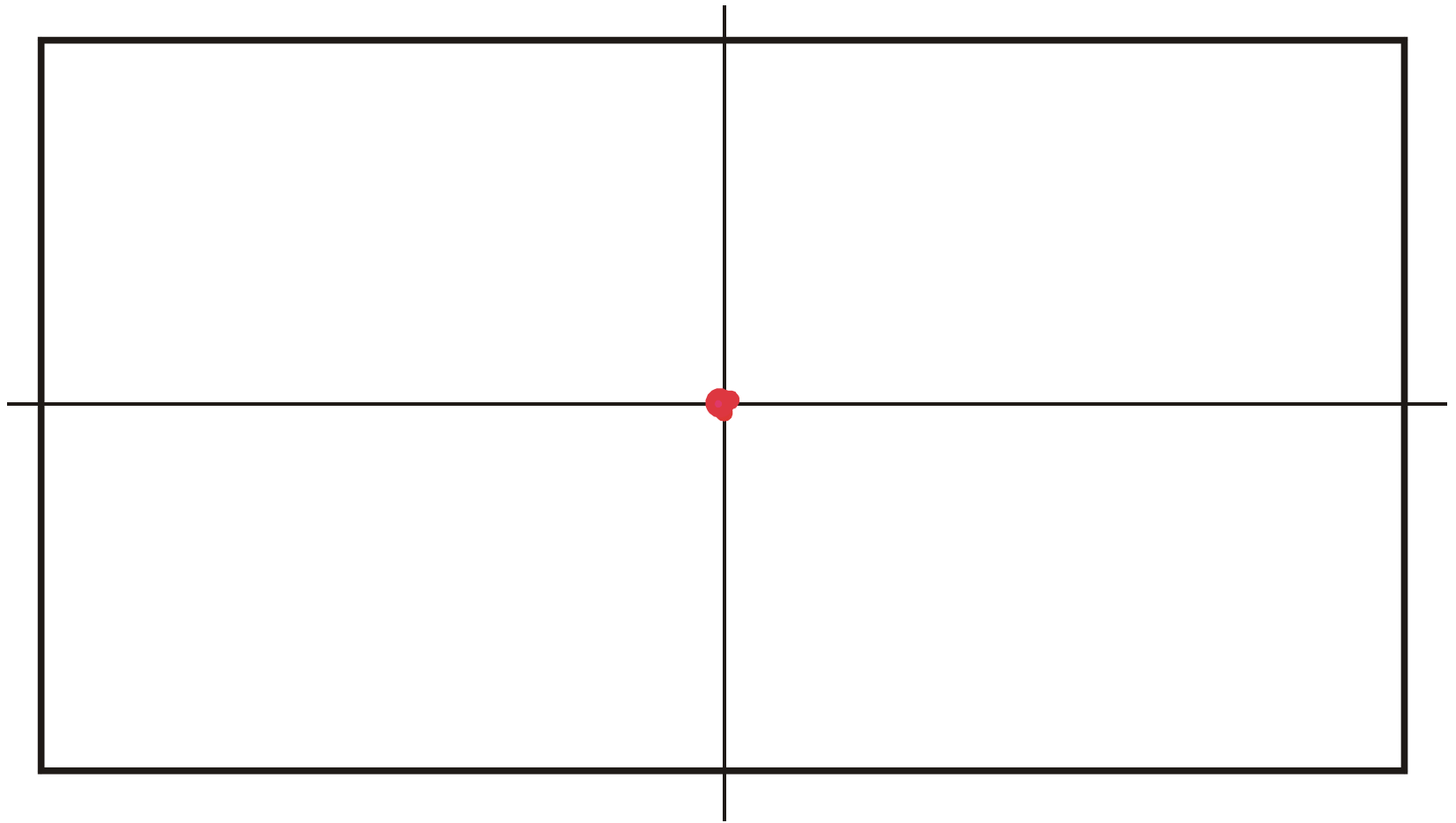


Evolutionary trajectory

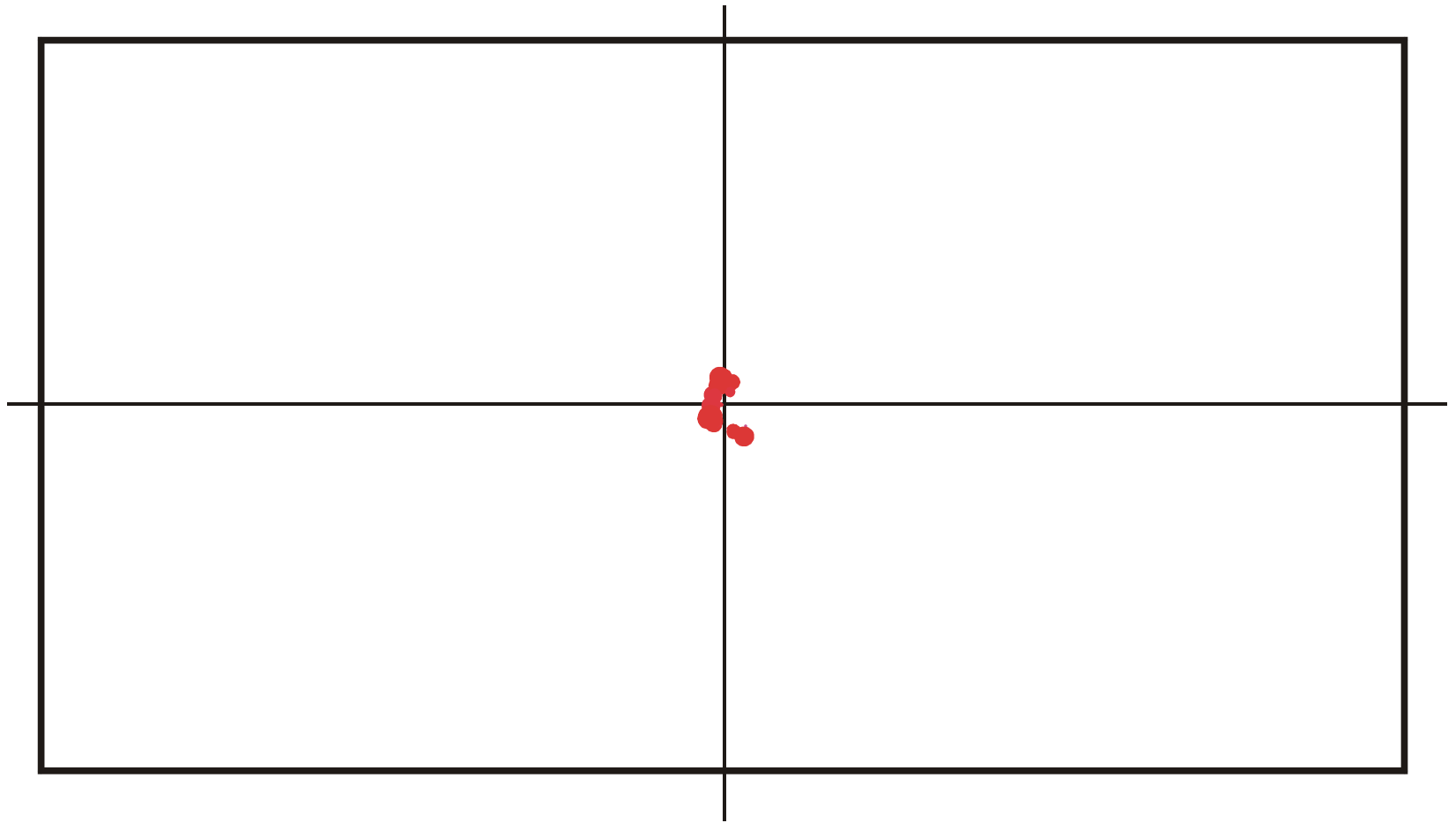
Spreading of the population on neutral networks

Drift of the population center in sequence space

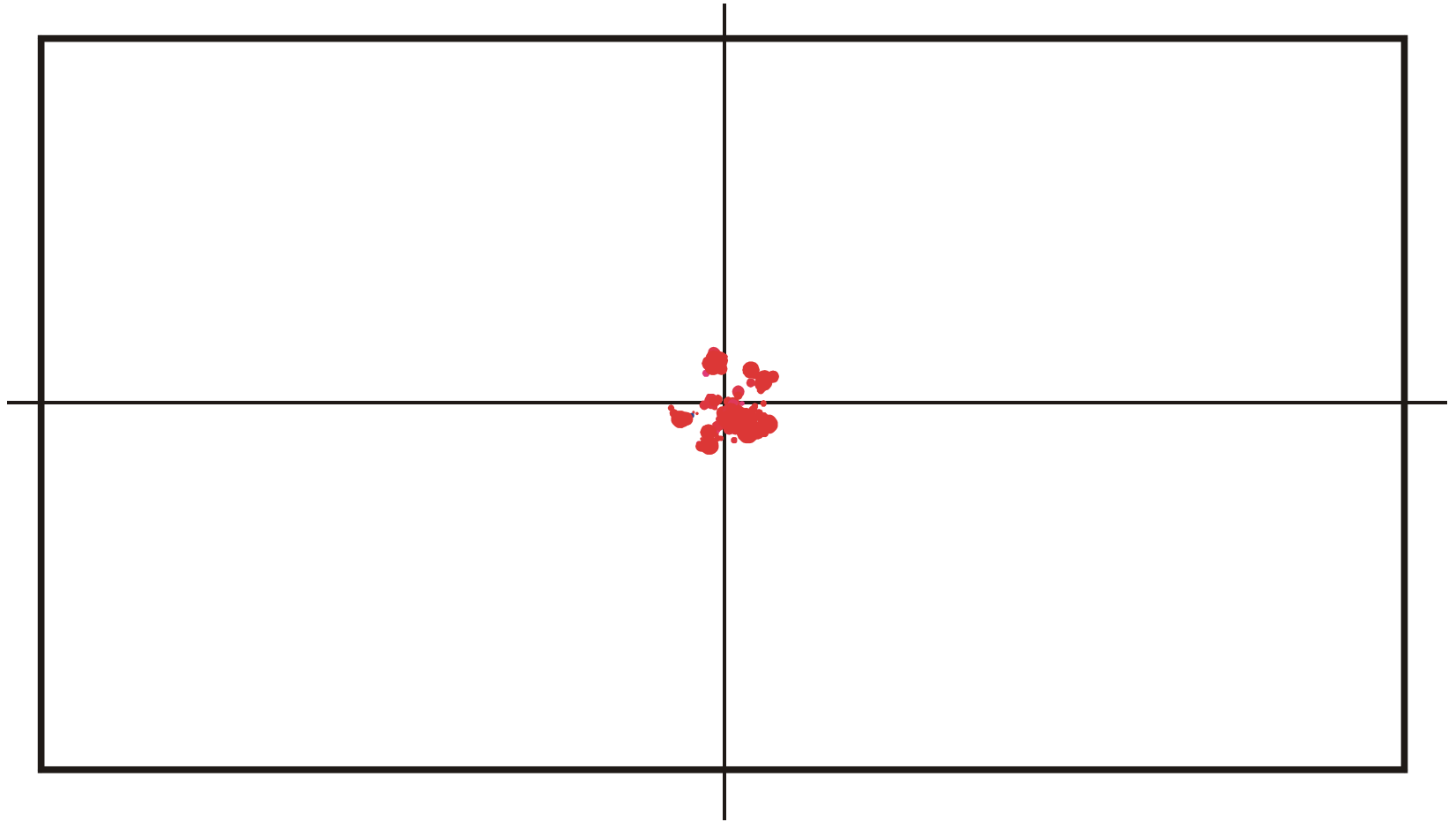




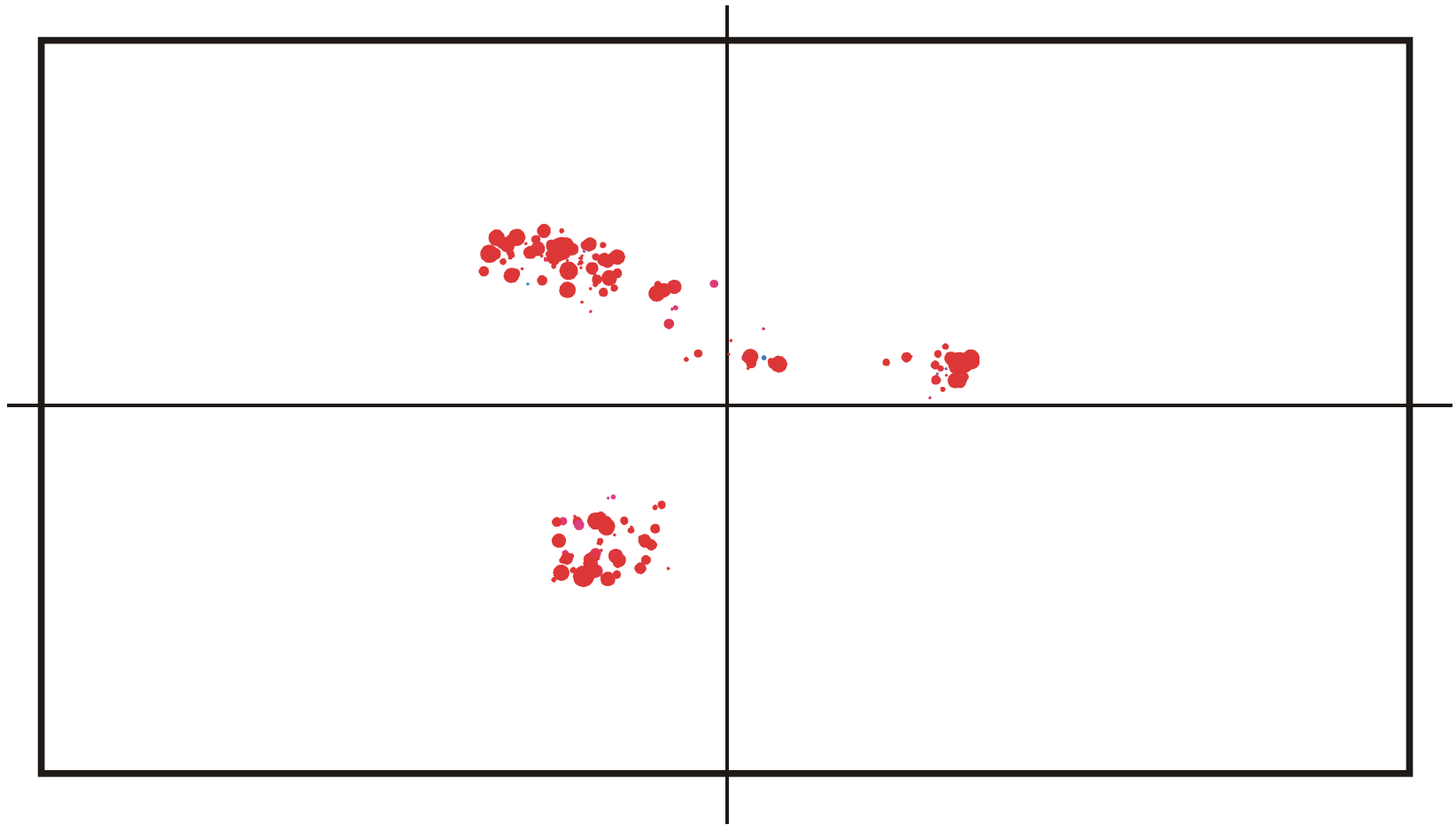
Spreading and evolution of a population on a neutral network:  $t = 150$



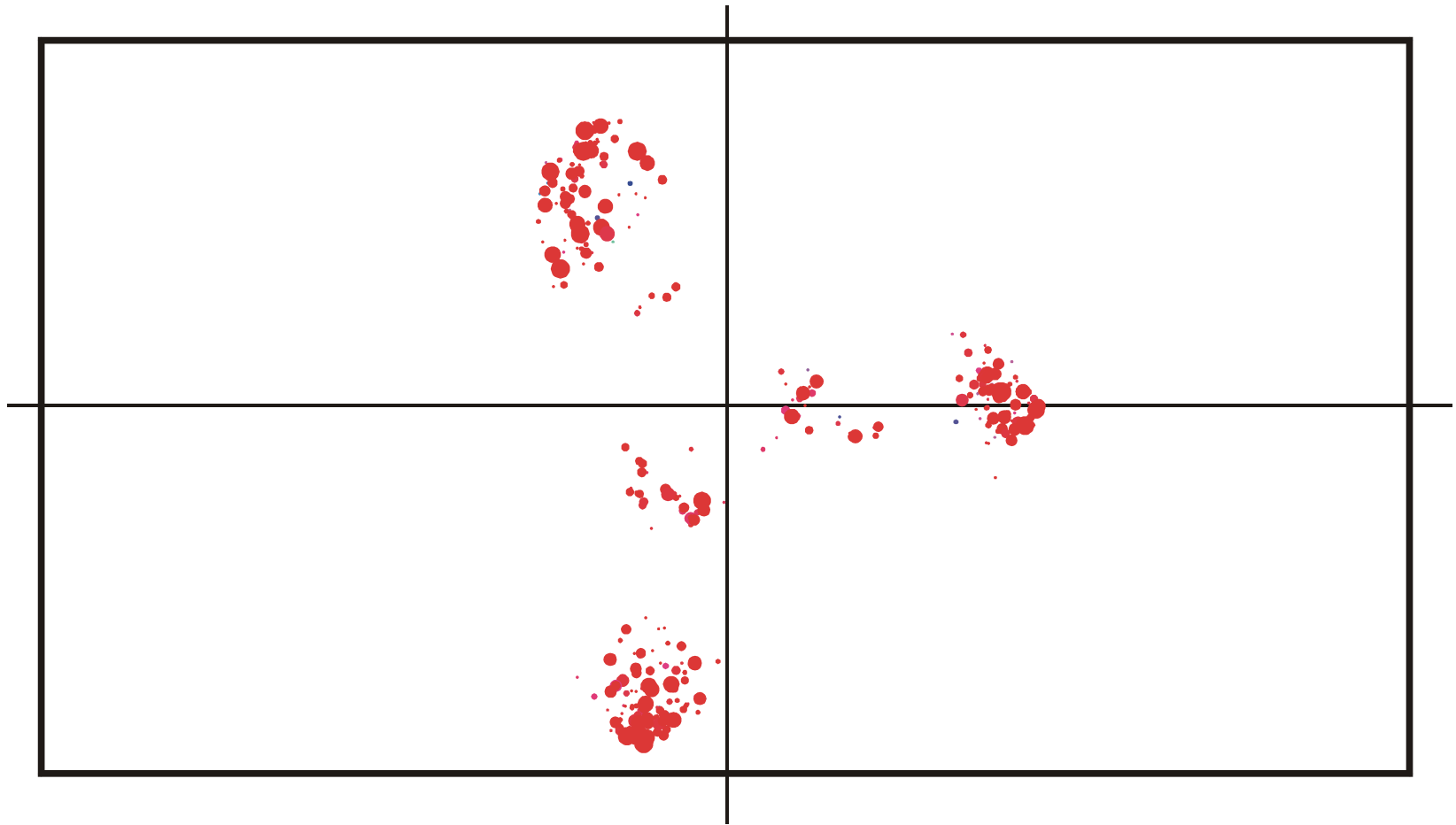
Spreading and evolution of a population on a neutral network :  $t = 170$



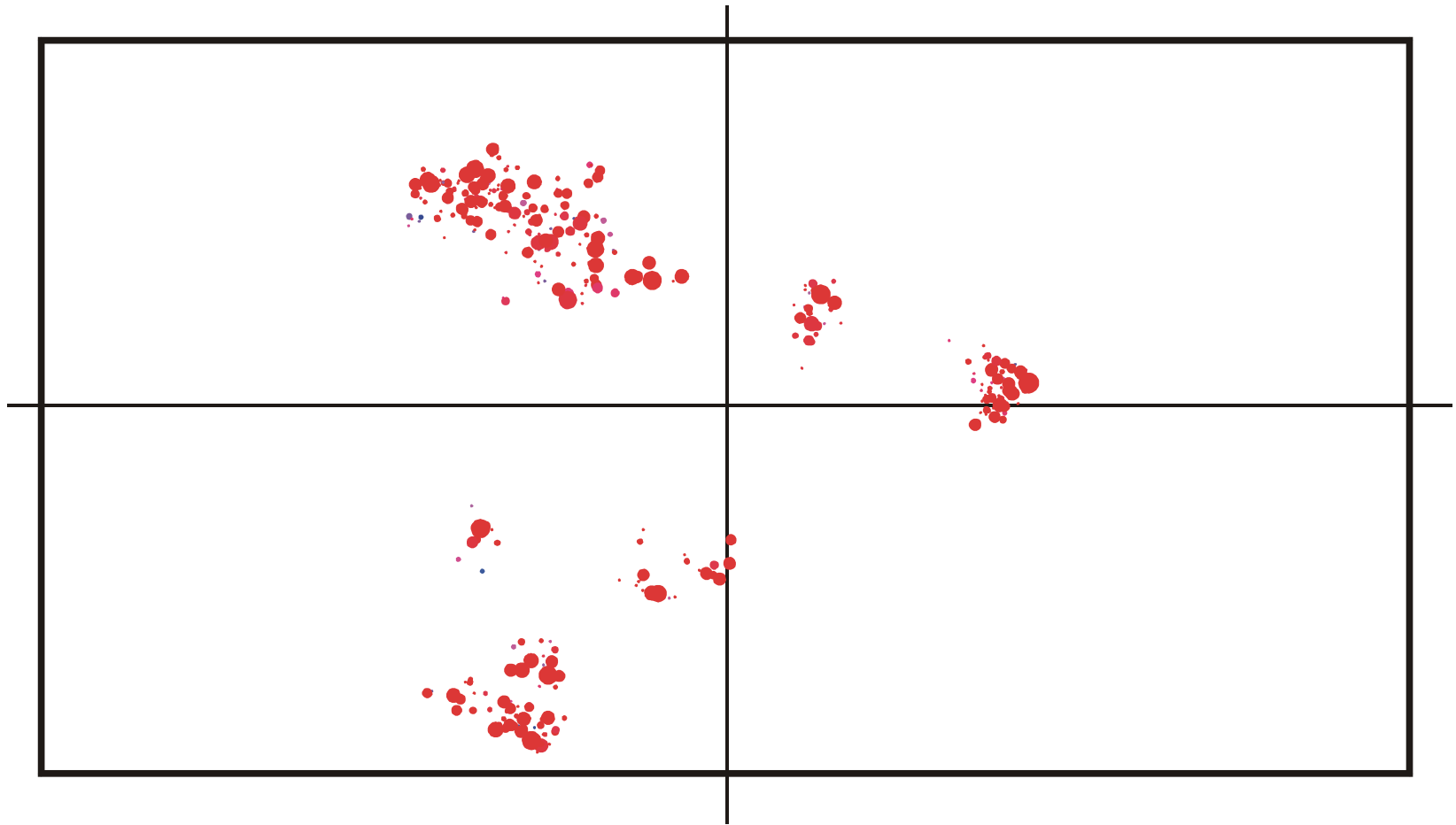
Spreading and evolution of a population on a neutral network :  $t = 200$



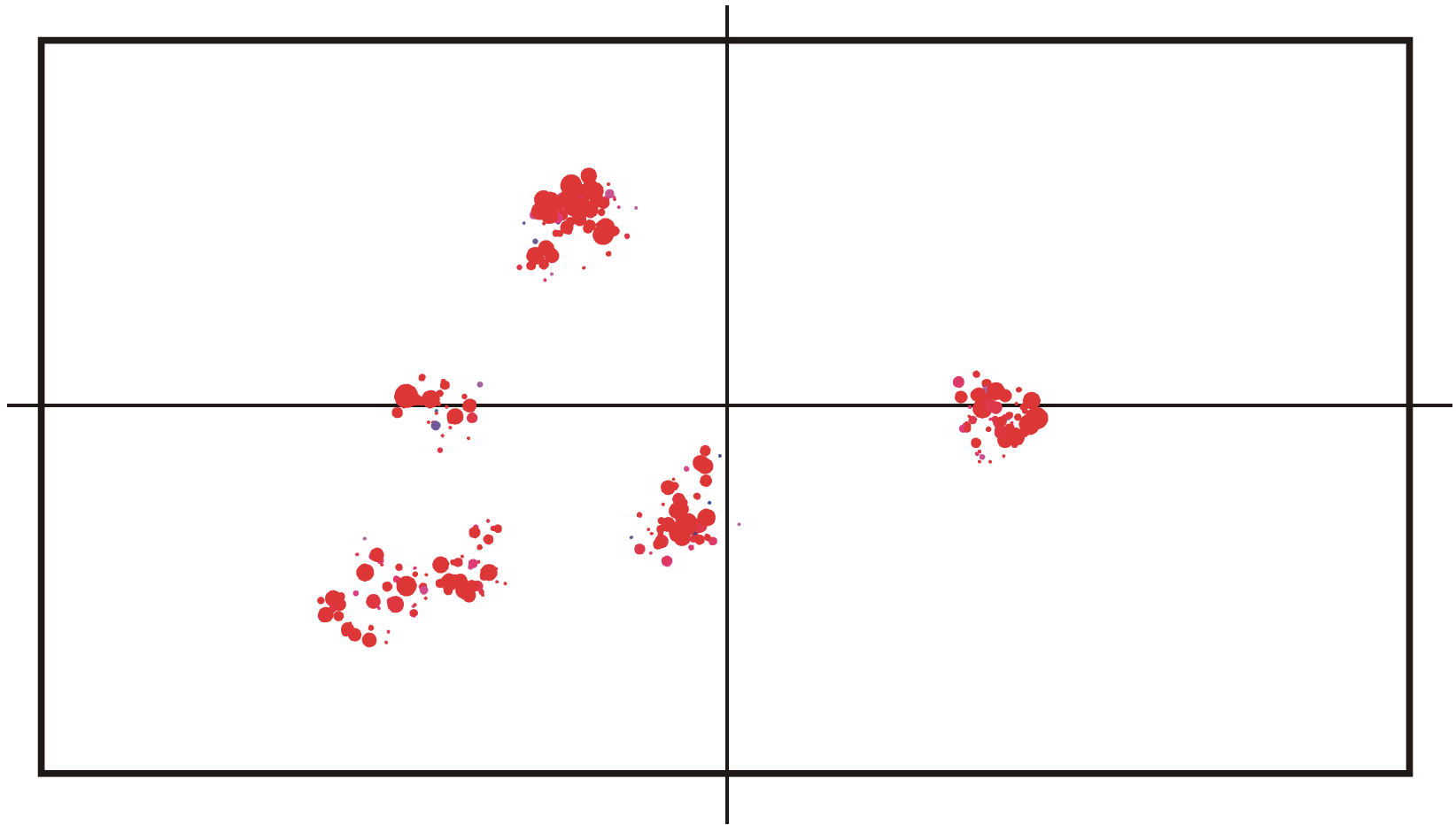
Spreading and evolution of a population on a neutral network :  $t = 350$



Spreading and evolution of a population on a neutral network :  $t = 500$

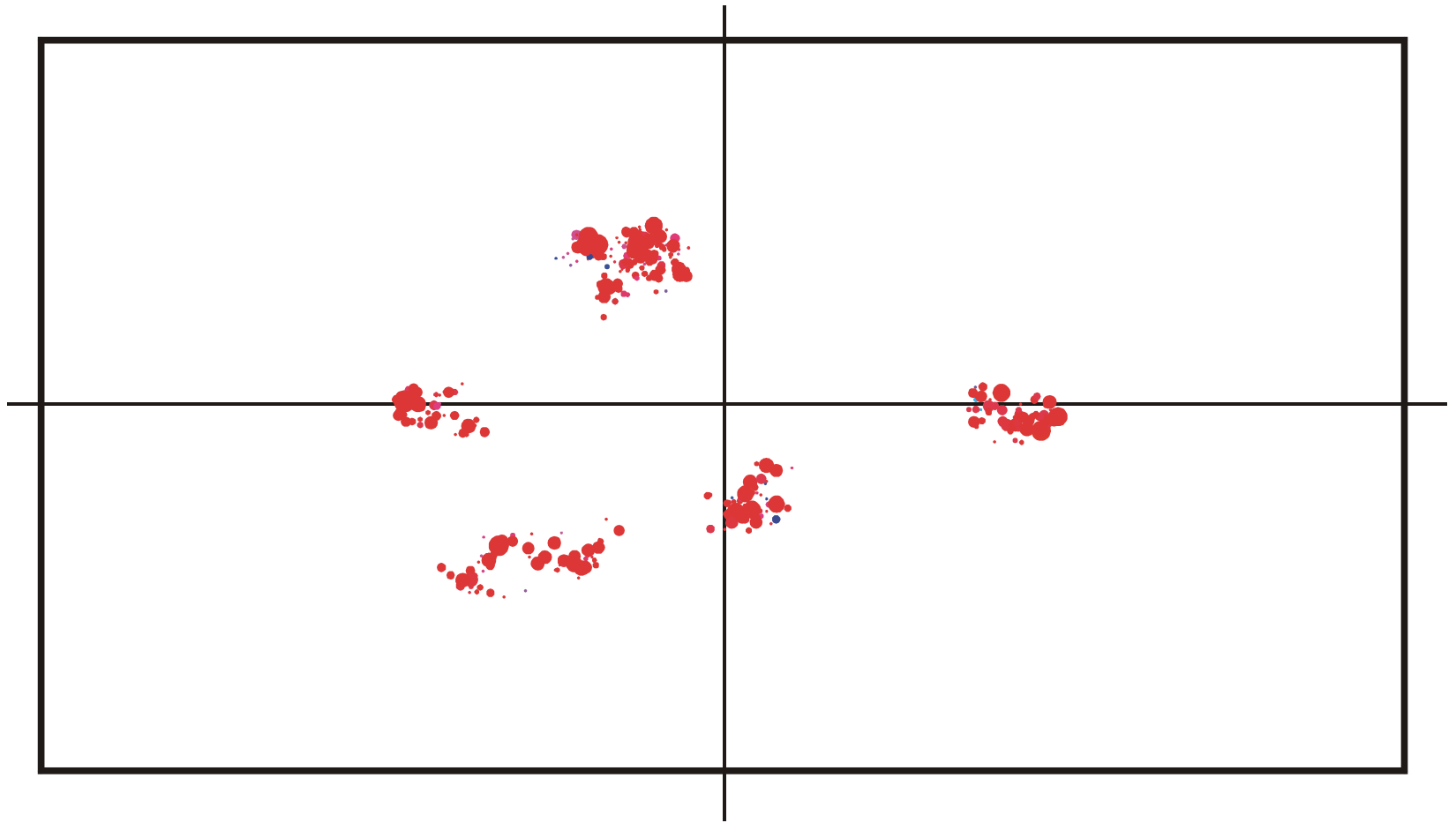


Spreading and evolution of a population on a neutral network :  $t = 650$

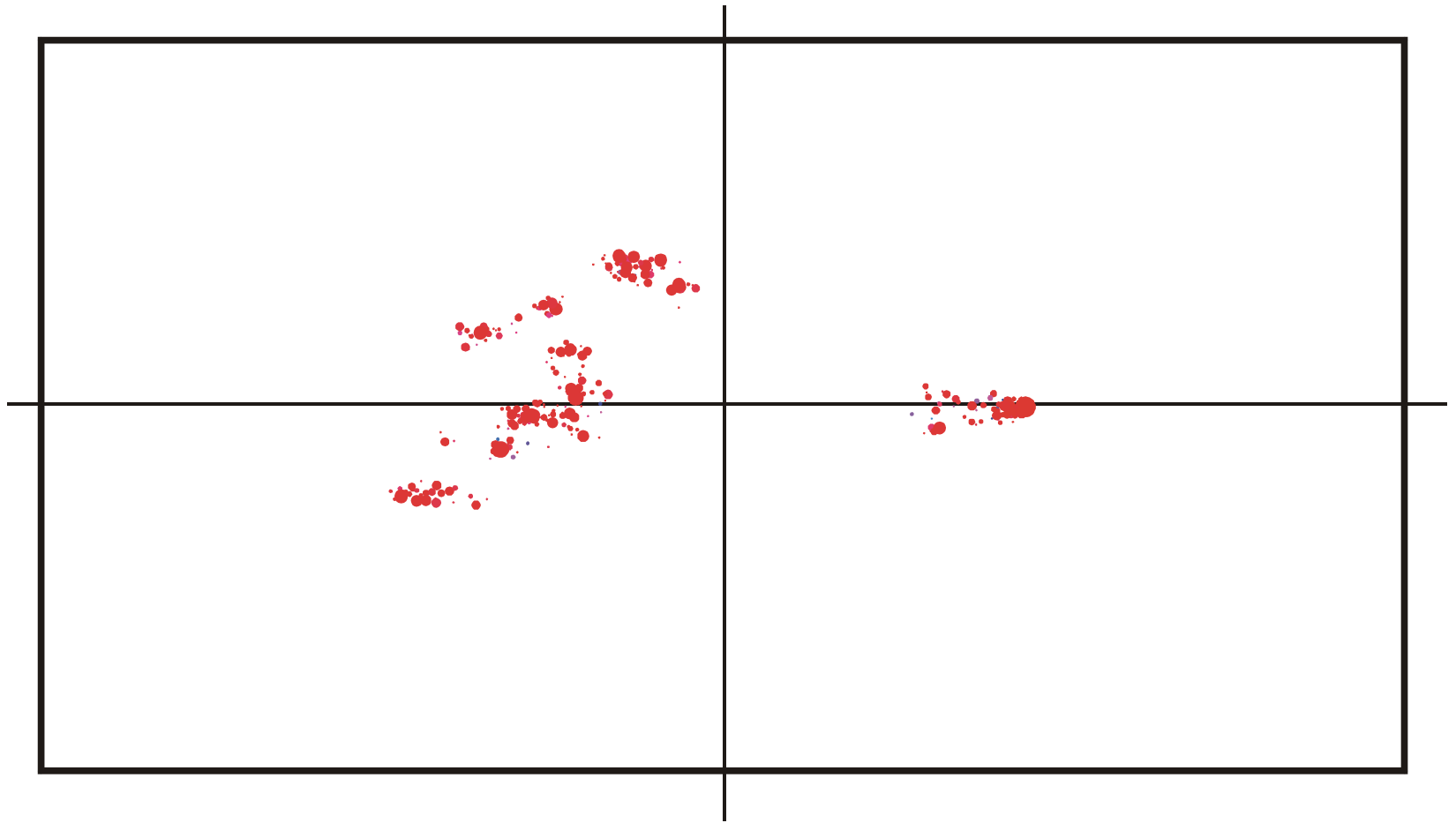


Spreading and evolution of a population on a neutral network :  $t = 820$

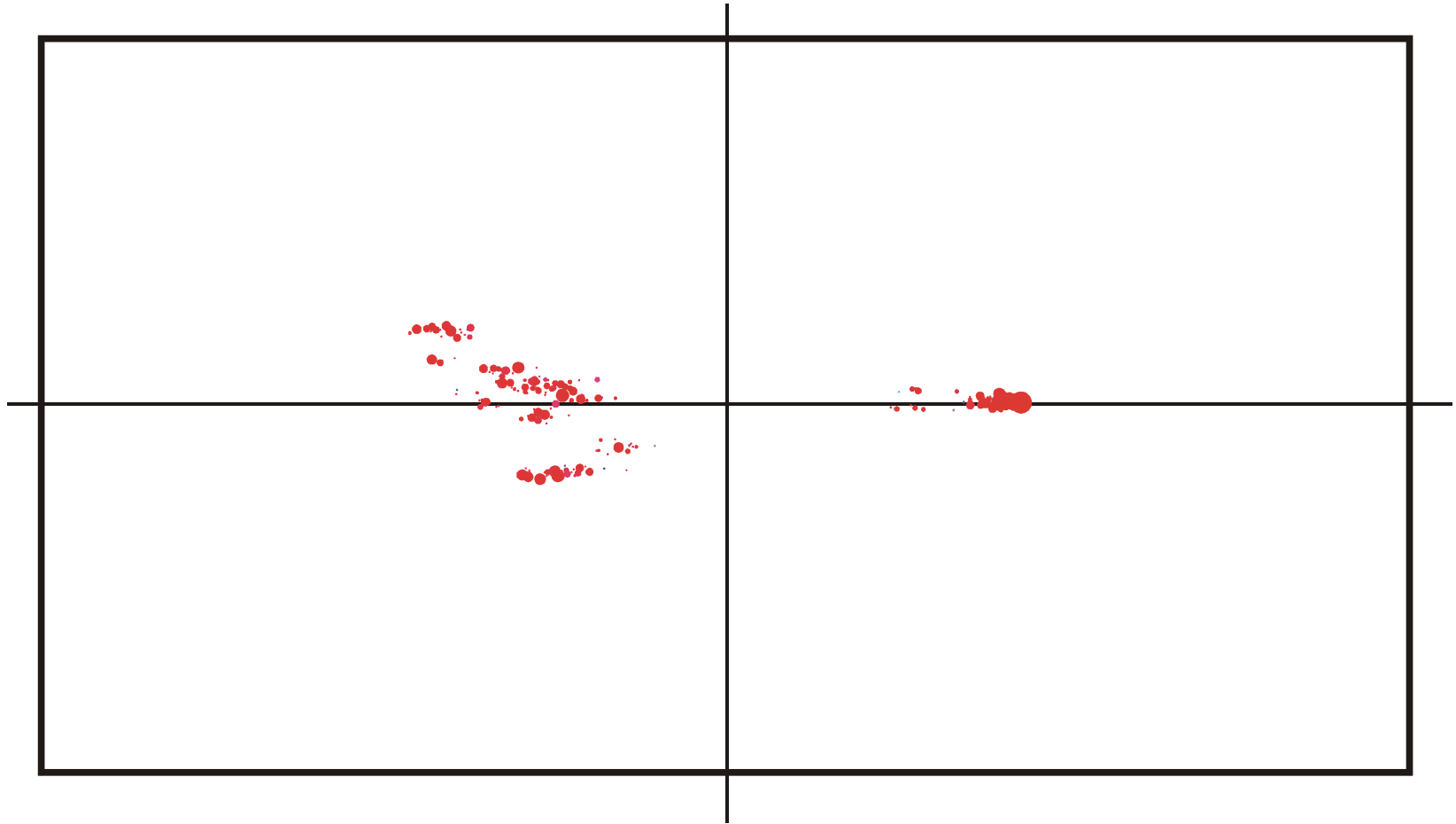




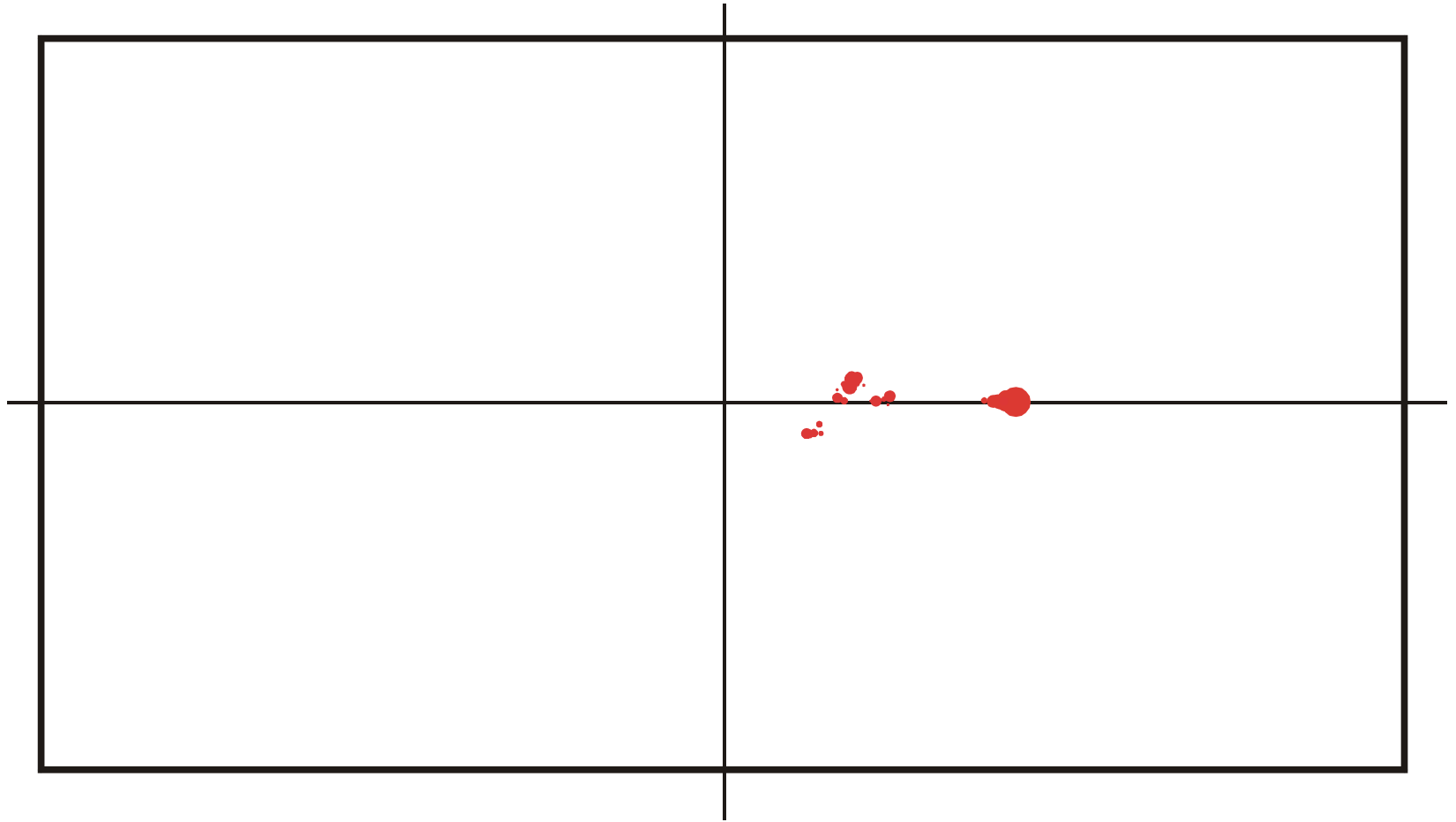
Spreading and evolution of a population on a neutral network :  $t = 825$



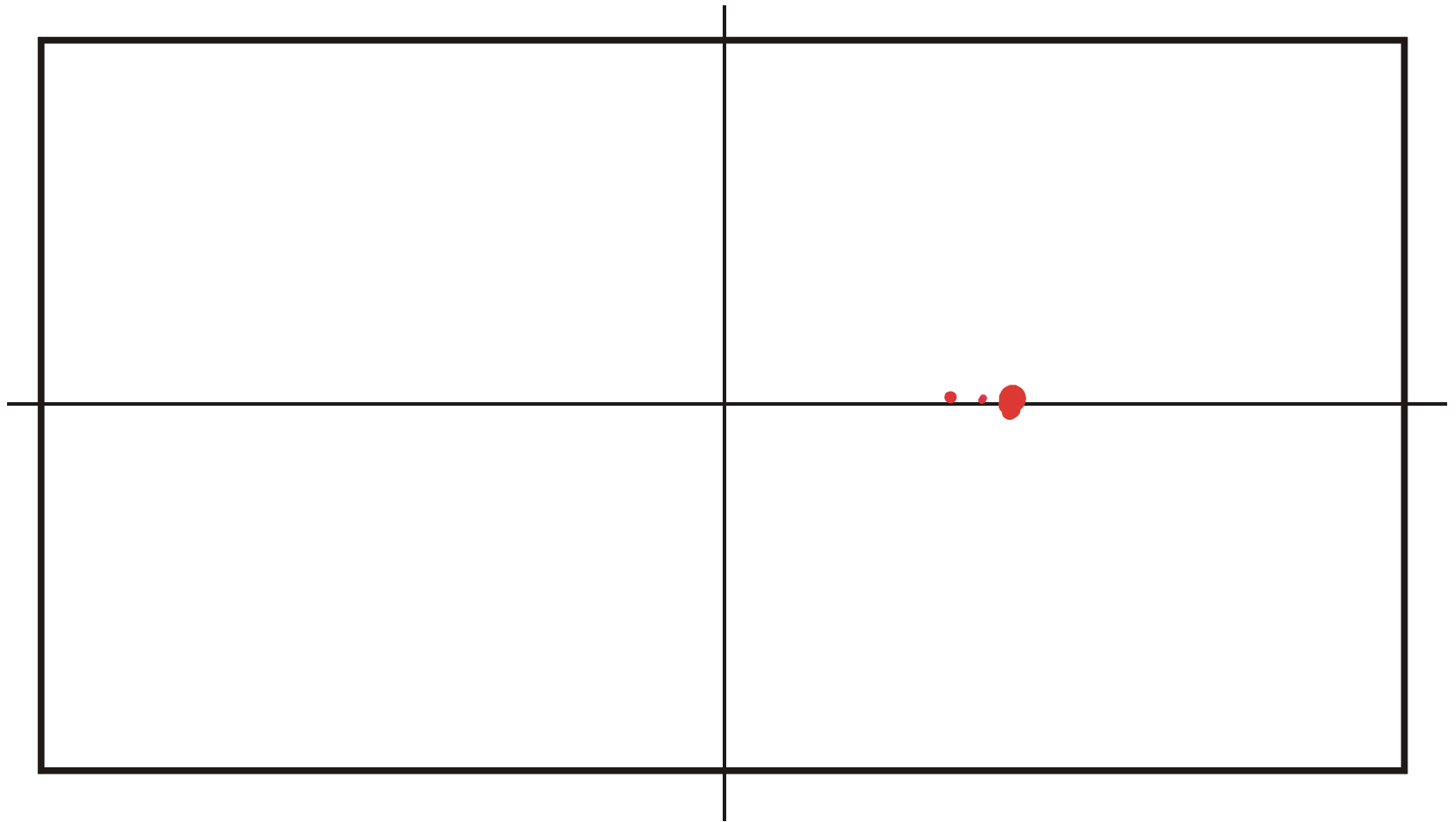
Spreading and evolution of a population on a neutral network :  $t = 830$



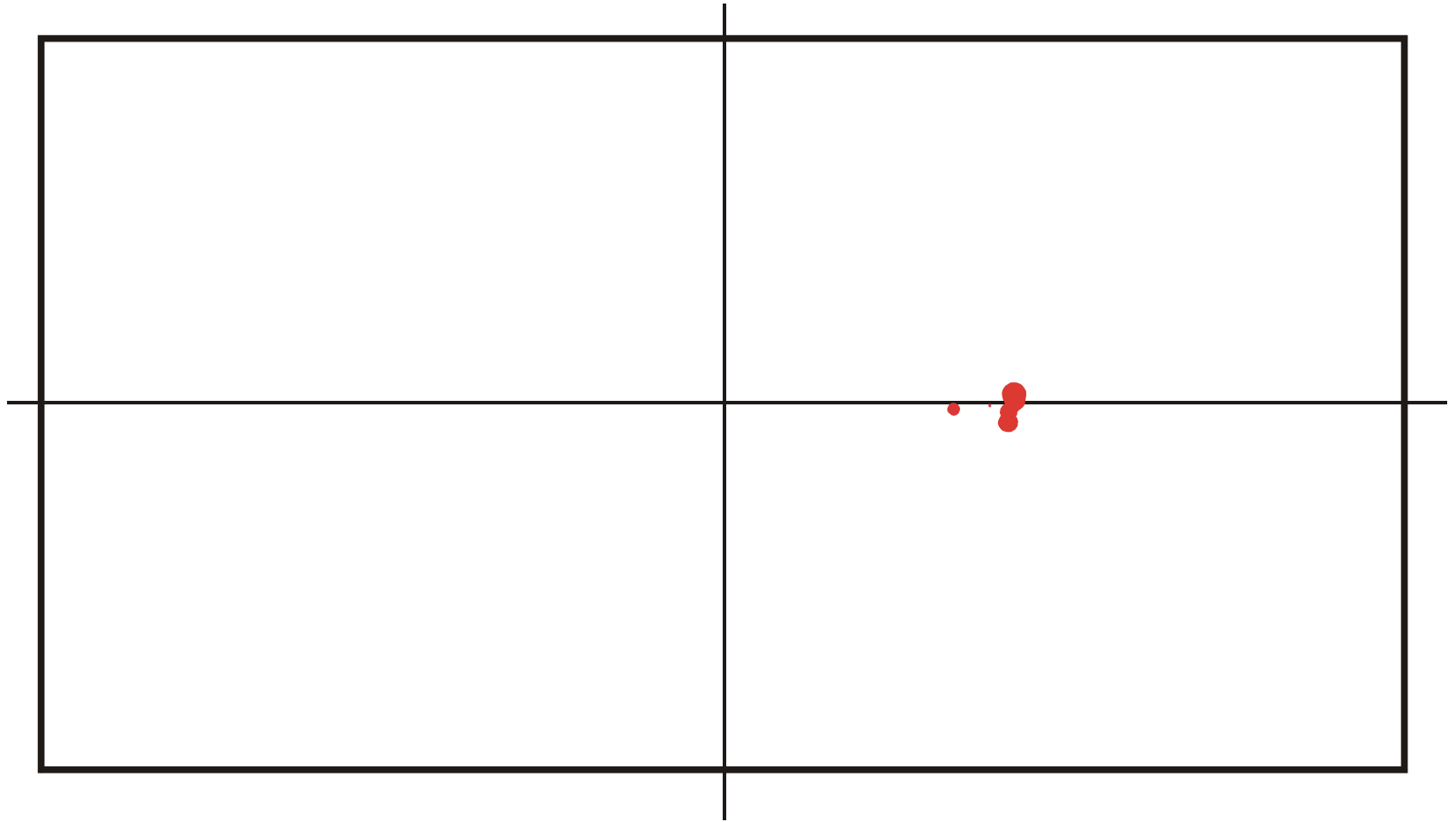
Spreading and evolution of a population on a neutral network :  $t = 835$



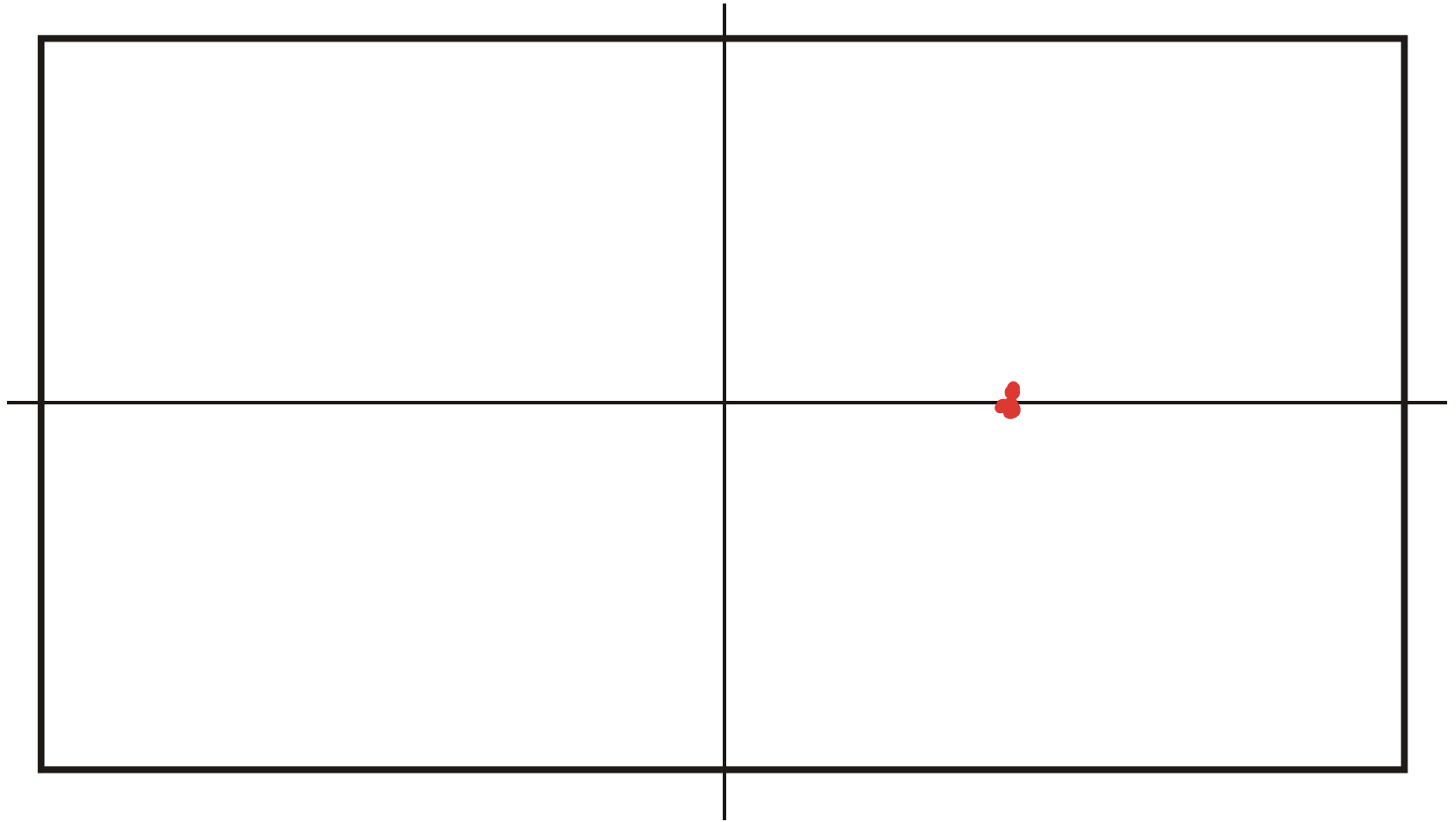
Spreading and evolution of a population on a neutral network :  $t = 840$



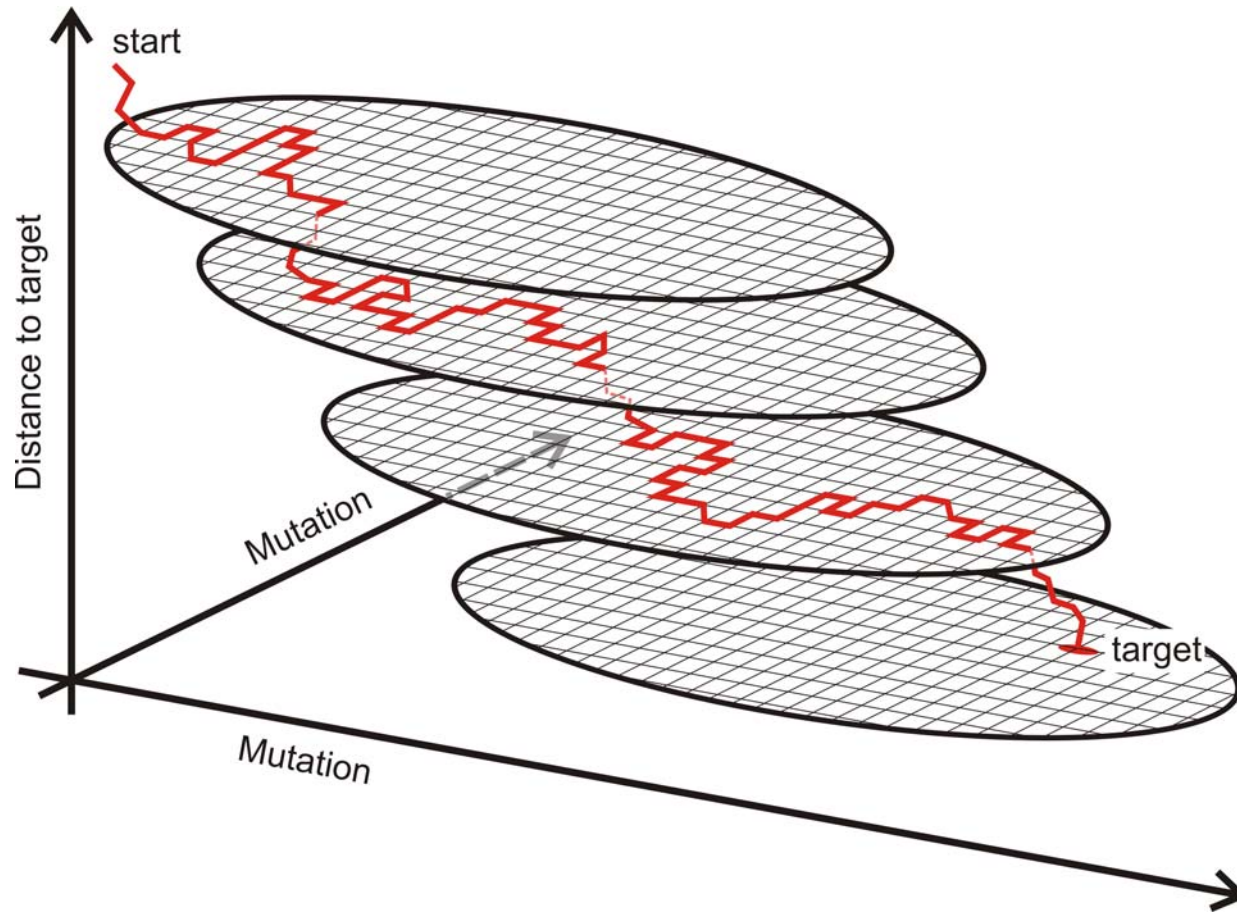
Spreading and evolution of a population on a neutral network :  $t = 845$



Spreading and evolution of a population on a neutral network :  $t = 850$



Spreading and evolution of a population on a neutral network :  $t = 855$



A sketch of optimization on neutral networks



**Table 8.** Statistics of the optimization trajectories. The table shows the results of sampled evolutionary trajectories leading from a random initial structure,  $S_I$ , to the structure of tRNA<sup>phe</sup>,  $S_T$ , as the target<sup>a</sup>. Simulations were performed with an algorithm introduced by Gillespie [55–57]. The time unit is here undefined. A mutation rate of  $p = 0.001$  per site and replication were used. The mean and standard deviation were calculated under the assumption of a log-normal distribution that fits well the data of the simulations.

Alphabet	Population size, $N$	Number of runs, $n_R$	Real time from start to target		Number of replications [ $10^7$ ]	
			Mean value	$\sigma$	Mean value	$\sigma$
<b>AUGC</b>	1 000	120	900	+1380 –542	1.2	+3.1 –0.9
	2 000	120	530	+880 –330	1.4	+3.6 –1.0
	3 000	1199	400	+670 –250	1.6	+4.4 –1.2
	10 000	120	190	+230 –100	2.3	+5.3 –1.6
	30 000	63	110	+97 –52	3.6	+6.7 –2.3
	100 000	18	62	+50 –28	–	–
<b>GC</b>	1 000	46	5160	+15700 –3890	–	–
	3 000	278	1910	+5180 –1460	7.4	+35.8 –6.1
	10 000	40	560	+1620 –420	–	–

<sup>a</sup> The structures  $S_I$  and  $S_T$  were used in the optimization:

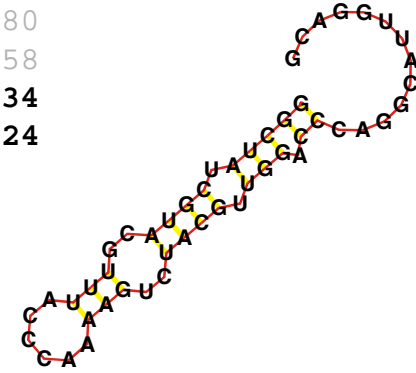
$S_I$ : ((.((((((((((((((((.....(((.....))).....)))))).)))))).))...(((.....)))

$S_T$ : ((((((...(((.....))))).((((.....))))).).....((((.....))))).))))).)....

Is the degree of neutrality in **GC** space much lower than in **AUGC** space ?

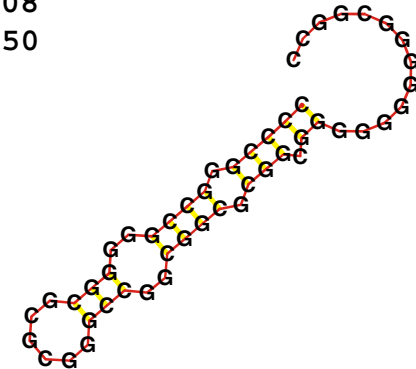
	<b>Number</b>	<b>Mean Value</b>	<b>Variance</b>	<b>Std.Dev.</b>
Total Hamming Distance:	150000	11.647973	23.140715	4.810480
Nonzero Hamming Distance:	99875	16.949991	30.757651	5.545958
Degree of Neutrality:	50125	<b>0.334167</b>	0.006961	<b>0.083434</b>
Number of Structures:	<b>1000</b>	<b>52.31</b>	85.30	<b>9.24</b>

1	(((((((((.....)))))))).)).....	<b>50125</b>	<b>0.334167</b>	
2	..(((((((((.....)))))))).)).....	2856	0.019040	
3	(((((((((((.....)))))))).)).....	2799	0.018660	
4	(((((((((((.....)))))))).)).....	2417	0.016113	
5	(((((((((((.....)))))))).)).....	2265	0.015100	
6	(((((((((((.....)))))))).)).....	2233	0.014887	



	<b>Number</b>	<b>Mean Value</b>	<b>Variance</b>	<b>Std.Dev.</b>
Total Hamming Distance:	50000	13.673580	10.795762	3.285691
Nonzero Hamming Distance:	45738	14.872054	10.821236	3.289565
Degree of Neutrality:	4262	<b>0.085240</b>	0.001824	<b>0.042708</b>
Number of Structures:	<b>1000</b>	<b>36.24</b>	6.27	<b>2.50</b>

1	(((((((((.....)))))))).)).....	<b>4262</b>	<b>0.085240</b>	
2	(((((((((((.....)))))))).)).....	1940	0.038800	
3	(((((((((((.....)))))))).)).....	1791	0.035820	
4	(((((((((((.....)))))))).)).....	1752	0.035040	
5	(((((((((((.....)))))))).)).....	1423	0.028460	



Shadow – Surrounding of an RNA structure in shape space – **AUGC** and **GC** alphabet

## Neutrality in evolution

Charles Darwin: *„ ... neutrality might exist ... ”*

Motoo Kimura: *„ ... neutrality is unavoidable and represents the main reason for changes in genotypes and leads to molecular phylogeny ... ”*

Current view: *„ ... neutrality is essential for successful optimization on rugged landscapes ... ”*

Proposed view: *„ ... neutrality provides the genetic reservoir in the rare and frequent mutation scenario ... ”*

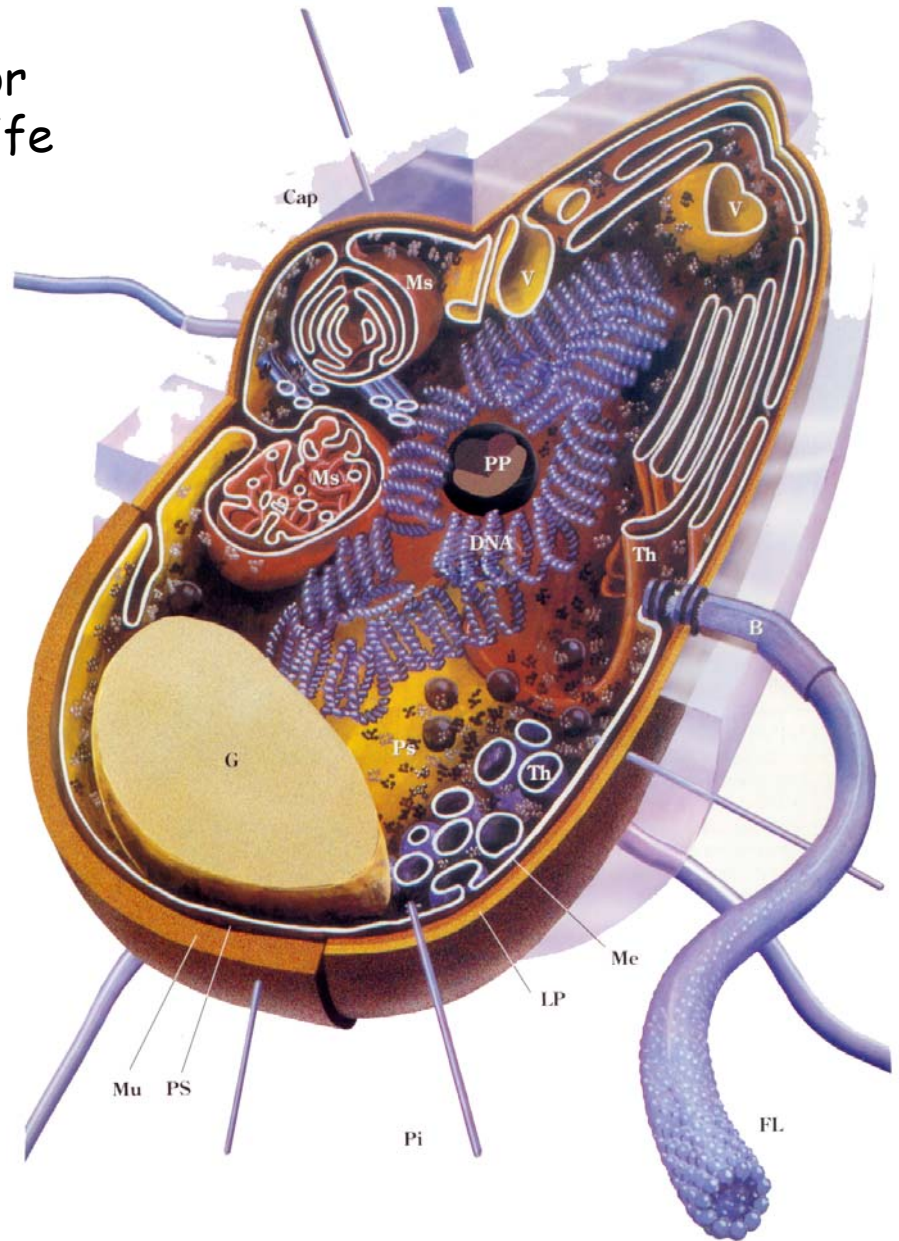
1. The chemistry of Darwinian evolution
2. RNA sequences and structures
3. Consequences of neutrality
4. Evolutionary optimization of RNA structure
5. **Complexity in biology**

The bacterial cell as an example for the simplest form of autonomous life

Escherichia coli genome:

4 million nucleotides

4460 genes



The structure of the bacterium *Escherichia coli*

---

---

EVOLUTIONARY TINKERING

*Blood . . . is the best possible thing to have coursing  
through one's veins.*

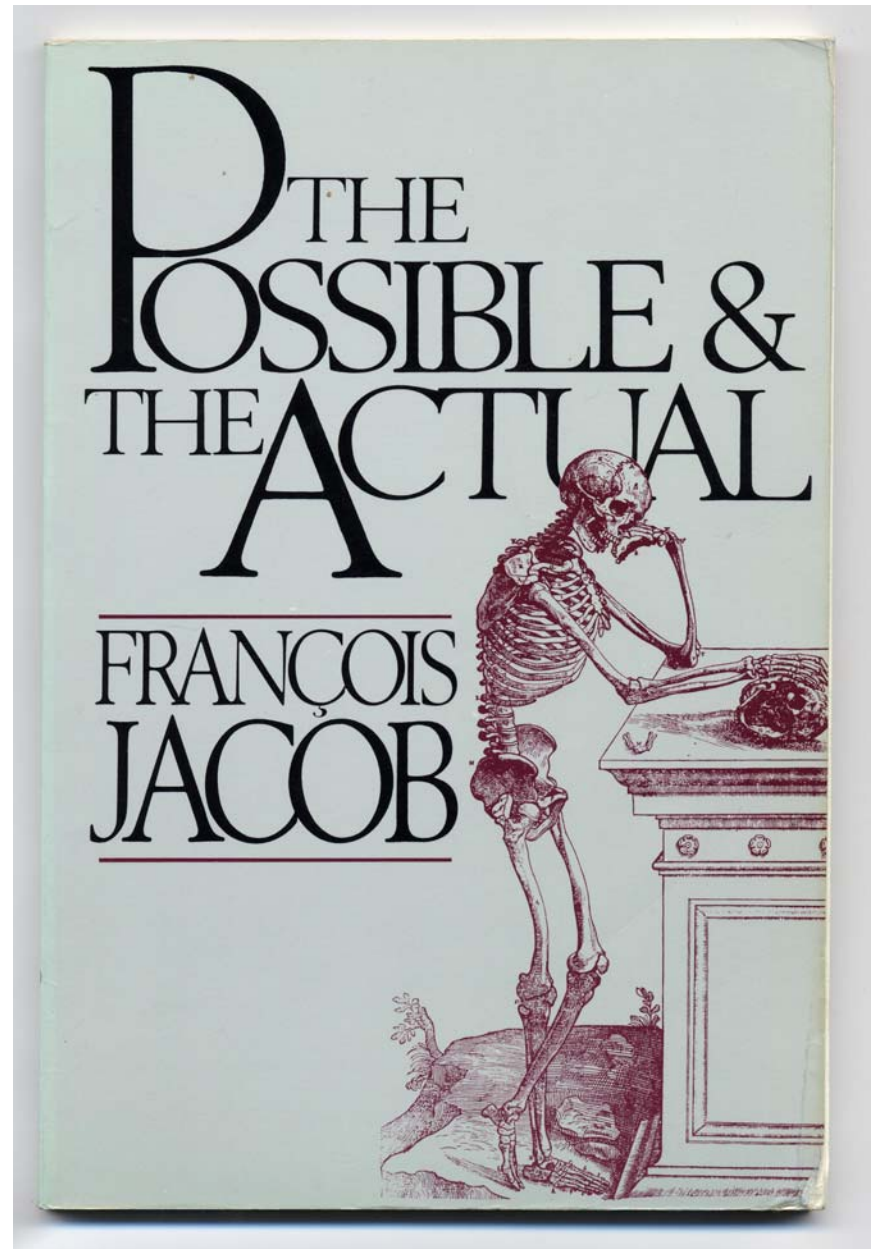
—Woody Allen, *Getting Even*

---

---

Evolution does not design with  
the eyes of an engineer,  
evolution works like a tinkerer.

Francois Jacob, Pantheon Books,  
New York 1982



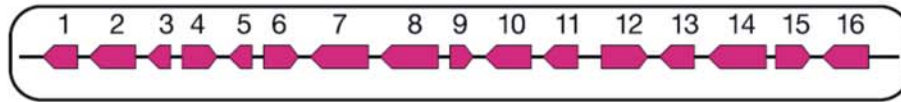
# The evolution of 'bricolage'

**DENIS DUBOULE** (denis.duboule@zoo.unige.ch)

**ADAM S. WILKINS** (edoffice@bioessays.demon.co.uk)

*The past ten years of developmental genetics have revealed that most of our genes are shared by other species throughout the animal kingdom. Consequently, animal diversity might largely rely on the differential use of the same components, either at the individual level through divergent functional recruitment, or at a more integrated level, through their participation in various genetic networks. Here, we argue that this inevitably leads to an increase in the interdependency between functions that, in turn, influences the degree to which novel variations can be tolerated. In this 'transitionist' scheme, evolution is neither inherently gradualist nor punctuated but, instead, progresses from one extreme to the other, together with the increased complexity of organisms.*

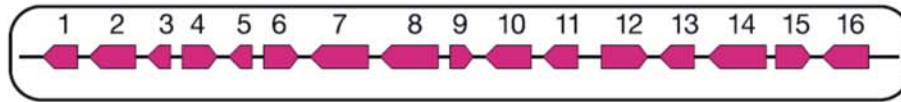
D. Duboule, A.S. Wilkins. 1998.  
The evolution of 'bricolage'.  
Trends in Genetics 14:54-59.



## A model for the genome duplication in yeast 100 million years ago

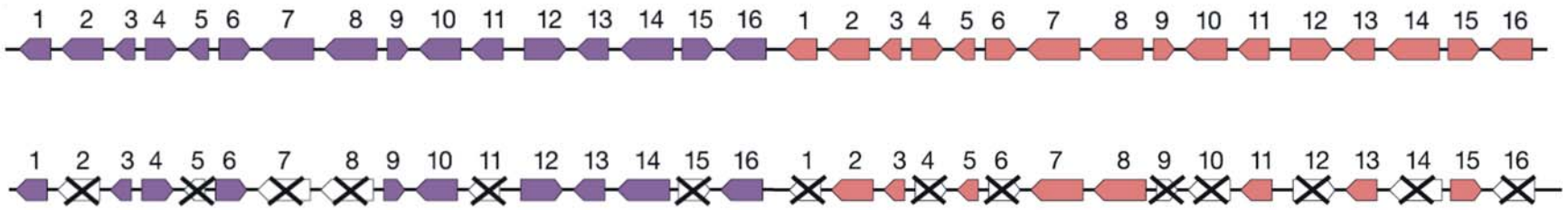
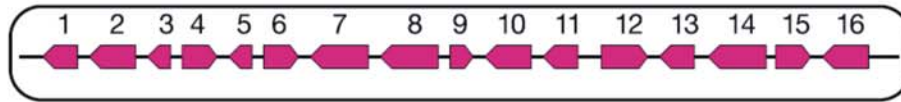
Manolis Kellis, Bruce W. Birren, and Eric S. Lander. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617-624, 2004





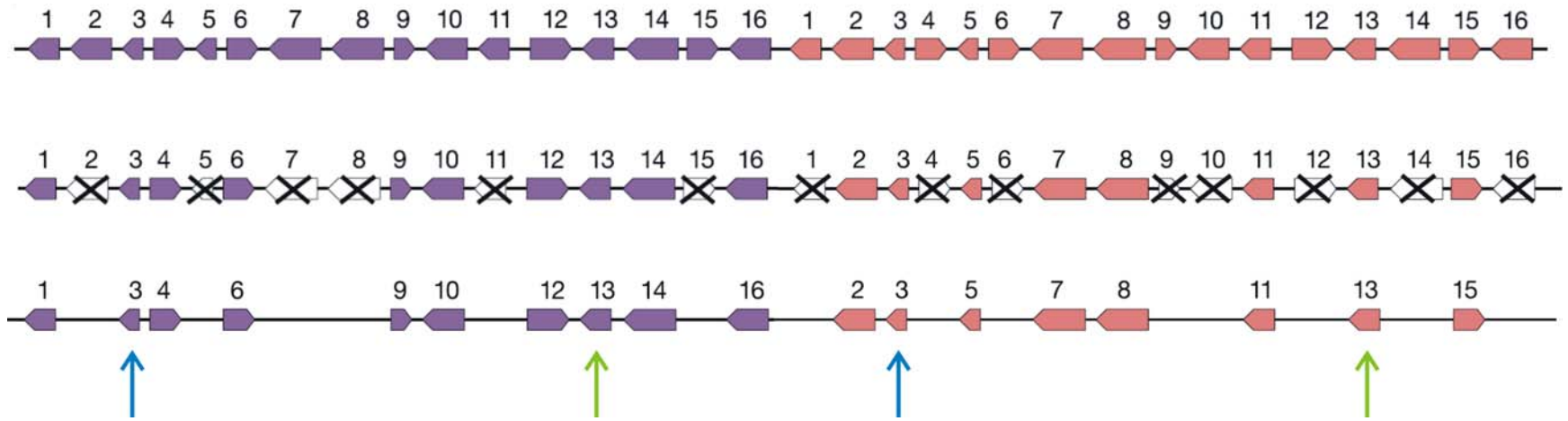
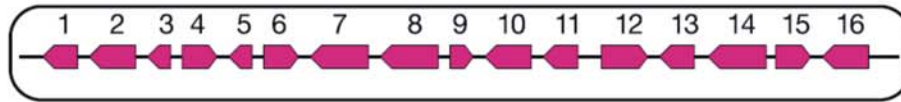
## A model for the genome duplication in yeast 100 million years ago

Manolis Kellis, Bruce W. Birren, and Eric S. Lander. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617-624, 2004



## A model for the genome duplication in yeast 100 million years ago

Manolis Kellis, Bruce W. Birren, and Eric S. Lander. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617-624, 2004



A model for the genome duplication in yeast 100 million years ago

Manolis Kellis, Bruce W. Birren, and Eric S. Lander. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617-624, 2004

# WHAT IS A GENE?

The idea of genes as beads on a DNA string is fast fading. Protein-coding sequences have no clear beginning or end and RNA is a key part of the information package, reports **Helen Pearson**.

'Gene' is not a typical four-letter word. It is not offensive. It is never bleeped out of TV shows. And where the meaning of most four-letter words is all too clear, that of gene is not. The more expert scientists become in molecular genetics, the less easy it is to be sure about what, if anything, a gene actually is.

Rick Young, a geneticist at the Whitehead Institute in Cambridge, Massachusetts, says that when he first started teaching as a young professor two decades ago, it took him about two hours to teach fresh-faced undergraduates what a gene was and the nuts and bolts of how it worked. Today, he and his colleagues need three months of lectures to convey the concept of the gene, and that's not because the students are any less bright. "It takes a whole semester to teach this stuff to talented graduates," Young says. "It used to be we could give a one-off definition and now it's much more complicated."

In classical genetics, a gene was an abstract concept — a unit of inheritance that ferried a characteristic from parent to child. As biochemistry came into its own, those characteristics were associated with enzymes or proteins, one for each gene. And with the advent of molecular biology, genes became real, physical things — sequences of DNA which when converted into strands of so-called messenger RNA could be used as the basis for building their associated protein piece by piece. The great coiled DNA molecules of the chromosomes were seen as long strings on which gene sequences sat like discrete beads.

This picture is still the working model for many scientists. But those at the forefront of genetic research see it as increasingly old-fashioned — a crude approximation that, at best, hides fascinating new complexities and, at worst, blinds its users to useful new paths of enquiry.

Information, it seems, is parceled out along chromosomes in a much more complex way than was originally supposed. RNA molecules are not just passive conduits through which the gene's message flows into the world but active regulators of cellular processes. In some cases, RNA may even pass information across generations — normally the sole preserve of DNA.

An eye-opening study last year raised the possibility that plants sometimes rewrite their DNA on the basis of RNA messages inherited from generations past<sup>1</sup>. A study on page 469 of this issue suggests that a comparable phenomenon might occur in mice, and by implication in other mammals<sup>2</sup>. If this type of phenomenon is indeed widespread, it "would have huge implications," says evolutionary geneticist

Laurence Hurst at the University of Bath, UK.

"All of that information seriously challenges our conventional definition of a gene," says molecular biologist Bing Ren at the University of California, San Diego. And the information challenge is about to get even tougher. Later this year, a glut of data will be released from the international Encyclopedia of DNA Elements (ENCODE) project. The pilot phase of ENCODE involves scrutinizing roughly 1% of the human genome in unprecedented detail; the aim is to find all the sequences that serve a useful purpose and explain what that purpose is. "When we started the ENCODE project I had a different view of what a gene was," says contributing researcher Roderic Guigo at the Center for Genomic Regulation in Barcelona. "The degree of complexity we've seen was not anticipated."

## Under fire

The first of the complexities to challenge molecular biology's paradigm of a single DNA sequence encoding a single protein was alternative splicing, discovered in viruses in 1977 (see 'Hard to track', overleaf). Most of the DNA sequences describing proteins in humans have a modular arrangement in which exons, which carry the instructions for making proteins, are interspersed with non-coding introns. In alternative splicing, the cell snips out introns and sews together the exons in various different orders, creating messages that can code for different proteins. Over the years geneticists have also documented overlapping genes, genes within genes and countless other weird arrangements (see 'Muddling over genes', overleaf).

Alternative splicing, however, did not in itself require a drastic reappraisal of the notion of a gene; it just showed that some DNA sequences could describe more than one protein. Today's assault on the gene concept is more far reaching, fuelled largely by studies that show the pre-

viously unimagined scope of RNA.

The one gene, one protein idea is coming under particular assault from researchers who are comprehensively extracting and analysing the RNA messages, or transcripts, manufactured by genomes, including the human and mouse genome. Researchers led by Thomas Gingeras at the company Affymetrix in Santa Clara, California, for example, recently studied all the transcripts from ten chromosomes across eight human cell lines and worked out

precisely where on the chromosomes each of the transcripts came from<sup>3</sup>.

The picture these studies paint is one of mind-boggling complexity. Instead of discrete genes dutifully mass-producing

identical RNA transcripts, a teeming mass of transcription converts many segments of the genome into multiple RNA ribbons of differing lengths. These ribbons can be generated from both strands of DNA, rather than from just one as was conventionally thought. Some of these transcripts come from regions of DNA previously identified as holding protein-coding genes. But many do not. "It's somewhat revolutionary," says Gingeras's colleague Phillip Kapranov. "We've come to the realization that the genome is full of overlapping transcripts."

Other studies, one by Guigo's team<sup>4</sup>, and one by geneticist Rotem Sorek<sup>5</sup>, now at Tel Aviv University, Israel, and his colleagues, have hinted at the reasons behind the mass of transcription. The two teams investigated occasional reports that transcription can start at a DNA sequence associated with one protein and run straight through into the gene for a completely different protein, producing a fused transcript. By delving into databases of human RNA transcripts, Guigo's team estimate that 4–5% of the DNA in regions conventionally recognized as genes is transcribed in this way. Producing fused transcripts could be one way for a cell to generate a greater variety of proteins from a limited number of exons, the researchers say.

Many scientists are now starting to think that the descriptions of proteins encoded in DNA know no borders — that each sequence reaches into the next and beyond. This idea will be one of the central points to emerge from the ENCODE project when its results are published later this year.

Kapranov and others say that they have documented many examples of transcripts in which protein-coding exons from one part of the genome combine with exons from another

**"We've come to the realization that the genome is full of overlapping transcripts."**

— Phillip Kapranov



Spools of DNA (above) still harbour surprises, with one protein-coding gene often overlapping the next.

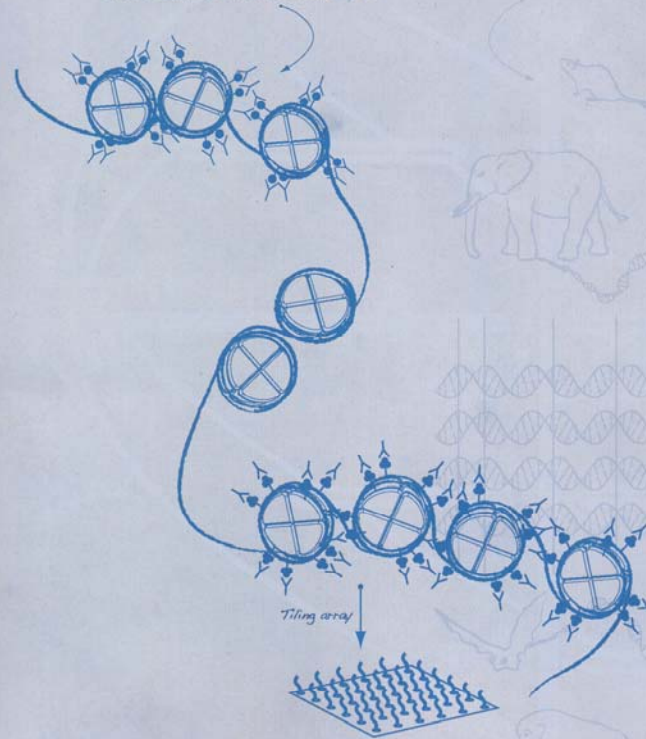
The difficulty to define the notion of „gene“.

Helen Pearson,  
*Nature* 441: 399-401, 2006

# nature

*Hi-stone-modification chromatin IP*

*Comparative syntenic alignment*



**MARS'S  
ANCIENT OCEAN**  
Polar wander  
solves an enigma

**THE DEPTHS OF  
DISGUST**  
Understanding the  
ugliest emotion

**MENTORING**  
How to be top

**NATUREJOBS**  
Contract  
research

## DECODING THE BLUEPRINT

The ENCODE pilot maps  
human genome function



ENCODE stands for  
**ENC**yclopedia **Of** **DNA** **E**lements.

ENCODE Project Consortium.  
Identification and analysis of functional  
elements in 1% of the human genome by  
the ENCODE pilot project.  
*Nature* **447**:799-816, 2007

## Acknowledgement of support

Fonds zur Förderung der wissenschaftlichen Forschung (FWF)  
Projects No. 09942, 10578, 11065, 13093  
13887, and 14898

Wiener Wissenschafts-, Forschungs- und Technologiefonds (WWTF)  
Project No. Mat05

Jubiläumsfonds der Österreichischen Nationalbank  
Project No. Nat-7813

European Commission: Contracts No. 98-0189, 12835 (NEST)

Austrian Genome Research Program – GEN-AU: Bioinformatics  
Network (BIN)

Österreichische Akademie der Wissenschaften

Siemens AG, Austria

Universität Wien and the Santa Fe Institute



Universität Wien

# Coworkers

**Peter Stadler, Bärbel M. Stadler**, Universität Leipzig, GE

**Paul E. Phillipson**, University of Colorado at Boulder, CO

**Heinz Engl, Philipp Kügler, James Lu, Stefan Müller**, RICAM Linz, AT

**Jord Nagel, Kees Pleij**, Universiteit Leiden, NL

**Walter Fontana**, Harvard Medical School, MA

**Christian Reidys, Christian Forst**, Los Alamos National Laboratory, NM

**Ulrike Göbel, Walter Grüner, Stefan Kopp, Jaqueline Weber**, Institut für  
Molekulare Biotechnologie, Jena, GE

**Ivo L.Hofacker, Christoph Flamm, Andreas Svrček-Seiler**, Universität Wien, AT

**Kurt Grünberger, Michael Kospach, Andreas Wernitznig, Stefanie Widder,  
Stefan Wuchty**, Universität Wien, AT

**Jan Cupal, Stefan Bernhart, Lukas Endler, Ulrike Langhammer, Rainer Machne,  
Ulrike Mückstein, Hakim Tafer, Thomas Taylor**, Universität Wien, AT



Universität Wien

## **Prediction of RNA secondary structures: from theory to models and real molecules**

**Peter Schuster**<sup>1,2</sup>

<sup>1</sup>Institut für Theoretische Chemie der Universität Wien, Währingerstraße 17, A-1090 Vienna, Austria

<sup>2</sup>The Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA

E-mail: [pbs@tbi.univie.ac.at](mailto:pbs@tbi.univie.ac.at)



Web-Page for further information:

<http://www.tbi.univie.ac.at/~pks>

