



# **How innovation occurs in evolution of molecules**

Peter Schuster

Institut für Theoretische Chemie und Molekulare  
Strukturbiologie der Universität Wien

Evolutionary innovation

Praha, 30.05.2002

**Darwinian principle is based on three functions:**

- **Reproduction** efficiency expressed by fitness of **phenotypes**.
- **Variation** of **genotypes** through imperfect copying and recombination.
- **Selection** of **phenotypes** based on differences in fitness.

**Two additional features are required:**

- Large reservoirs of genotypes and sufficiently rich repertoires of phenotypes.
- Mapping of genotypes into phenotypes with suitable properties.

The **genotypes** or **genomes** of individuals are DNA or RNA sequences. They are changing from generation to generation through mutation and recombination. Species are reproductively related ensembles of individuals.

Genotypes unfold into **phenotypes**, being molecular structures, viruses or organisms, which are the targets of the evolutionary selection process.

The most common mutations are **point mutations**, which consist of single nucleotide exchanges. The **Hamming distance** of two sequences is the minimal number of single nucleotide exchanges that mutually converts the two sequence into each other.



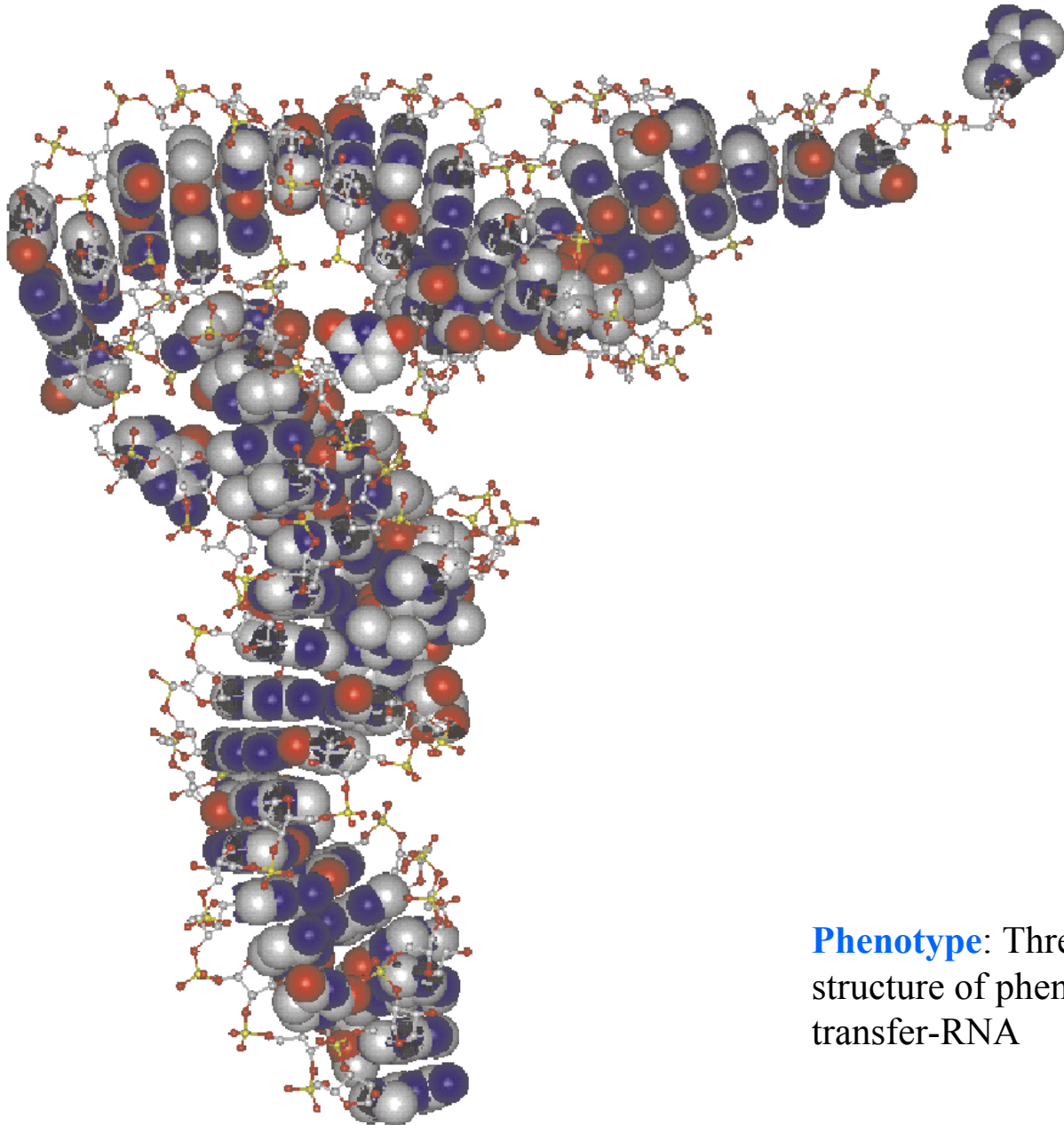
**A** = adenylate

**U** = uridylate

**C** = cytidylate

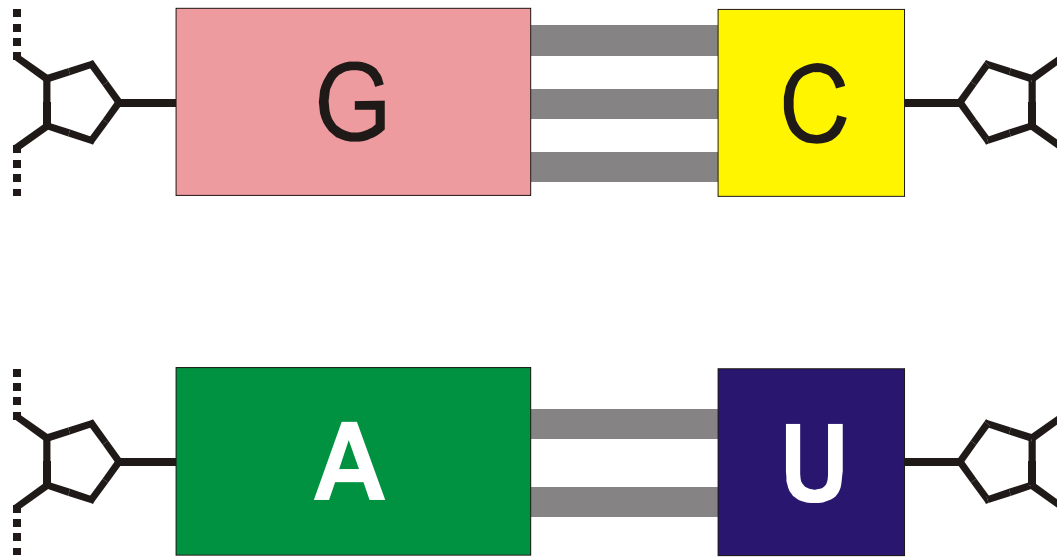
**G** = guanylate

**Genotype:** The sequence of an RNA molecule consisting of monomers chosen from four classes, A, U, G, and C.

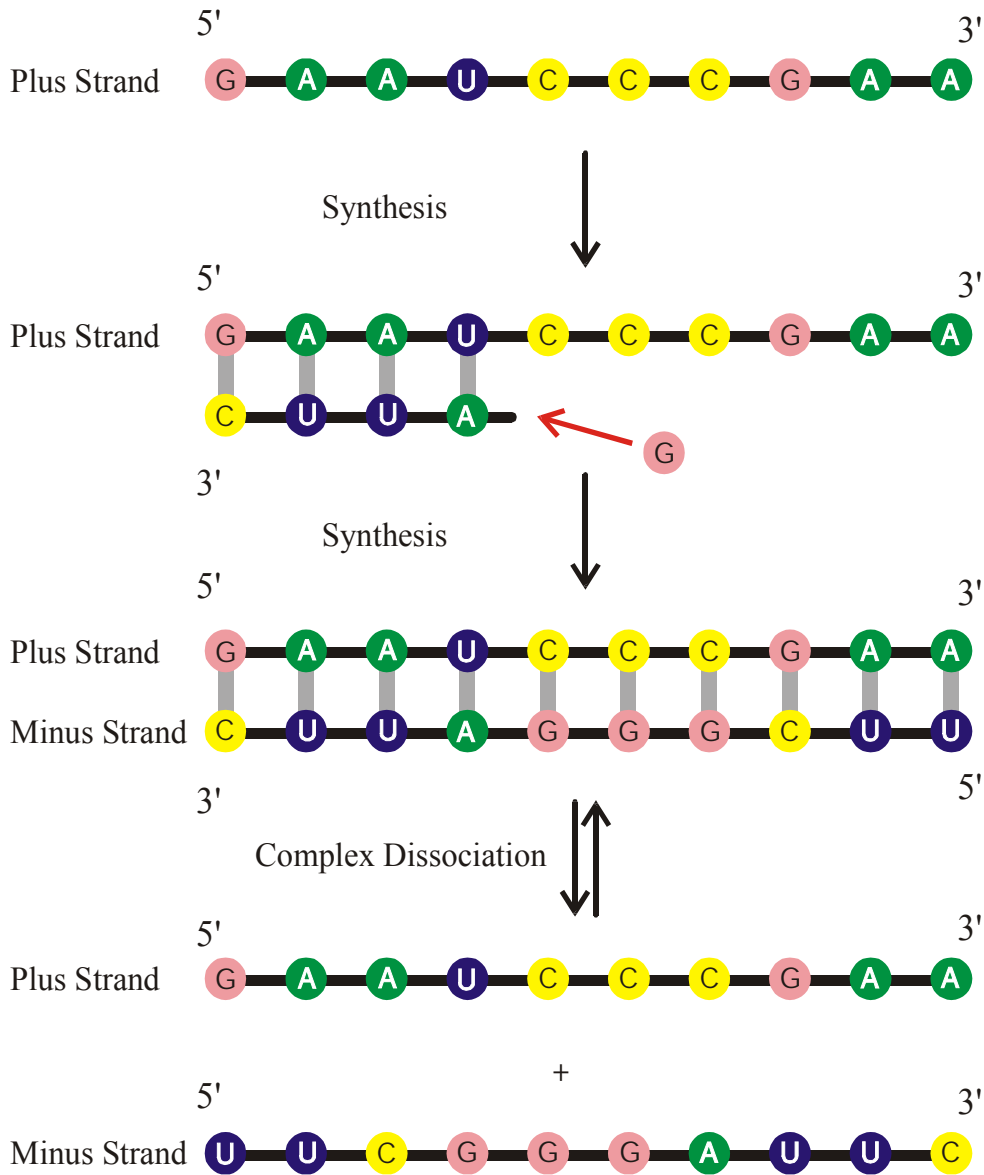


**Phenotype:** Three-dimensional structure of phenylalanyl transfer-RNA

## Hydrogen bonds

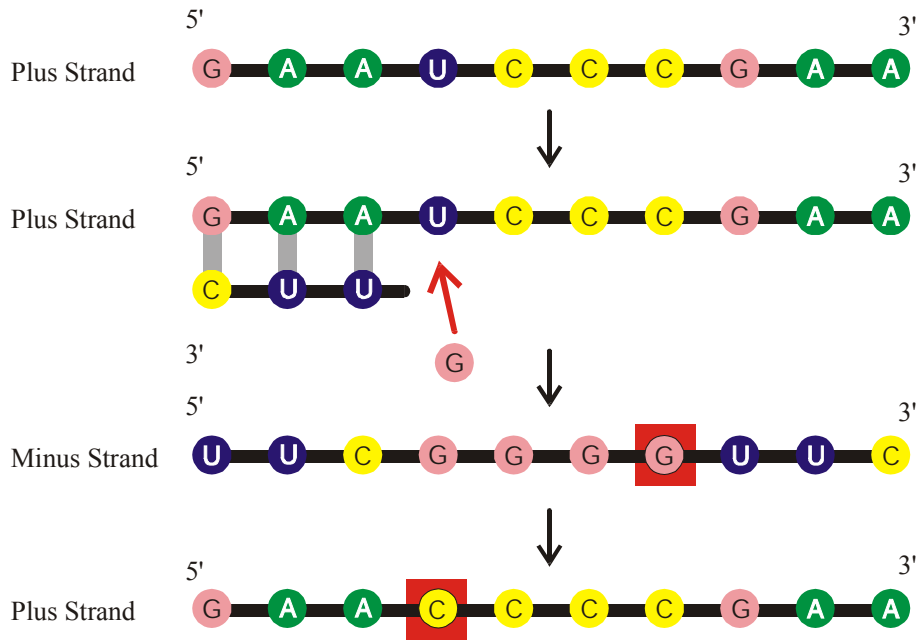


Hydrogen bonding between nucleotide bases is the principle of template action of RNA and DNA.



Complementary replication as the simplest copying mechanism of RNA

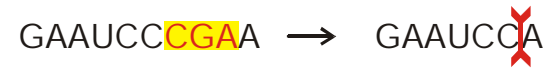




**Point Mutation**



**Insertion**



**Deletion**

Mutations represent the mechanism of variation in nucleic acids.

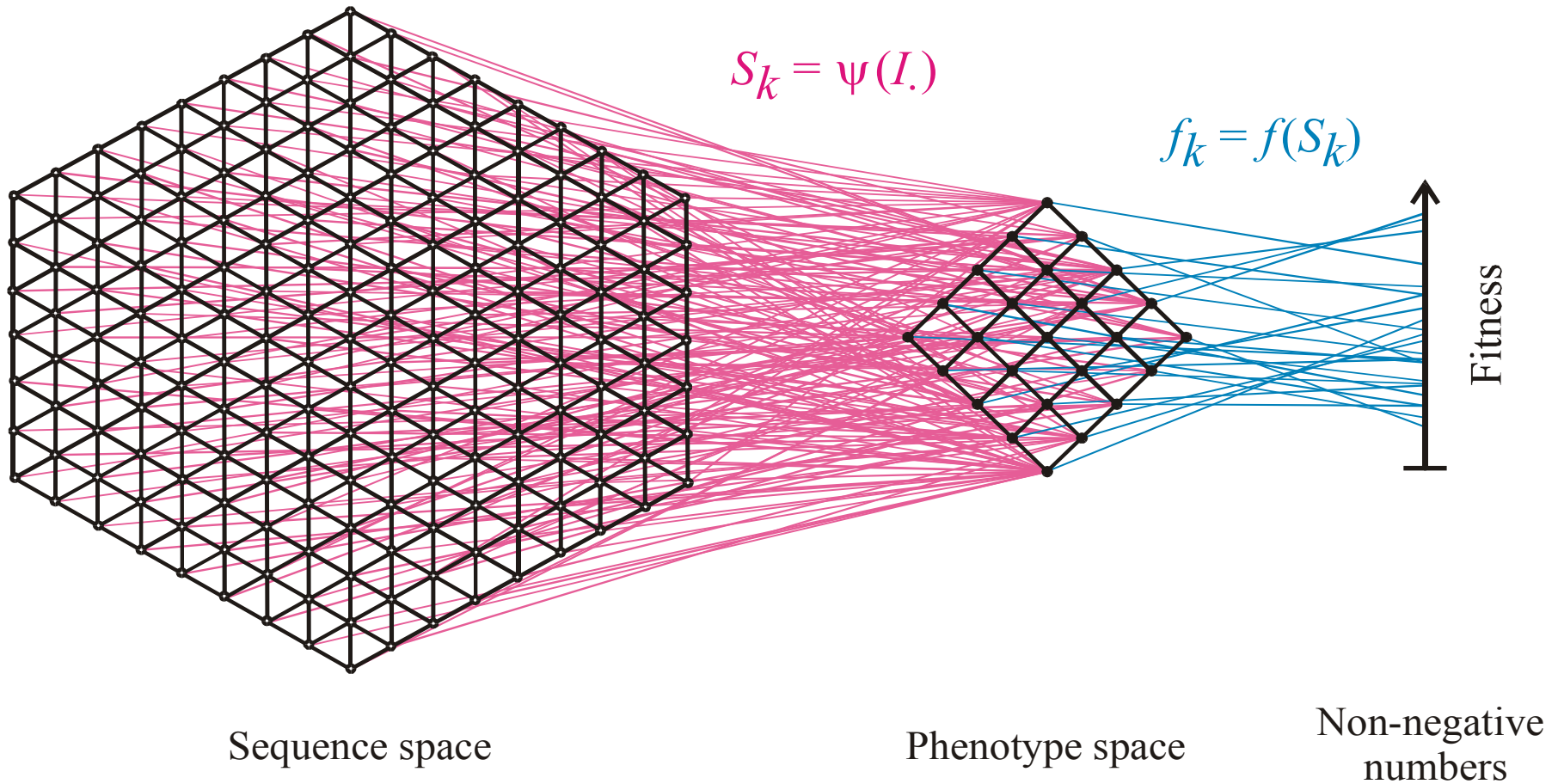


$4^{27} = 1.801 \times 10^{16}$  possible different sequences

Combinatorial diversity of sequences:  $N = 4^{\{$

- A = adenylate
- U = uridylate
- C = cytidylate
- G = guanylate

Combinatorial diversity of heteropolymers illustrated by means of an RNA aptamer that binds to the antibiotic tobramycin



Mapping from sequence space into phenotype space and into fitness values

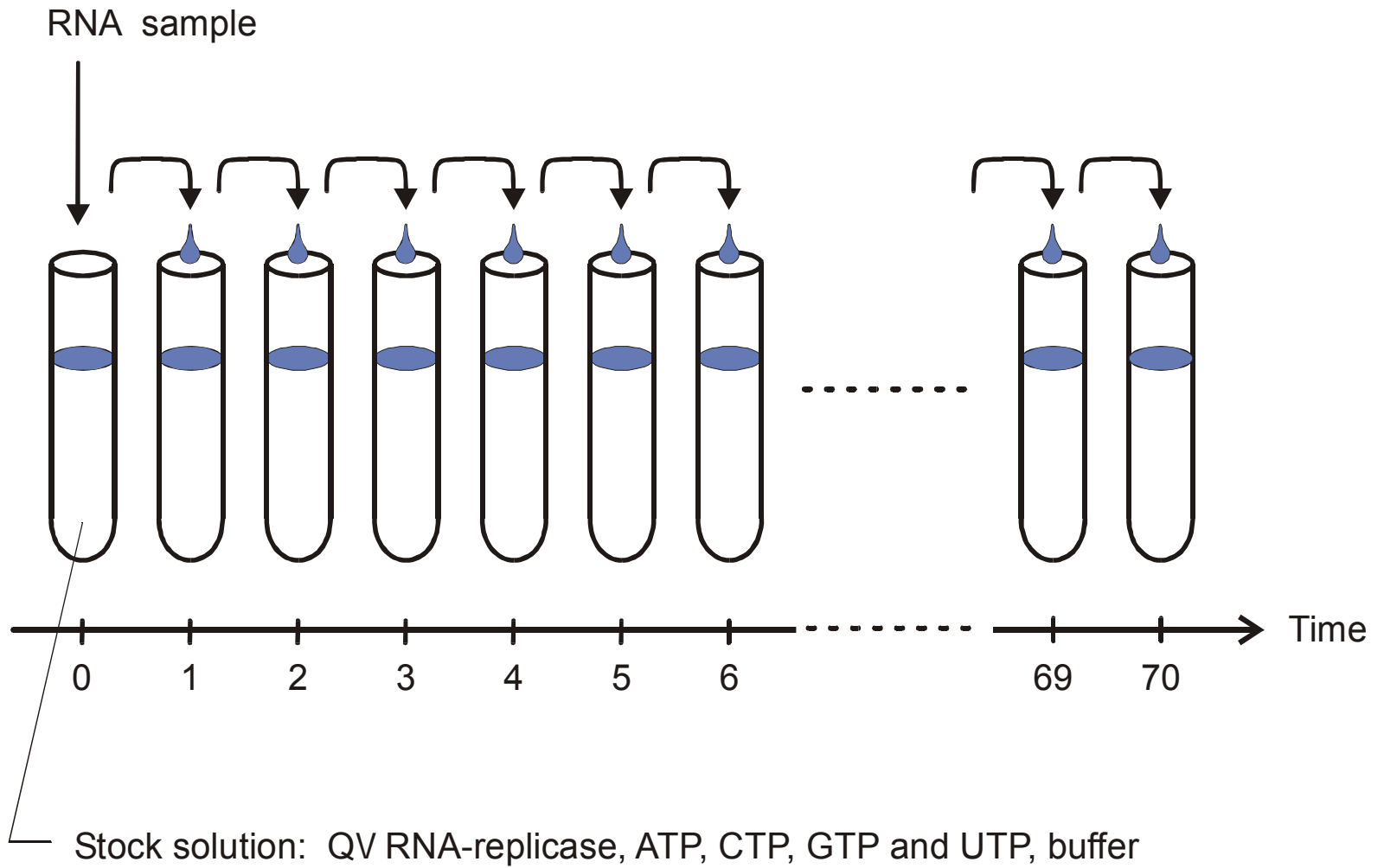
## Evolution of RNA molecules based on Q $\beta$ phage

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

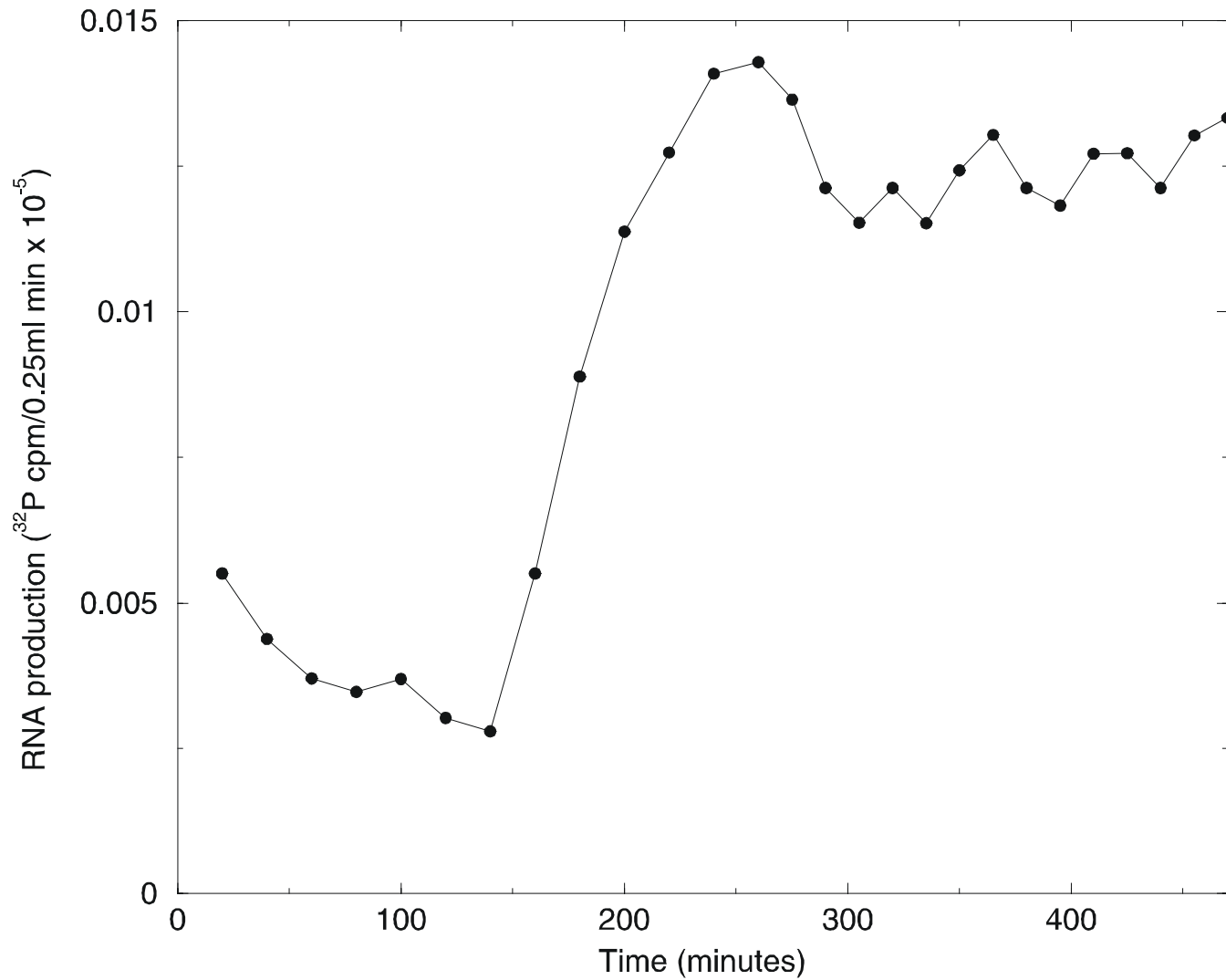
S.Spiegelman, *An approach to the experimental analysis of precellular evolution*. Quart.Rev.Biophys. **4** (1971), 213-253

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

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



The serial transfer technique applied to RNA evolution *in vitro*



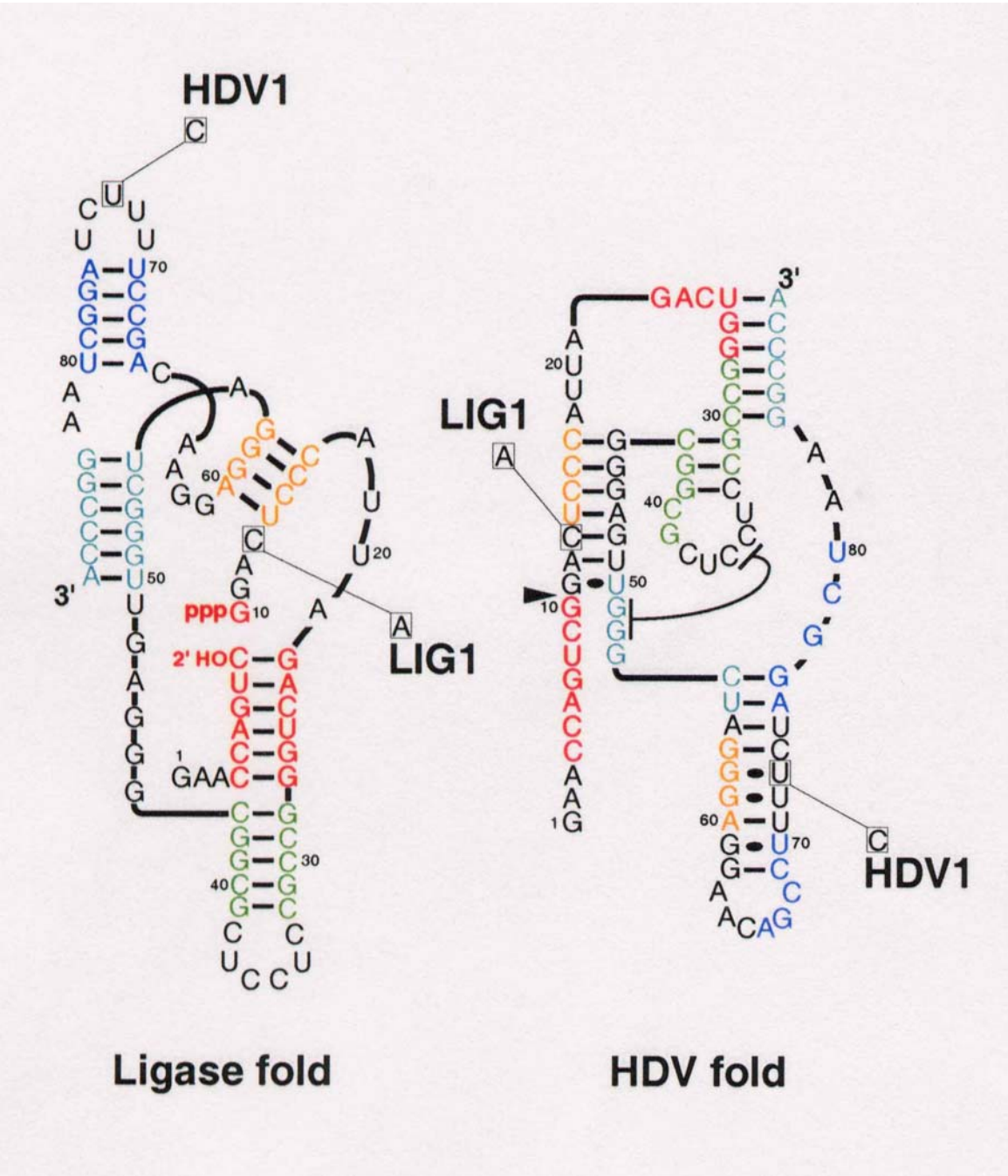
The increase in RNA production rate during a serial transfer experiment

## **A ribozyme switch**

E.A.Schultes, D.B.Bartel, *One sequence, two ribozymes: Implication for the emergence of new ribozyme folds*. Science **289** (2000), 448-452







The sequence at the *intersection*:

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



S0092-8240(96)00089-4

## GENERIC PROPERTIES OF COMBINATORY MAPS: NEUTRAL NETWORKS OF RNA SECONDARY STRUCTURES<sup>1</sup>

■ CHRISTIAN REIDYS\*, †, PETER F. STADLER\*, ‡  
 and PETER SCHUSTER\*, ‡, §, ¶

\*Santa Fe Institute,  
 Santa Fe, NM 87501, U.S.A.

†Los Alamos National Laboratory,  
 Los Alamos, NM 87545, U.S.A.

‡Institut für Theoretische Chemie der Universität Wien,  
 A-1090 Wien, Austria

§Institut für Molekulare Biotechnologie,  
 D-07708 Jena, Germany

(E.mail: pks@tbi.univie.ac.at)

Random graph theory is used to model and analyse the relationships between sequences and secondary structures of RNA molecules, which are understood as mappings from sequence space into shape space. These maps are non-invertible since there are always many orders of magnitude more sequences than structures. Sequences folding into identical structures form *neutral networks*. A neutral network is embedded in the set of sequences that are *compatible* with the given structure. Networks are modeled as graphs and constructed by random choice of vertices from the space of compatible sequences. The theory characterizes neutral networks by the mean fraction of neutral neighbors ( $\lambda$ ). The networks are connected and percolate sequence space if the fraction of neutral nearest neighbors exceeds a threshold value ( $\lambda > \lambda^*$ ). Below threshold ( $\lambda < \lambda^*$ ), the networks are partitioned into a largest “giant” component and several smaller components. Structures are classified as “common” or “rare” according to the sizes of their pre-images, i.e. according to the fractions of sequences folding into them. The neutral networks of any pair of two different common structures almost touch each other, and, as expressed by the conjecture of *shape space covering* sequences folding into almost all common structures, can be found in a small ball of an arbitrary location in sequence space. The results from random graph theory are compared to data obtained by folding large samples of RNA sequences. Differences are explained in terms of specific features of RNA molecular structures. © 1997 Society for Mathematical Biology

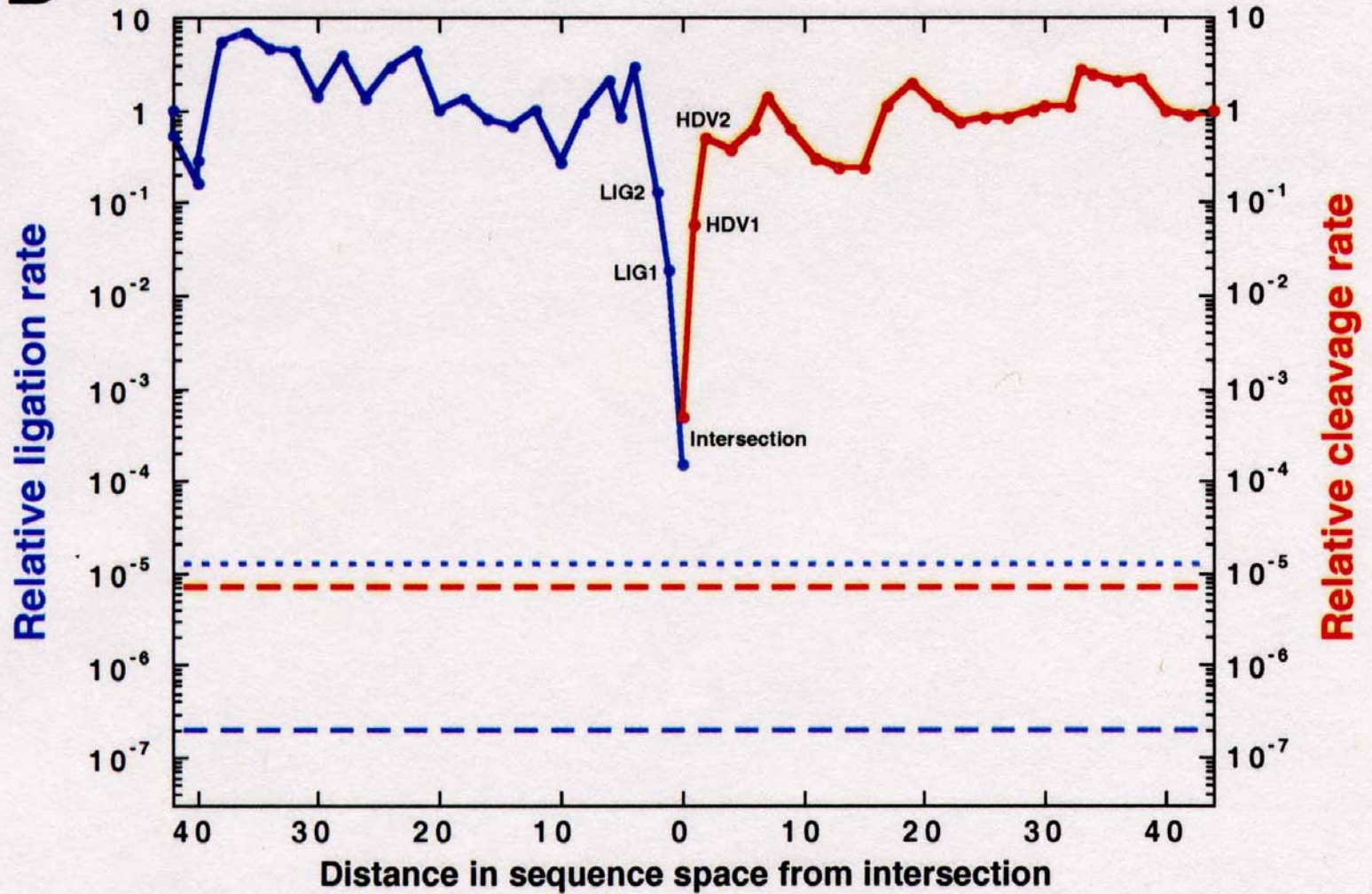
**THEOREM 5. INTERSECTION-THEOREM.** *Let  $s$  and  $s'$  be arbitrary secondary structures and  $C[s], C[s']$  their corresponding compatible sequences. Then,*

$$C[s] \cap C[s'] \neq \emptyset.$$

*Proof.* Suppose that the alphabet admits only the complementary base pair  $[XY]$  and we ask for a sequence  $x$  compatible to both  $s$  and  $s'$ . Then  $f(s, s') \cong D_m$  operates on the set of all positions  $\{x_1, \dots, x_n\}$ . Since we have the operation of a dihedral group, the orbits are either cycles or chains and the cycles have even order. A constraint for the sequence compatible to both structures appears only in the cycles where the choice of bases is not independent. It remains to be shown that there is a valid choice of bases for each cycle, which is obvious since these have even order. Therefore, it suffices to choose an alternating sequence of the pairing partners  $X$  and  $Y$ . Thus, there are at least two different choices for the first base in the orbit. ■

*Remark.* A generalization of the statement of theorem 5 to three different structures is false.

Reference for the definition of the intersection and the proof of the *intersection theorem*

**B**

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



# From sequences to shapes and back: a case study in RNA secondary structures

PETER SCHUSTER<sup>1,2,3</sup>, WALTER FONTANA<sup>3</sup>, PETER F. STADLER<sup>2,3</sup>  
AND IVO L. HOFACKER<sup>2</sup>

<sup>1</sup> Institut für Molekulare Biotechnologie, Beutenbergstrasse 11, PF 100813, D-07708 Jena, Germany

<sup>2</sup> Institut für Theoretische Chemie, Universität Wien, Austria

<sup>3</sup> Santa Fe Institute, Santa Fe, U.S.A.

## SUMMARY

RNA folding is viewed here as a map assigning secondary structures to sequences. At fixed chain length the number of sequences far exceeds the number of structures. Frequencies of structures are highly non-uniform and follow a generalized form of Zipf's law: we find relatively few common and many rare ones. By using an algorithm for inverse folding, we show that sequences sharing the same structure are distributed randomly over sequence space. All common structures can be accessed from an arbitrary sequence by a number of mutations much smaller than the chain length. The sequence space is percolated by extensive neutral networks connecting nearest neighbours folding into identical structures. Implications for evolutionary adaptation and for applied molecular evolution are evident: finding a particular structure by mutation and selection is much simpler than expected and, even if catalytic activity should turn out to be sparse in the space of RNA structures, it can hardly be missed by evolutionary processes.

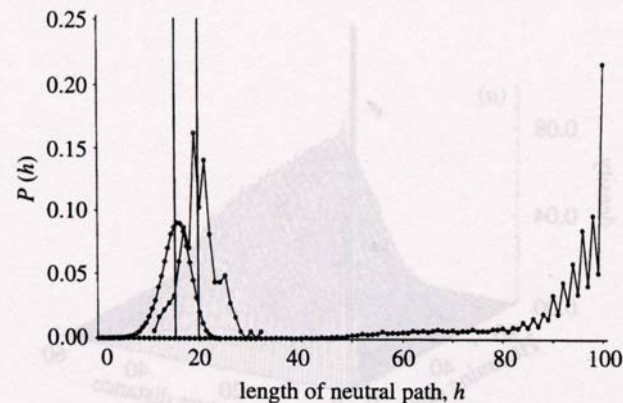


Figure 4. Neutral paths. A neutral path is defined by a series of nearest neighbour sequences that fold into identical structures. Two classes of nearest neighbours are admitted: neighbours of Hamming distance 1, which are obtained by single base exchanges in unpaired stretches of the structure, and neighbours of Hamming distance 2, resulting from base pair exchanges in stacks. Two probability densities of Hamming distances are shown that were obtained by searching for neutral paths in sequence space: (i) an upper bound for the closest approach of trial and target sequences (open circles) obtained as endpoints of neutral paths approaching the target from a random trial sequence (185 targets and 100 trials for each were used); (ii) a lower bound for the closest approach of trial and target sequences (open diamonds) derived from secondary structure statistics (Fontana *et al.* 1993a; see this paper, §4); and (iii) longest distances between the reference and the endpoints of monotonously diverging neutral paths (filled circles) (500 reference sequences were used).

No new principle will declare  
itself from below a heap of  
facts.

Sir Peter Medawar, 1985

$$dx_j / dt = \sum_i k_i x_i - x_j \Phi$$

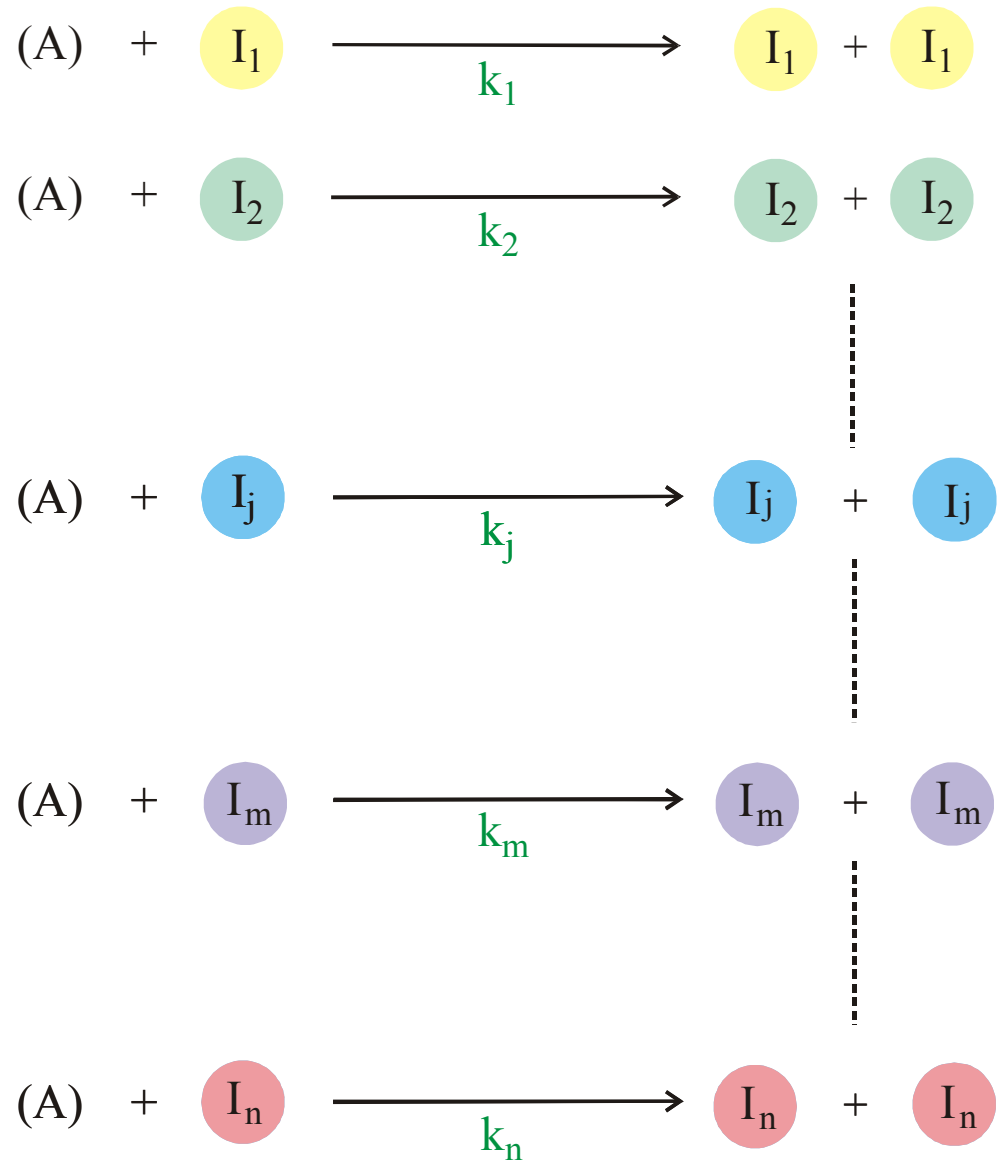
$$\Phi = \sum_i k_i x_i ; \quad \sum_i x_i = 1$$

$$[A] = a = \text{constant}$$

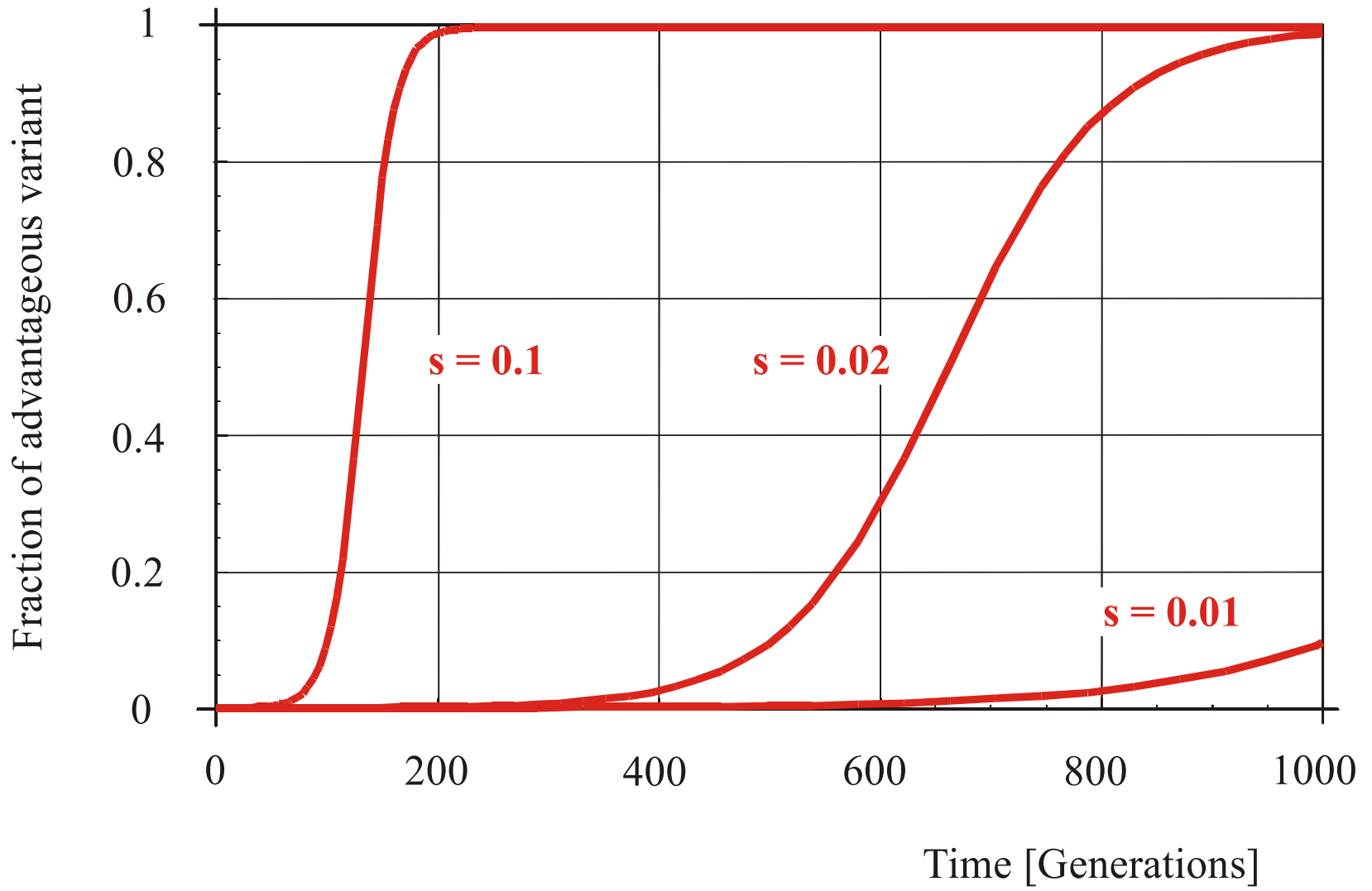
$$k_m = \max \{k_j; j=1,2,\dots,n\}$$

$$x_m(t) \approx 1 \text{ for } t \gg \tau$$

$$s = (k_{m+1} - k_m) / k_m$$



Selection of the „fittest“ or fastest replicating species



Selection of advantageous mutants in populations of  $N = 10\,000$  individuals



# Theory of molecular evolution

M.Eigen, *Self-organization of matter and the evolution of biological macromolecules*. Naturwissenschaften **58** (1971), 465-526

M.Eigen, P.Schuster, *The hypercycle. A principle of natural self-organization. Part A: Emergence of the hypercycle*. Naturwissenschaften **58** (1977), 465-526

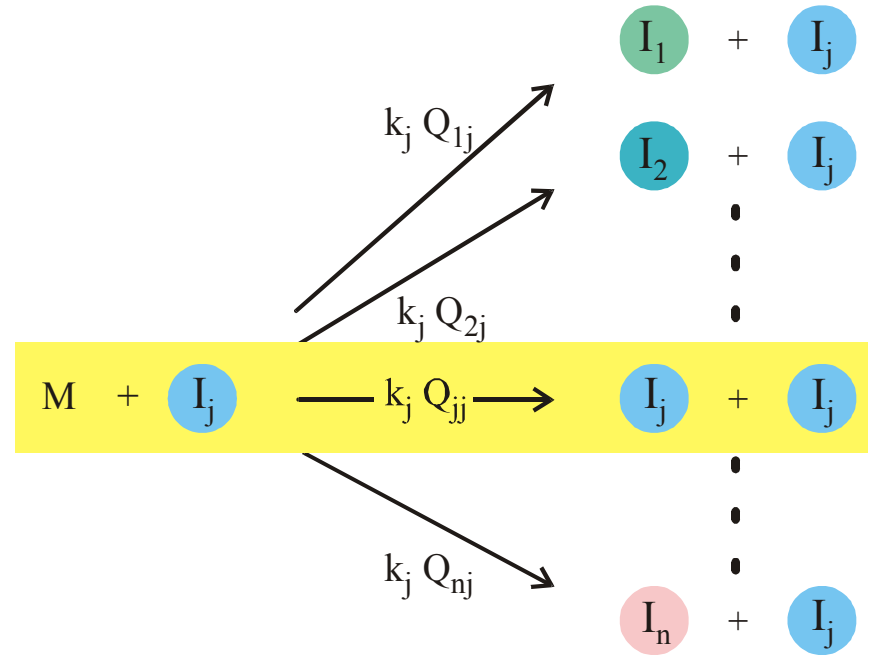
M.Eigen, P.Schuster, *The hypercycle. A principle of natural self-organization. Part B: The abstract hypercycle*. Naturwissenschaften **65** (1978), 7-41

M.Eigen, P.Schuster, *The hypercycle. A principle of natural self-organization. Part C: The realistic hypercycle*. Naturwissenschaften **65** (1978), 341-369

J.S.McCaskill, *A localization threshold for macromolecular quasi-species from continuously distributed replication rates*. J.Chem.Phys. **80** (1984), 5194-5205

M.Eigen, J.McCaskill, P.Schuster, *The molecular quasispecies*. Adv.Chem.Phys. **75** (1989), 149-263

C. Reidys, C.Forst, P.Schuster, *Replication and mutation on neutral networks*. Bull.Math.Biol. **63** (2001), 57-94



$$\sum_i Q_{ij} = 1$$

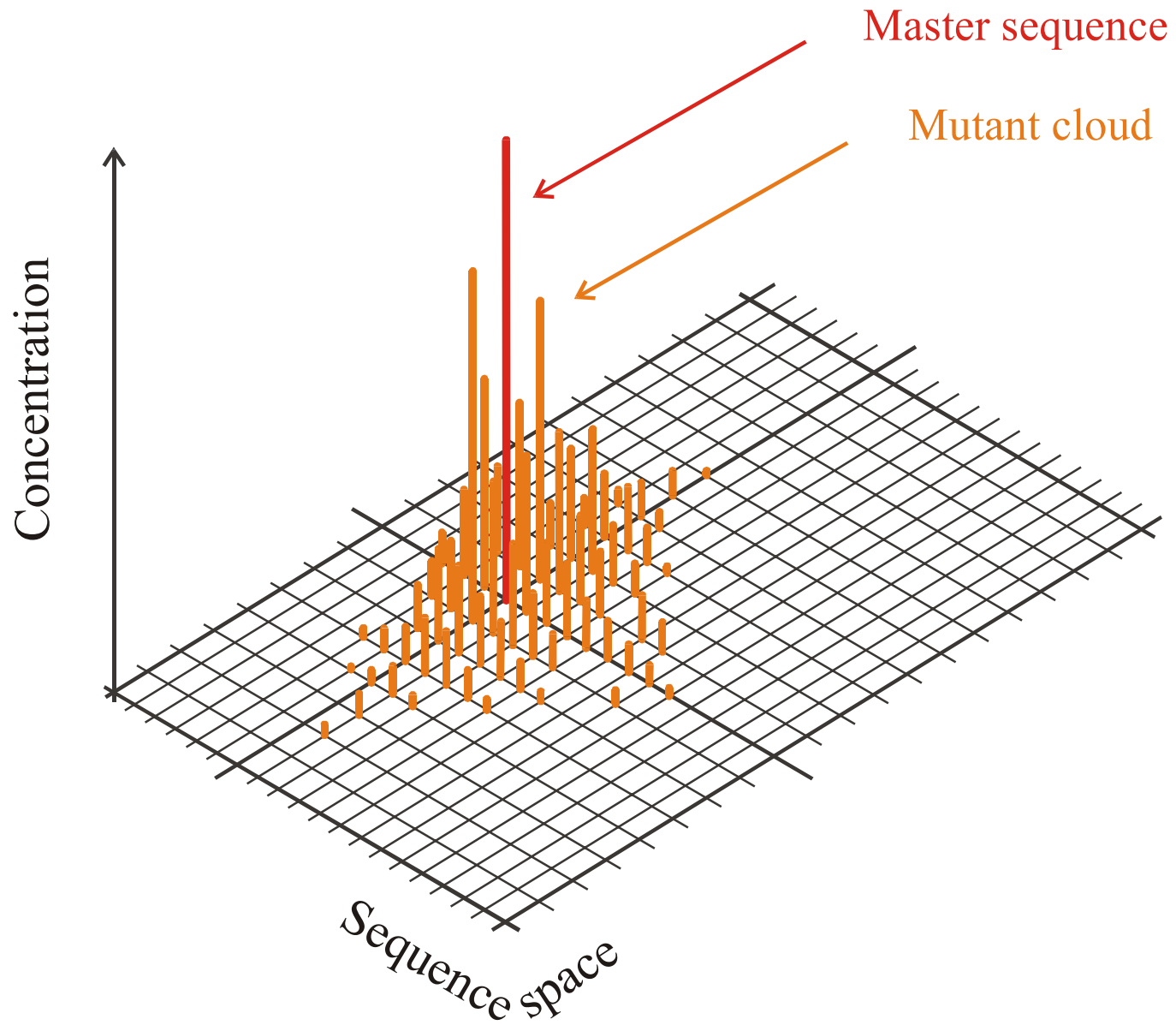
$$Q_{ij} = (1-p)^{n-d(i,j)} p^{d(i,j)} ; \quad p \text{ ..... error rate per digit}$$

$d(i,j)$  ..... Hamming distance between  $I_i$  and  $I_j$

$$dx_j / dt = \sum_i k_i Q_{ji} x_i - x_j \Phi$$

$$\Phi = \sum_i k_i x_i ; \quad \sum_i x_i = 1$$

Chemical kinetics of replication  
and mutation as parallel reactions



The molecular quasispecies in sequence space

The **RNA model** considers **RNA sequences** as **genotypes** and simplified **RNA structures**, called **secondary structures**, as **phenotypes**.

Variation is restricted to **point mutations**.

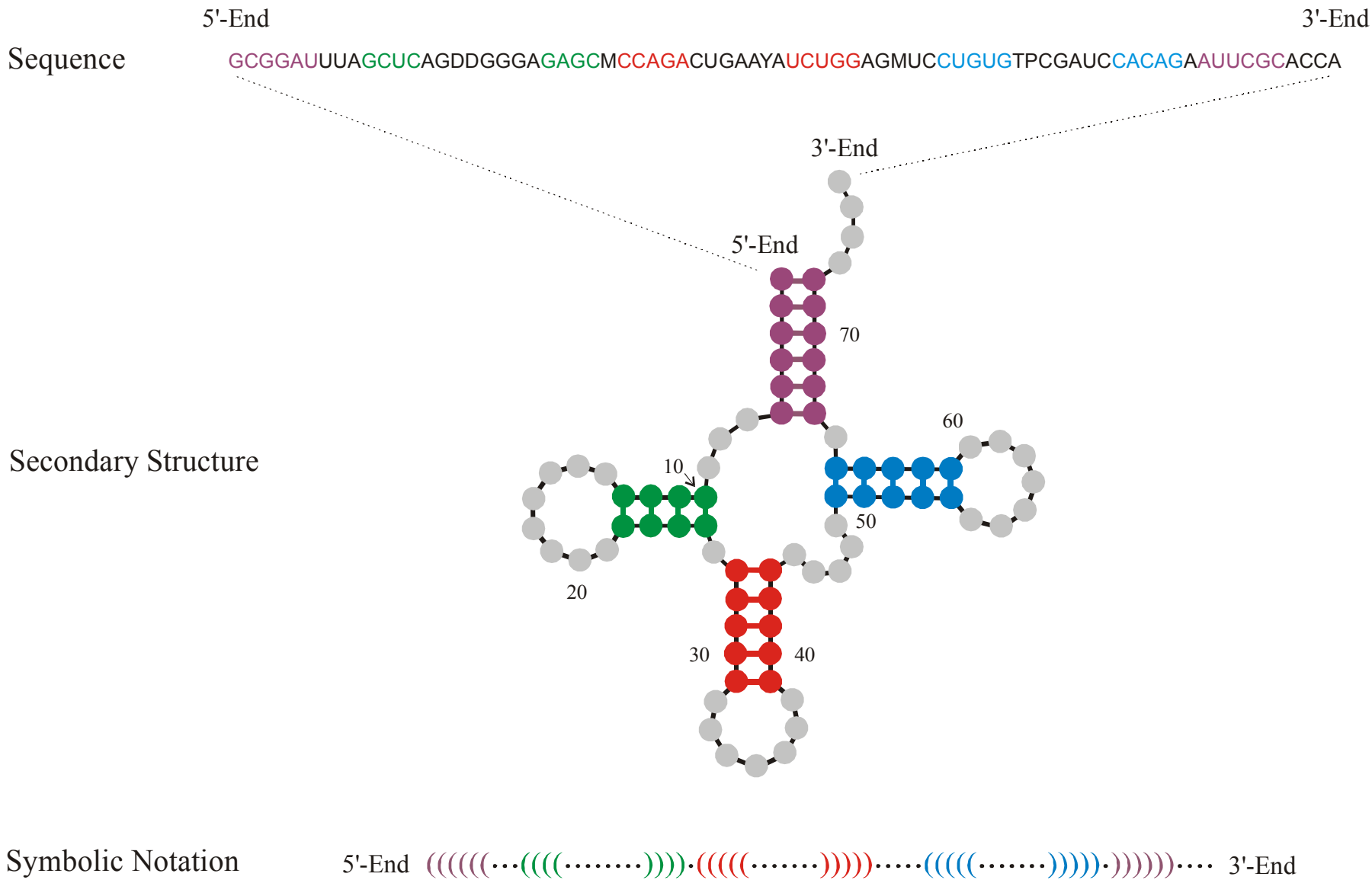
The **mapping** from genotypes into phenotypes is many-to-one. Hence, it is redundant and not invertible.

Genotypes, i.e. RNA sequences, which are mapped onto the same phenotype, i.e. the same RNA secondary structure, form **neutral networks**. Neutral networks are represented by graphs in sequence space.

## **RNA secondary structures and their properties**

RNA secondary structures are **listings of** Watson-Crick and GU wobble **base pairs**, which are free of knots and pseudoknots. Secondary structures are **folding intermediates** in the formation of full three-dimensional structures.

D.Thirumalai, N.Lee, S.A.Woodson, and D.K.Klimov.  
*Annu.Rev.Phys.Chem.* **52**:751-762 (2001)



Definition and formation of the secondary structure of phenylalanyl-tRNA

## RNA minimum free energy structures

Efficient algorithms based on dynamical programming are available for computation of secondary structures for given sequences. Inverse folding algorithms compute sequences for given secondary structures.

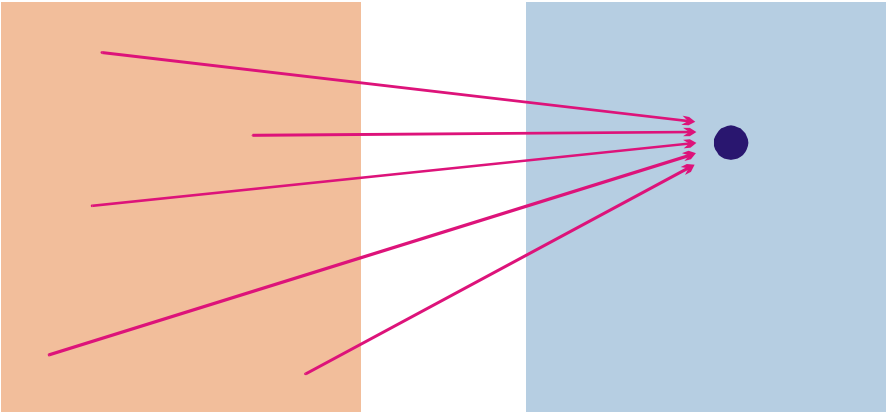
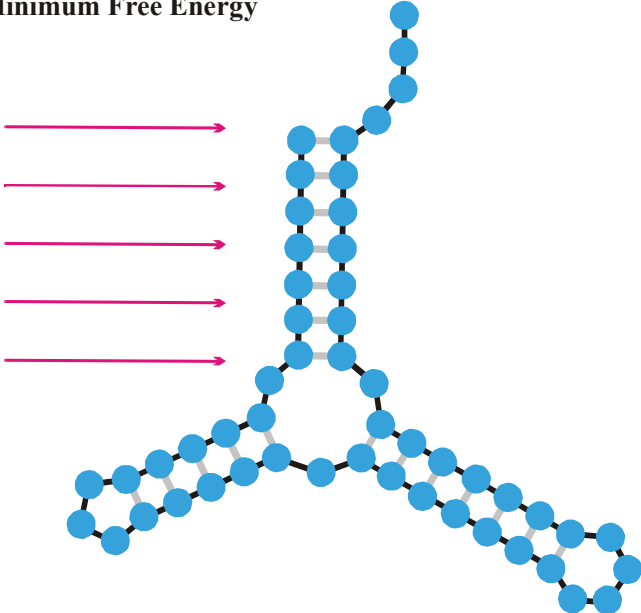
M.Zuker and P.Stiegler. *Nucleic Acids Res.* **9**:133-148 (1981)

**Vienna RNA Package:** <http://www.tbi.univie.ac.at> (includes **inverse folding**, **suboptimal structures**, **kinetic folding**, etc.)

I.L.Hofacker, W. Fontana, P.F.Stadler, L.S.Bonhoeffer,  
M.Tacker, and P. Schuster. *Mh.Chem.* **125**:167-188 (1994)

**Criterion of  
Minimum Free Energy**

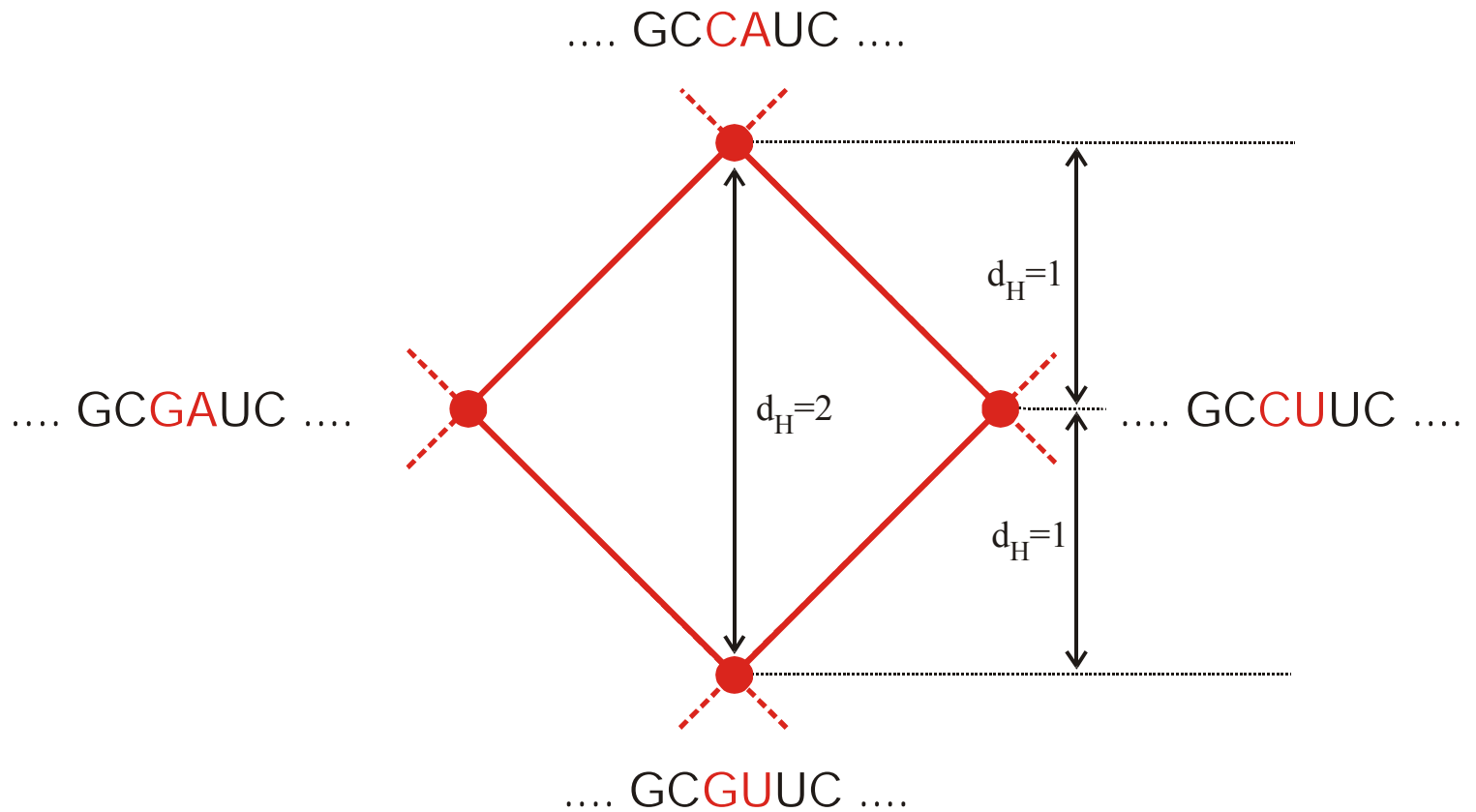
UUUAGCCAGCGCGAGUCGUGCGGACGGGGUUAUCUCUGUCGGGCUAGGGCGC  
GUGAGCGCGGGGCACAGUUUCUCAAGGAUGUAAGUUUUUGCCGUUUUUCUGG  
UUAGCGAGAGAGGAGGCUUCUAGACCCAGCUCUCUGGGUCGUUGCUGAUGCG  
CAUUGGUGCUAAUGAUUUAGGGCUGUAUJCCUGUAUAGCGAUCAGUGUCCG  
GUAGGCCCUUUGACAUAAGAUUUUUCCAUGGGUGGAGAUGGCCAUUGCAG



Sequence Space

Shape Space





Point mutations as moves in sequence space

$S_1$ : CGTCGTTACAATTTA **G**GTTATGTGCGAATTC **A**CAAATT **G**AAAA **T**ACAAGAG . . . . .  
 $S_2$ : CGTCGTTACAATTTA **A**GTTATGTGCGAATTC **C**CAAATT **A**AAAA **C**ACAAGAG . . . . .

Hamming distance  $d_H(S_1, S_2) = 4$

- (i)  $d_H(S_1, S_1) = 0$
- (ii)  $d_H(S_1, S_2) = d_H(S_2, S_1)$
- (iii)  $d_H(S_1, S_3) < d_H(S_1, S_2) + d_H(S_2, S_3)$

The Hamming distance induces a metric in sequence space

## Mutant class

0

1

2

3

4

5

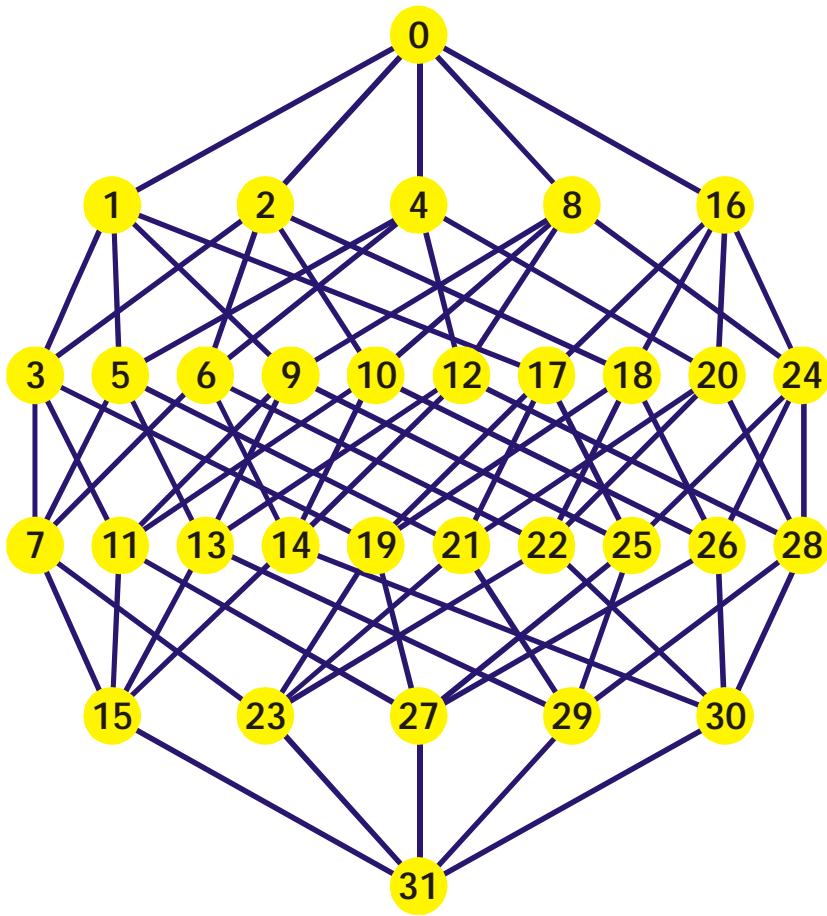
Binary sequences are encoded by their decimal equivalents:

C = 0 and G = 1, for example,

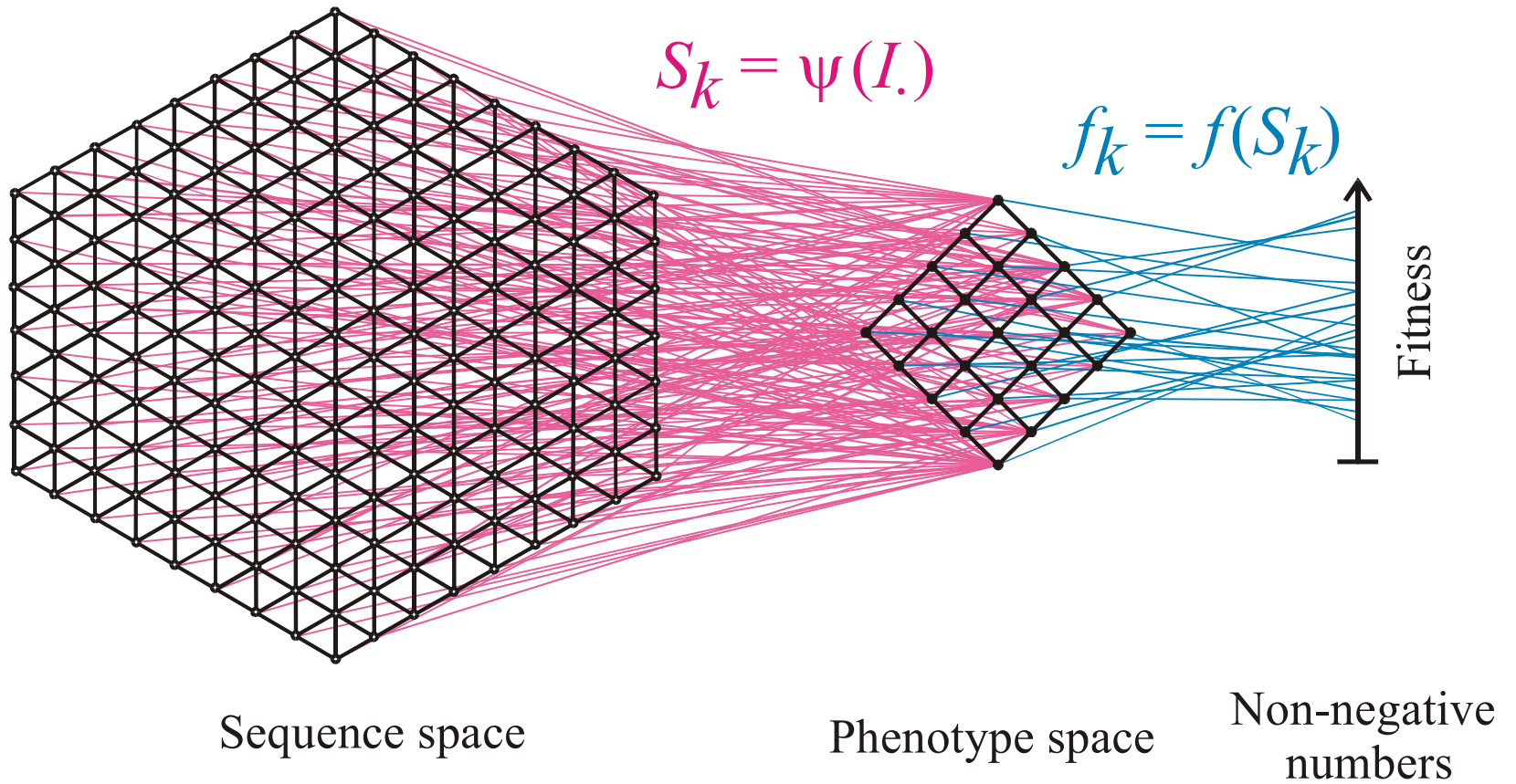
"0"  $\equiv$  00000 = CCCCC,

"14"  $\equiv$  01110 = CGGGC,

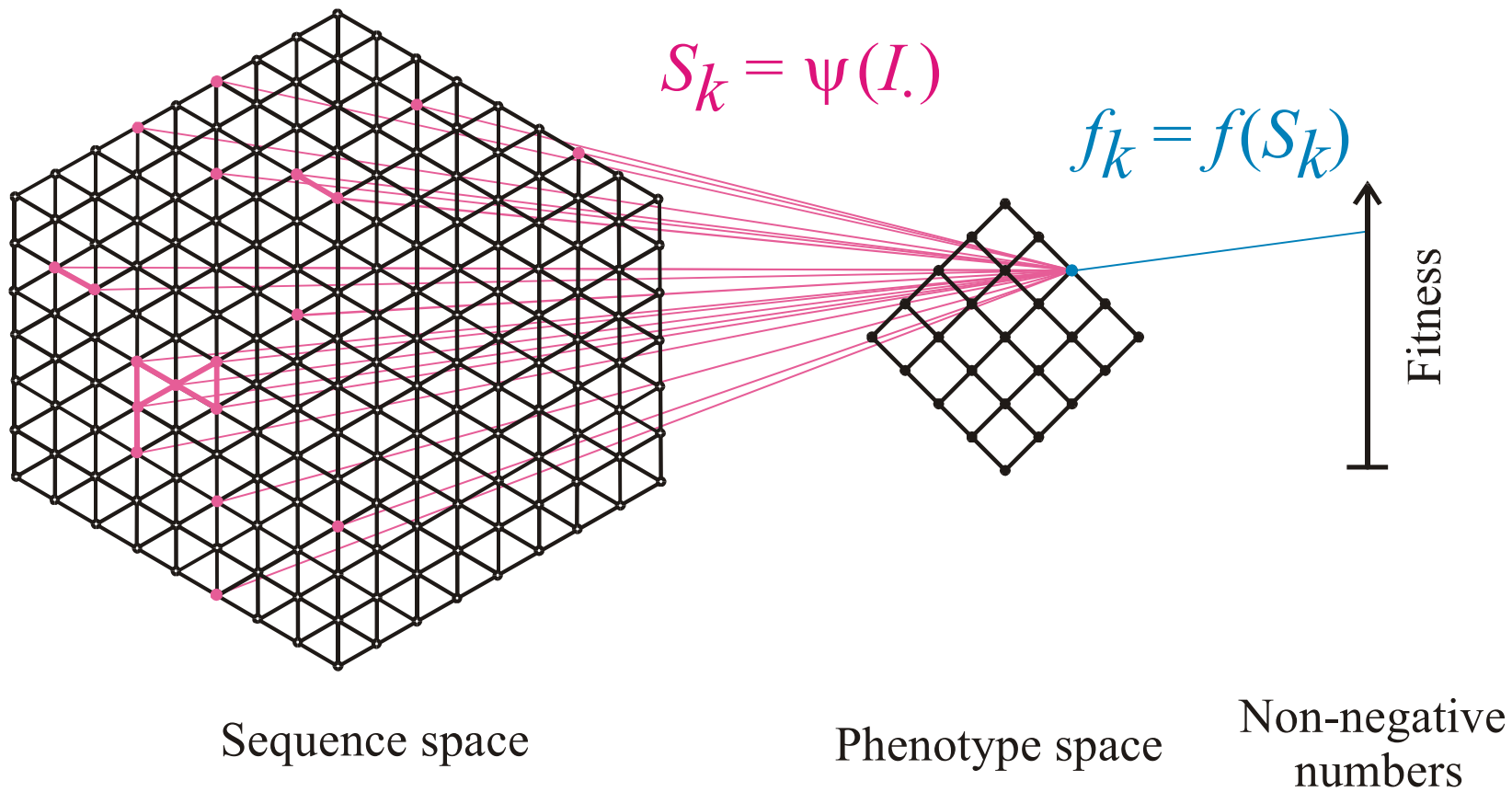
"29"  $\equiv$  11101 = GGGCG, etc.

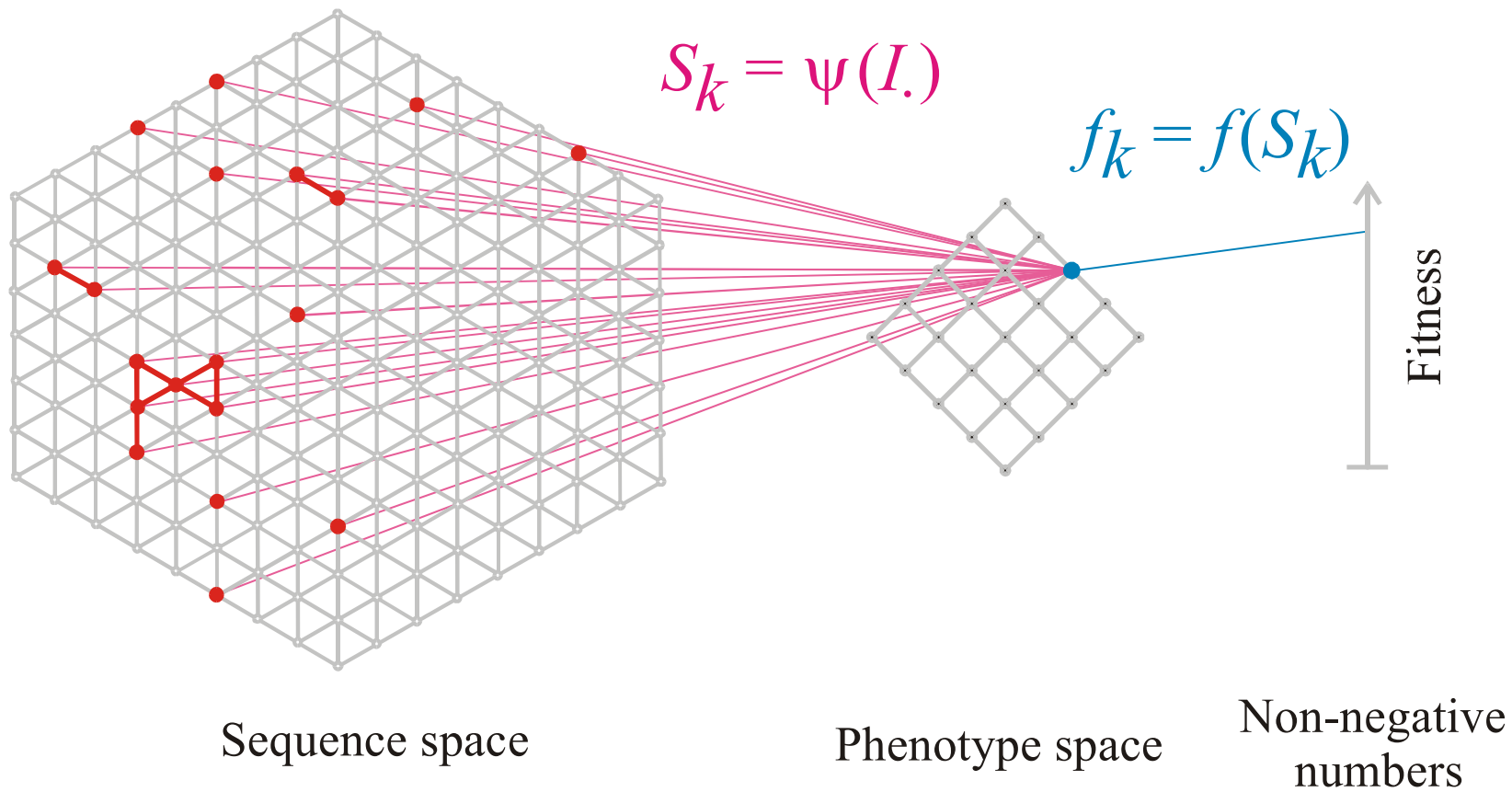


Sequence space of binary sequences of chain length  $n=5$



Mapping from sequence space into phenotype space and into fitness values



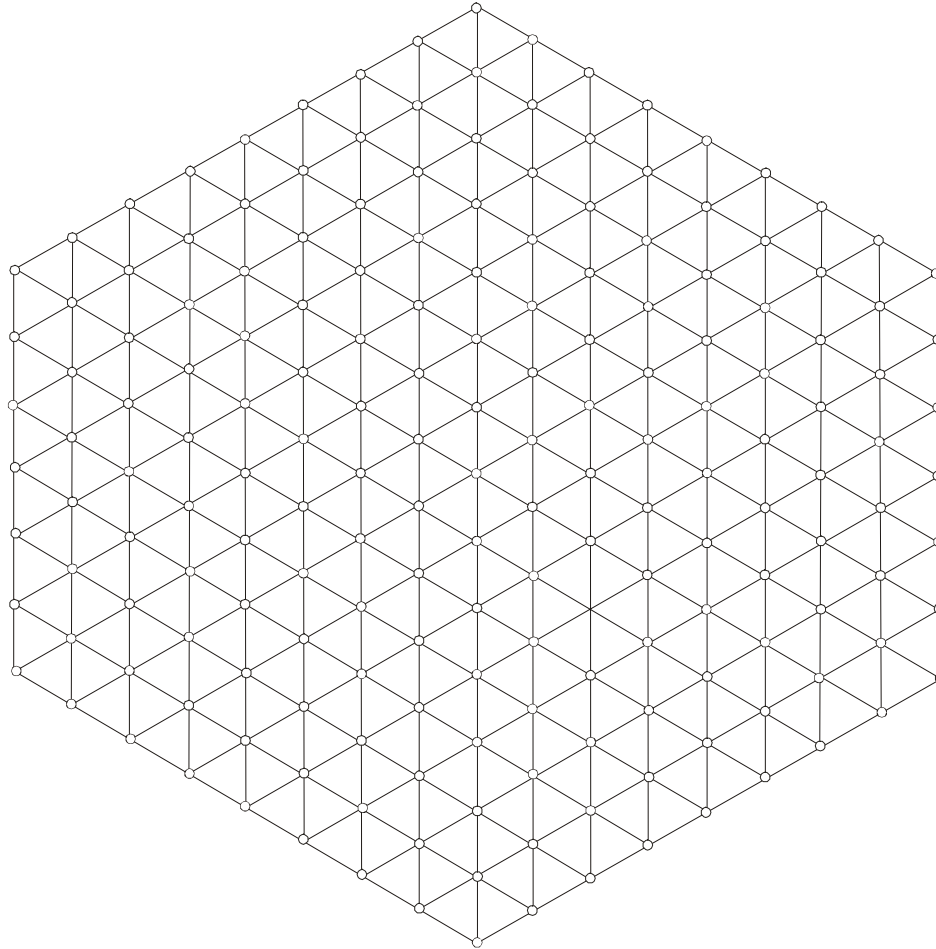


Neutral networks of small RNA molecules can be computed by exhaustive folding of complete sequence spaces, i.e. all RNA sequences of a given chain length. This number,  $N=4^l$ , becomes very large with increasing length, and is prohibitive for numerical computations.

Neutral networks can be modelled by **random graphs** in sequence space. In this approach, nodes are inserted randomly into sequence space until the size of the pre-image, i.e. the number of neutral sequences, matches the neutral network to be studied.

Step 00

Sketch of sequence space

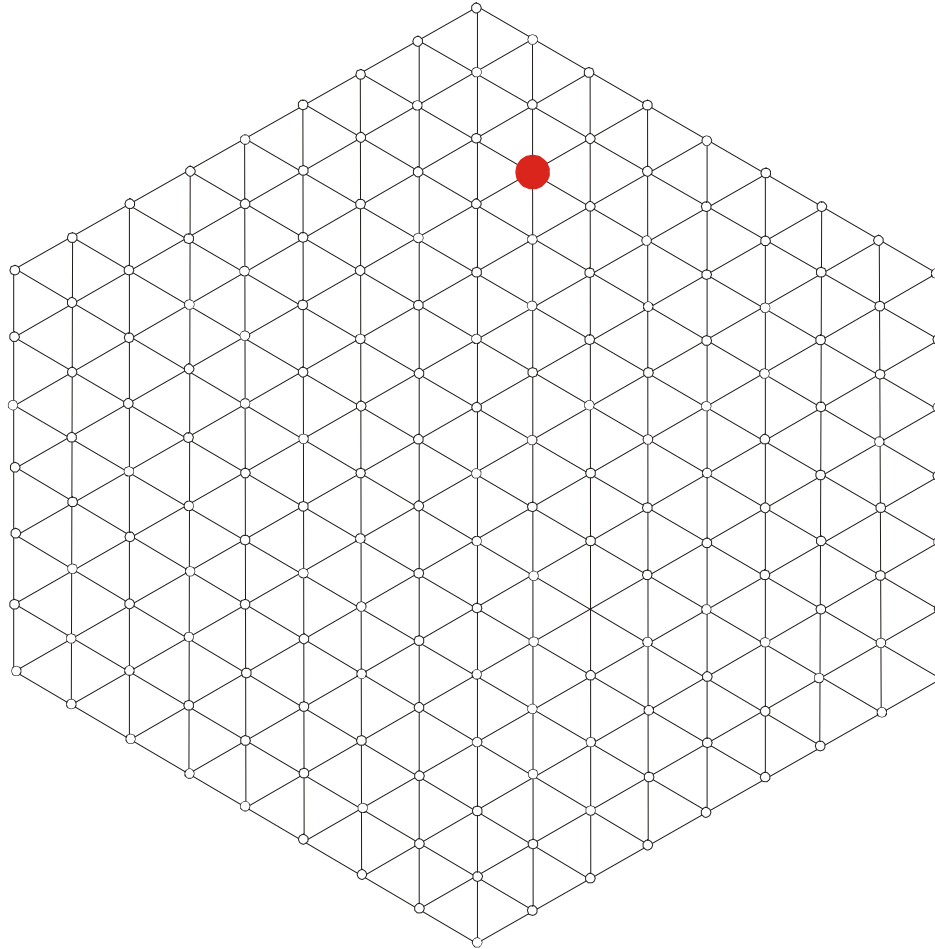


Random graph approach to neutral networks



Step 01

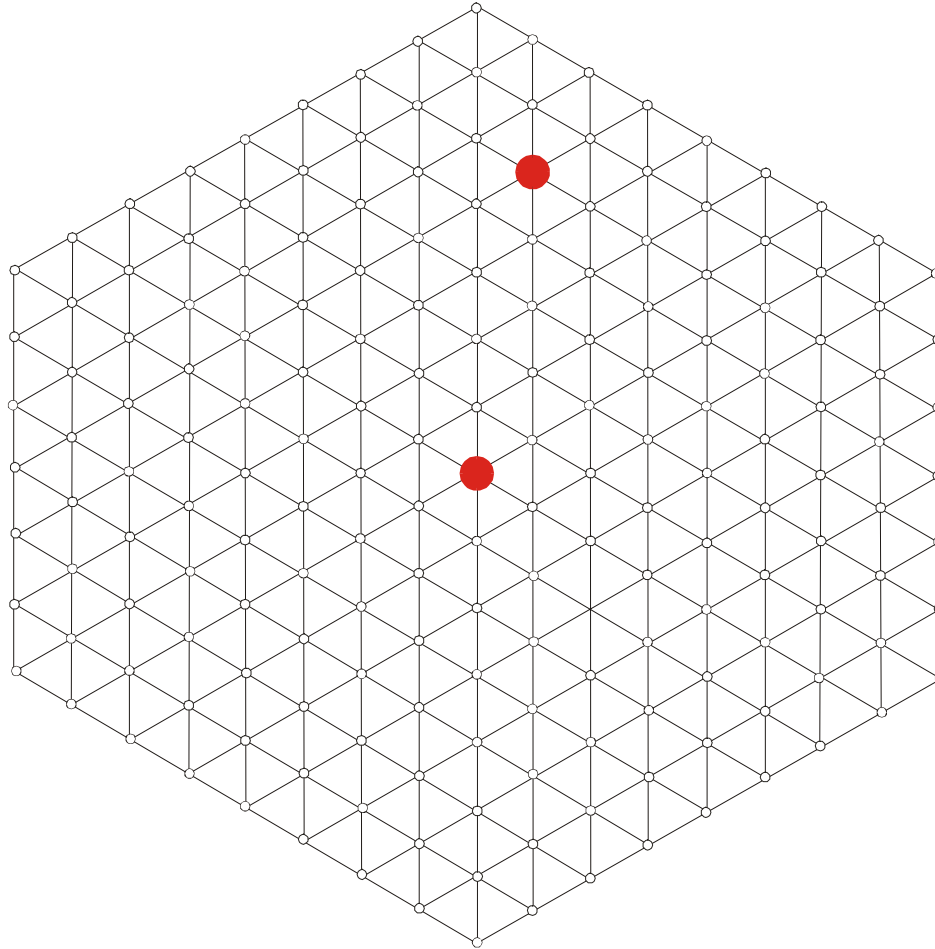
Sketch of sequence space



Random graph approach to neutral networks

Step 02

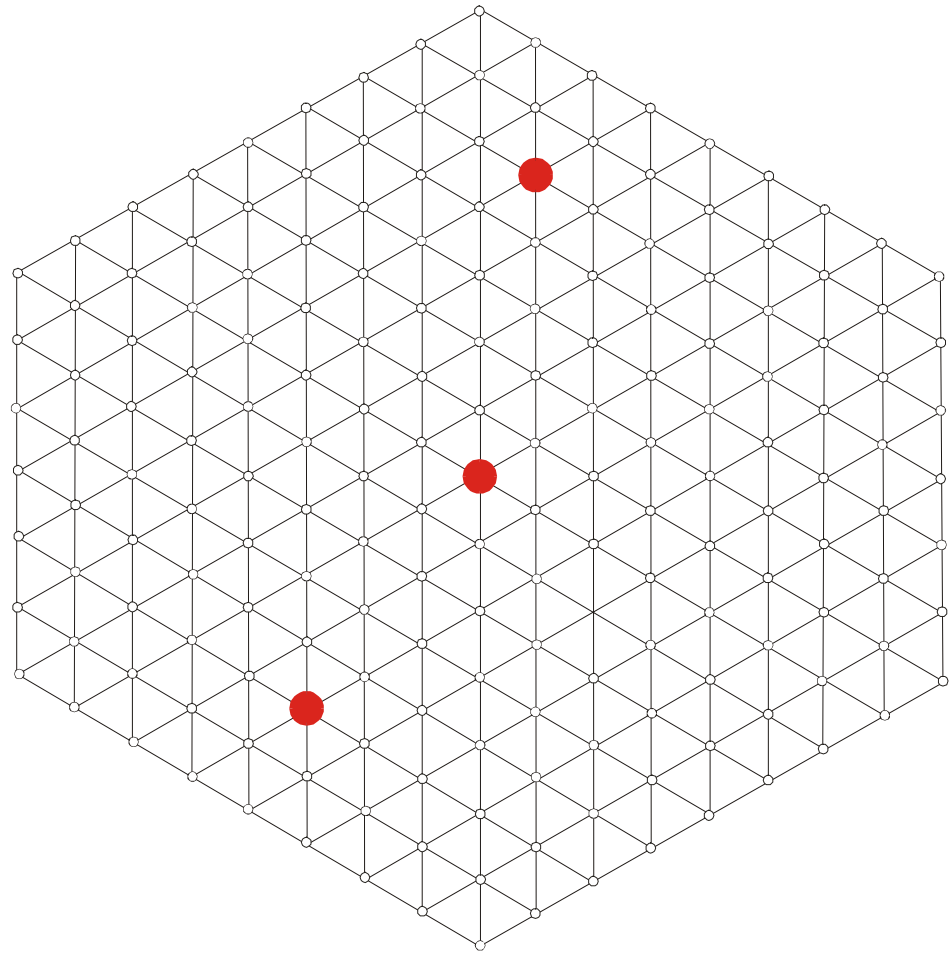
Sketch of sequence space



Random graph approach to neutral networks

Step 03

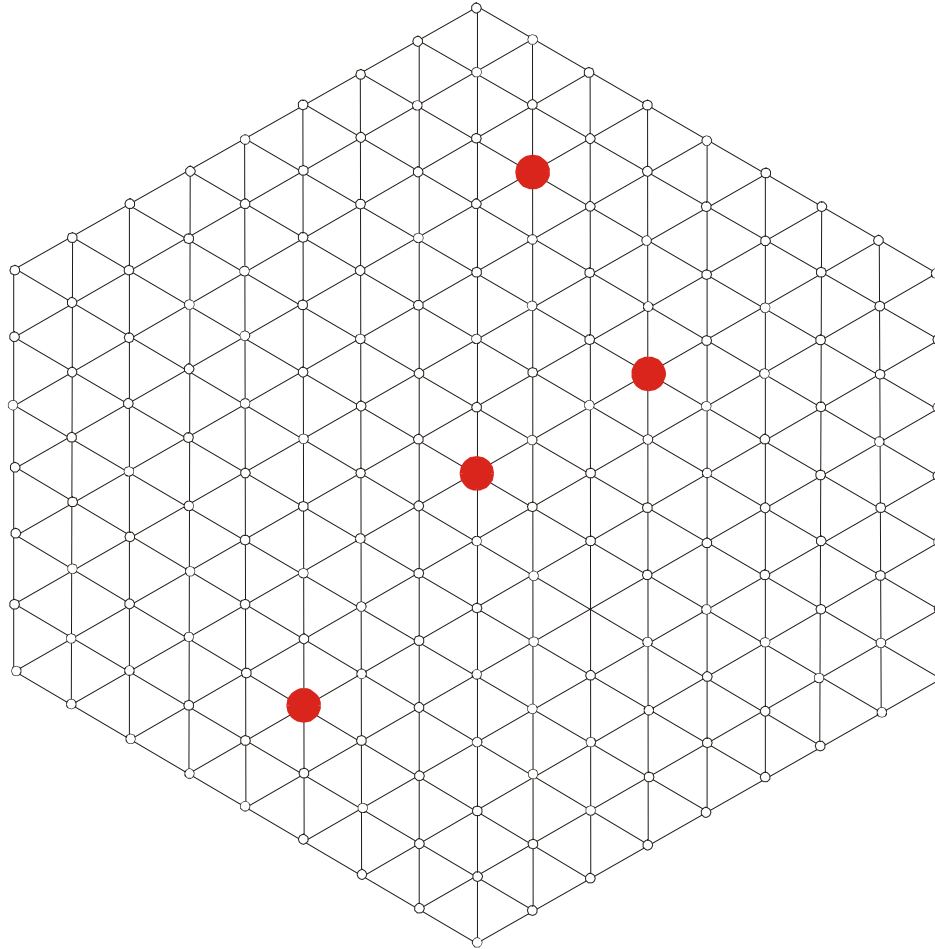
Sketch of sequence space



Random graph approach to neutral networks

Step 04

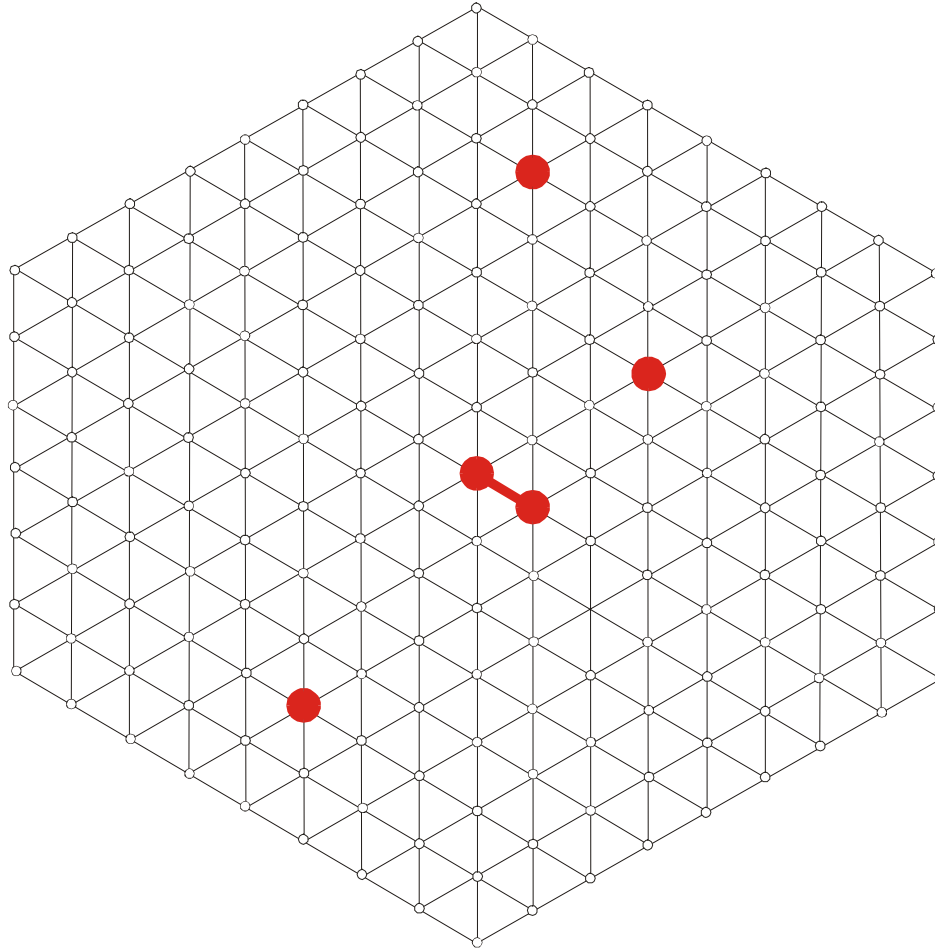
Sketch of sequence space



Random graph approach to neutral networks

Step 05

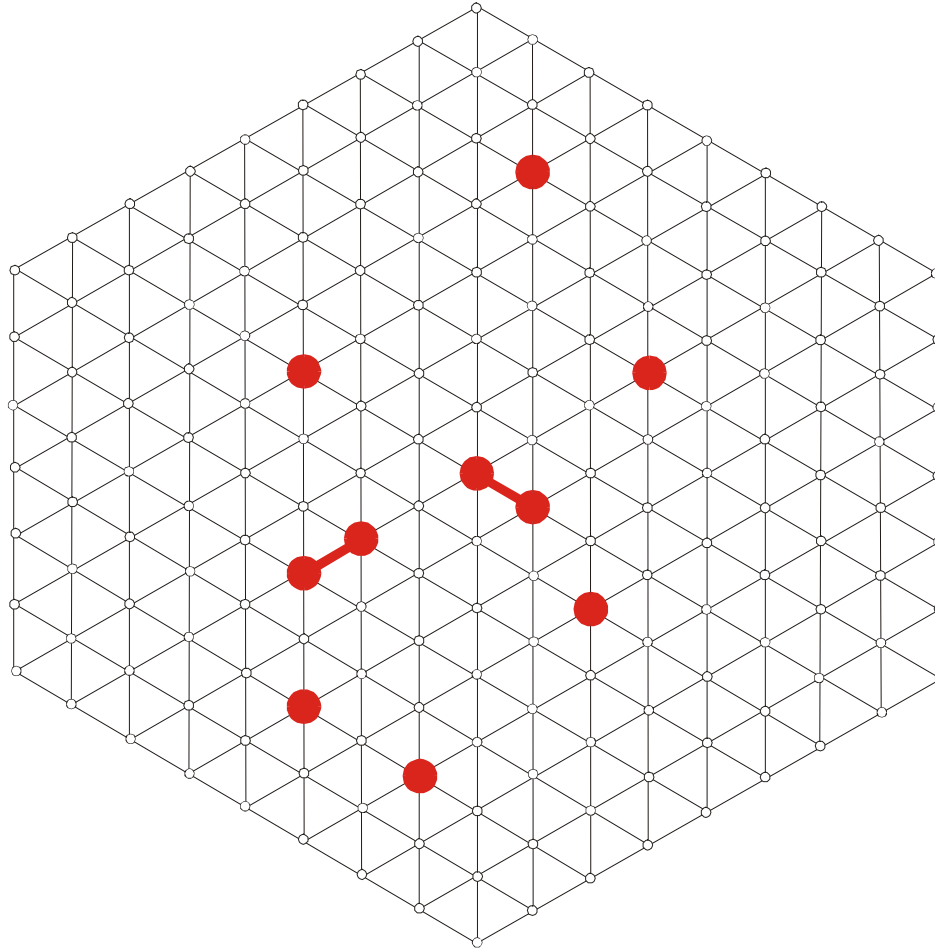
Sketch of sequence space



Random graph approach to neutral networks

Step 10

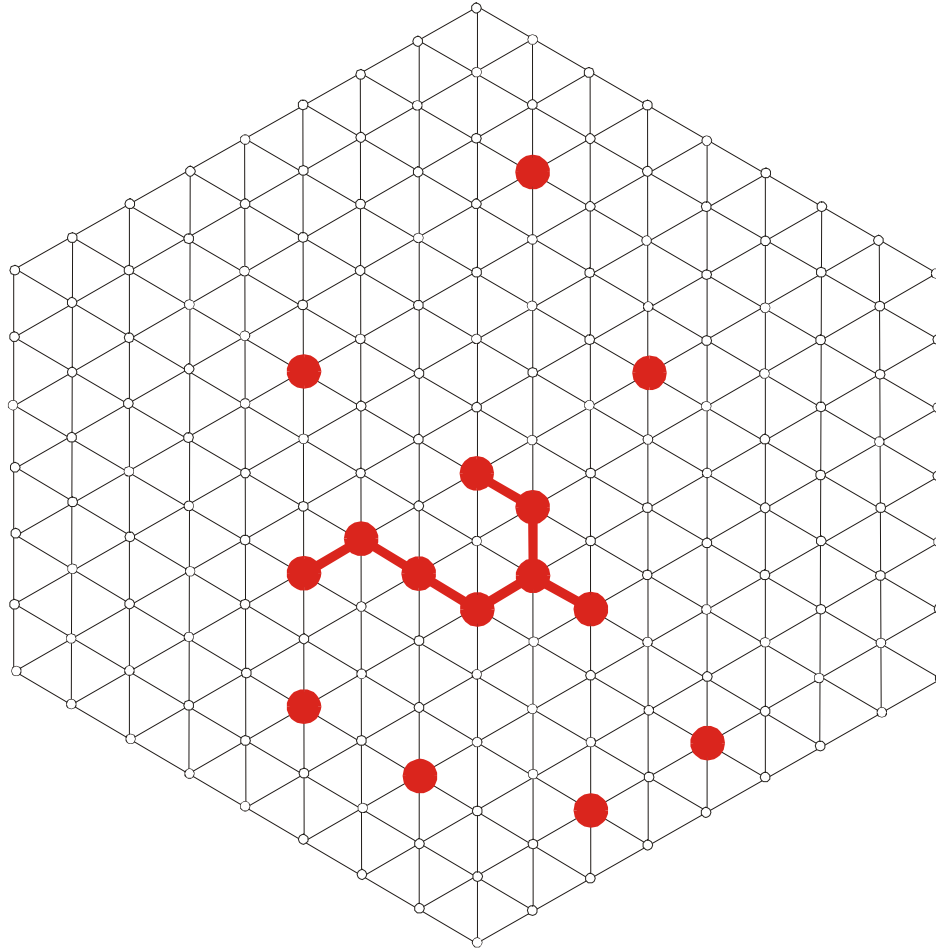
Sketch of sequence space



Random graph approach to neutral networks

Step 15

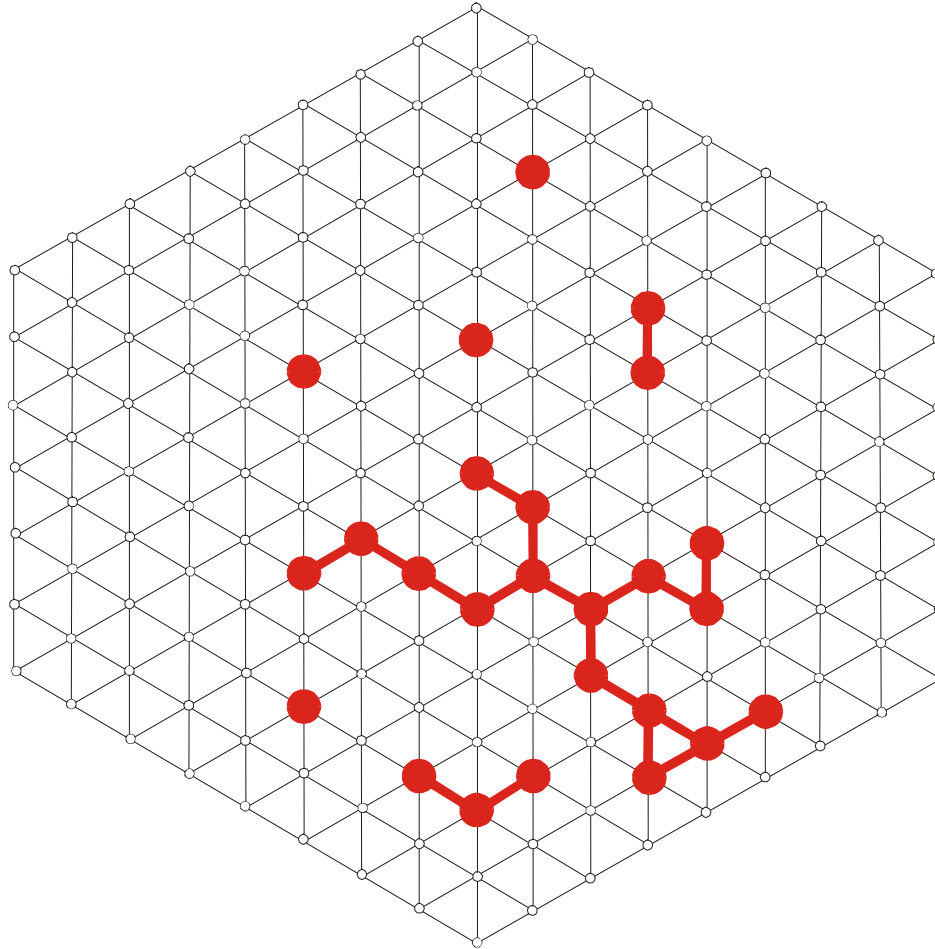
Sketch of sequence space



Random graph approach to neutral networks

Step 25

Sketch of sequence space

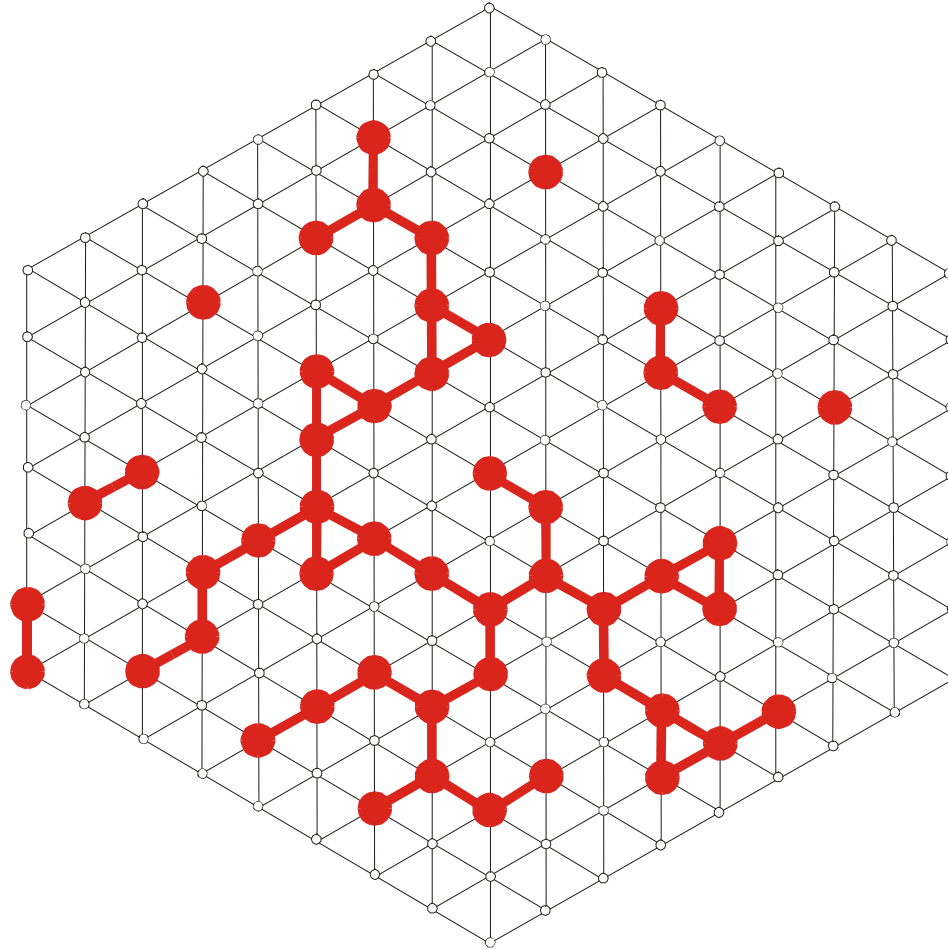


Random graph approach to neutral networks



Step 50

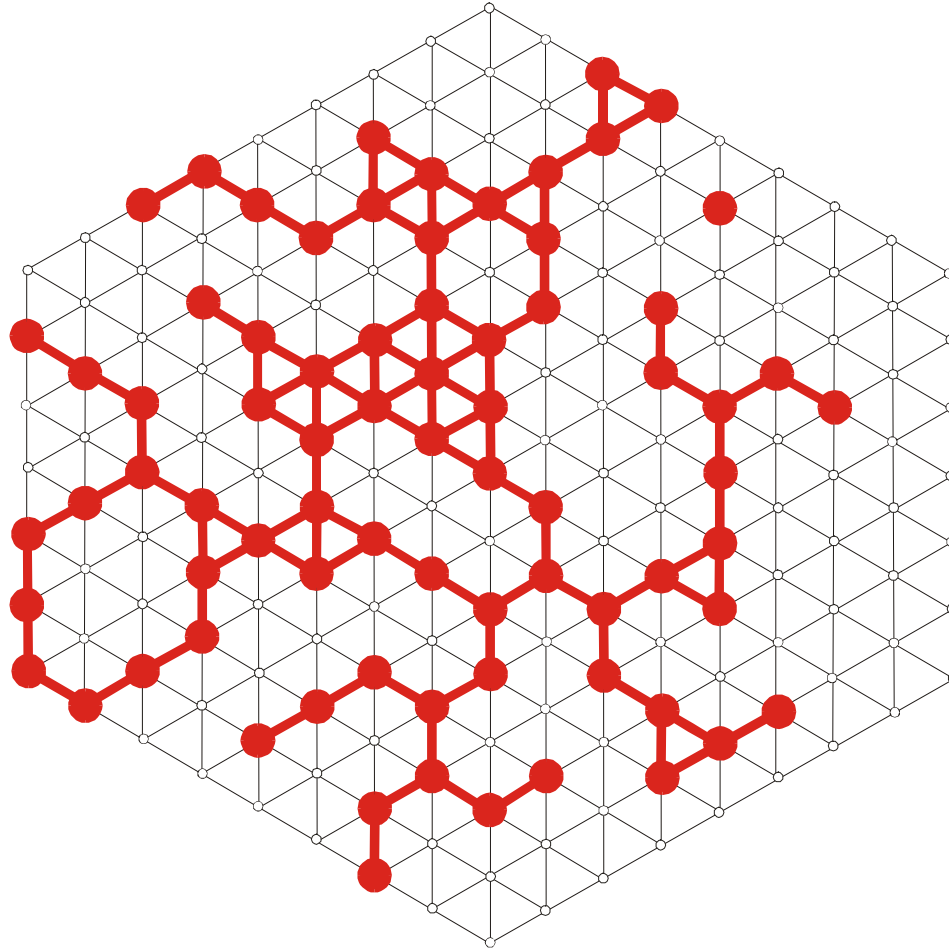
Sketch of sequence space



Random graph approach to neutral networks

Step 75

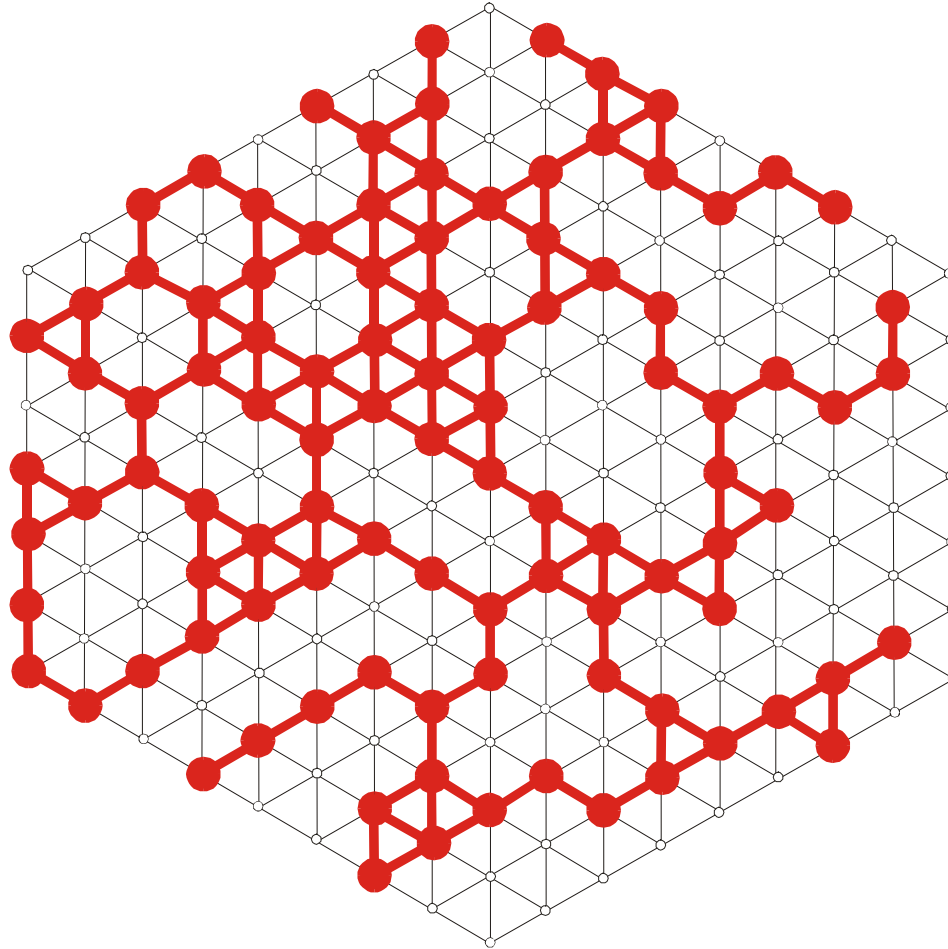
Sketch of sequence space



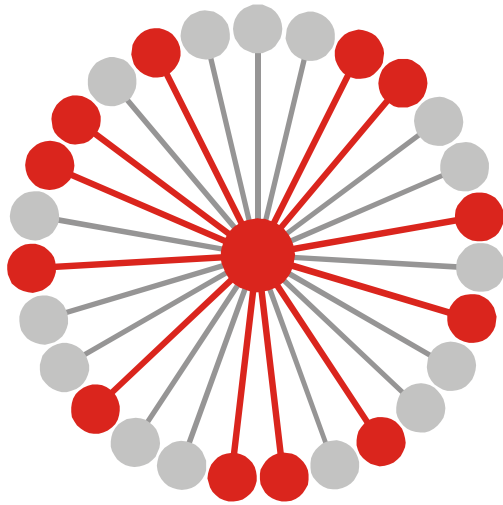
Random graph approach to neutral networks

Step 100

Sketch of sequence space



Random graph approach to neutral networks



$$G_k = m^{-1}(S_k) \cup \{I_j \mid m(I_j) = S_k\} \cup q$$

$$\lambda_j = 12 / 27, \quad \bar{\lambda}_k = \frac{\sum_{j \in |G_k|} \hat{\lambda}_j(k)}{|G_k|}$$

Connectivity threshold:  $\lambda_{cr} = 1 - \kappa^{-1}/(\kappa-1)$

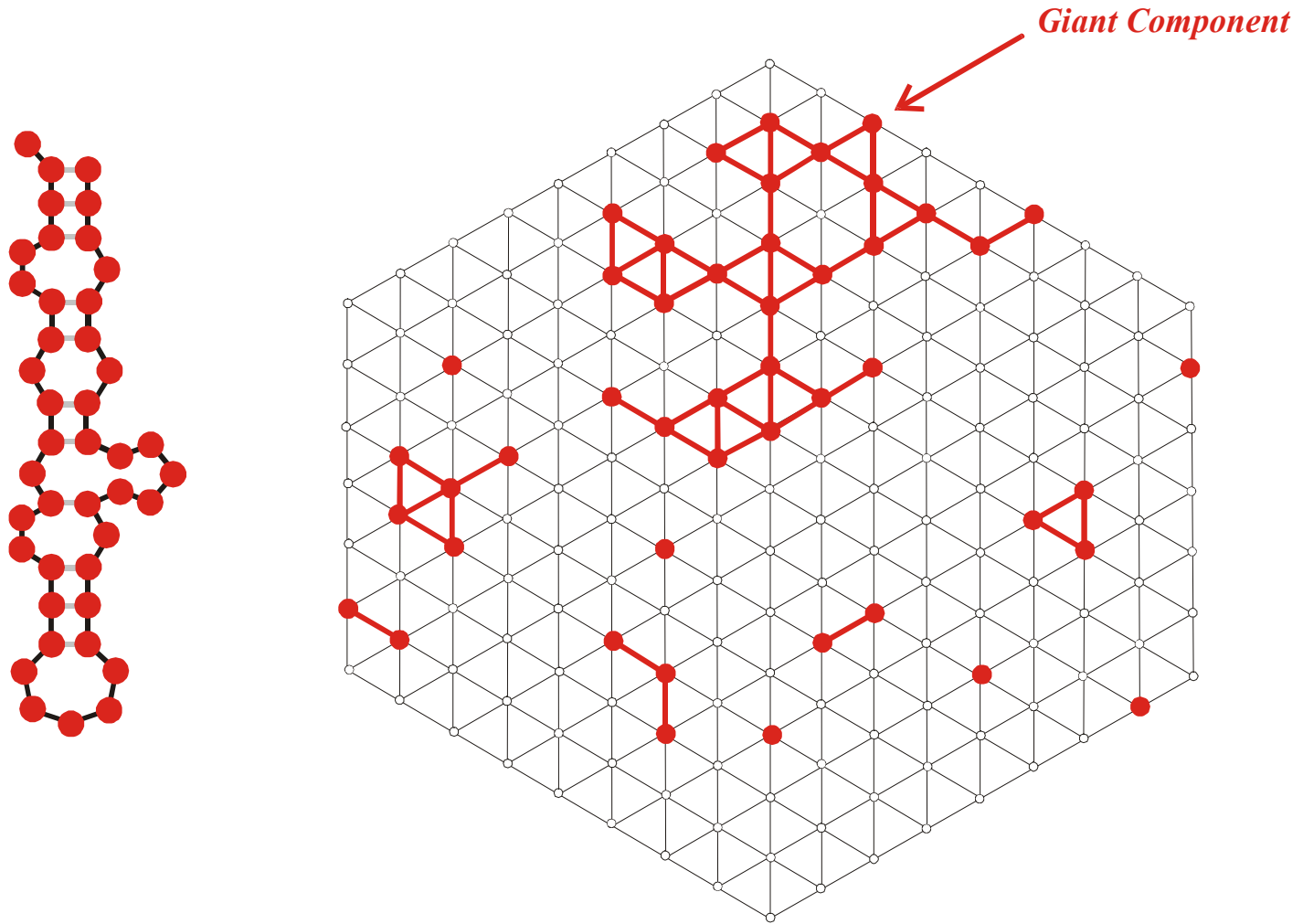
Alphabet size  $\kappa$ : **AUGC**  $\kappa = 4$

$\bar{\lambda}_k > \lambda_{cr}$  ..... network  $G_k$  is connected

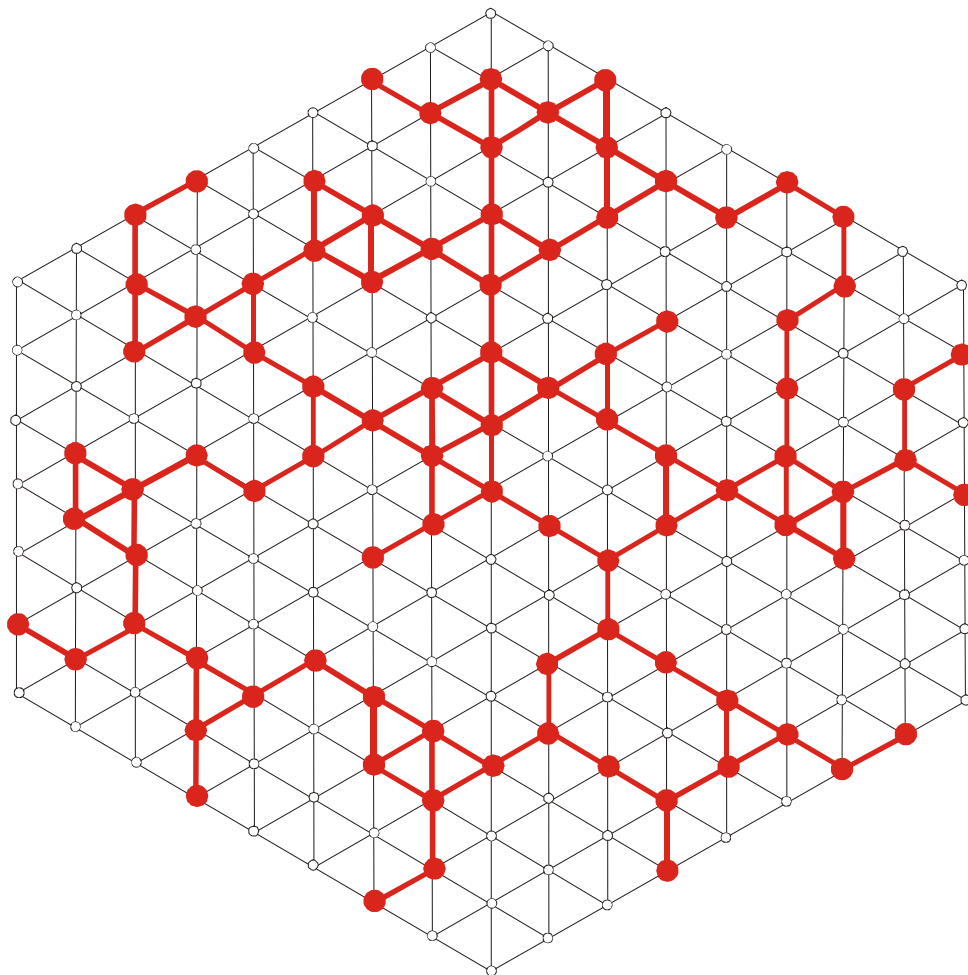
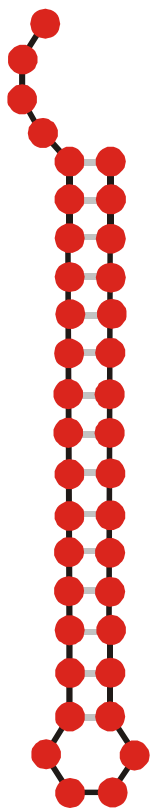
$\bar{\lambda}_k < \lambda_{cr}$  ..... network  $G_k$  is **not** connected

$\kappa$	$\lambda_{cr}$
2	0.5
3	0.4226
4	0.3700

Mean degree of neutrality and connectivity of neutral networks



A multi-component neutral network



A connected neutral network

## **Optimization of RNA molecules *in silico***

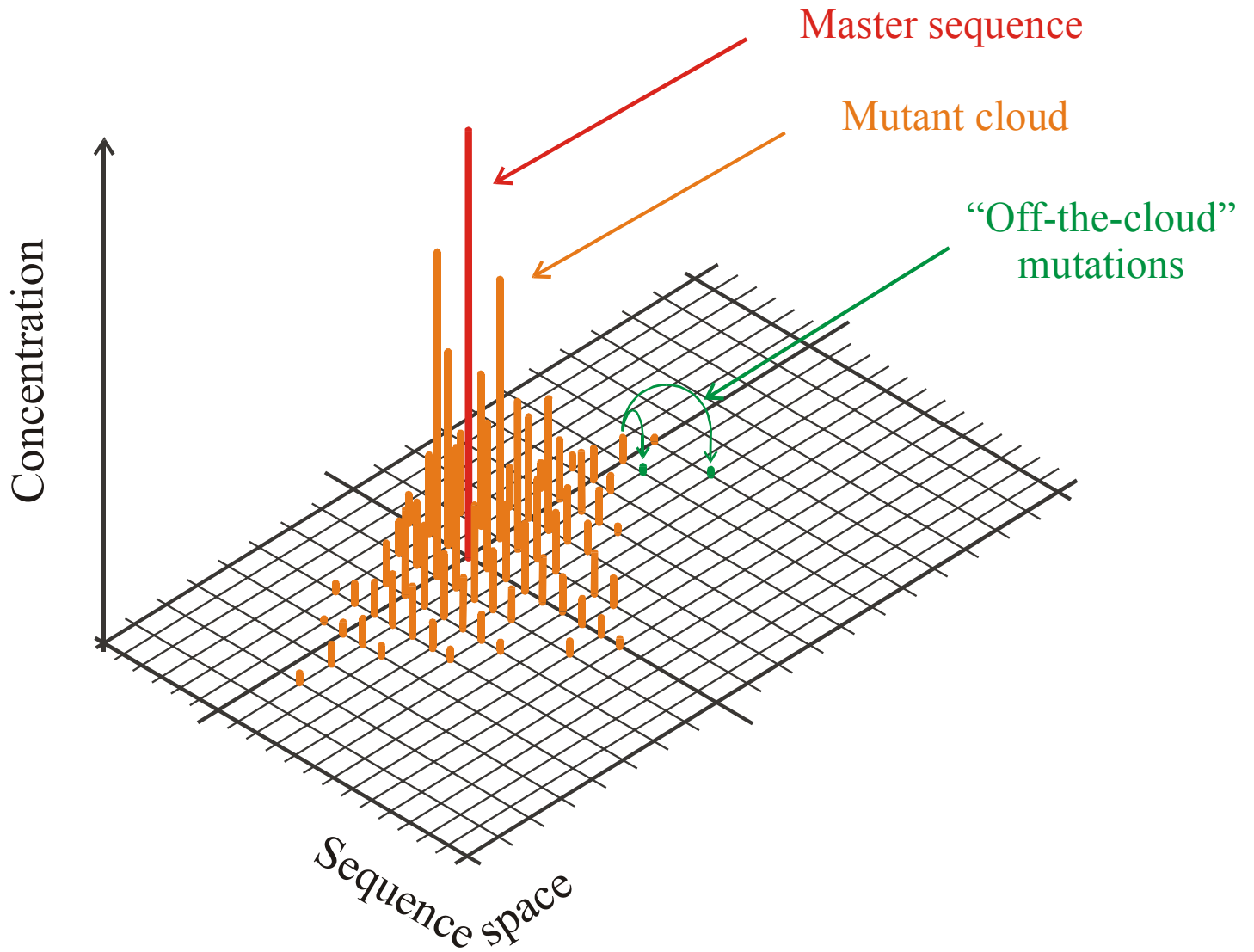
W.Fontana, P.Schuster, *A computer model of evolutionary optimization*. Biophysical Chemistry **26** (1987), 123-147

W.Fontana, W.Schnabl, P.Schuster, *Physical aspects of evolutionary optimization and adaptation*. Phys.Rev.A **40** (1989), 3301-3321

M.A.Huynen, W.Fontana, P.F.Stadler, *Smoothness within ruggedness. The role of neutrality in adaptation*. Proc.Natl.Acad.Sci.USA **93** (1996), 397-401

W.Fontana, P.Schuster, *Continuity in evolution. On the nature of transitions*. Science **280** (1998), 1451-1455

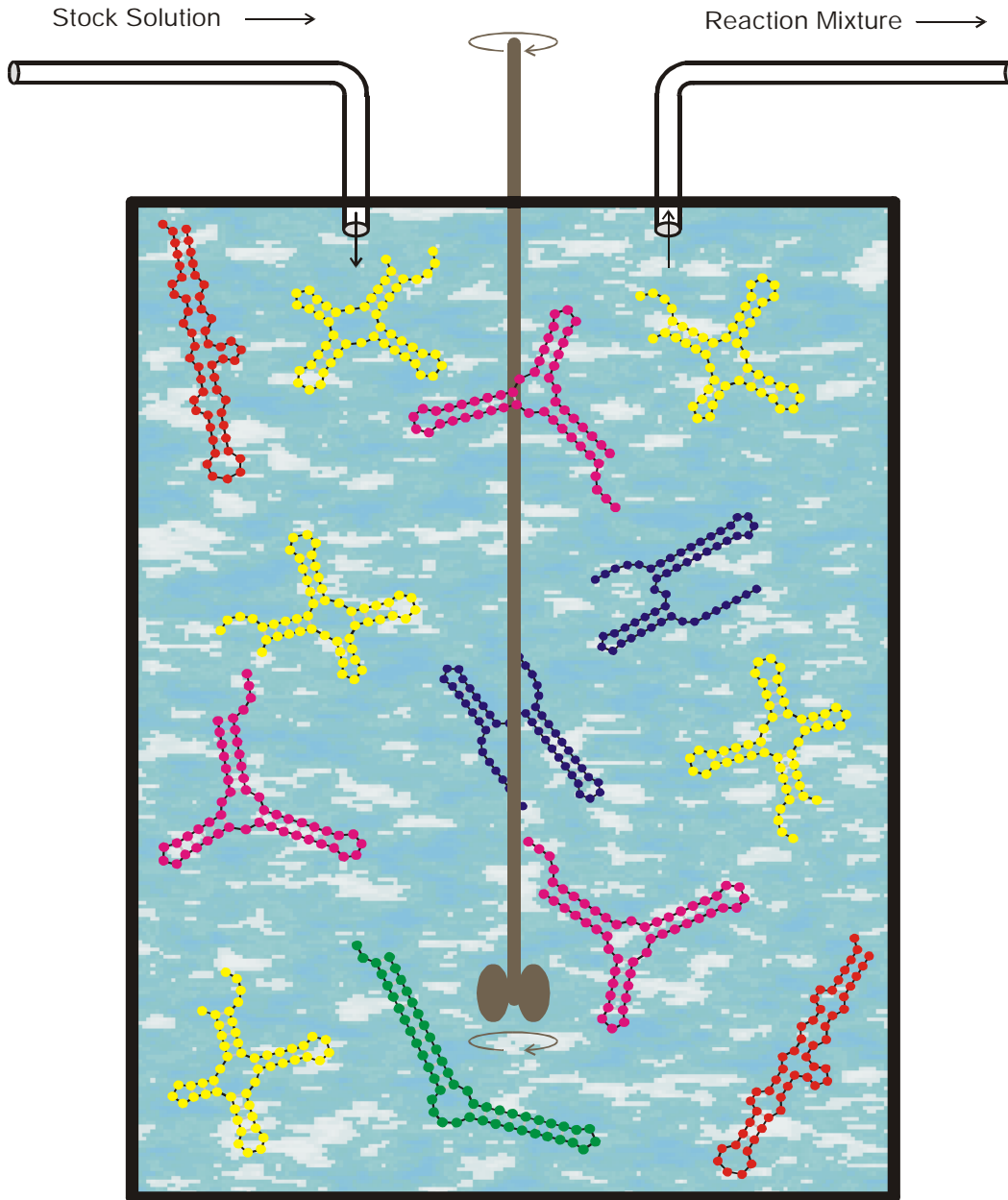
W.Fontana, P.Schuster, *Shaping space. The possible and the attainable in RNA genotype-phenotype mapping*. J.Theor.Biol. **194** (1998), 491-515



The molecular quasispecies  
in sequence space





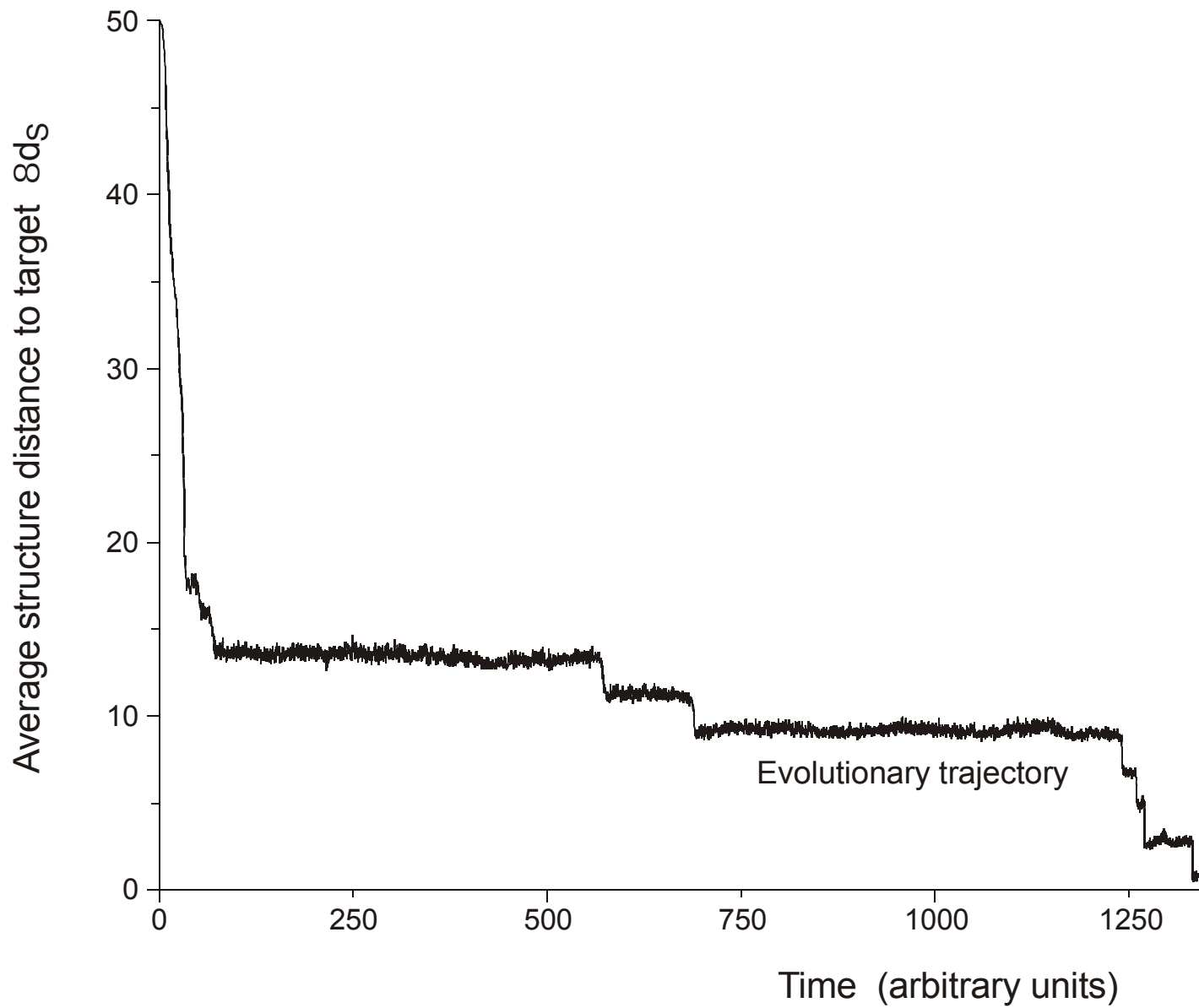


Fitness function for optimization in the flow reactor:

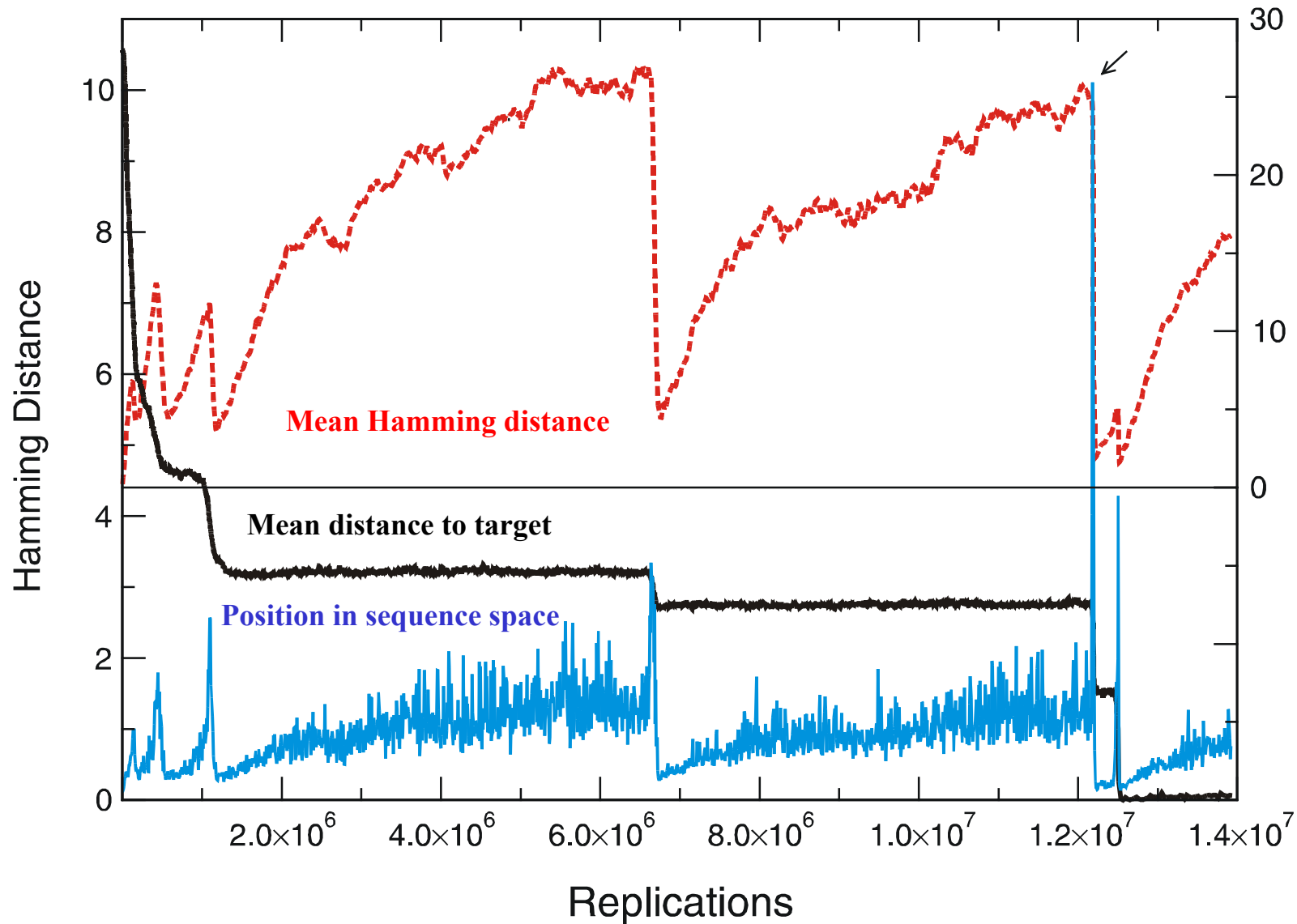
$$f_k = [ / [U + \delta d_S^{(k)}]$$

$$\delta d_S^{(k)} = d^s(I_k, I_h)$$

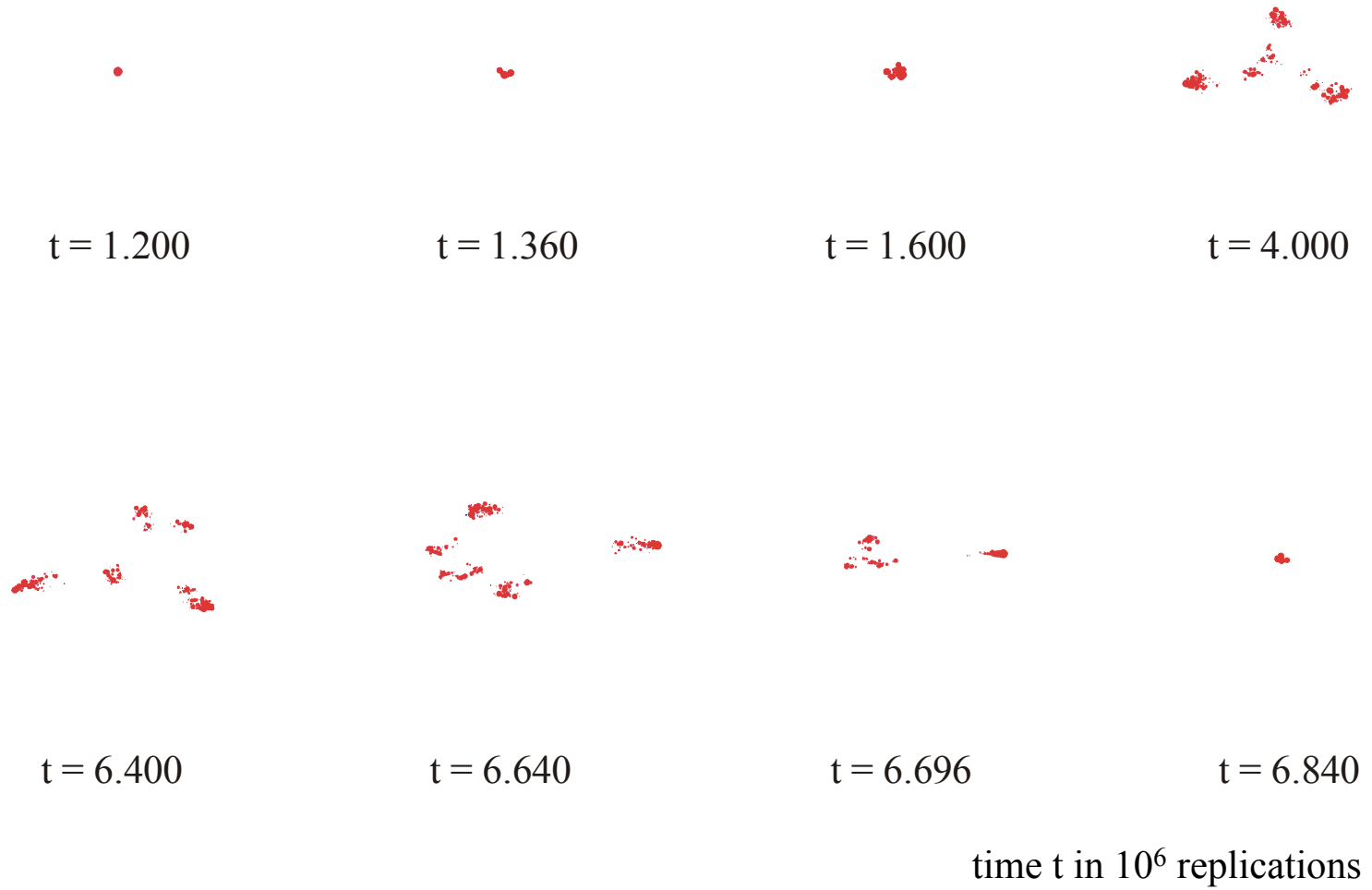
The flowreactor as a device for studies of evolution *in vitro* and *in silico*



*In silico* optimization in the flow reactor: Trajectory

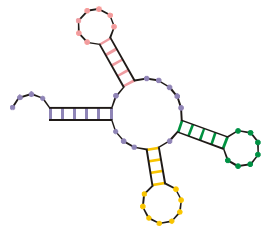
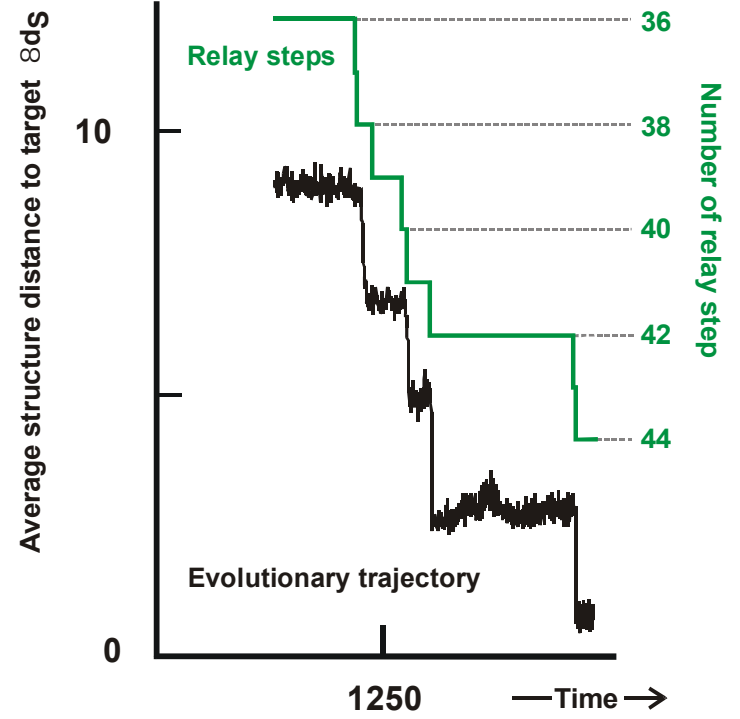
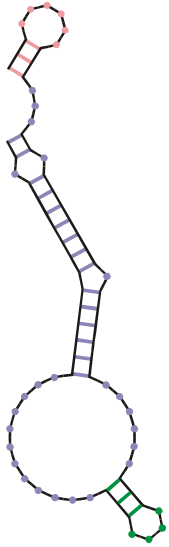


Variation of the population in genotype space during optimization of phenotypes



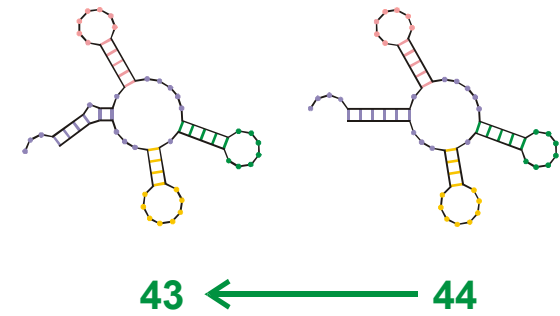
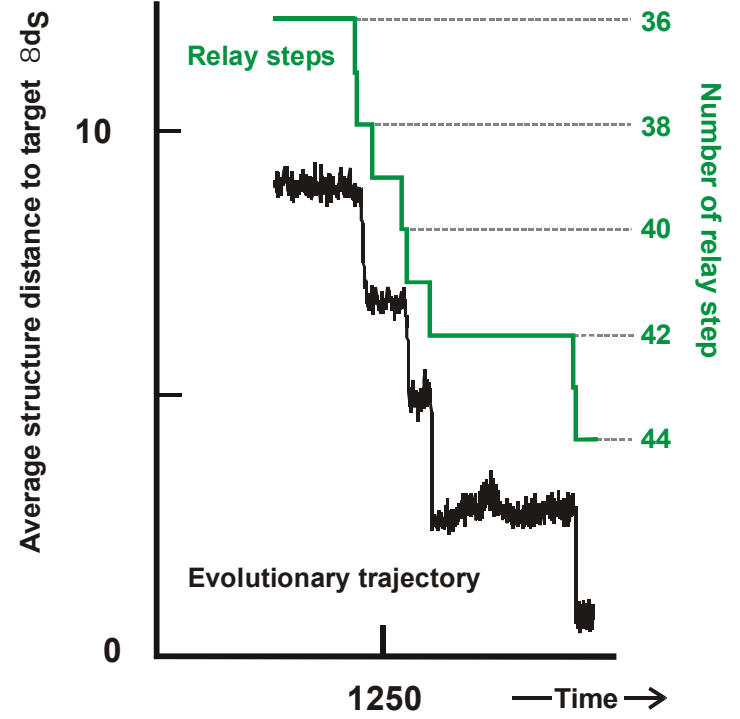
Spreading of a population during neutral evolution on a fitness plateau

00

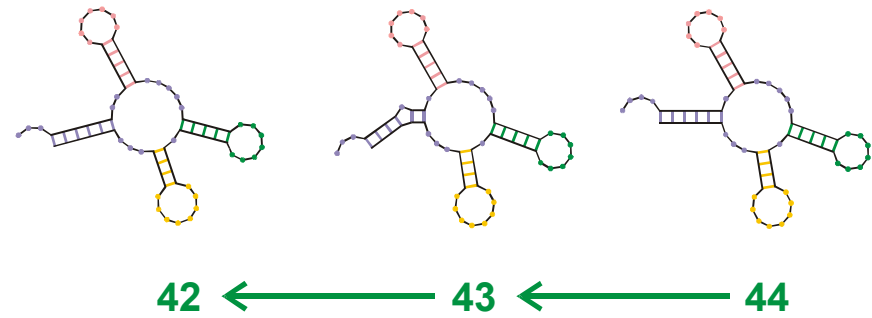
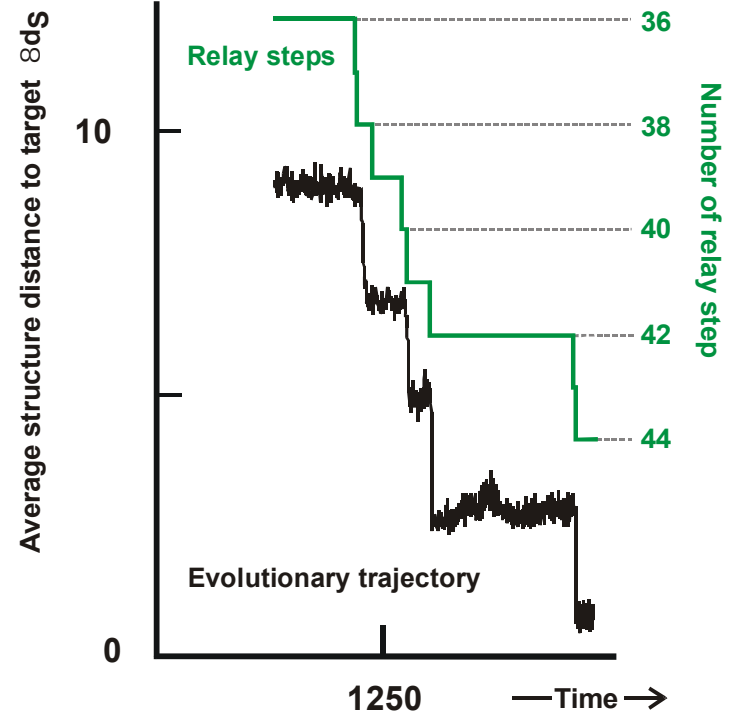


44

Initial structure and final conformation of the optimization process

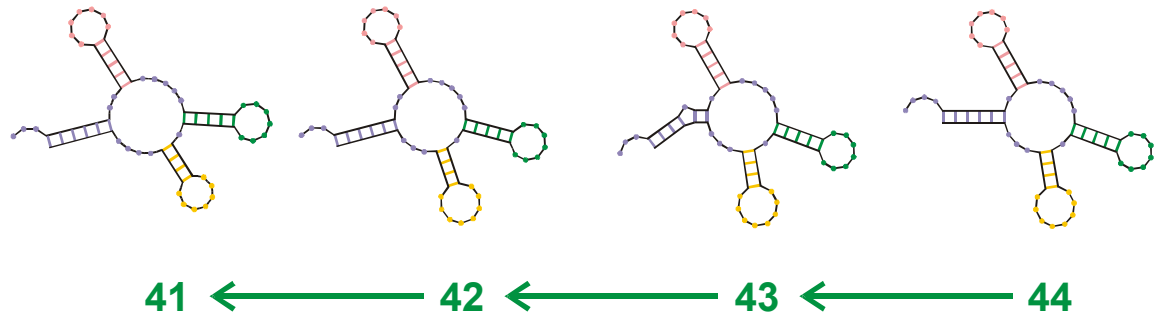
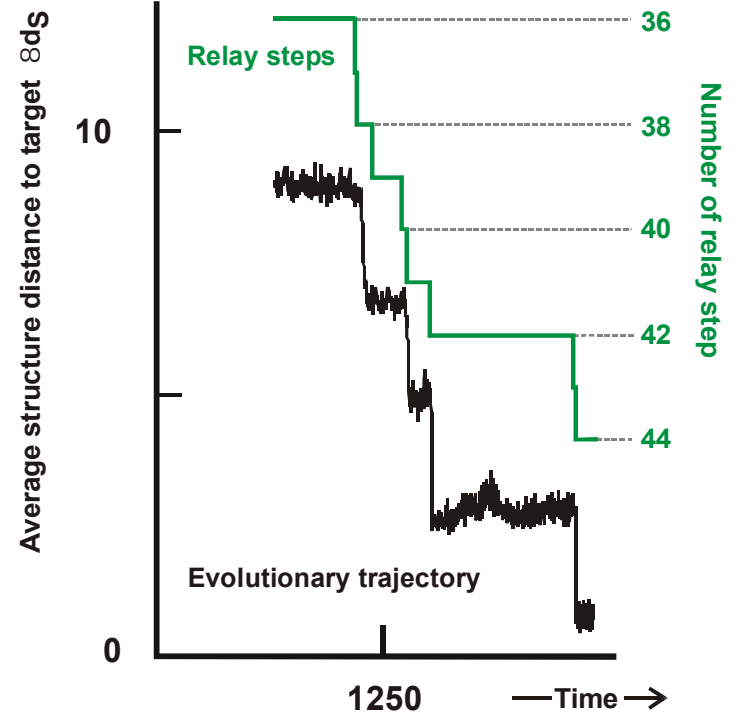


Reconstruction of the last step 43  $\dot{\simeq}$  44

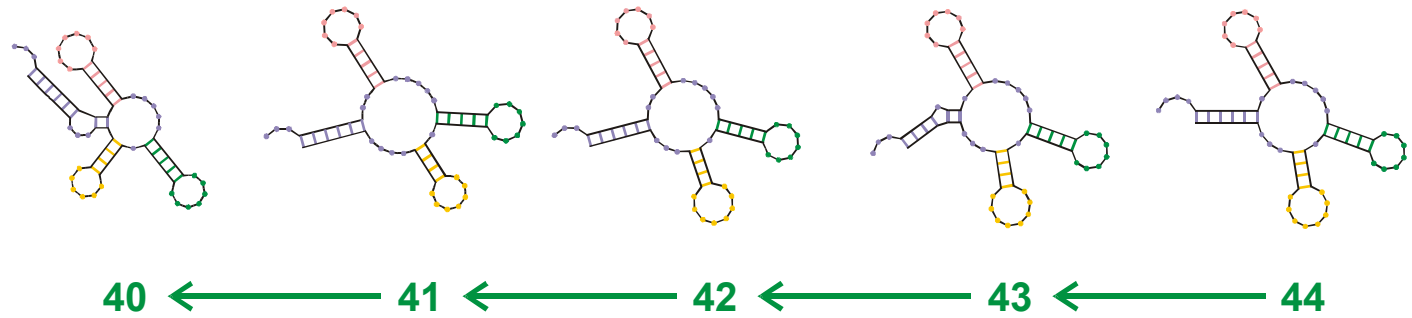
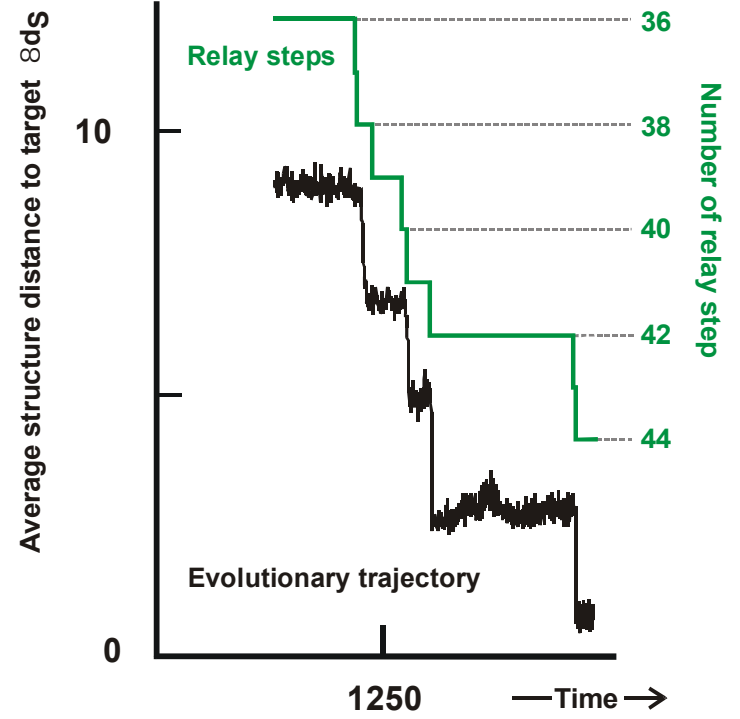


Reconstruction of last-but-one step 42  $\checkmark$  43 ( $\checkmark$  44)

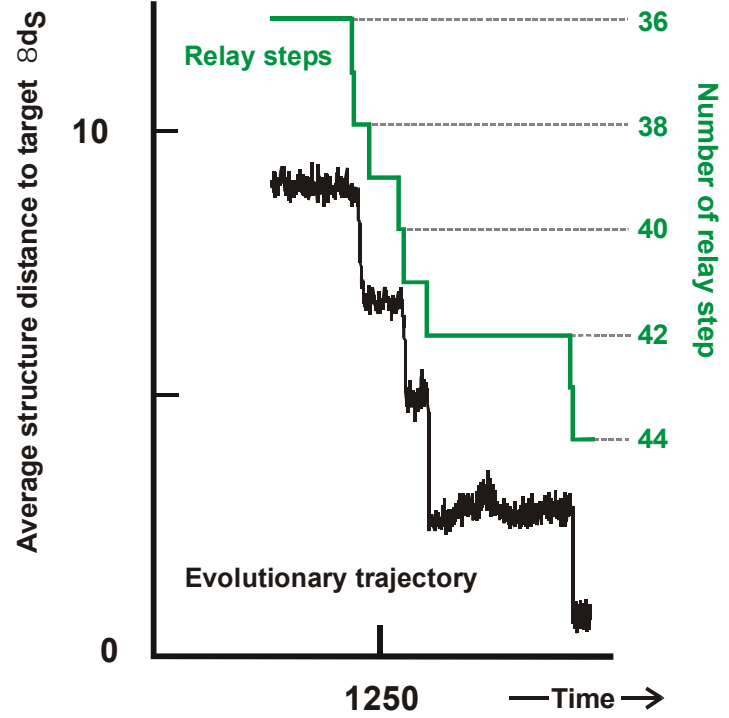




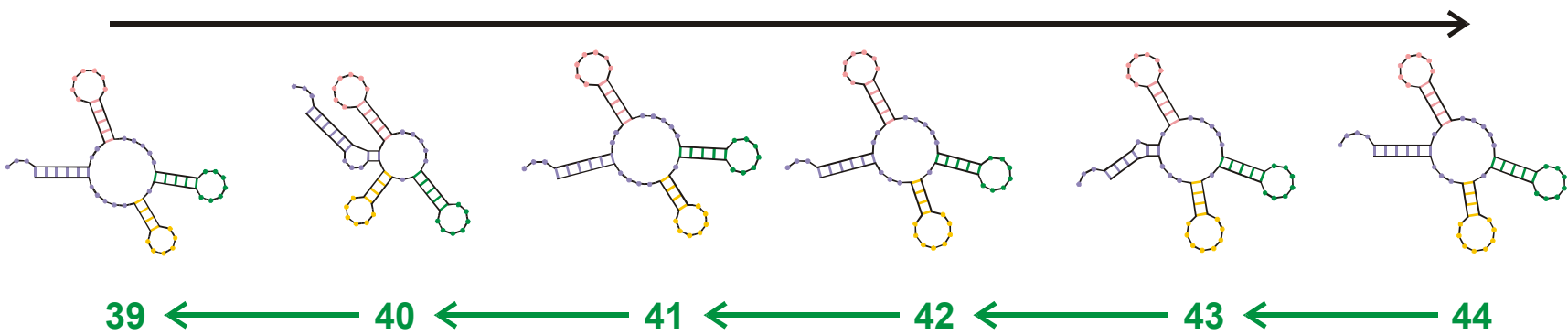
Reconstruction of step 41 š 42 (š 43 š 44)



Reconstruction of step 40 š 41 (š 42 š 43 š 44)



Evolutionary process



Reconstruction

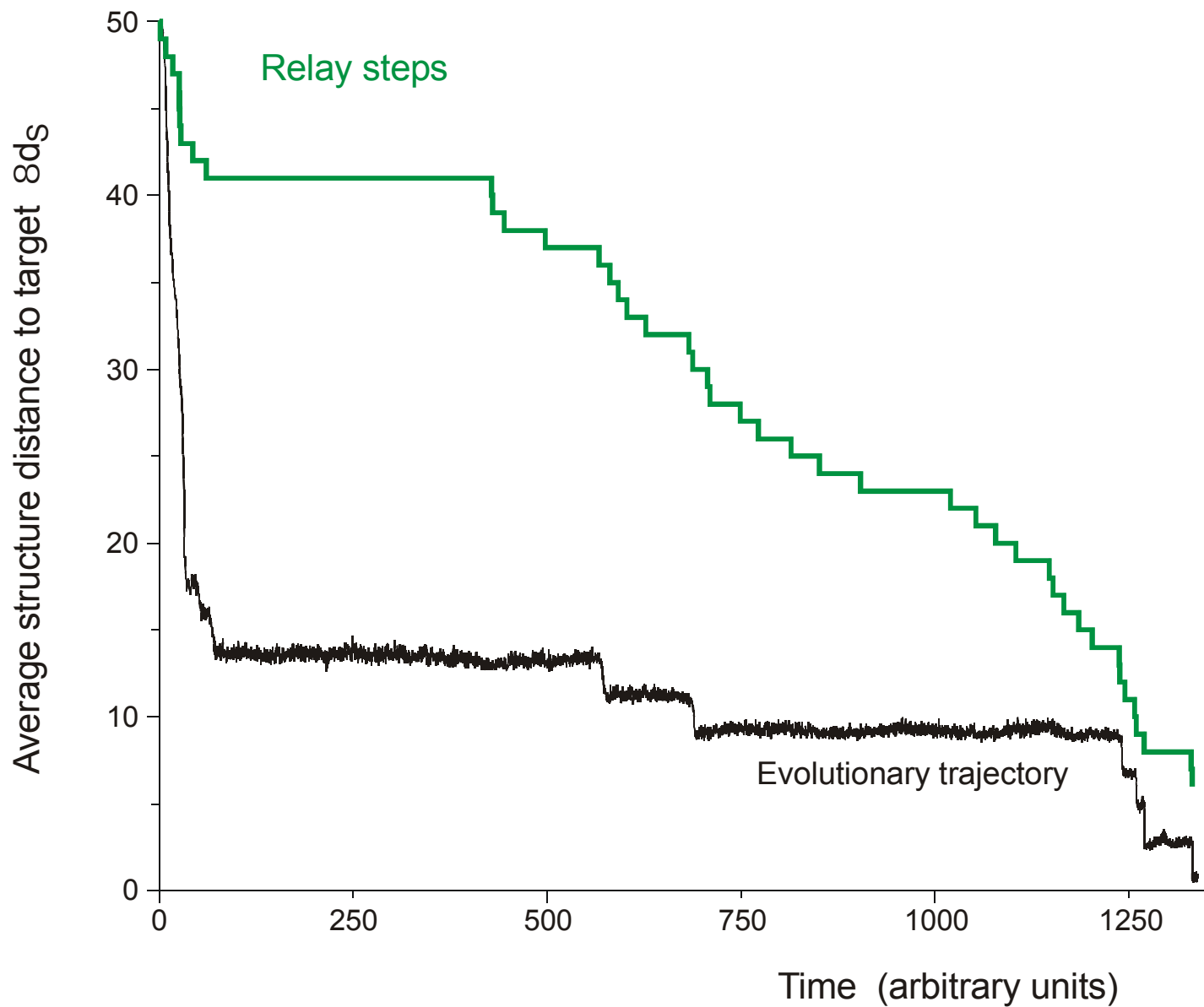
Reconstruction of the relay series

entry 39 GGGAUACAUGUGGCCCCUCAAGGCC**C**UAGCGAAACUGCUGCUGAAACCGUGUGAAUAAUCCGCACCCUGUCCCGA  
 ((((((.....(((.....))))).(((.....))))).(((.....))))).(((.....))))).(((.....))))).  
 exit GGGAUAUAC**G**AGGCCCGUCAAGGCC**G**UAGCGAA**C**CGACUGUUGAAACUGUG**C**GAAUAAUCCGCACCCUGUCCCG**G**  
 entry 40 GGGAUAUAC**G**GGGCCCGUCAAGGCC**G**UAGCGAAACCGACUGUUGAAACUGUGCGAAUAAUCCGCACCCUGUCCCGGG  
 ((((((.....(((.....))))).(((.....))))).(((.....))))).(((.....))))).(((.....))))).  
 exit GGGAUAUACGGG**G**CCCGUCAAGGCC**G**UAGCGAAACCGACUGUUGA**G**ACUGUGCGAAUAAUCCGCACCCUGUCCCGGG  
 entry 41 GGGAUAUACGGG**G**CCCGUCAAGGCC**G**UAGCGAAACCGACUGUUGAGACUGUGCGAAUAAUCCGCACCCUGUCCCGGG  
 ((((((.....(((.....))))).(((.....))))).(((.....))))).(((.....))))).(((.....))))).  
 exit GGGAUAUACGGGCCCC**U**UCAAG**G**CC**A**UAGCGAAACCGACUGUUGA**A**ACUGUGCGAAUAAUCCGCACCCUGUCCCG**G**  
 entry 42 GGGAUAUACGGGCCCC**U**UCAAG**C**CAUAGCGAAACCGACUGUUGAAACUGUGCGAAUAAUCCGCACCCUGUCCCGGA  
 ((((((.....(((.....))))).(((.....))))).(((.....))))).(((.....))))).(((.....))))).  
 exit GGGAU**G**AUAGGG**C**GUGU**G**AUAGCCCAUAGCGAAAC**CC**CC**C**GUGA**G**CUUGUGCGA**CG**UUUGUGCACCUGUCCCG**C**  
 entry 43 GGG**A**GAUAGGG**C**GUGUGAUAGCCCAUAGCGAAAC**CC**CC**C**CGUGAGCUUGUGCGACGUUUGUGCACCUGUCCCGCU  
 ((((((.....(((.....))))).(((.....))))).(((.....))))).(((.....))))).(((.....))))).  
 exit GGG**A**GAUAGGG**C**GUGUGAUAGCCCAUAGCGAAAC**CC**CC**C**CGUGAGCUUGUGCGACGUUUGUGCACCUGUCCCGCU  
 entry 44 GGG**C**AGAUAGGG**C**GUGUGAUAGCCCAUAGCGAAAC**CC**CC**C**CGUGAGCUUGUGCGACGUUUGUGCACCUGUCCCGCU  
 ((((((.....(((.....))))).(((.....))))).(((.....))))).(((.....))))).(((.....))))).

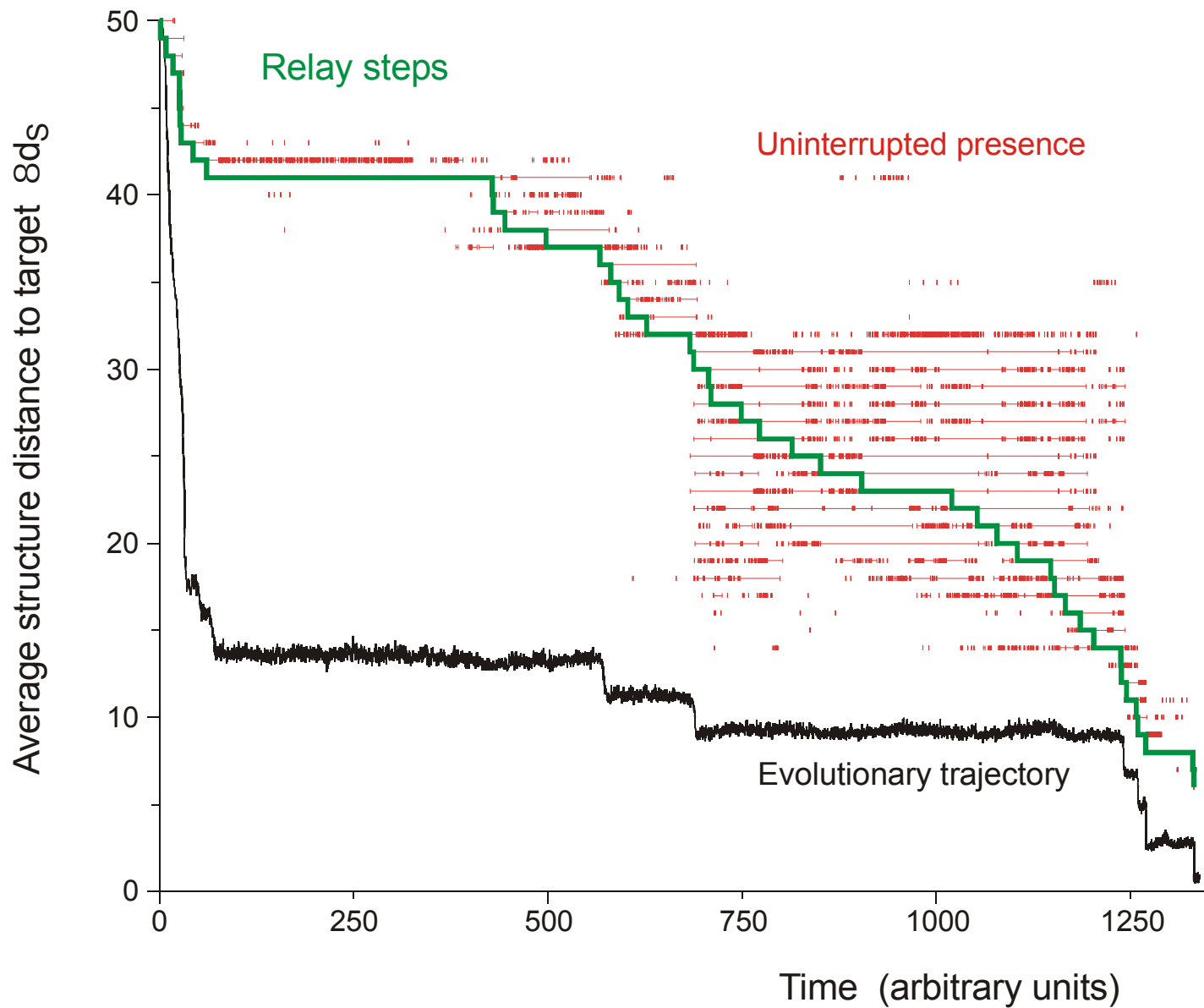
**Transition inducing point mutations**

**Neutral point mutations**

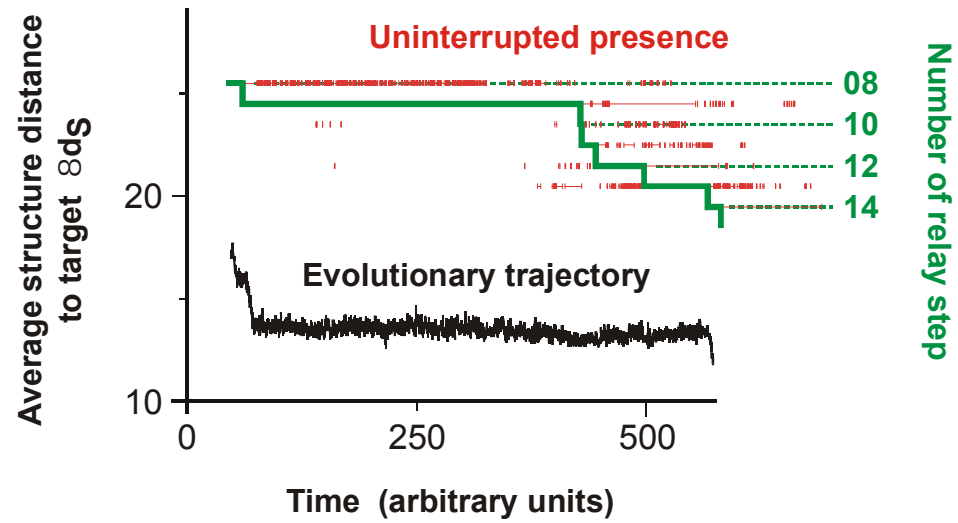
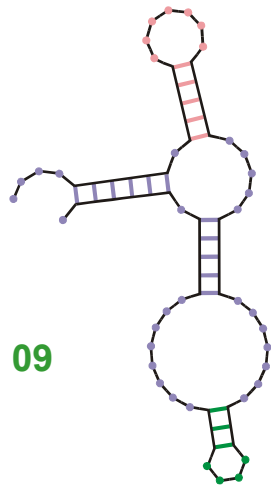
Change in RNA sequences during the final five relay steps 39 § 44



*In silico* optimization in the flow reactor: Trajectory and relay steps



*In silico* optimization in the flow reactor: Uninterrupted presence



entry 8 GGUAUGGGCGUUGAAUAGUAGGGUUUAAACCAAUCGGCAACGAUCUCGUGUGCGCAUUUCAUAUCCCGUACAGAA  
 .((((((((((((.....(((.....)).....)))))).....((((.....)))))))))).....  
 exit 8 GGUAUGGGCGUUGAAUAUAJAGGGUUUAAACCAAUCGGCCAACGAUCUCGUGUGCGCAUUUCAUAUCCAUAACAGAA  
 entry 9 GGUAUGGGCGUUGAAUAAUAGGGUUUAAACCAAUCGGCCAACGAUCUCGUGUGCGCAUUUCAUAUAACCAUACAGAA  
 .(((((((.....((((.....)).....)))))).....((((.....)))))).....  
 exit 9 UGGAUGGACGUUGAAUAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACACCGUCCCAAG  
 entry 10 UGGAUGGACGUUGAAUAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACACCGUCCCAAG  
 .(((((((.....((((.....)).....)))))).....((((.....)))))).....  
 exit 10 UGGAUGGACGUUGAAUAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACAGCGUCCCAAG

**Transition inducing point mutations**

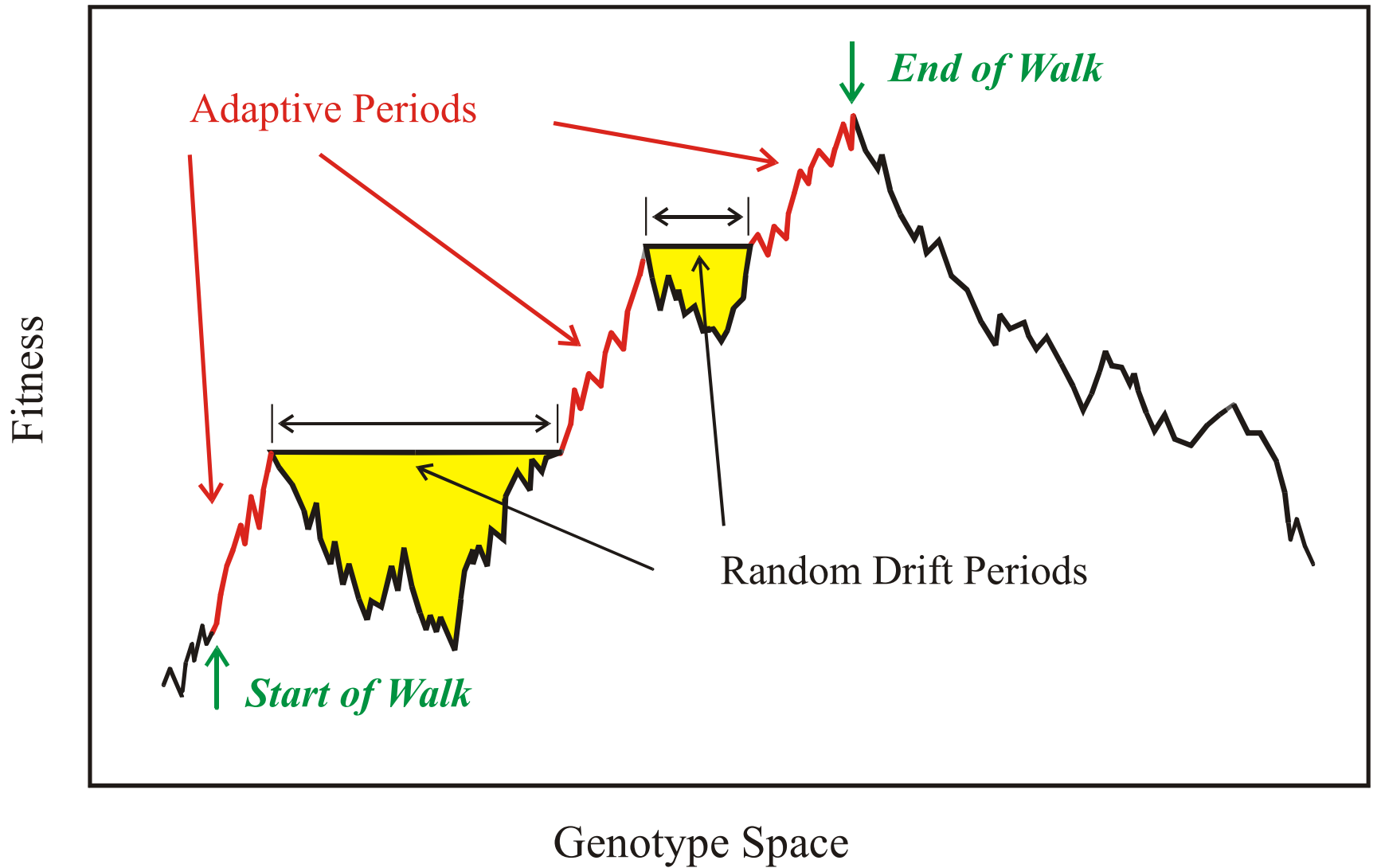
**Neutral point mutations**

Neutral genotype evolution during phenotypic stasis

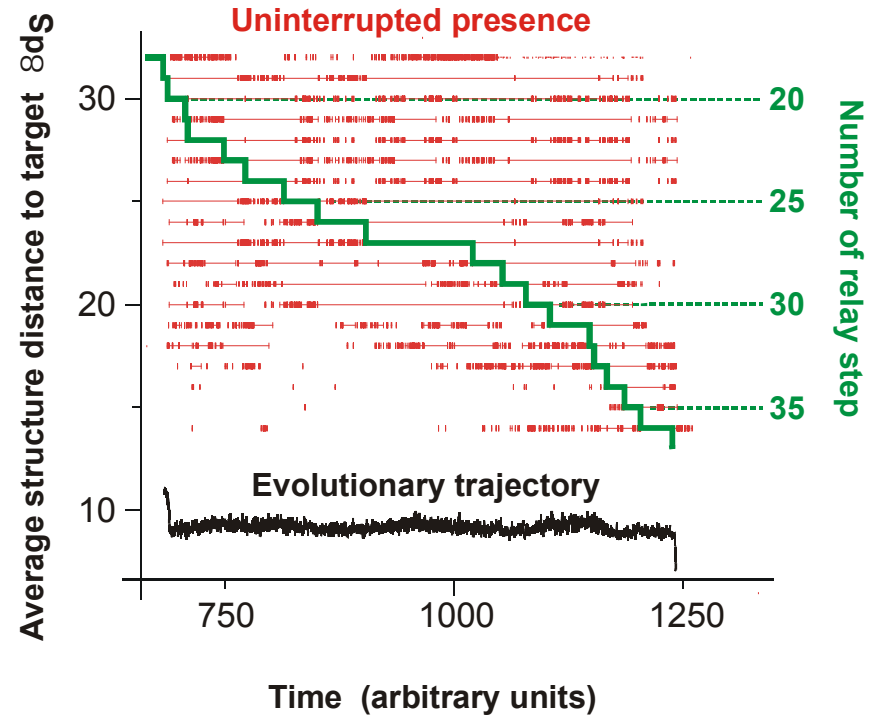
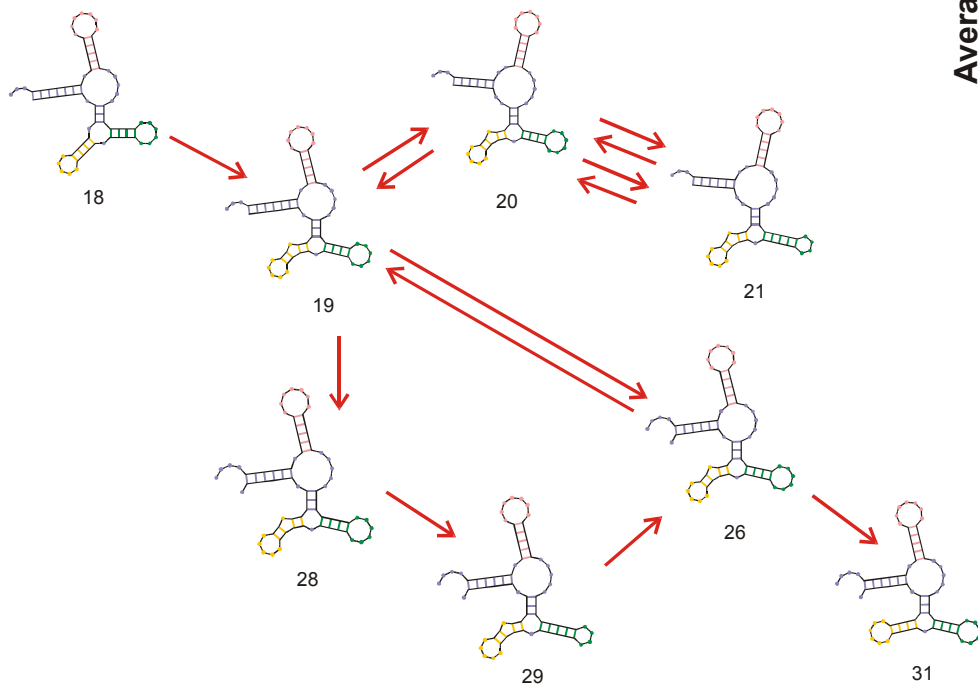
„...Variations neither useful not injurious would not be affected by natural selection, and would be left either a fluctuating element, as perhaps we see in certain polymorphic species, or would ultimately become fixed, owing to the nature of the organism and the nature of the conditions.  
...“  
...

Charles Darwin, Origin of species (1859)

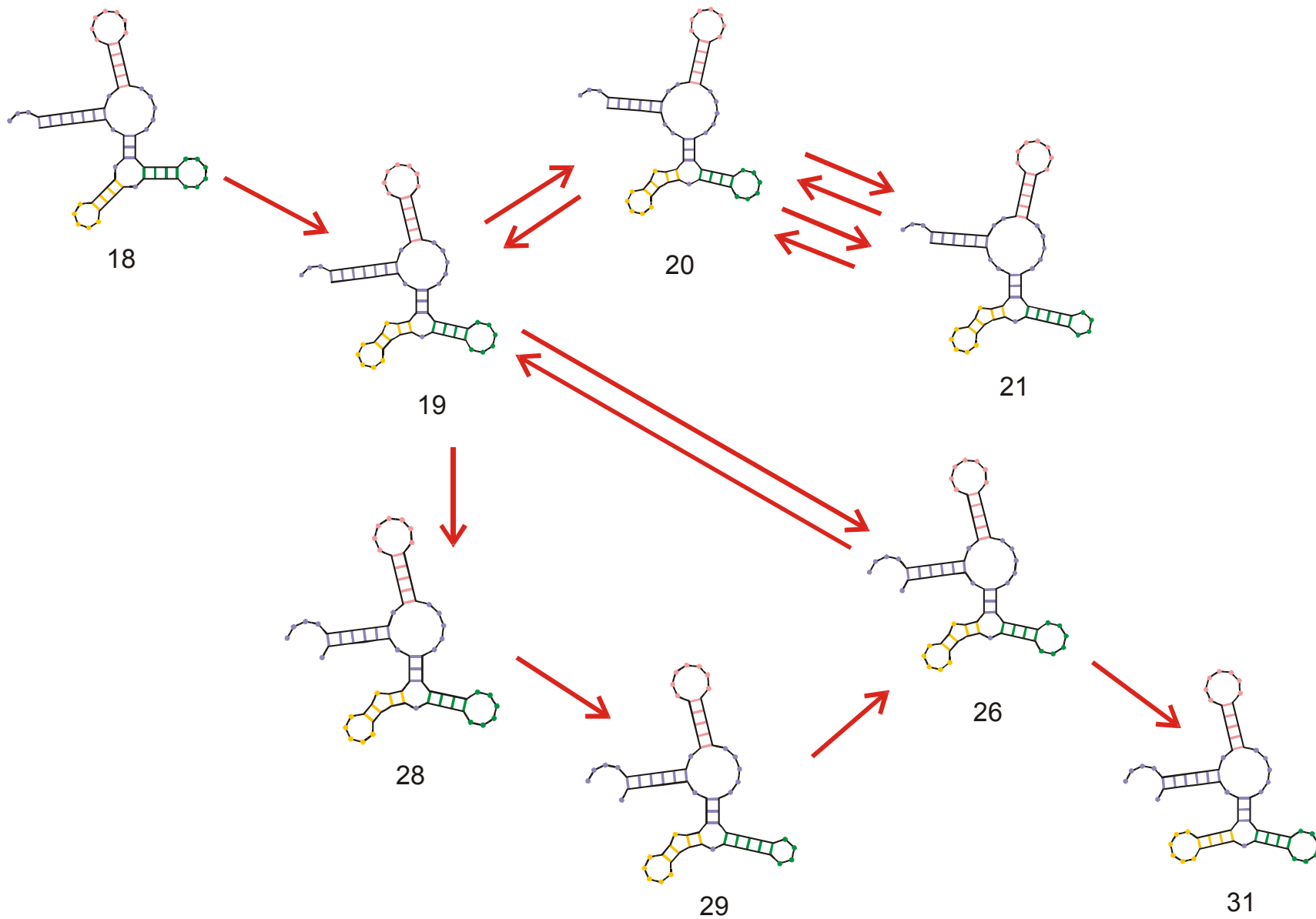




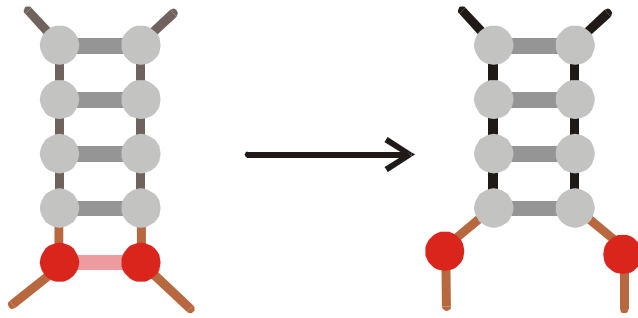
Evolution in genotype space sketched as a non-descending walk in a fitness landscape



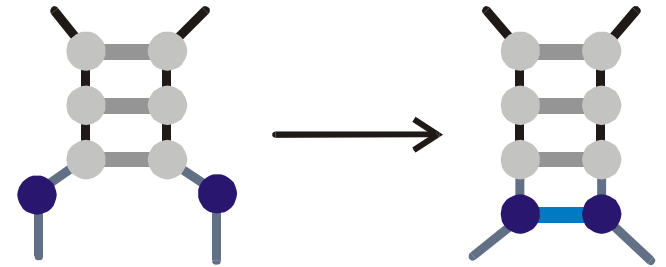
A random sequence of **minor** or continuous **transitions** in the relay series



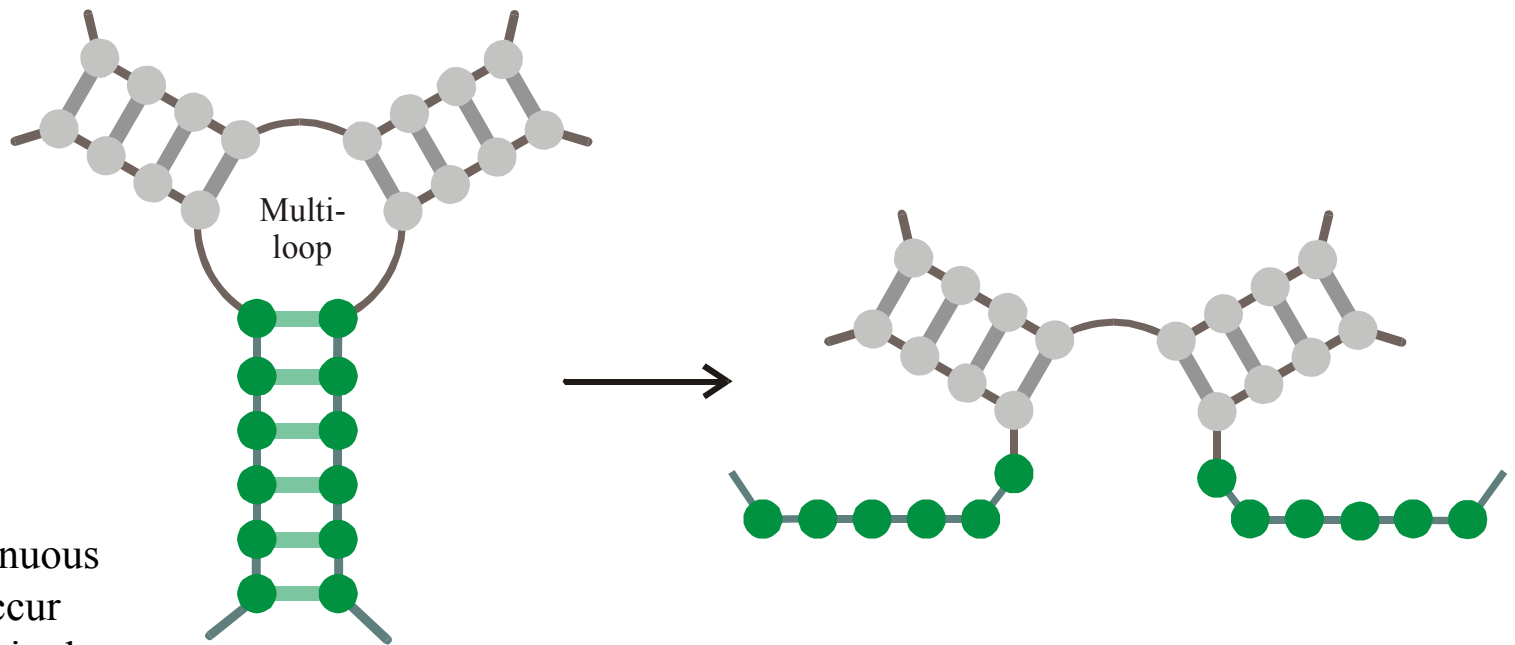
A random sequence of **minor** or continuous **transitions** in the relay series



Shortening of Stacks

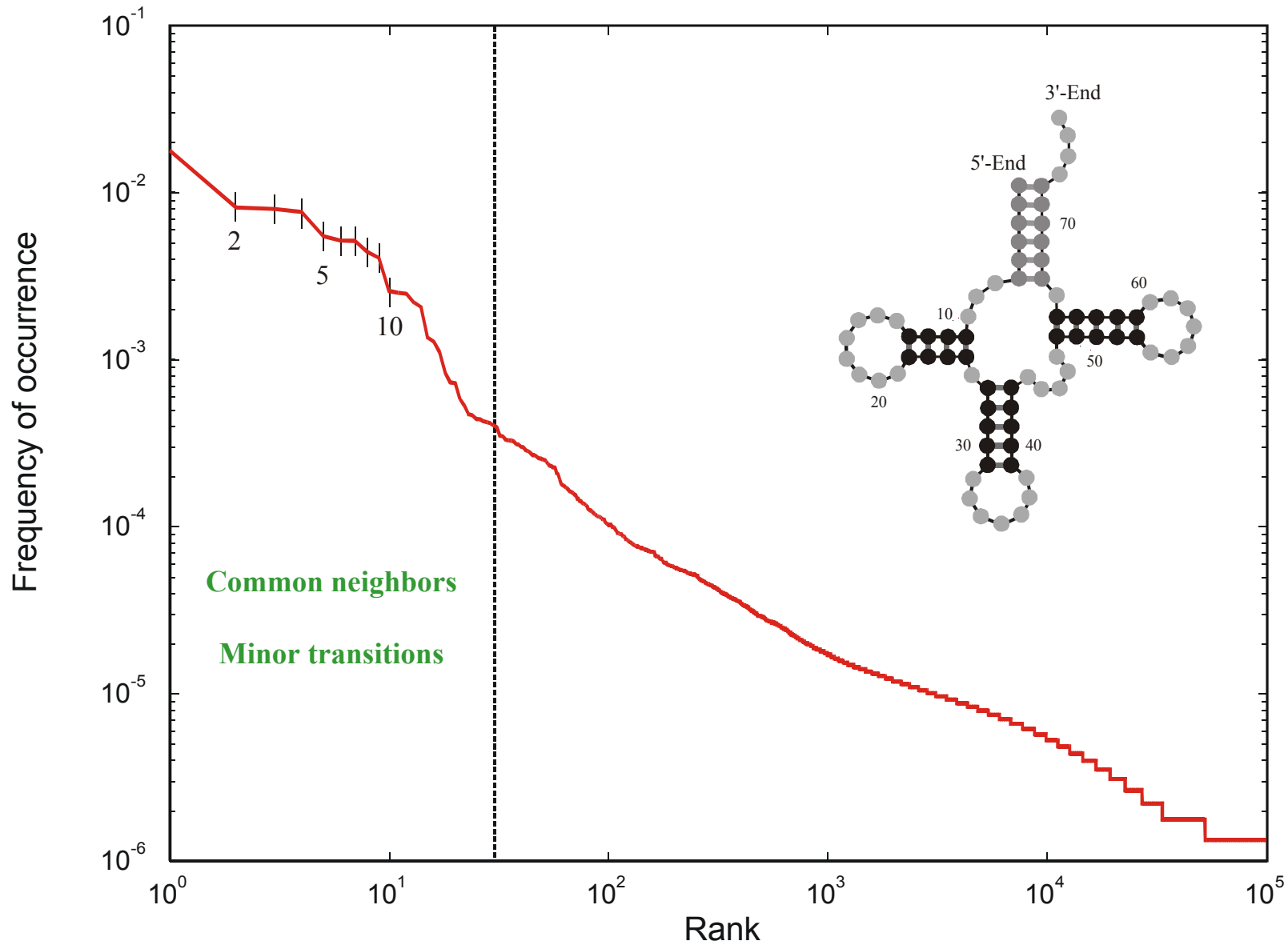


Elongation of Stacks

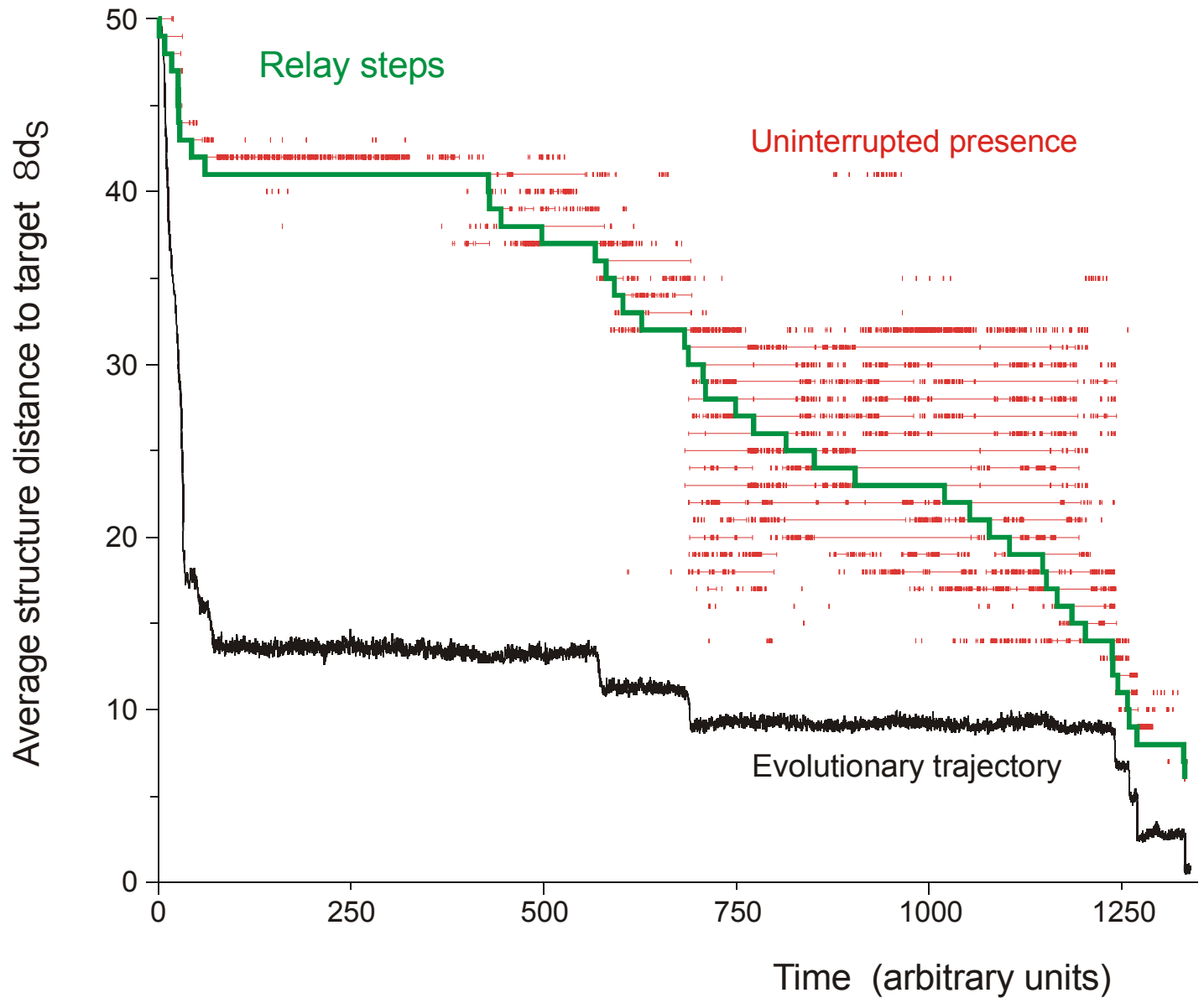


Opening of Constrained Stacks

**Minor** or continuous **transitions**: Occur **frequently** on single point mutations

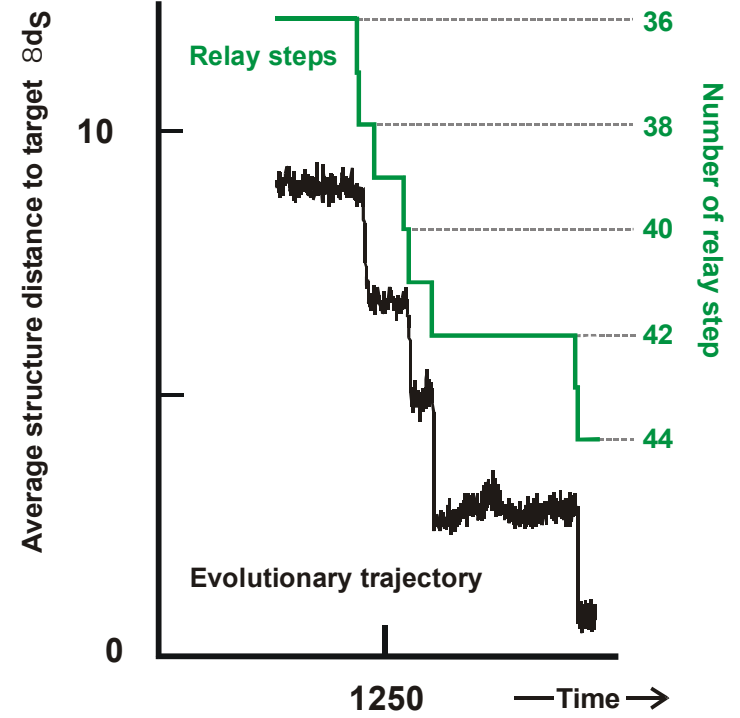
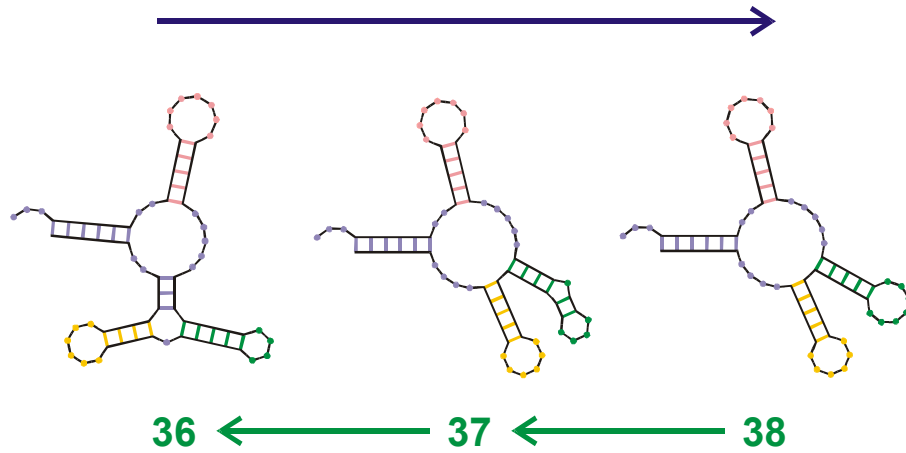


Probability of occurrence of different structures in the mutational neighborhood of tRNA<sup>phe</sup>



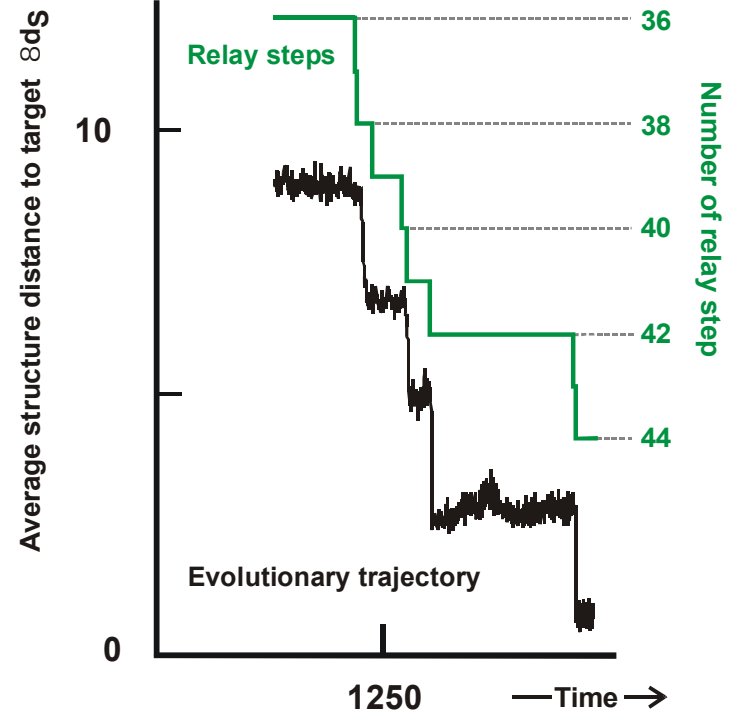
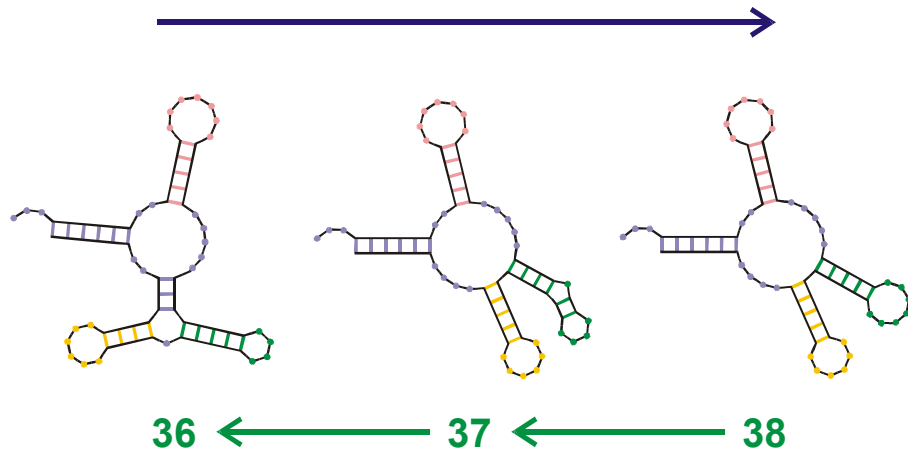
*In silico* optimization in the flow reactor: **Uninterrupted presence**

## Major transition leading to clover leaf

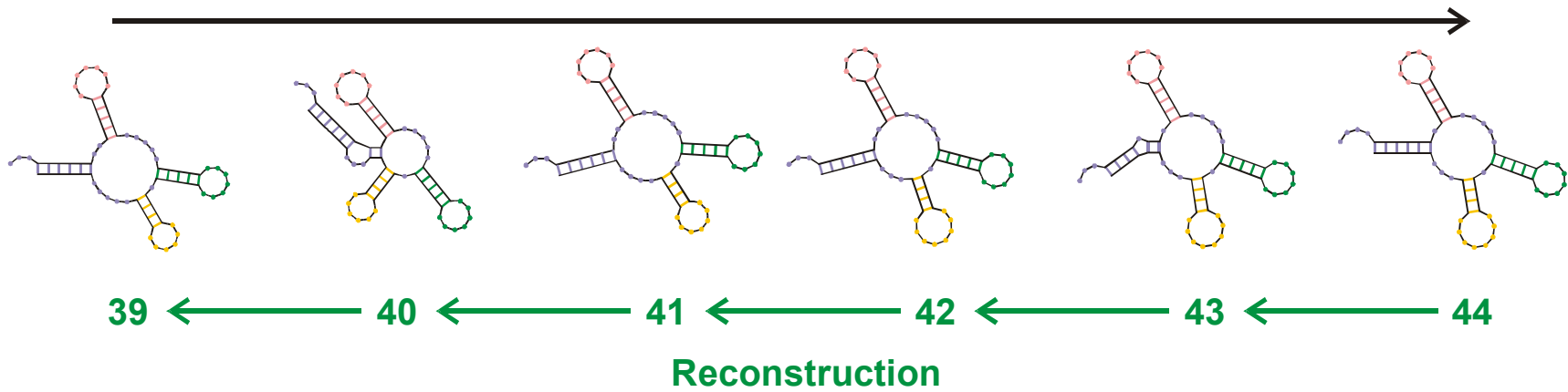


Reconstruction of a **major transitions** 36  $\rightarrow$  37 ( $\rightarrow$  38)

### Major transition leading to clover leaf

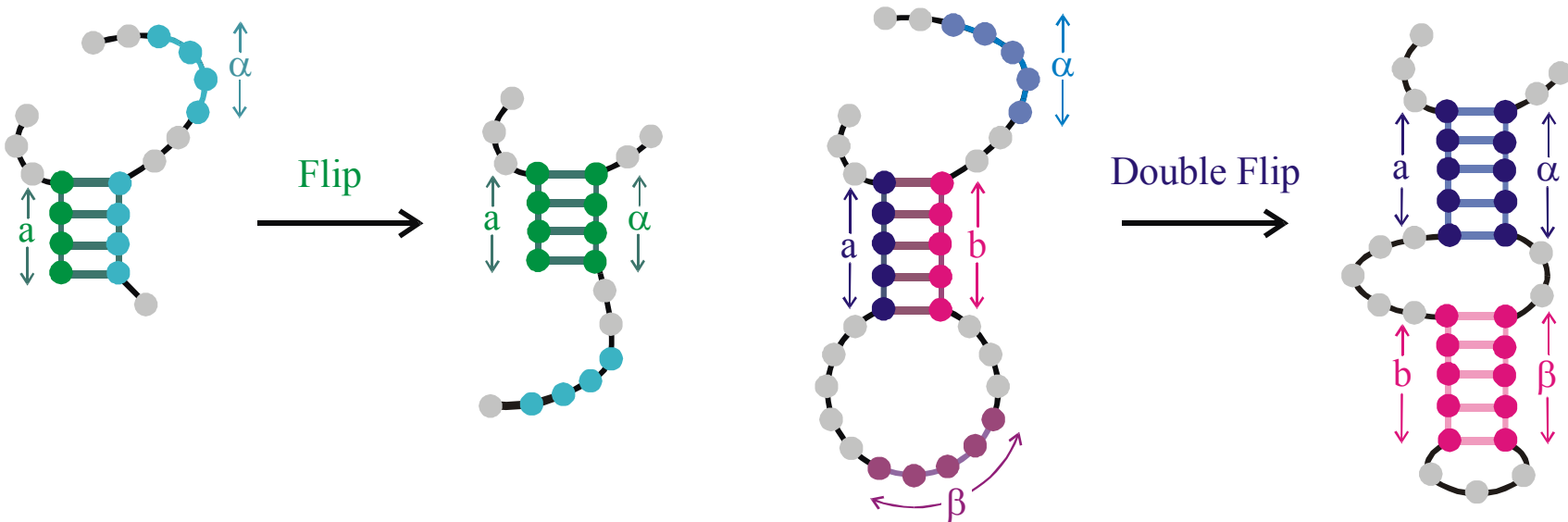
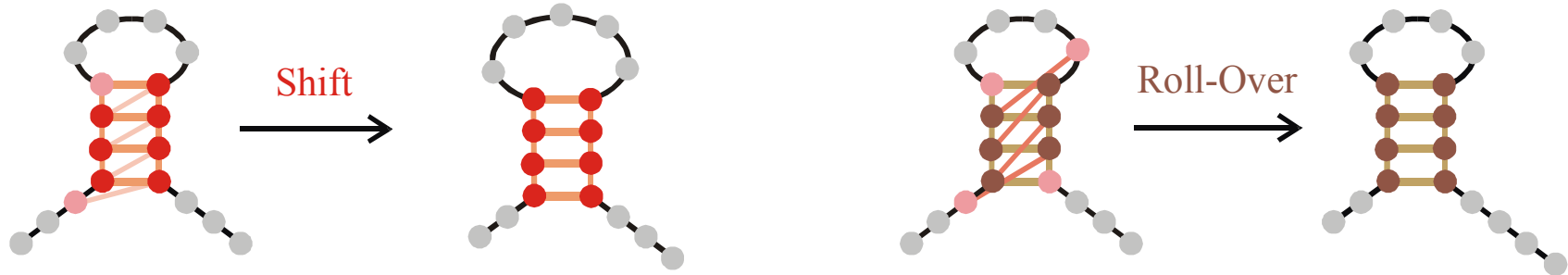


### Evolutionary process

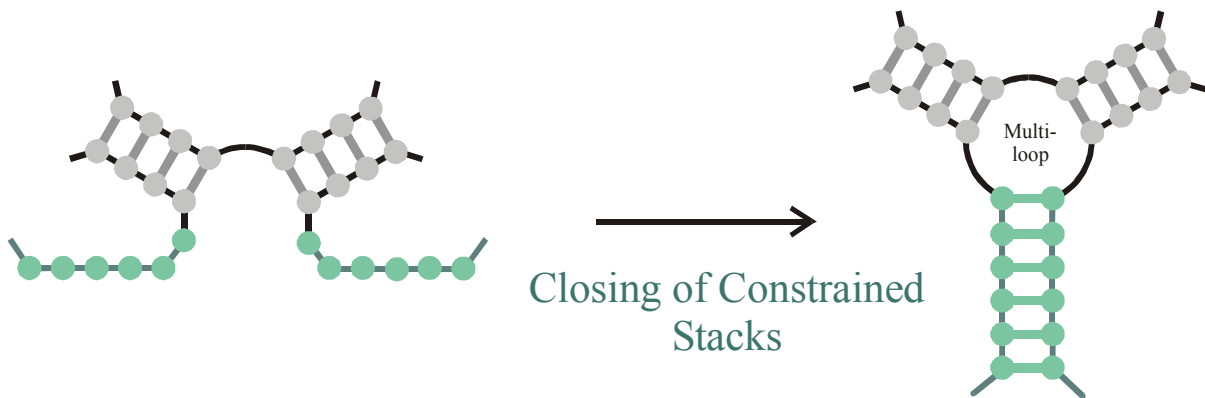


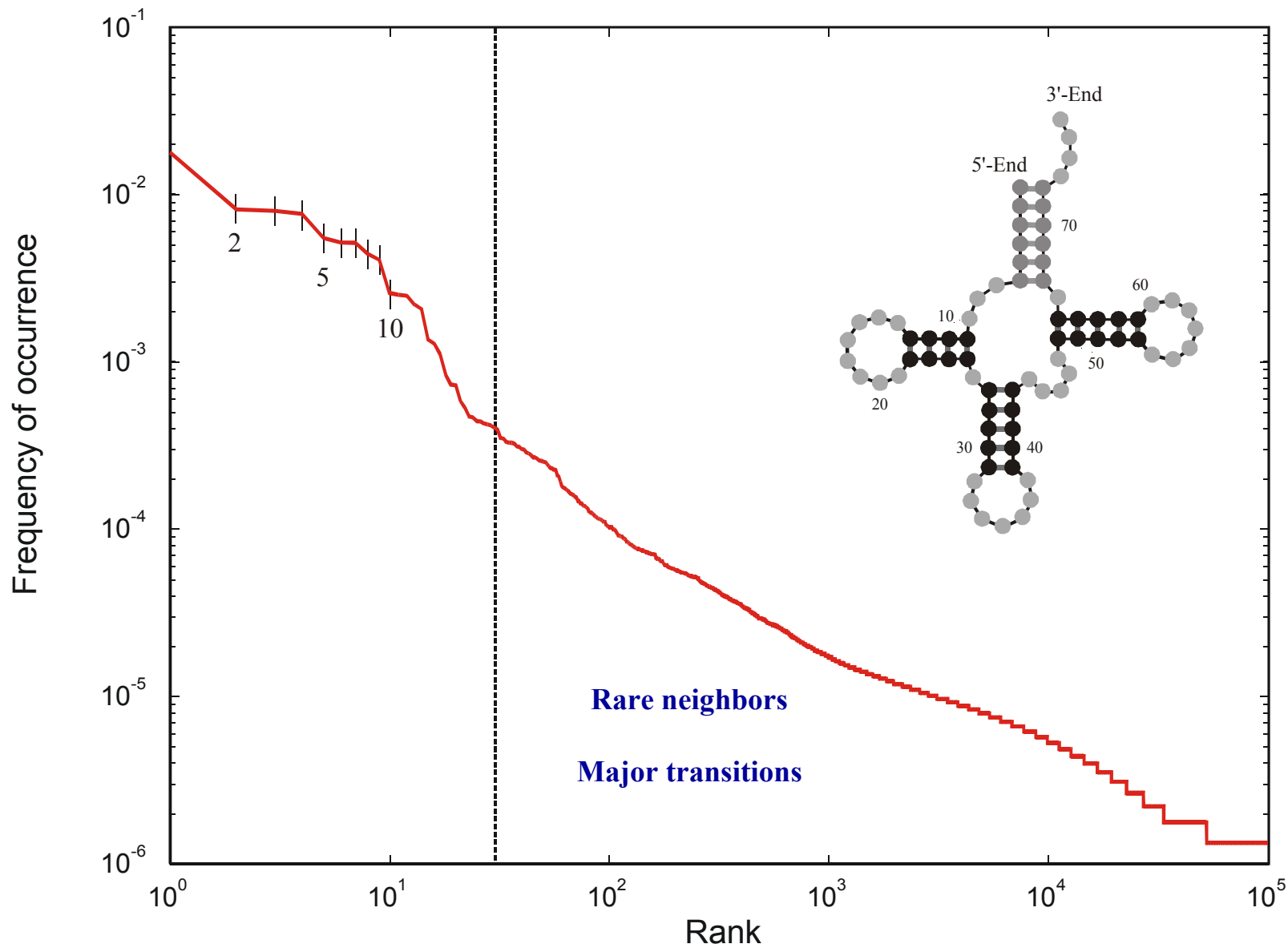
Final reconstruction 36  $\hat{S}$  44



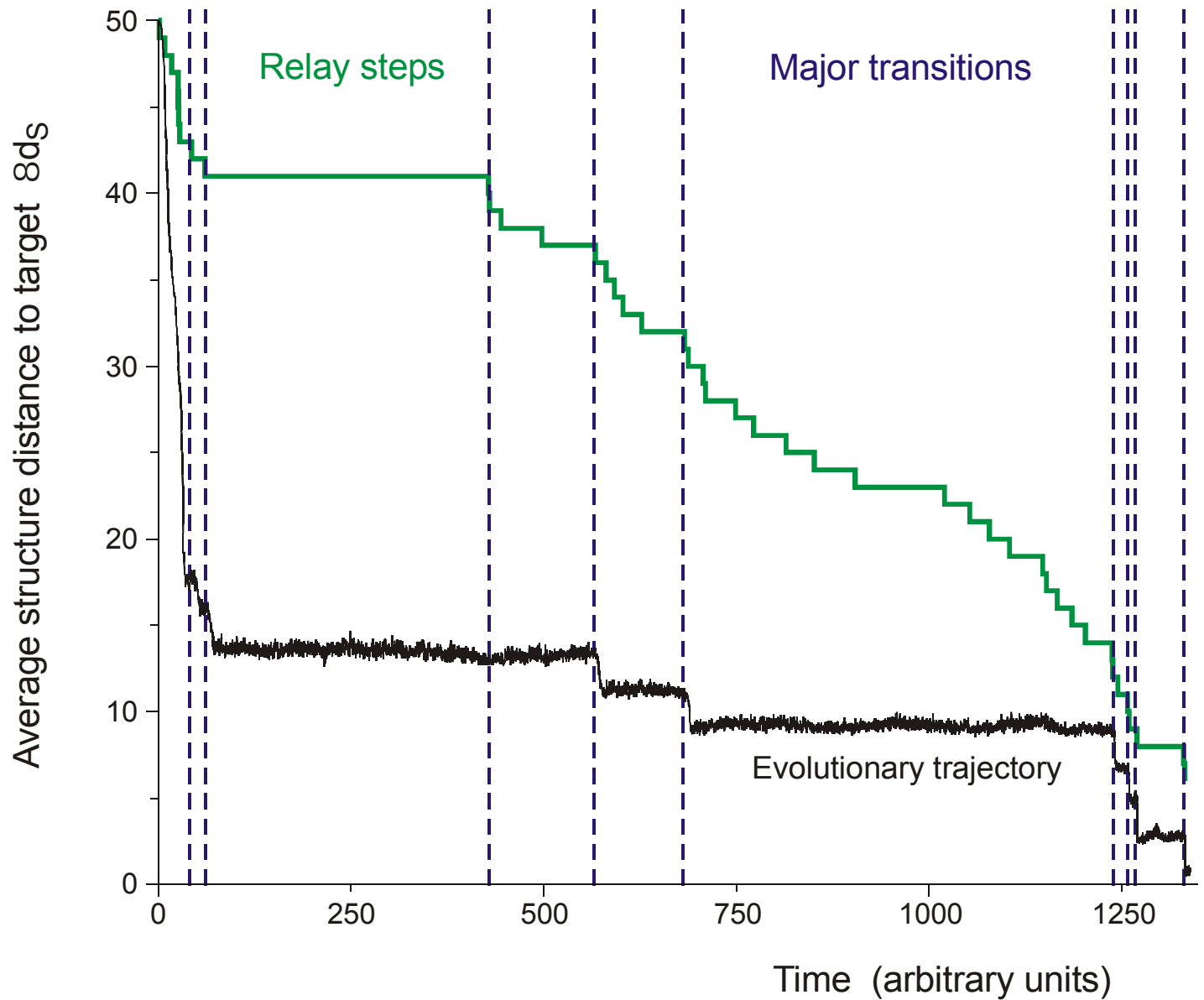


**Major** or discontinuous transitions: **Structural innovations**, occur **rarely** on single point mutations

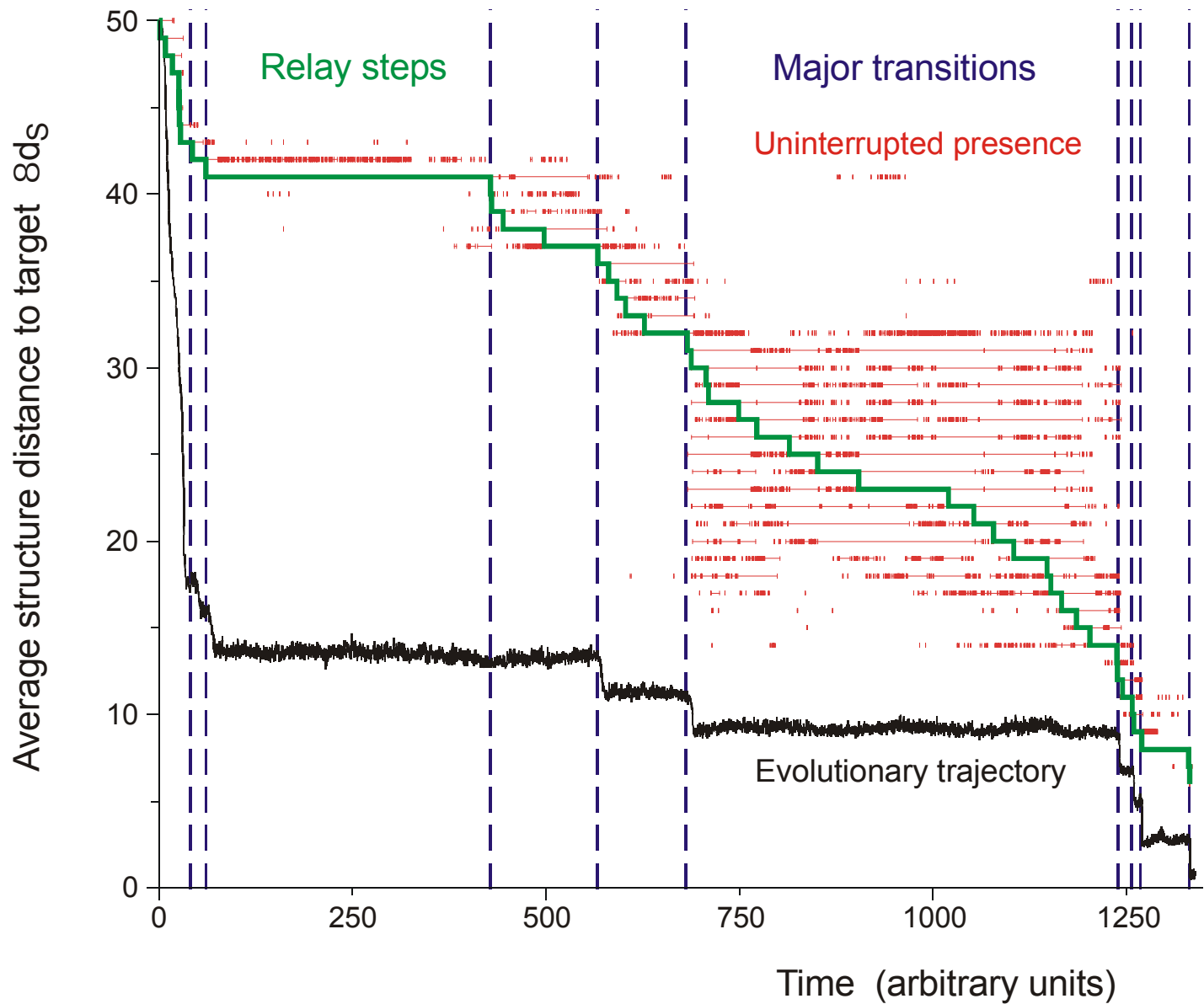




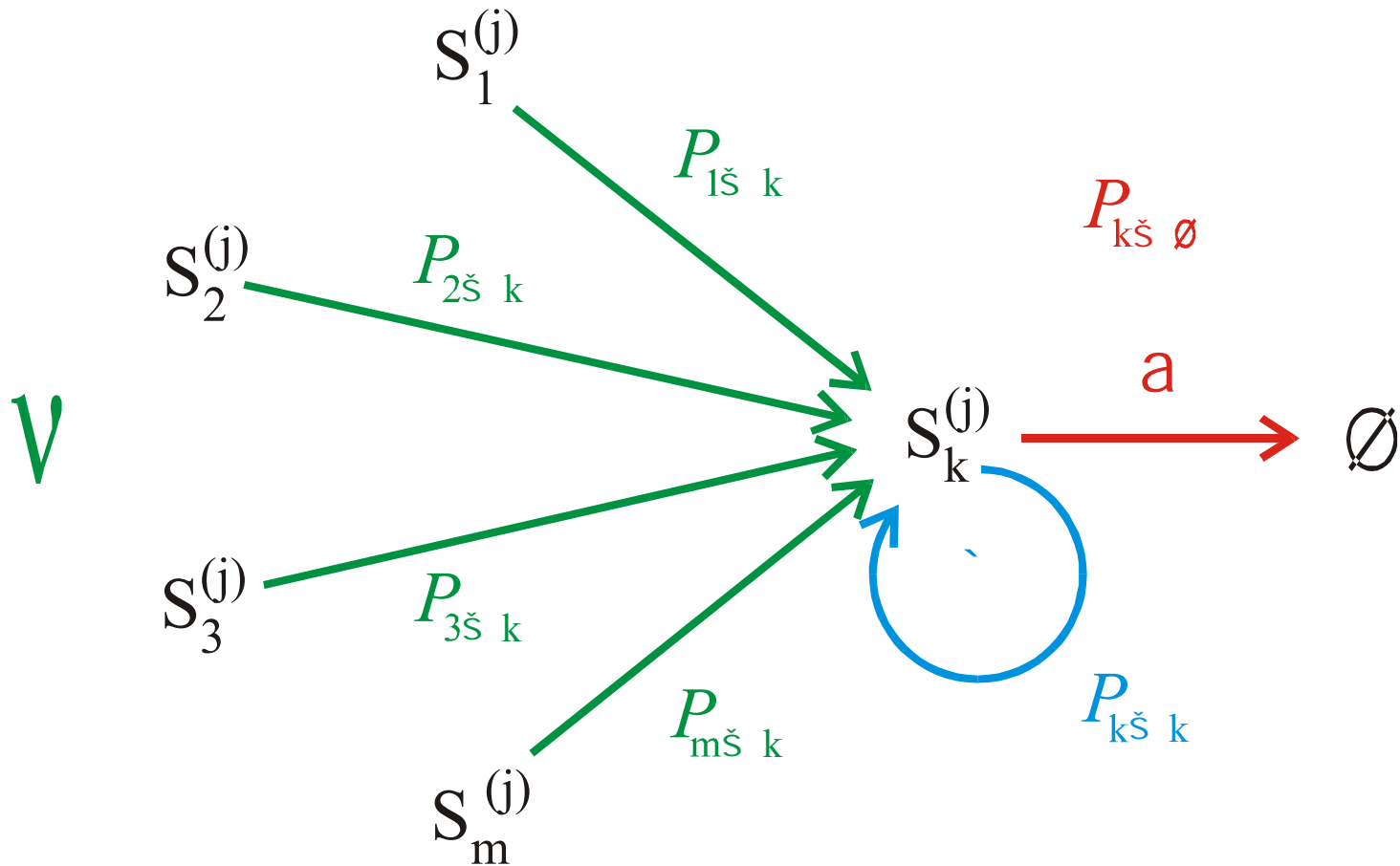
Probability of occurrence of different structures in the mutational neighborhood of tRNA<sup>phe</sup>



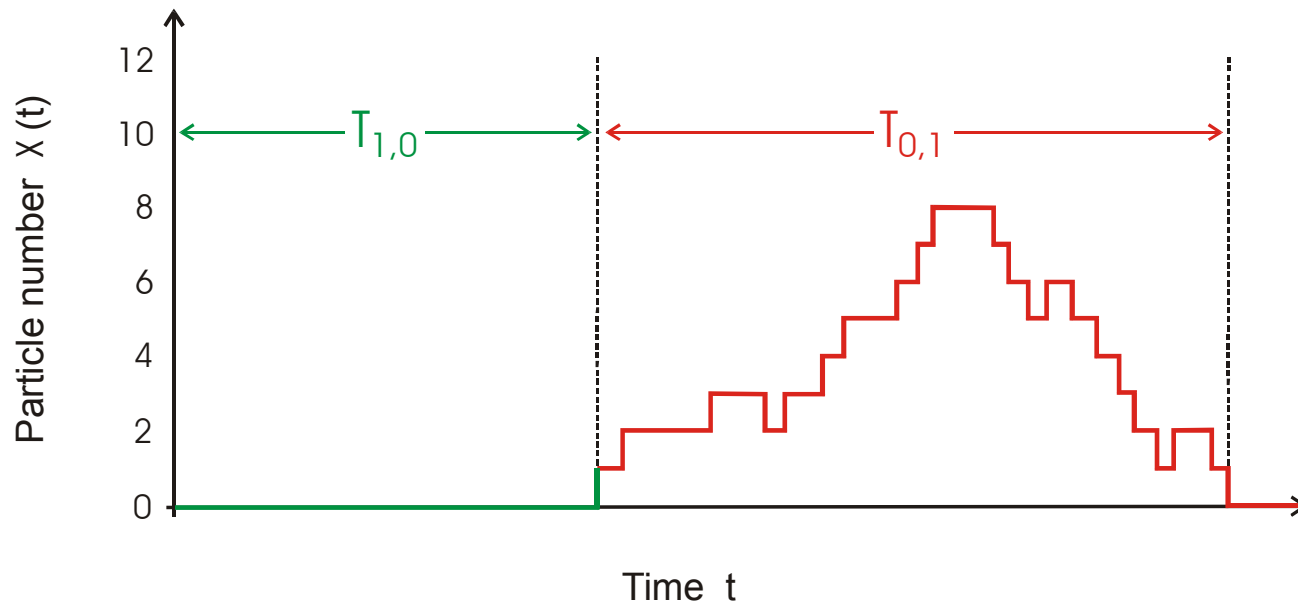
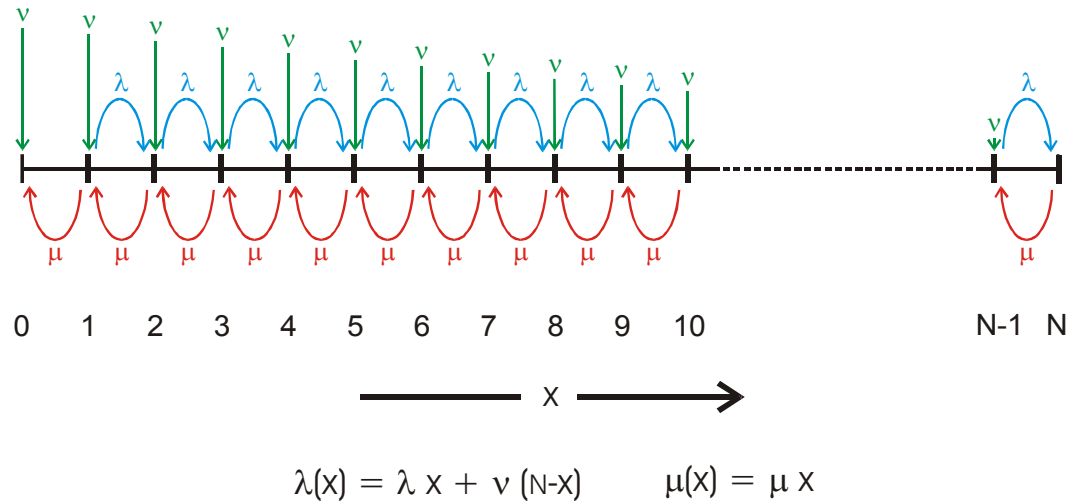
*In silico* optimization in the flow reactor: **Major transitions**



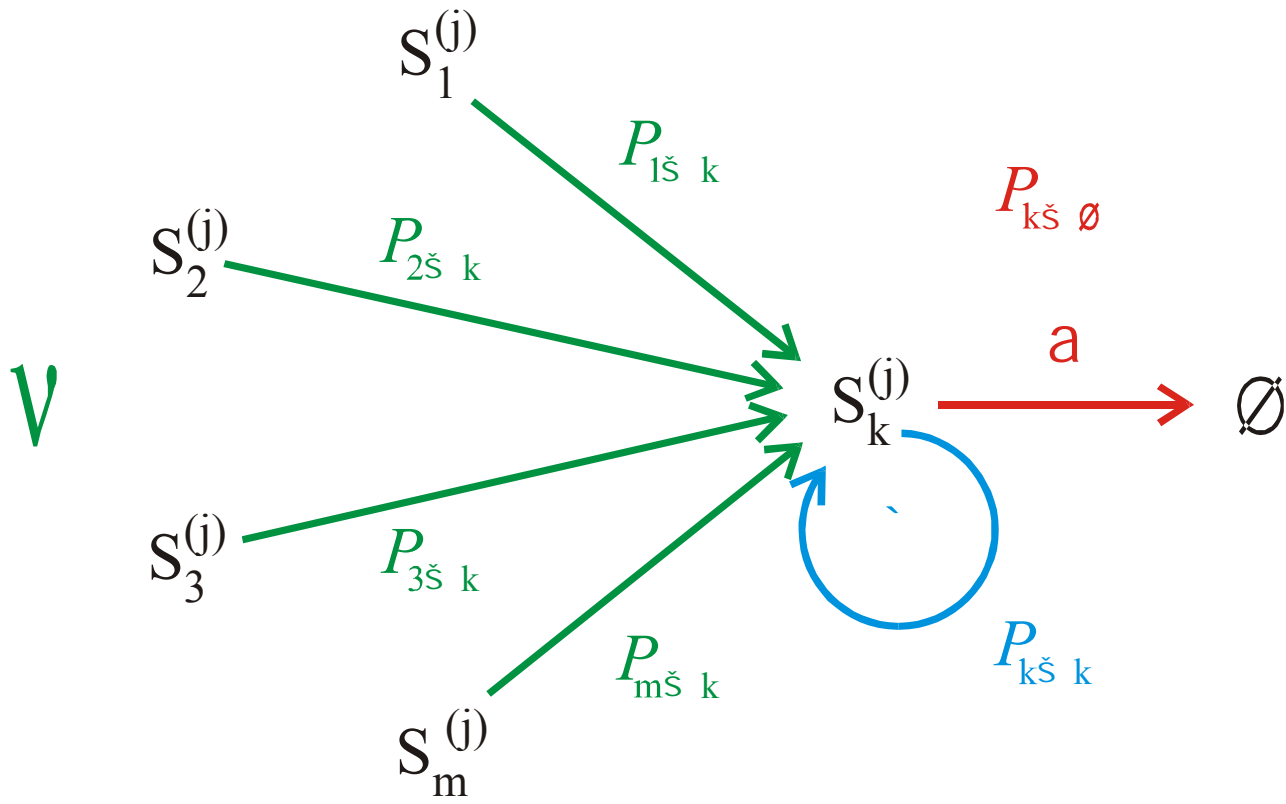
*In silico* optimization in the flow reactor



Transition probabilities determining the presence of phenotype  $S_k^{(j)}$  in the population



Calculation of transition probabilities by means of a birth-and-death process with immigration

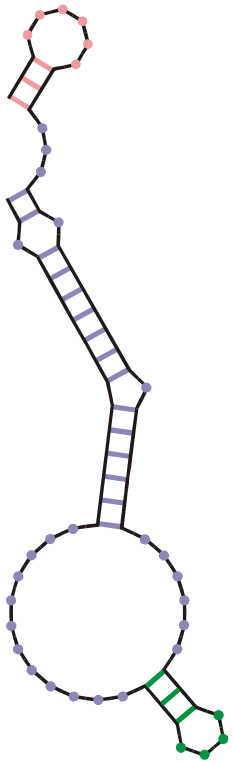


$$N_{\text{sat}}^{(j)} = \frac{1}{p \cdot l \cdot \langle f^{(j)} \rangle}$$

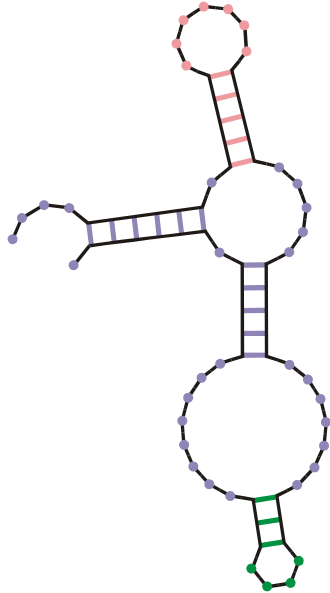
## Statistics of evolutionary trajectories

Population Size $N$	Number of Replications $\langle n_{\text{rep}} \rangle$	Number of Transitions $\langle n_{\text{tr}} \rangle$	Number of Major Transitions $\langle n_{\text{dtr}} \rangle$	Epochal Phase $\langle d_{\tau}^s(t_{\text{ep}}) \rangle$
1 000	$(5.5 \pm [6.9, 3.1]) \times 10^7$	$92.7 \pm [80.3, 43.0]$	$8.8 \pm [2.4, 1.9]$	$23.7 \pm [5.0, 4.1]$
2 000	$(6.0 \pm [11.1, 3.9]) \times 10^7$	$55.7 \pm [30.7, 19.8]$	$8.9 \pm [2.8, 2.1]$	$22.2 \pm [5.1, 4.2]$
3 000	$(6.6 \pm [21.0, 5.0]) \times 10^7$	$44.2 \pm [25.9, 16.3]$	$8.1 \pm [2.3, 1.8]$	$20.9 \pm [2.4, 2.2]$
10 000	$(1.2 \pm [1.3, 0.6]) \times 10^8$	$35.9 \pm [10.3, 8.0]$	$10.3 \pm [2.6, 2.1]$	$18.4 \pm [2.3, 2.1]$
20 000	$(1.5 \pm [1.4, 0.7]) \times 10^8$	$28.8 \pm [5.8, 4.8]$	$9.0 \pm [2.8, 2.2]$	$17.5 \pm [2.5, 2.2]$
30 000	$(2.2 \pm [3.1, 1.3]) \times 10^8$	$29.8 \pm [7.3, 5.9]$	$8.7 \pm [2.4, 1.9]$	$16.7 \pm [2.0, 1.8]$
100 000	$(3 \pm [2, 1]) \times 10^8$	$24 \pm [6, 5]$	$9 \pm 2$	$17 \pm 1$

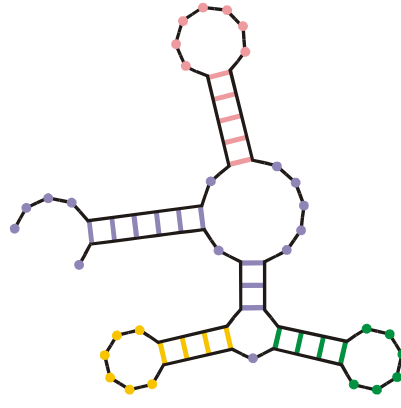




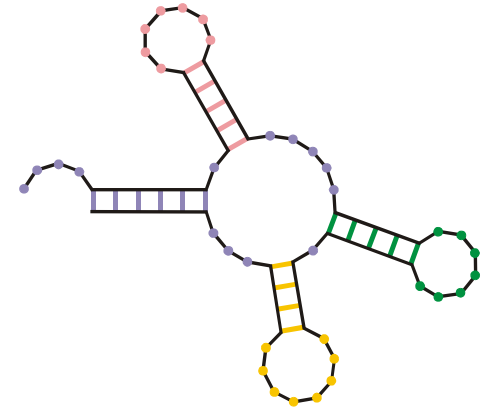
00



09



31

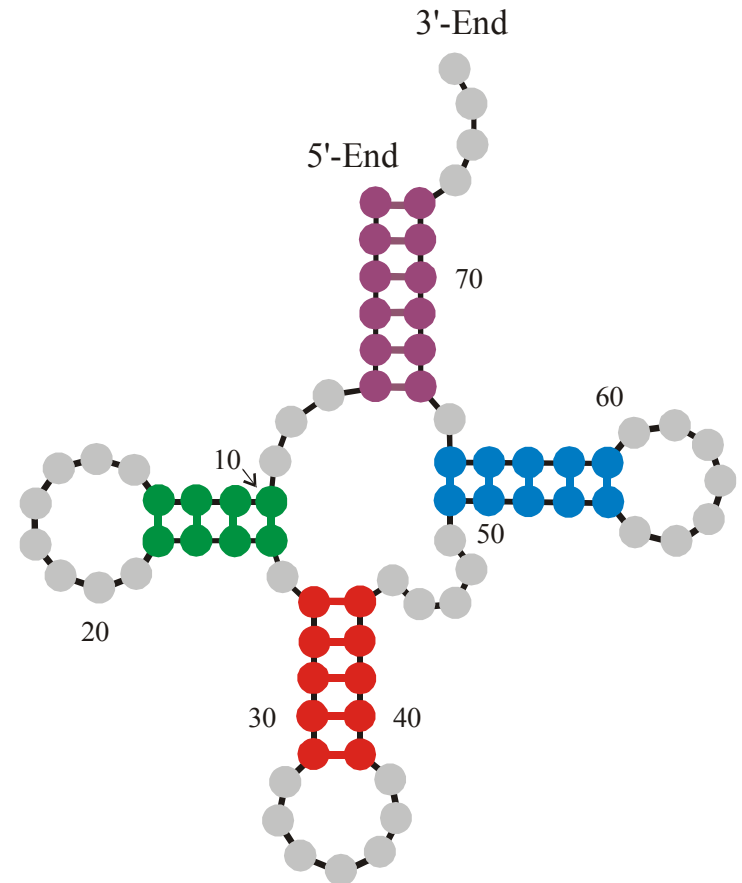


44

Three important steps in the formation of the tRNA clover leaf from a randomly chosen initial structure

Stable tRNA clover leaf structures built from binary, **GC**-only, sequences exist. The corresponding sequences are readily found through inverse folding. Optimization by mutation and selection in the flow reactor has so far always been unsuccessful.

The neutral network of the tRNA clover leaf in **GC** sequence space is not connected, whereas to the corresponding neutral network in **AUGC** sequence space is very close to the critical connectivity threshold,  $\lambda_{cr}$ . Here, both inverse folding and optimization in the flow reactor are successful.



**The success of optimization depends on the connectivity of neutral networks.**

## Main results of computer simulations of molecular evolution

- Individual trajectories are not reproducible. The sequences of the target structures obtained and the relay series were different. Nevertheless, solutions of comparable or the same quality are almost always achieved.
- **Transitions** between molecular phenotypes represented by RNA structures can be classified with respect to the induced structural changes. **Minor transitions** of high probability of occurrence are opposed by **major transitions** of low probability.
- **Major transitions** represent the relevant **structural innovations** in the course of molecular evolution.
- The number of **minor transitions** decreases with increasing population size.
- The number of **major transitions** or **structural innovations** is approximately constant for given start and stop structures.
- Not all structures are accessible through evolution in the flow reactor.

# **Coworkers**

**Walter Fontana**, Santa Fe Institute, NM

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

**Peter F. Stadler**, Universität Wien, AT

**Ivo L. Hofacker**

**Christoph Flamm**

**Bärbel Stadler, Andreas Wernitznig**, Universität Wien, AT

**Michael Kospach, Ulrike Langhammer, Ulrike Mückstein, Stefanie Widder,  
Jan Cupal, Kurt Grünberger, Andreas Svrček-Seiler, Stefan Wuchty**

**Ulrike Göbel**, Institut für Molekulare Biotechnologie, Jena, GE

**Walter Grüner, Stefan Kopp, Jaqueline Weber**

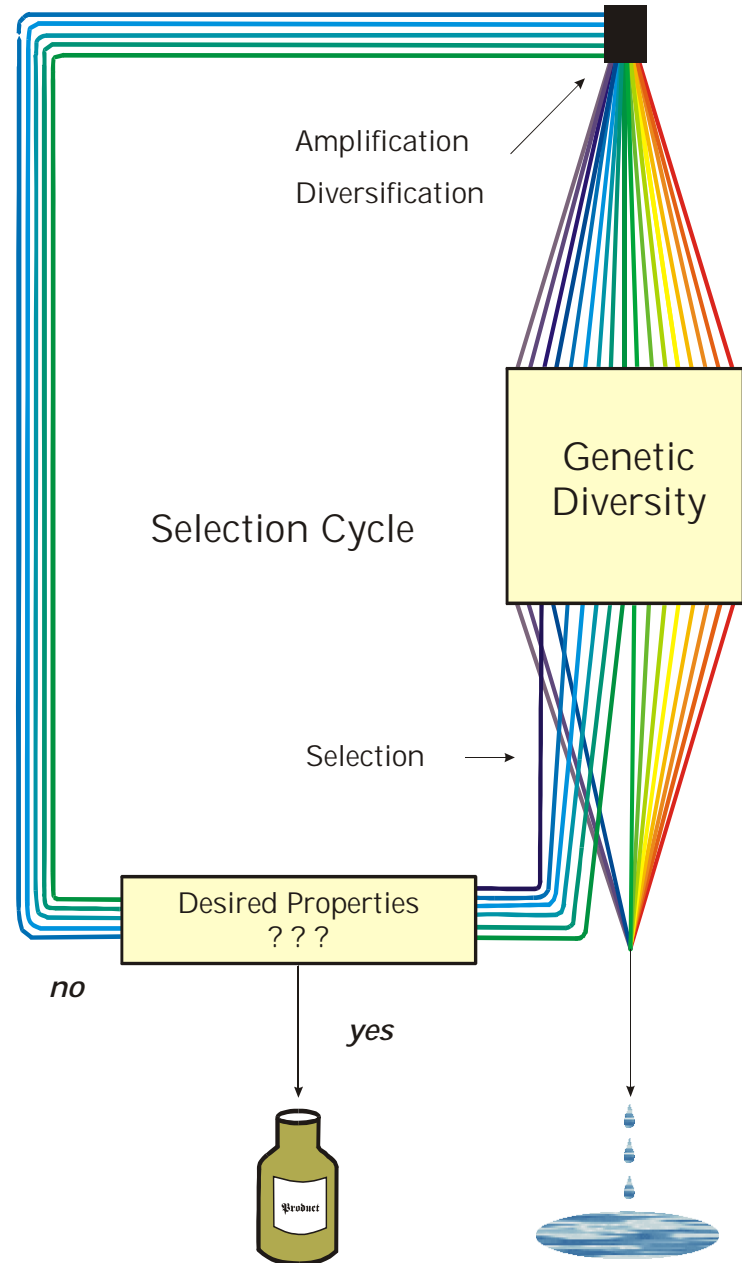
## Evolutionary design of RNA molecules

D.B.Bartel, 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

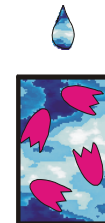
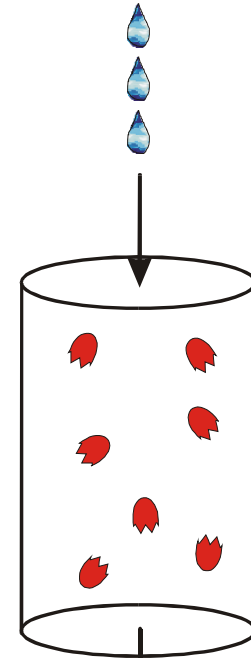
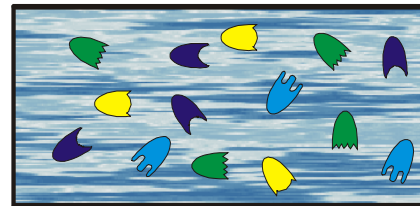
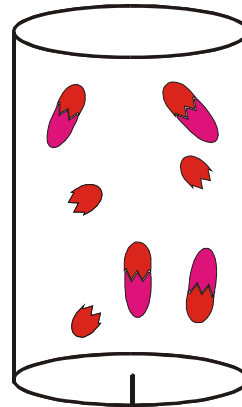
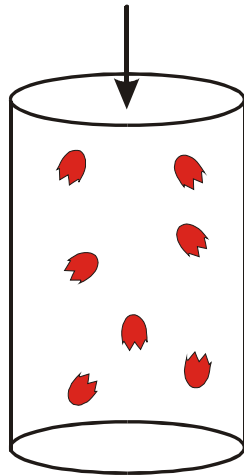
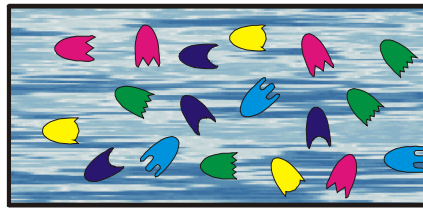


Selection cycle used in applied molecular evolution to design molecules with predefined properties

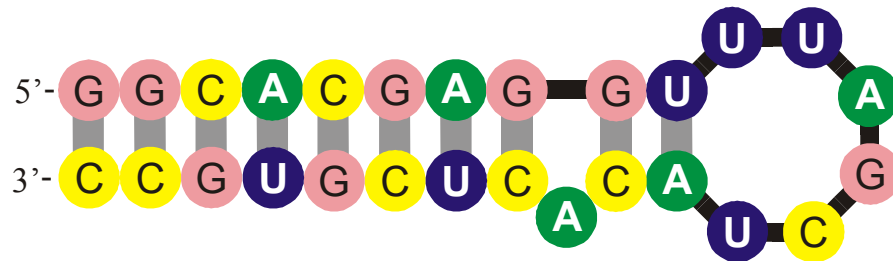
## Retention of binders

## Elution of binders

Chromatographic column



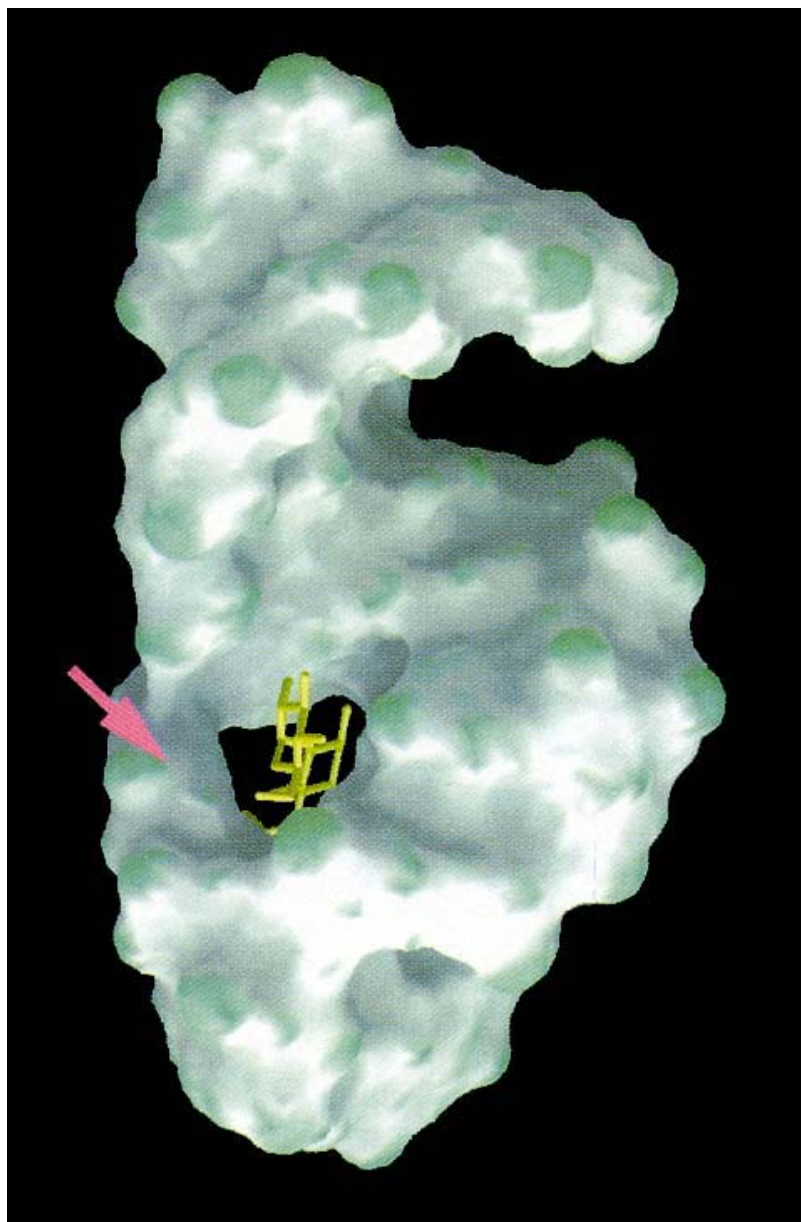
The SELEX technique for the evolutionary design of *aptamers*



Formation of secondary structure of the tobramycin binding RNA aptamer

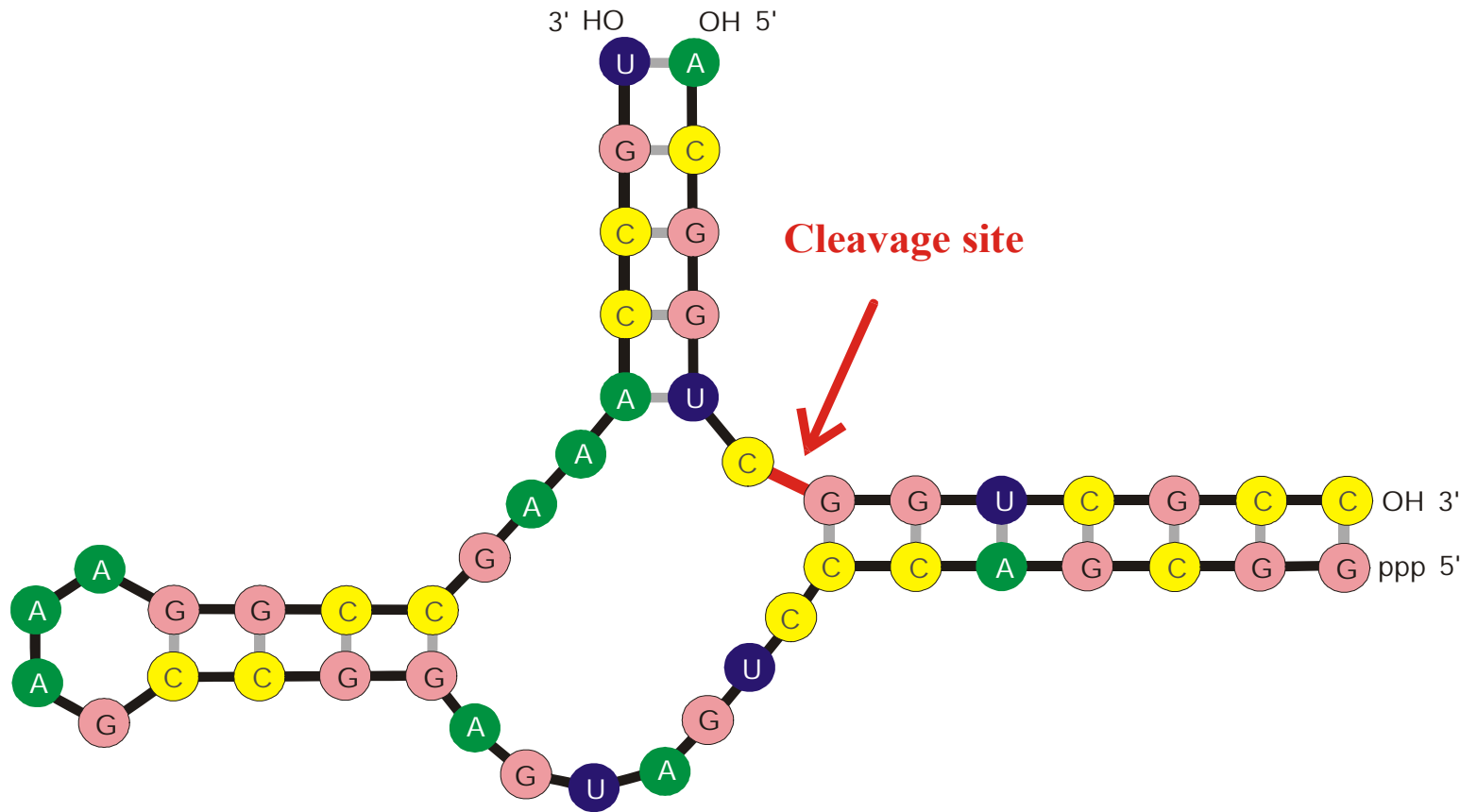
L. Jiang, A. K. Suri, R. Fiala, D. J. Patel, *Chemistry & Biology* 4:35-50 (1997)





The three-dimensional structure of the  
tobramycin aptamer complex

L. Jiang, A. K. Suri, R. Fiala, D. J. Patel,  
*Chemistry & Biology* 4:35-50 (1997)



The "hammerhead" ribozyme

The smallest known  
catalytically active  
RNA molecule