

Complex behavior from simple molecular systems

Peter Schuster

Institut für Theoretische Chemie, Universität Wien, Austria

and

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



Franqui Symposium 2008 in honor of Pierre Gaspard

Complex Systems: A fundamental Science Perspective

Brussels, 03.– 06.09.2008

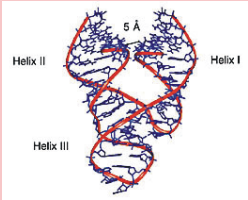
Web-Page for further information:

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

Review article:

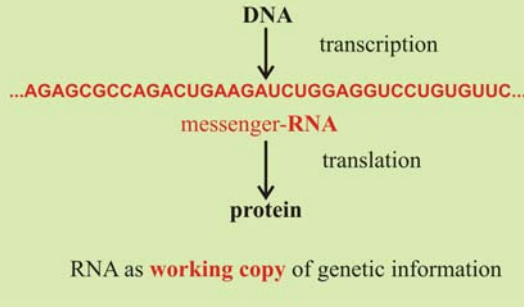
Peter Schuster. Prediction of RNA molecules: From theory to models and real molecules. *Rep.Prog.Phys.* **69**:1419-1477, 2006.

RNA as catalyst

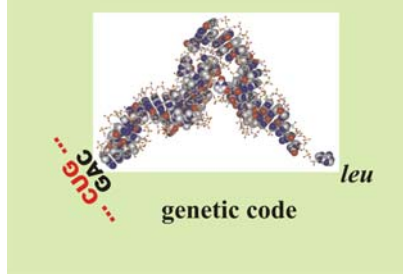


Ribozyme

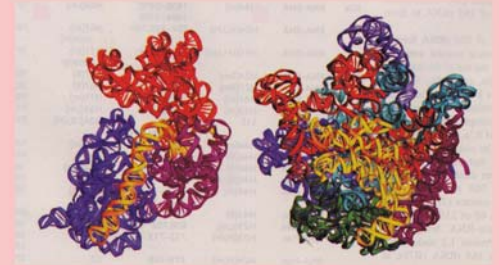
RNA as transmitter of genetic information



RNA as adapter molecule



RNA is the catalytic subunit in supramolecular complexes



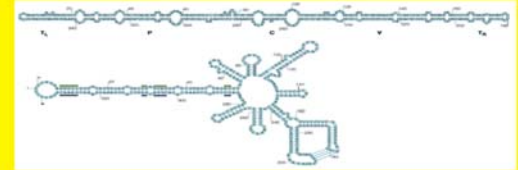
The **ribosome** is a **ribozyme** !

RNA

The RNA world as a precursor of the current DNA + protein biology

RNA is modified by epigenetic control

RNA editing, alternative splicing



Viroids

RNA as carrier of genetic information

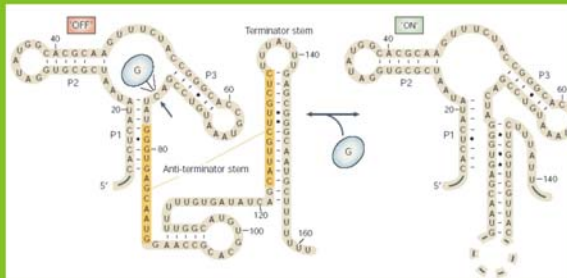
RNA viruses and retroviruses

RNA evolution *in vitro*

Evolutionary biotechnology

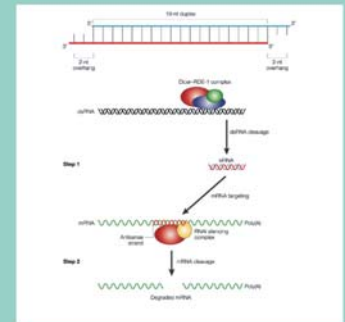
RNA aptamers, artificial ribozymes, allosteric ribozymes

Allosteric control of transcribed RNA



Riboswitches controlling transcription and translation through **metabolites**

RNA as regulator of gene expression



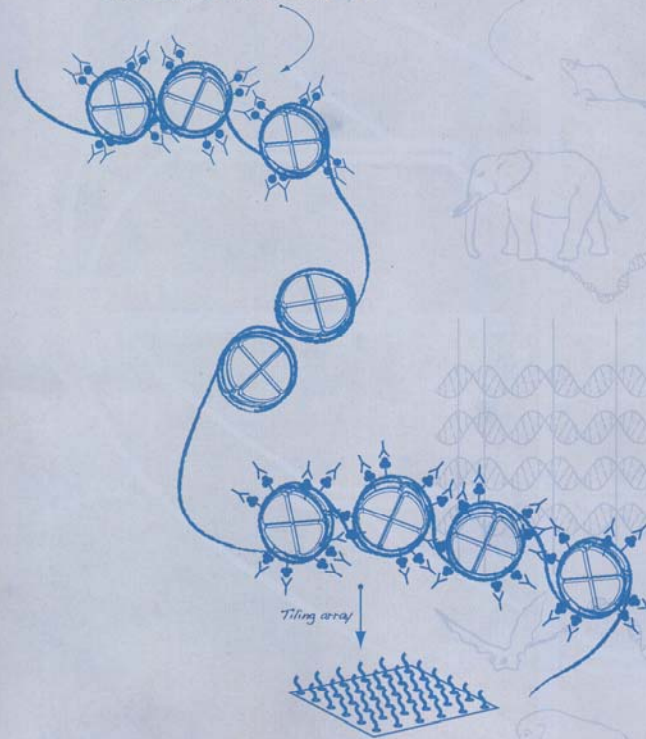
Gene regulation by small interfering RNAs

RNA – The magic molecule

nature

Hi-stone-modification chromatin IP

Comparative genomics 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

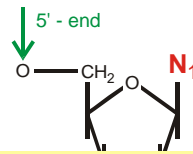


ENCODE stands for
ENCyclopedia **Of** **DNA** **E**lements.

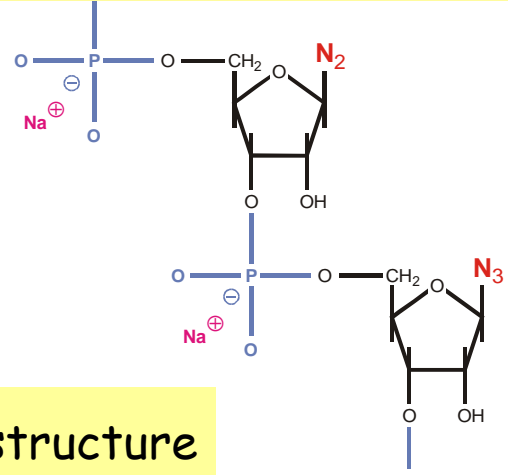
ENCODE Project Consortium.
Identification and analysis of functional
elements in 1% of the human genome by
the ENCODE pilot project.
Nature **447**:799-816, 2007

1. Minimum free energy structures of RNA
2. Suboptimal structures of RNA
3. Kinetic folding and RNA switches
4. Chemistry of Darwinian evolution
5. Consequences of neutrality
6. Evolutionary optimization of RNA structure

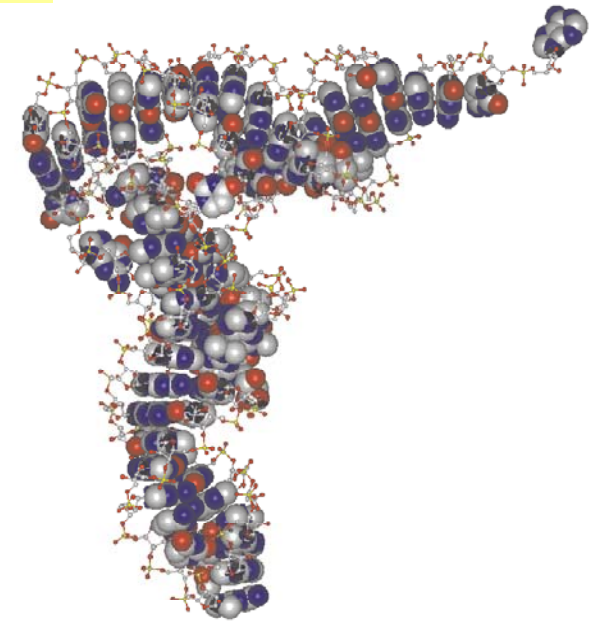
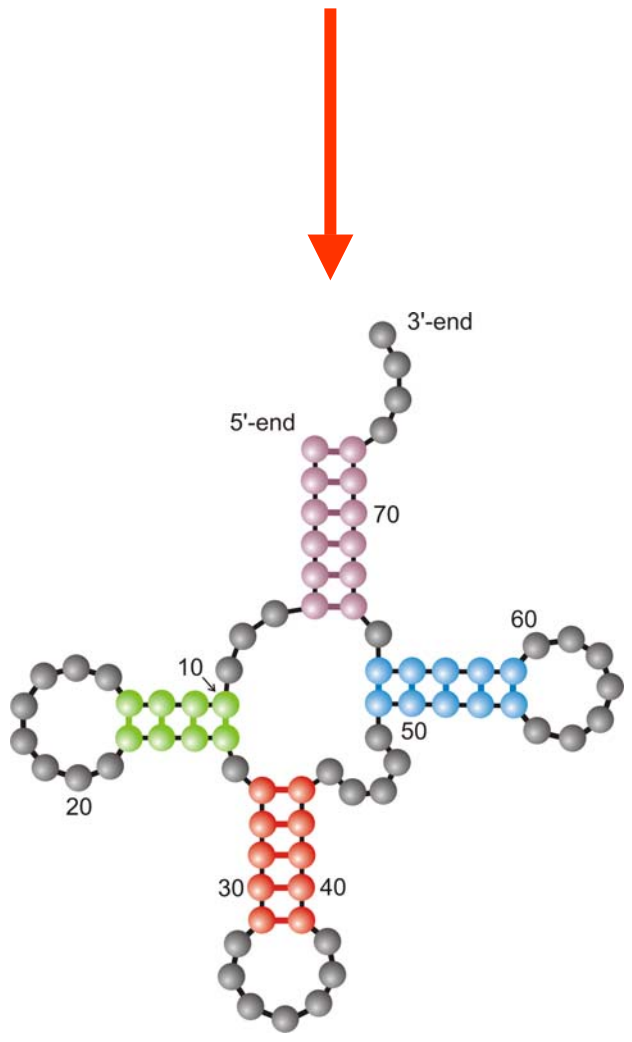
1. **Minimum free energy structures of RNA**
2. Suboptimal structures of RNA
3. Kinetic folding and RNA switches
4. Chemistry of Darwinian evolution
5. Consequences of neutrality
6. Evolutionary optimization of RNA structure

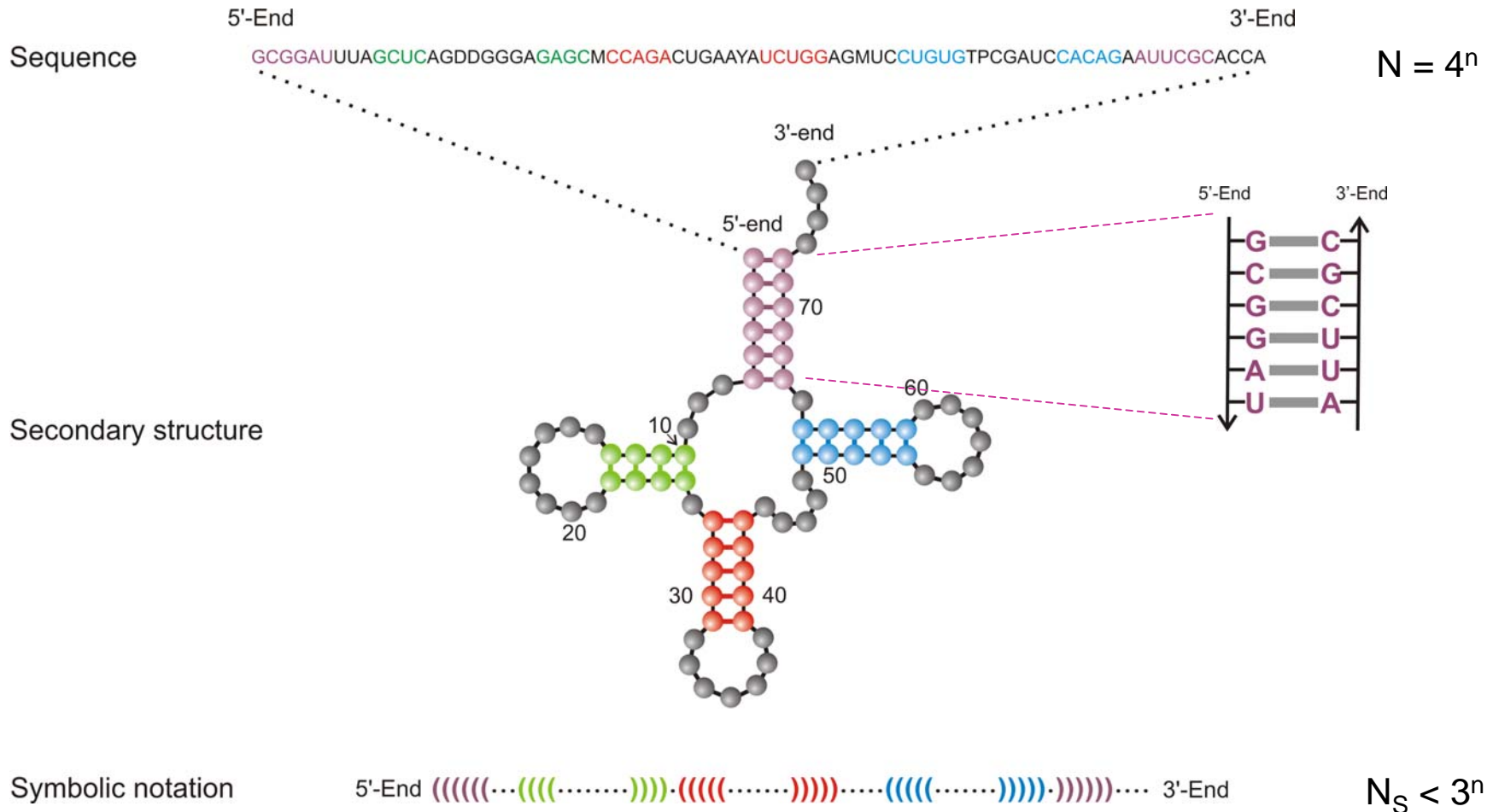


5'-end **GCGGAUUUAGCUC**AGUUGGGAGAG**CGCCAGACUGAAGAUCUGG**AGGUC**CUGUGUUCGAUCCACAGAAUUCGCACCA** 3'-end



Definition of RNA structure

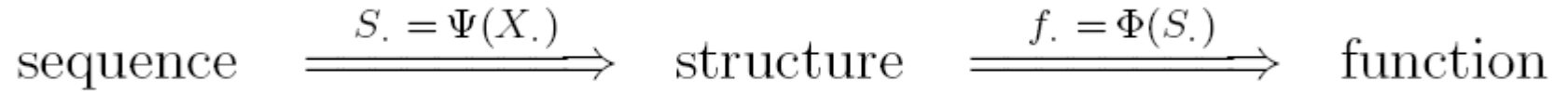




Criterion: Minimum free energy (mfe)

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

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



The paradigm of structural biology

What is neutrality ?

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

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

RNA sequence

GUAUCGAAAUACGUAGCGUAUGGGGAUGCUGGACGGUCCCAUCGGUACUCCA

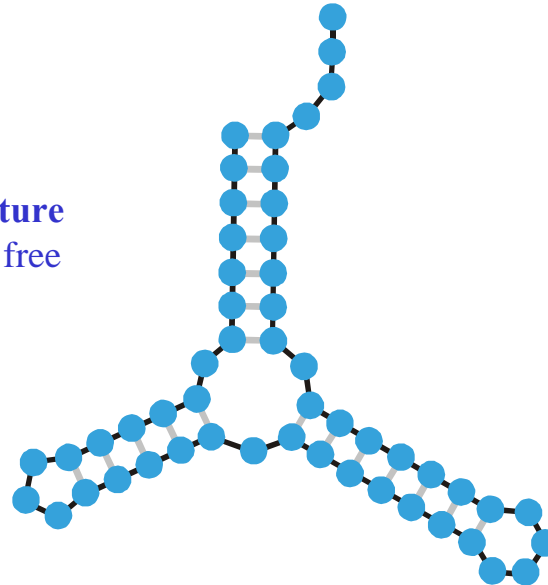
RNA folding:
Structural biology,
spectroscopy of
biomolecules,
understanding
molecular function

Biophysical chemistry:
thermodynamics and
kinetics



Empirical parameters

RNA structure
of minimal free
energy

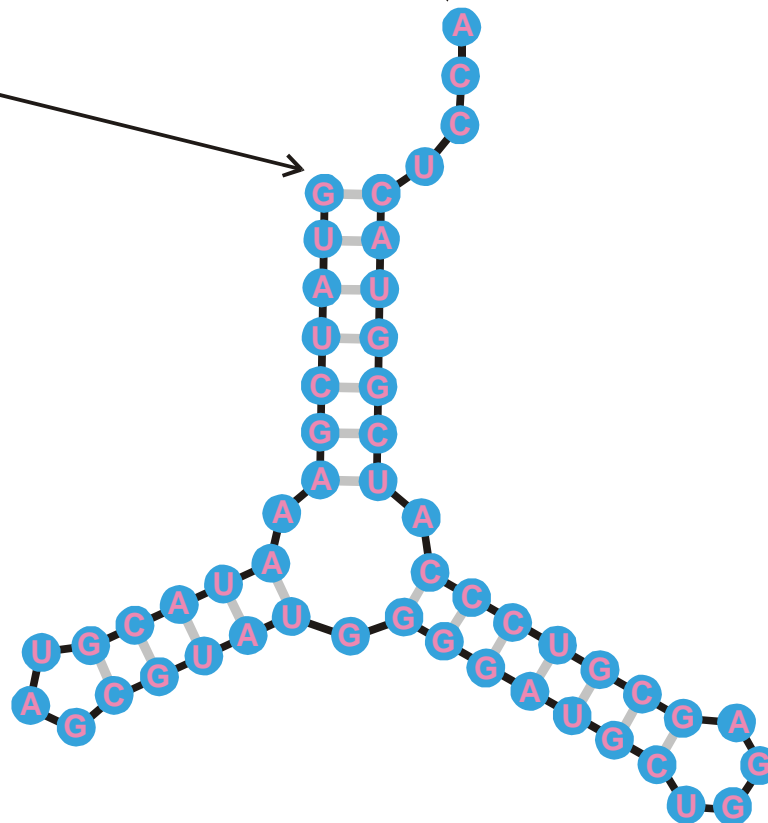
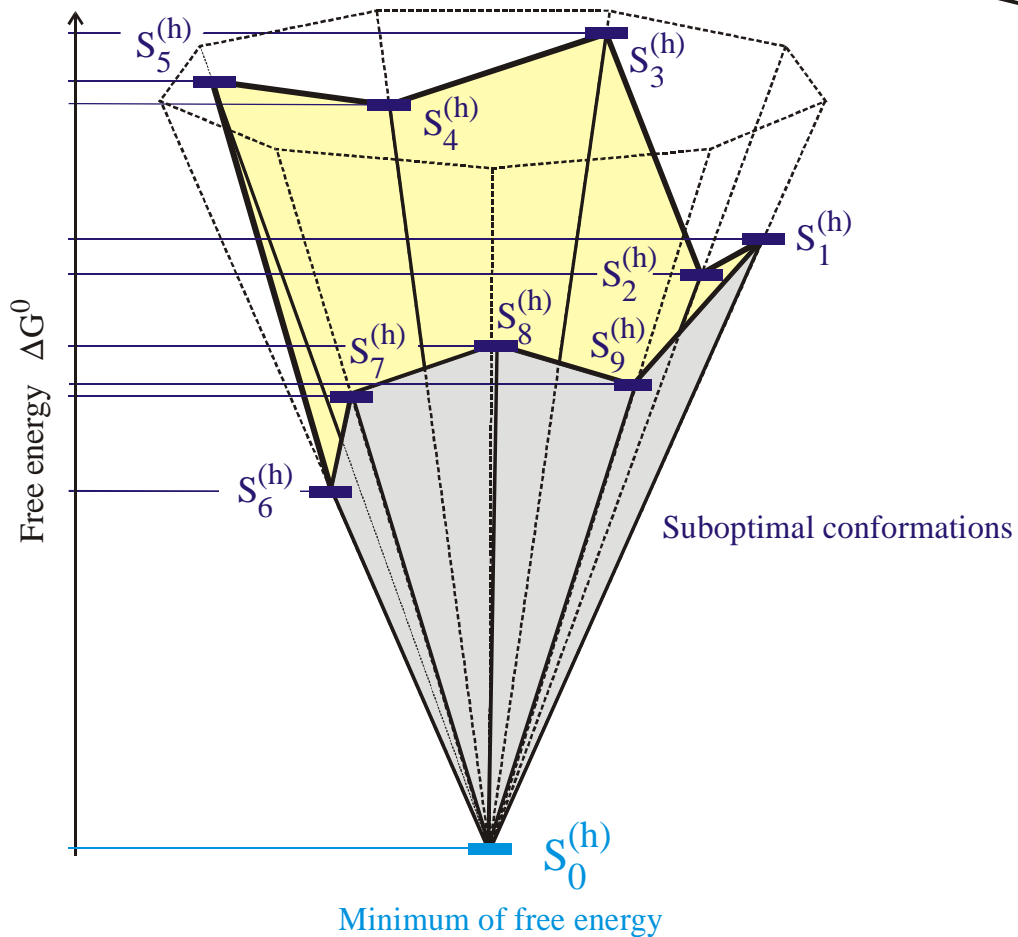


Sequence, structure, and design

5'-end

3'-end

GUAUCGAAUACGUAGCGUAUGGGGAUGCUGGACGGUCCCAUCGGUACUCCA



The minimum free energy structures on a discrete space of conformations



Minimum free energy structure

Extension of the notion of structure

RNA sequence

GUAUCGAAAUACGUAGCGUAUGGGGAUGCUGGACGGUCCCAUCGGUACUCCA

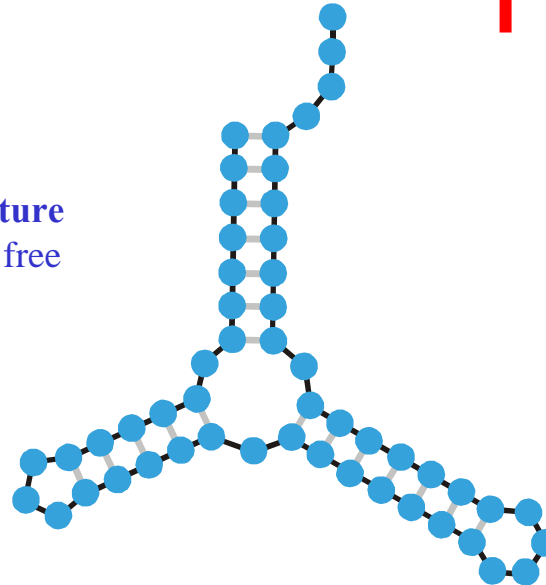
RNA folding:
Structural biology,
spectroscopy of
biomolecules,
understanding
molecular function

Iterative determination
of a sequence for the
given secondary
structure

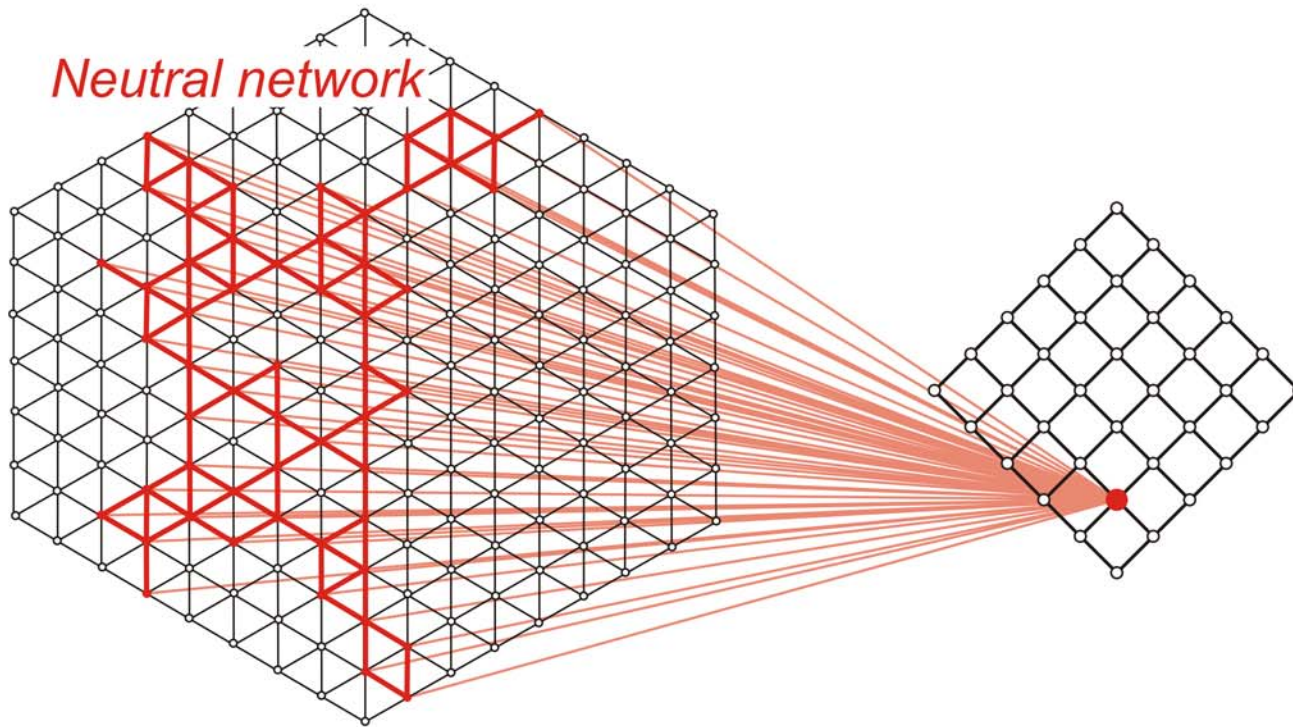
**Inverse Folding
Algorithm**

Inverse folding of RNA:
Biotechnology,
design of biomolecules
with predefined
structures and functions

RNA structure
of minimal free
energy



Sequence, structure, and design



Sequence space

Structure space

many sequences

⇒

one structure

Space of genotypes: $Q = \{X_1, X_2, X_3, X_4, \dots, X_N\}$; Hamming metric

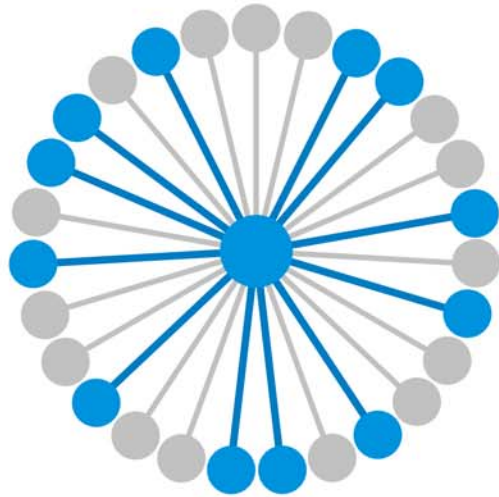
Space of phenotypes: $S = \{S_1, S_2, S_3, S_4, \dots, S_M\}$; metric (not required)

$$N \gg M$$

$$\psi(X_j) = S_k$$

$$G_k = \psi^{-1}(S_k) \doteq \{ X_j \mid \psi(X_j) = S_k \}$$

A mapping ψ and its inversion



$$\lambda_j = 12 / 27 = 0.444$$

$$\mathbf{G}_k = \psi^{-1}(\mathbf{S}_k) \doteq \{ I_j \mid \psi(I_j) = \mathbf{S}_k \}$$

$$\bar{\lambda}_k = \frac{\sum_{j \in |\mathbf{G}_k|} \lambda_j(k)}{|\mathbf{G}_k|}$$

Alphabet size κ :

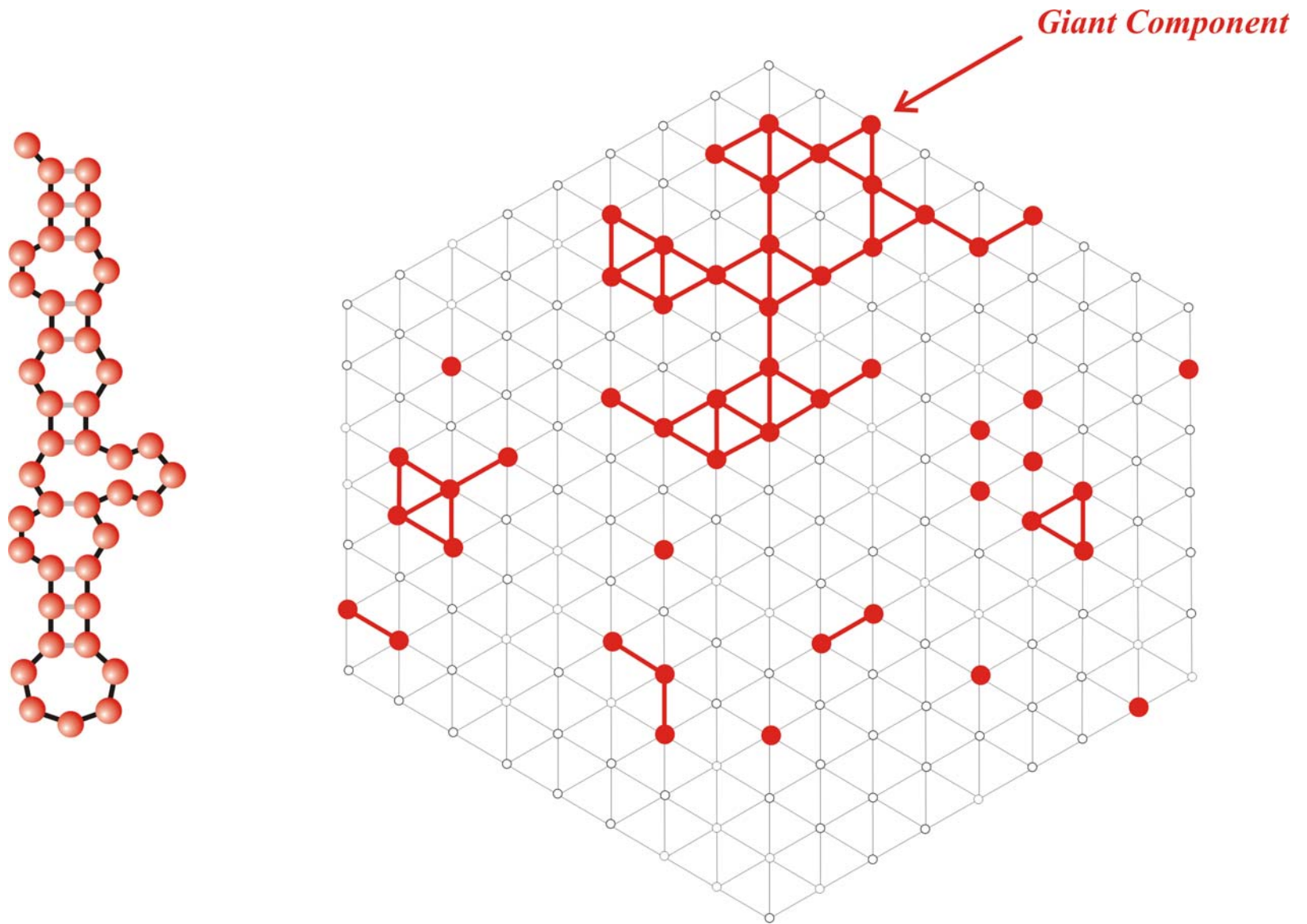
κ	λ_{cr}	
2	0.5	AU,GC,DU
3	0.423	AUG , UGC
4	0.370	AUGC

$\bar{\lambda}_k > \lambda_{cr}$ network \mathbf{G}_k is connected

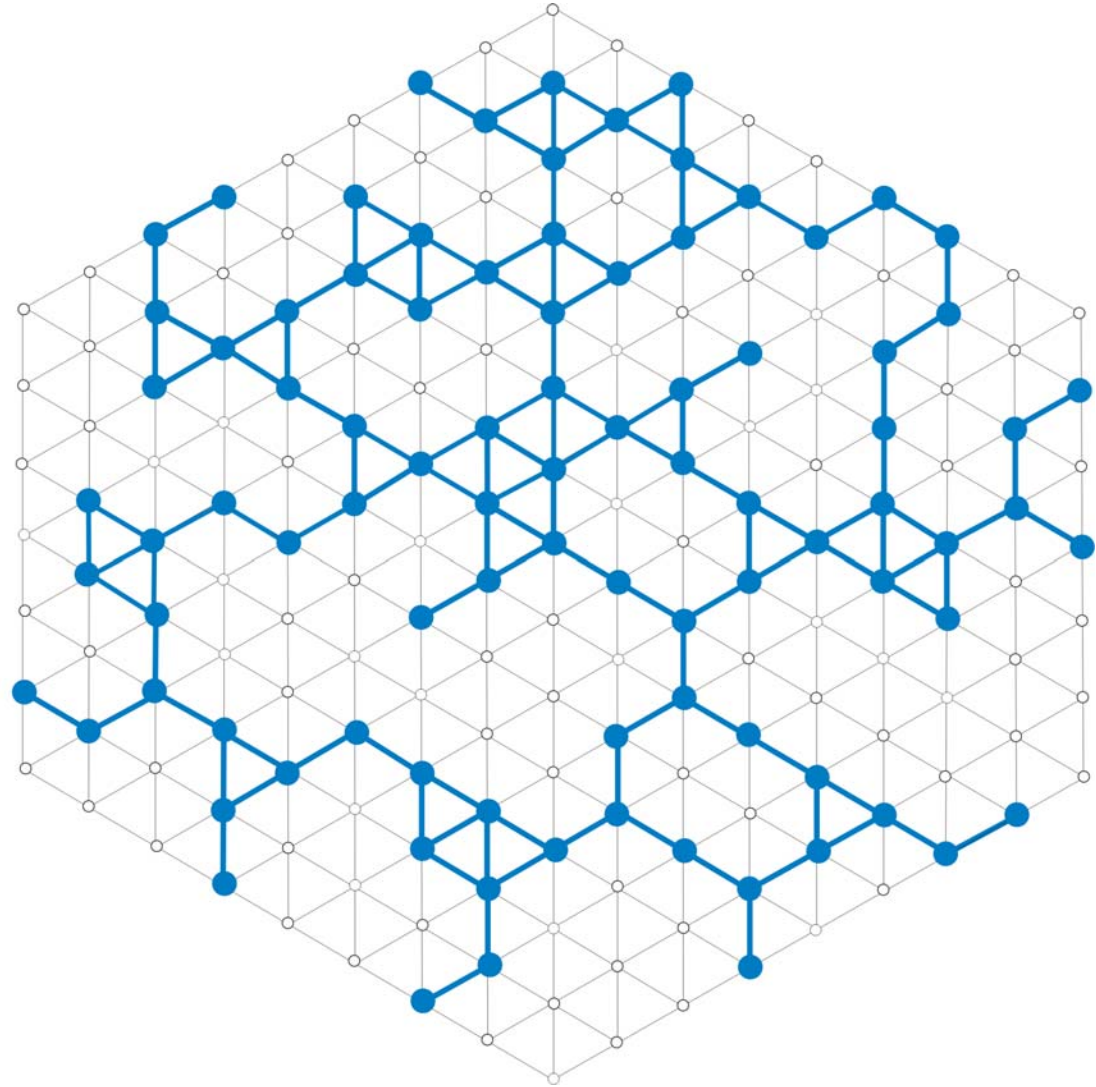
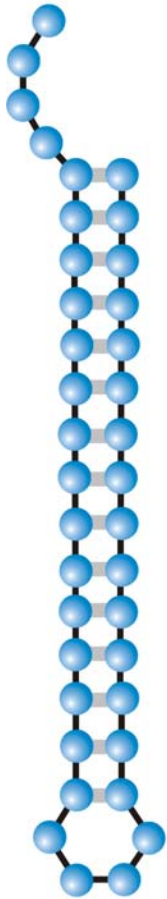
$\bar{\lambda}_k < \lambda_{cr}$ network \mathbf{G}_k is **not** connected

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

Degree of neutrality of neutral networks and the connectivity threshold

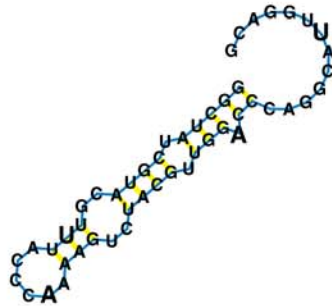


A multi-component neutral network formed by a rare structure: $\lambda < \lambda_{cr}$



A connected neutral network formed by a common structure: $\lambda > \lambda_{\text{cr}}$

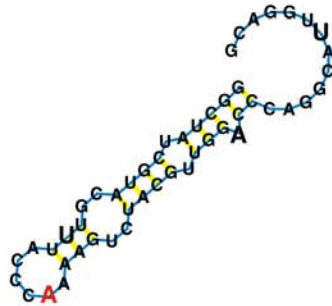
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG



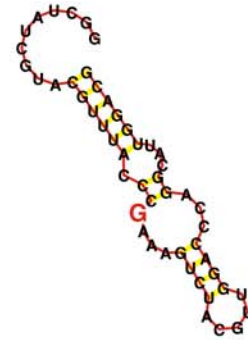
One error neighborhood – Surrounding of an RNA molecule of chain length $n=50$ in sequence and shape space

GGCUAUCGUACGUUUACCCGAAAGUCUACGUUGGACCCAGGCAUUGGACG

GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG

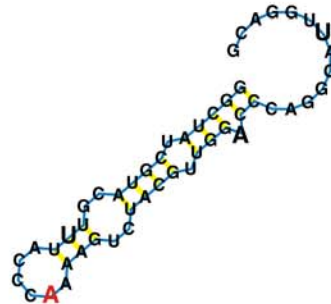


One error neighborhood – Surrounding of an RNA molecule of chain length $n=50$ in sequence and shape space

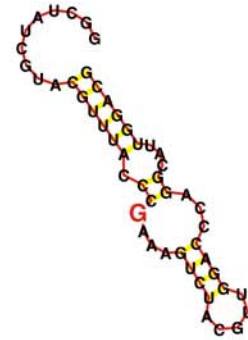


GGCUAUCGUACGUUUACCCGAAAGUCUACGUUGGACCCAGGCAUUGGACG

GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG



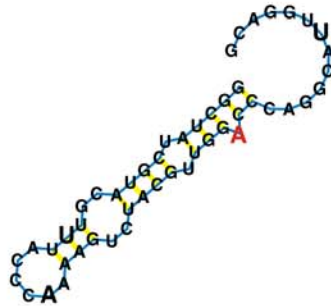
One error neighborhood – Surrounding of an RNA molecule of chain length $n=50$ in sequence and shape space



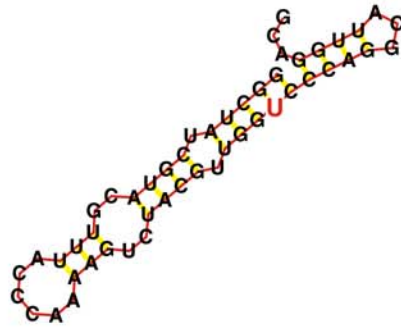
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGUCCAGGCAUUGGACG

GGCUAUCGUACGUUUACCCGAAAGUCUACGUUGGACCCAGGCAUUGGACG

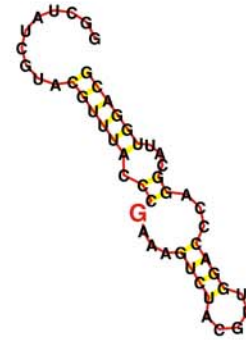
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCAGGCAUUGGACG



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

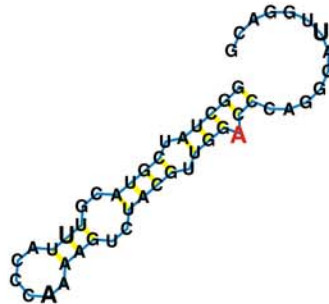


GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGG**U**CCAGGCAUUGGACG



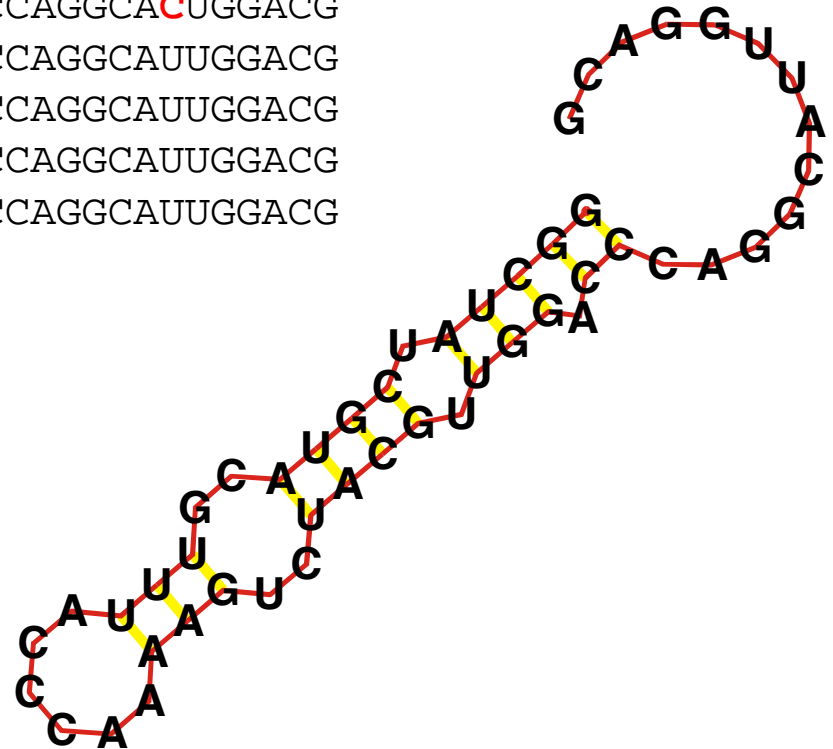
GGCUAUCGUACGUUUACCC**G**AAAGUCUACGUUGGACCCAGGCAUUGGACG

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



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

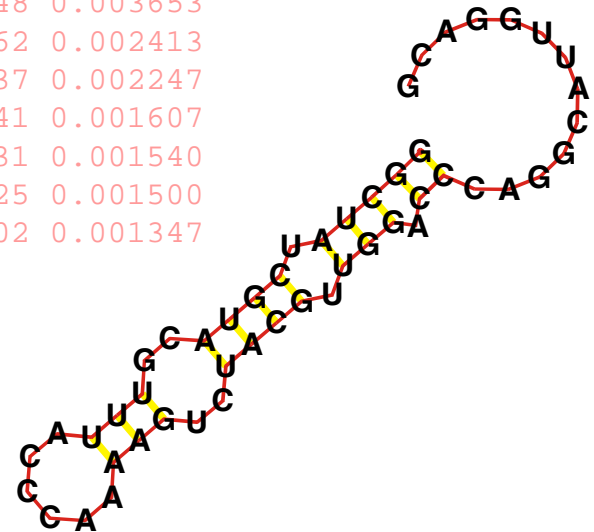
GGCUAUCGUAU**U**GUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUA**A**GACG
GGCUAUCGUACGUUUAC**U**CAAAGUCUACGUUGGACCCAGGCAUUGGACG
GGCUAUCGUACG**C**UUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG
GGC**C**AUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG
GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG
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	0.334167	0.006961	0.083434
Number of Structures:	1000	52.31	85.30	9.24

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

Results from RNA minimum free energy structures:

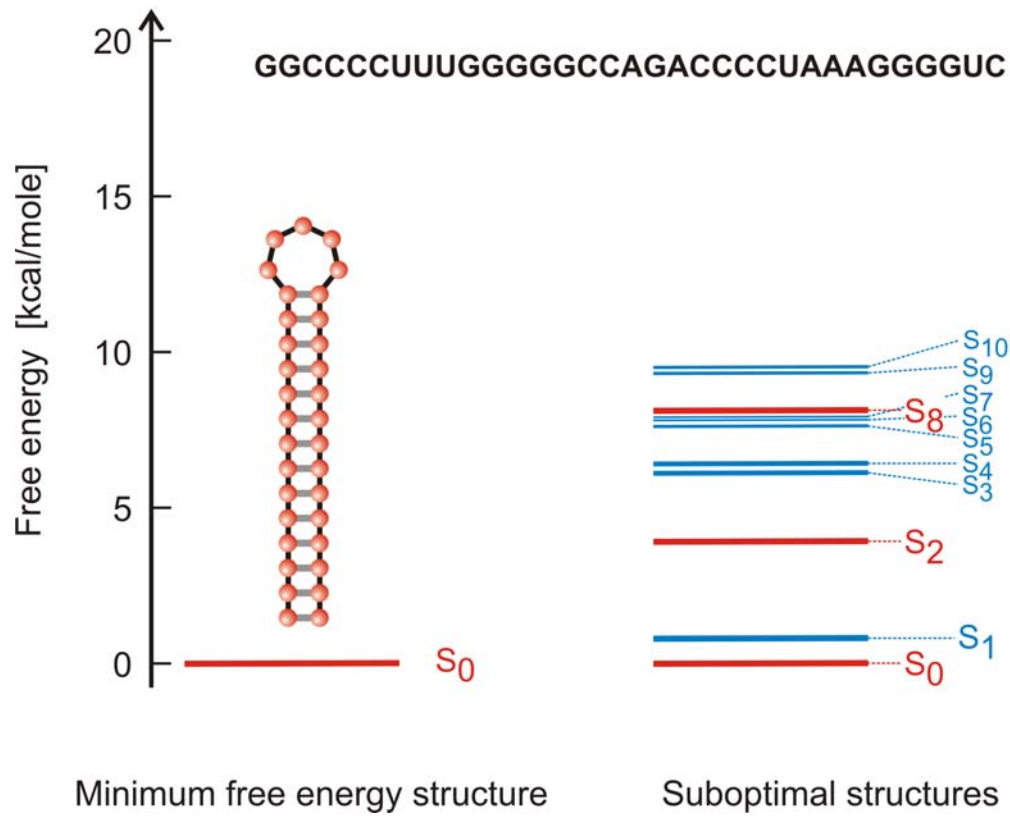
- RNA minimum free energy structures show neutrality: Many sequences fold into the same (secondary) structure.
- The single base mutation neighborhood contains structures from neutral sequences **and** a great variety of other structures: Biopolymer landscapes are **rugged**.

1. Minimum free energy structures of RNA
2. **Suboptimal structures of RNA**
3. Kinetic folding and RNA switches
4. Chemistry of Darwinian evolution
5. Consequences of neutrality
6. Evolutionary optimization of RNA structure



Minimum free energy structure

Extension of the notion of structure

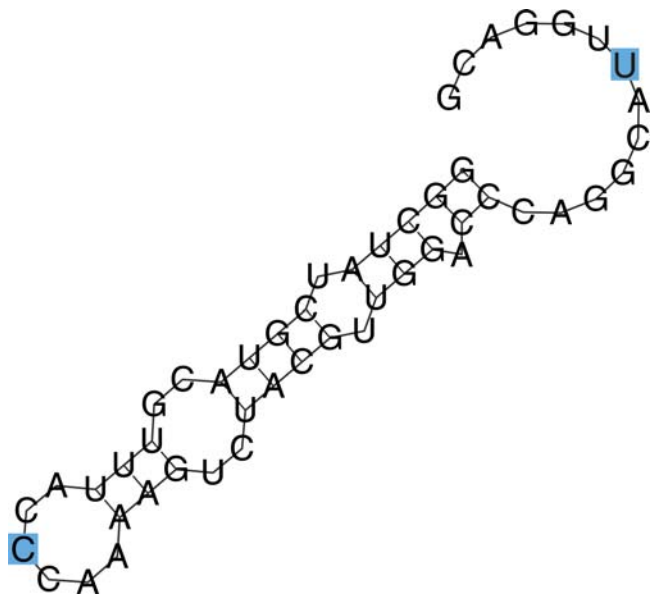


Extension of the notion of structure

GGCUAUCGUACGUUUACCCAAAAGUCUACGUUGGACCCAGGCAUUGGACG

(((((((((.....)))))))).))..... -7.30

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



GGCUAUCGUACGUUUACAAAAGUCUACGUUGGACCCAGGCAUUGGACG

(((((((((.....)))))))).))..... -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

At equilibrium and temperature T the conformations form a Boltzmann ensemble that contains S_j with the Boltzmann weight $\gamma_j(T) = g_j \exp(-(\varepsilon_j - \varepsilon_0)/RT)/Q(T)$, where R is the Boltzmann constant for 1 mole, $R = N_L \cdot k_B$, and $Q(T)$ is the partition function

$$Q(T) = \sum_i g_i \exp\left(-(\varepsilon_i - \varepsilon_0)/RT\right).$$

$$\gamma_j(T) = g_j \exp\left(-(\varepsilon_j - \varepsilon_0)/RT\right)/Q(T)$$

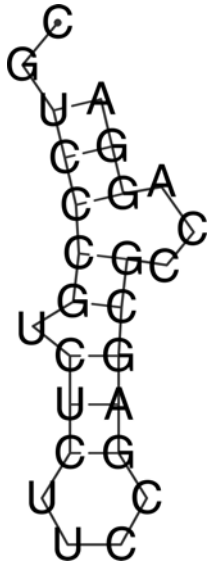
$$P(X, T) = \sum_k \gamma_k(T) A(S_k) \quad \text{or} \quad p_{ij}(X, T) = \sum_k \gamma_k(T) a_{ij}(S_k)$$

$A(S_k)$... adjacency matrix of structure S_k

$p_{ij}(X, T)$... base pairing probability

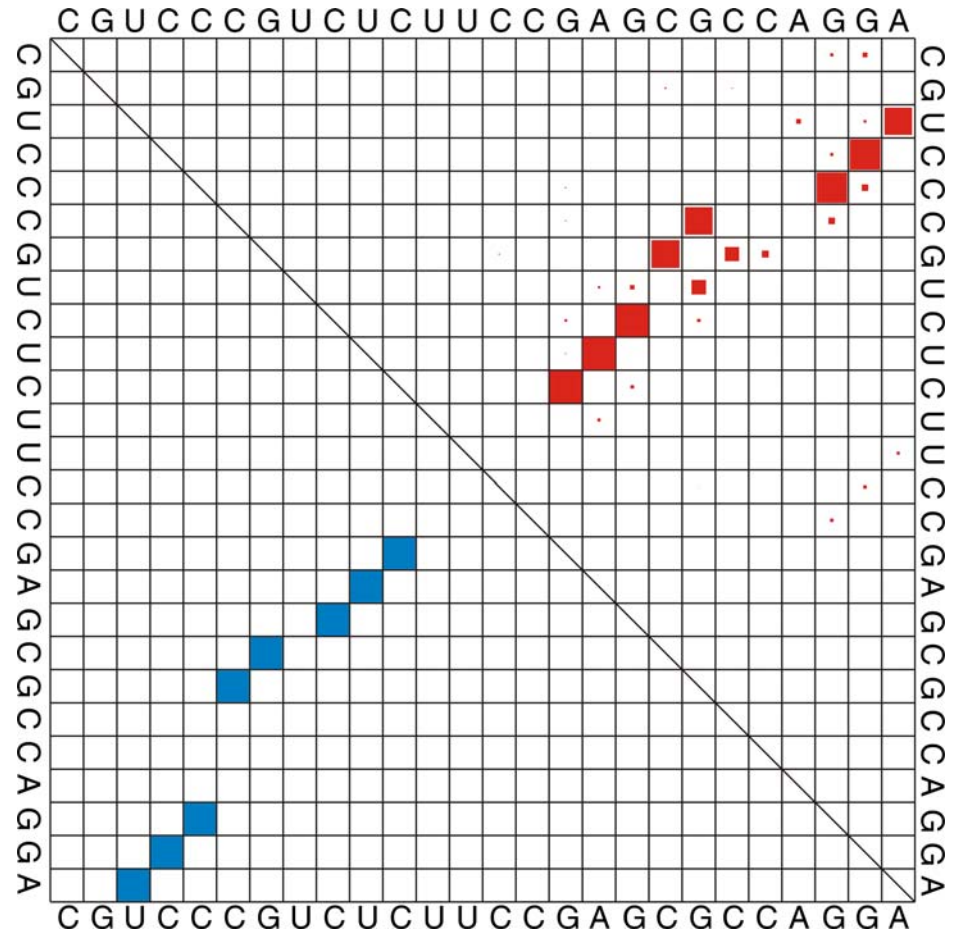
X ... sequence

Usage of the partition function to analyze the spectrum of suboptimal states



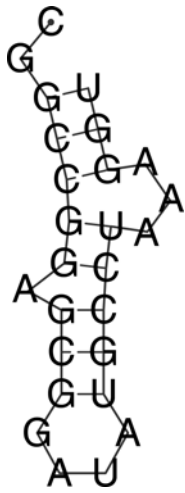
CGUCCCGUCUCUCCGAGCGCCAGGA

- ..((((((.((((.....)))))).....)) -4.50
- ..(((.(((((((.....))))).))..)) -3.70
- ...((((.((((.....)))))).....). -3.60
-((.((((.....))))))..... -3.00
- ...((.(((((((.....))))).))..). -2.80
- ..(((.(.((((.....))))..)).)) -2.60
- (.((..((.((((.....))))))..)).) -2.50



mfe-weight: 0.46336

Suboptimal structures and partition function of a small RNA molecule: $n = 26$



CGGCCGGAGCGGAUAUGCCUAAAGGU

..((((((.((((.....)))))).....))) -3.70

..(((.....(((.....)))).....))) -3.60

..(((.(.(((.....)))..)).....))) -3.50

..(((..(((.....)))..)).....))) -3.30

..(((..(((.....)))..)).....))) -3.30

..(((.(.(((.....)))).....))) -3.10

(.(((.....)).).....(((.....))) -2.90

..(((.....(((.....))).....))) -2.90

..(((.....))).....(((.....))) -2.90

..(((((((.....)))).....))) -2.70

..(((.....(((.....))).....))) -2.60

..(((.....))).....(((.....))) -2.60

..(((.(.(((.....)))..)).....))) -2.50

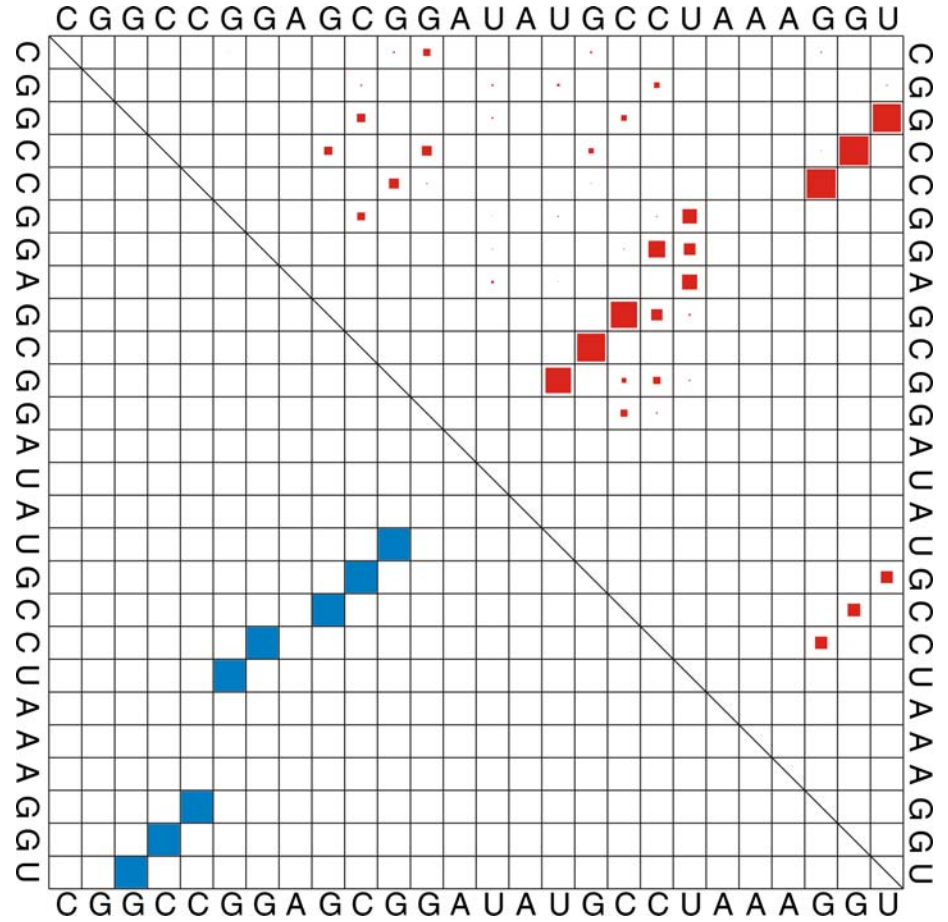
..(((..(((.....)))..)).....))) -2.50

..(((.....)))..... -2.30

..(((..(((.....)))..)).....))) -2.30

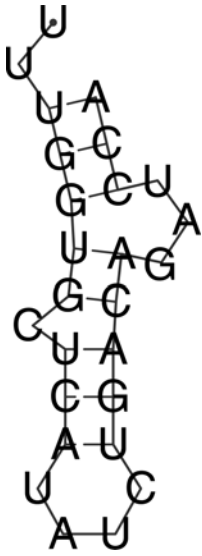
..(((..(((.....)))..)).....))) -2.30

.....(((.....)))..... -2.20



mfe-weight: 0.13642

Suboptimal structures and partition function of a small RNA molecule: $n = 26$



UUUGGUGCUCAUAUCUGACAGAUCCA

..((((((.....))))))...)) -1.10

..(((.....(((.....)))).....)) -1.00

...(((.....(((.....)))).....). -1.00

...(((.....(((.....)))).....). -0.90

..((((((.....(((.....)))).....)) -0.70

..(((.....(((.....)))).....)) -0.60

...(((.....(((.....)))).....). -0.60

...(((.....(((.....)))).....). -0.50

.....(((.....(((.....)))).....) -0.20

..(((.....(((.....)))).....)) -0.10

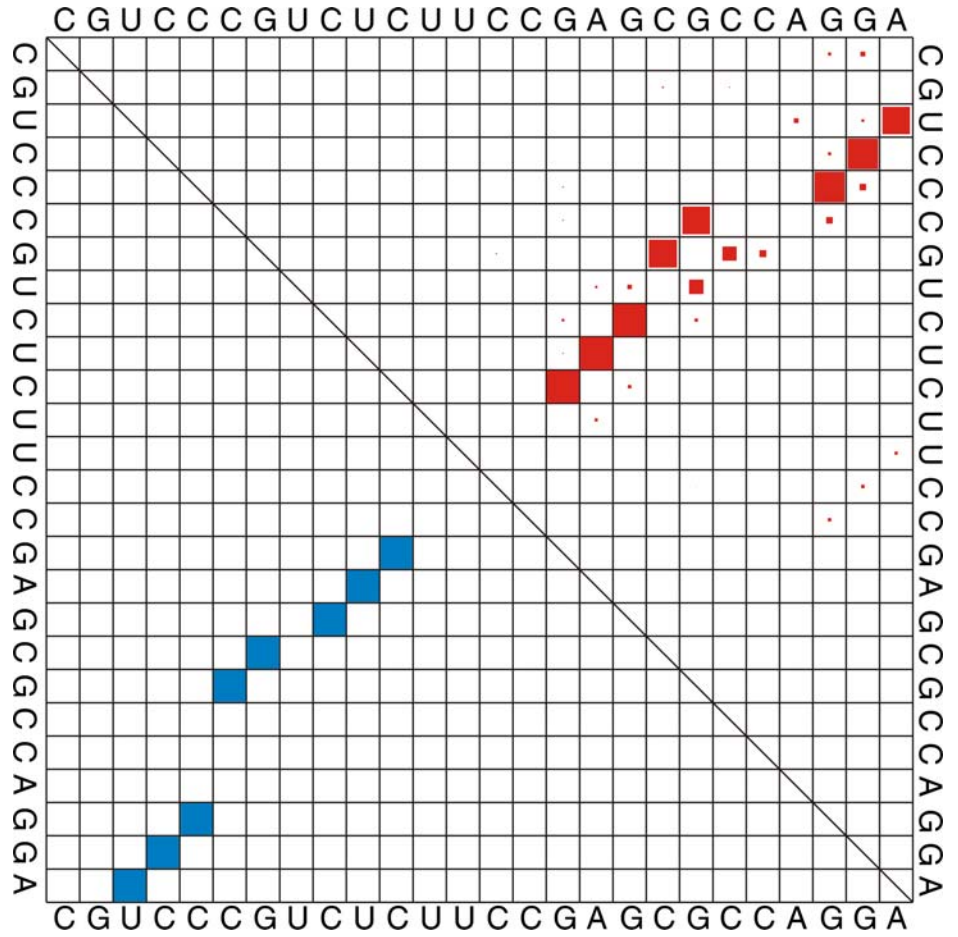
..(((.....(((.....)))).....). -0.10

(((((.....(((.....)))).....)).....) 0.00

...(((.....(((.....)))).....). 0.00

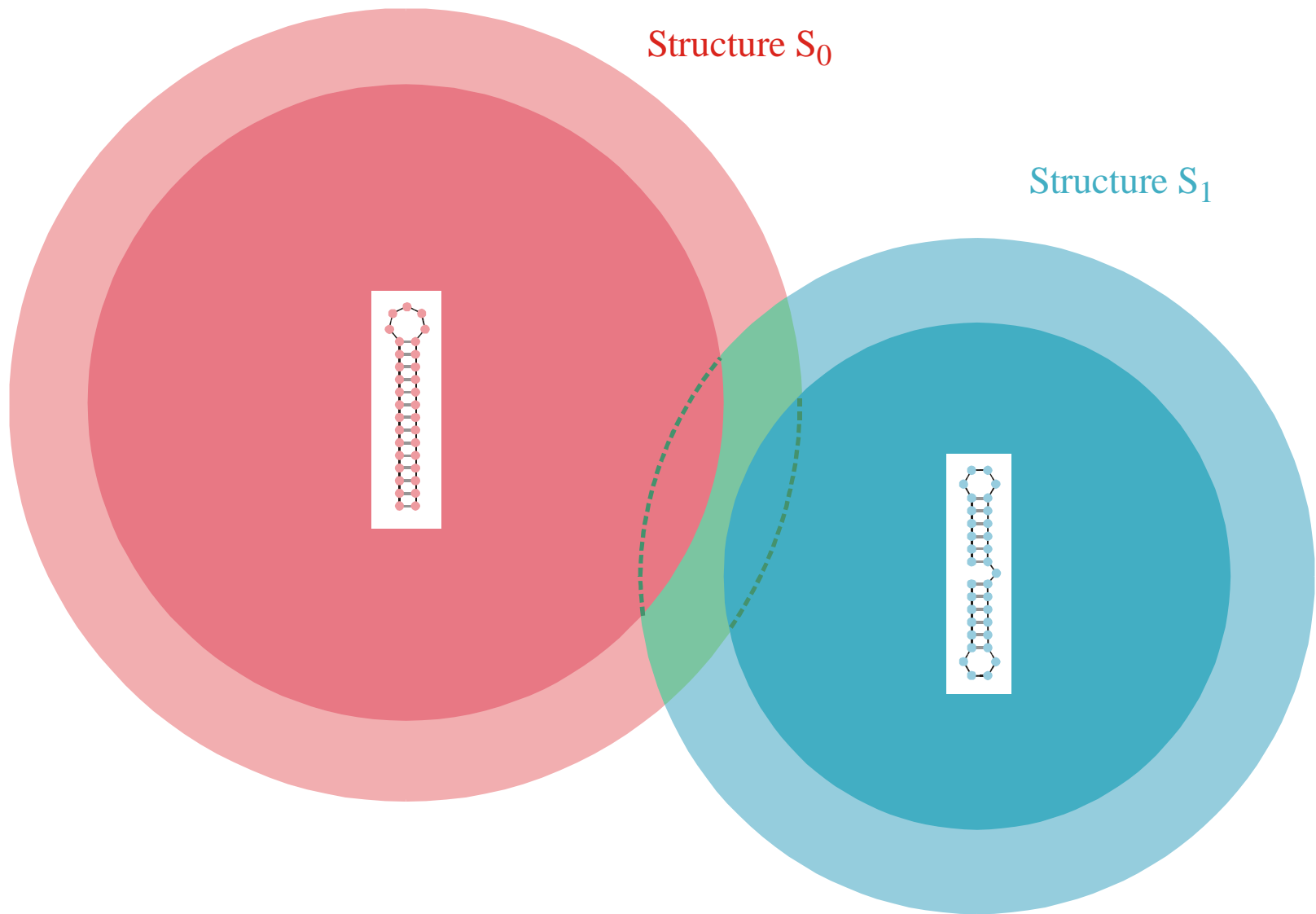
...(((.....(((.....)))).....). 0.00

..... 0.00



mfe-weight: 0.09514

Suboptimal structures and partition function of a small RNA molecule: n = 26



Intersection of two compatible sets: $C_0 \cap C_1$

The intersection of two compatible sets is always non empty: $C_0 \cap C_1 \neq \emptyset$



S0092-8240(96)00089-4

GENERIC PROPERTIES OF COMBINATORY MAPS: NEUTRAL NETWORKS OF RNA SECONDARY STRUCTURES¹

■ CHRISTIAN REIDYS*, †, PETER F. STADLER*, ‡
 and PETER SCHUSTER*, ‡, §, ¶²

*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 (λ). 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. ■

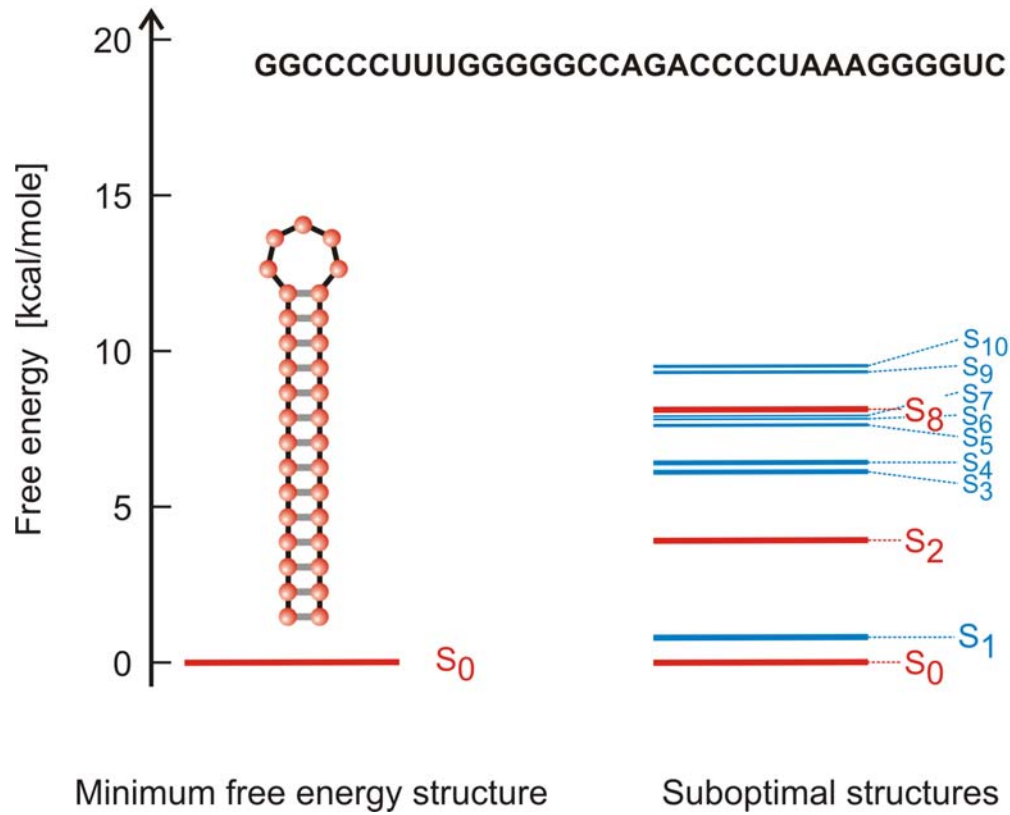
Remark. A generalization of the statement of theorem 5 to three different structures is false.

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

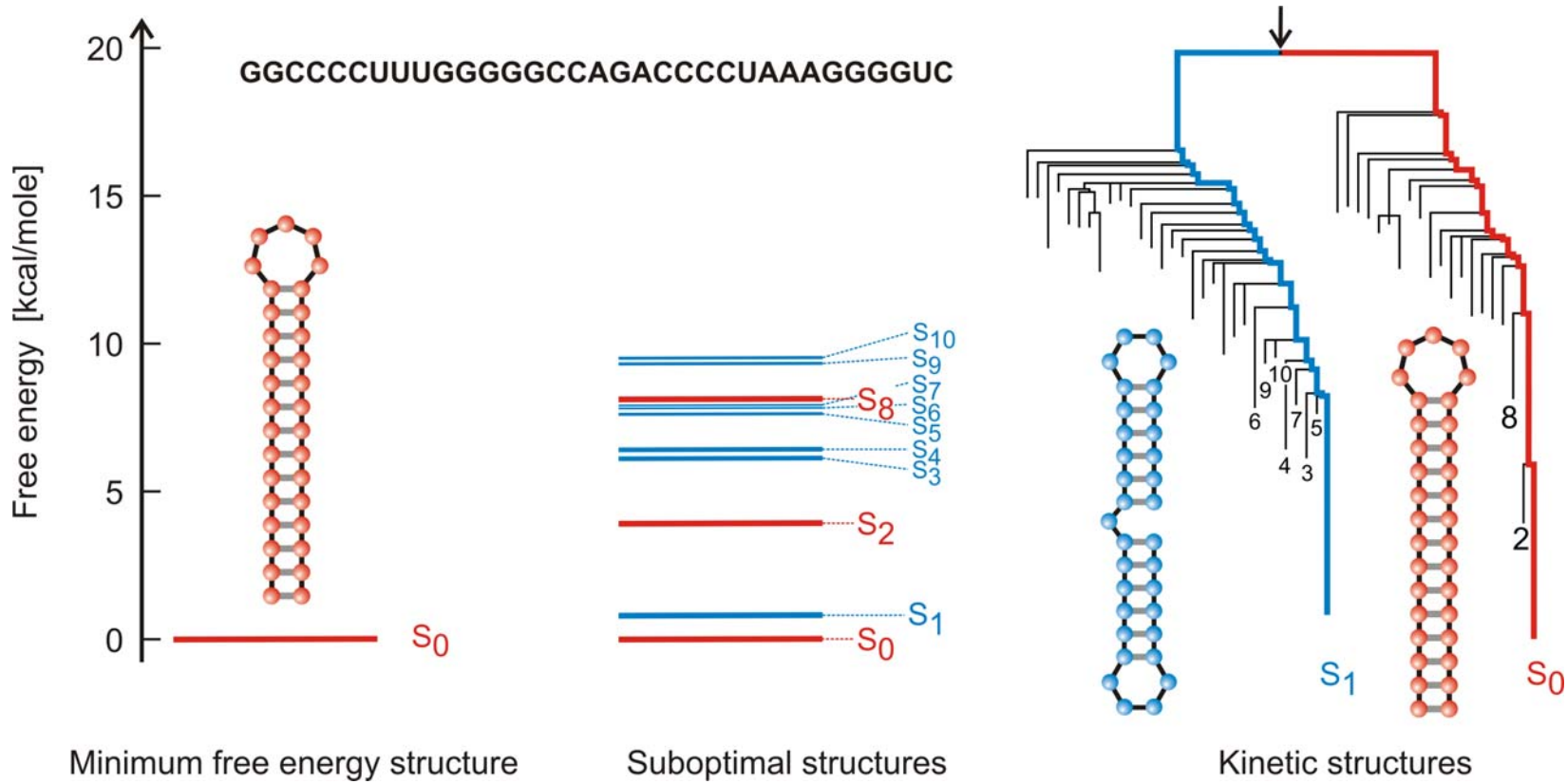
Results from RNA suboptimal structures:

- Neutral RNA sequences differ with respect to their spectra of suboptimal structures.
- Suboptimal RNA structures with low free energies contribute substantially to the partition function.
- Nature selects for stable structures in the sense that the contribution of the mfe structure to the partition function is large.
- For every pair of structures it is possible to find a sequence that can form both. This is not (always) true for three structures.

1. Minimum free energy structures of RNA
2. Suboptimal structures of RNA
- 3. Kinetic folding and RNA switches**
4. Chemistry of Darwinian evolution
5. Consequences of neutrality
6. Evolutionary optimization of RNA structure



Extension of the notion of structure



Extension of the notion of structure

The Folding Algorithm

A sequence \mathbf{I} specifies an energy ordered set of compatible structures $\mathfrak{S}(\mathbf{I})$:

$$\mathfrak{S}(\mathbf{I}) = \{\mathbf{S}_0, \mathbf{S}_1, \dots, \mathbf{S}_m, \mathbf{O}\}$$

A trajectory $\mathfrak{Z}_k(\mathbf{I})$ is a time ordered series of structures in $\mathfrak{S}(\mathbf{I})$. A folding trajectory is defined by starting with the open chain \mathbf{O} and ending with the global minimum free energy structure \mathbf{S}_0 or a metastable structure \mathbf{S}_k which represents a local energy minimum:

$$\mathfrak{Z}_0(\mathbf{I}) = \{\mathbf{O}, \mathbf{S}(1), \dots, \mathbf{S}(t-1), \mathbf{S}(t), \mathbf{S}(t+1), \dots, \mathbf{S}_0\}$$

$$\mathfrak{Z}_k(\mathbf{I}) = \{\mathbf{O}, \mathbf{S}(1), \dots, \mathbf{S}(t-1), \mathbf{S}(t), \mathbf{S}(t+1), \dots, \mathbf{S}_k\}$$

Formulation of kinetic RNA folding as a stochastic process

Master equation

$$\frac{dP_k}{dt} = \sum_{i=0}^{m+1} (P_{ik}(t) - P_{ki}(t)) = \sum_{i=0}^{m+1} k_{ik} P_i - P_k \sum_{i=0}^{m+1} k_{ki}$$

$$k = 0, 1, \dots, m+1$$

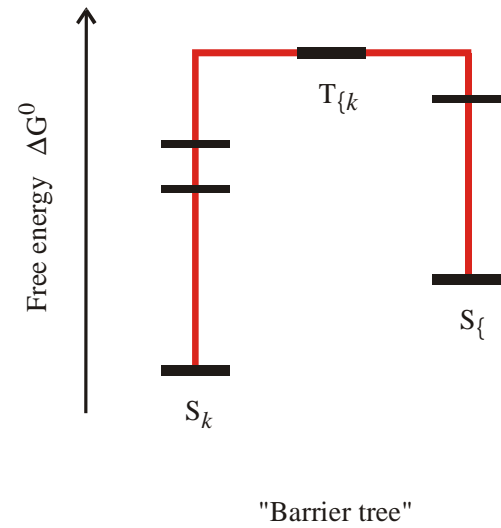
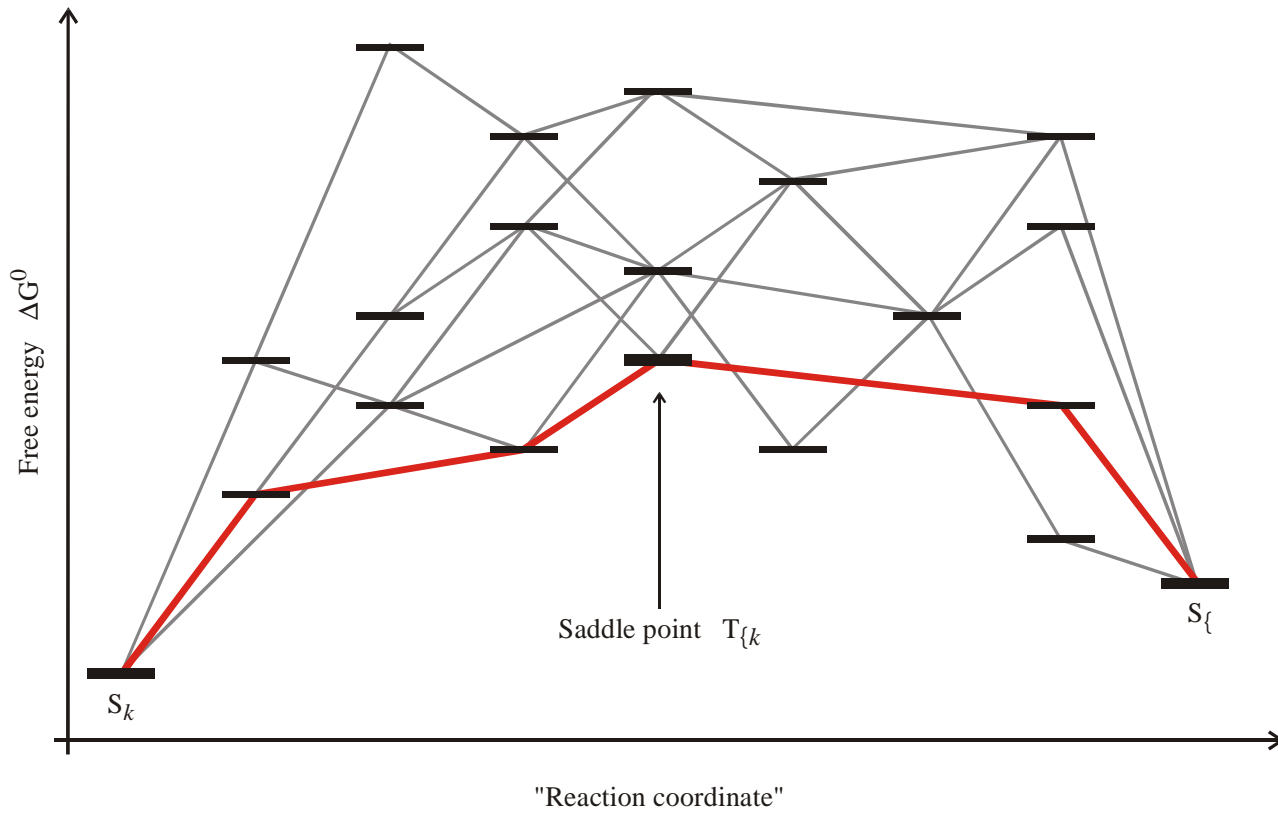
Transition probabilities $P_{ij}(t) = \text{Prob}\{\mathbf{S}_i \rightarrow \mathbf{S}_j\}$ are defined by

$$P_{ij}(t) = P_i(t) k_{ij} = P_i(t) \exp(-\Delta G_{ij}/2RT) / \Sigma_i$$

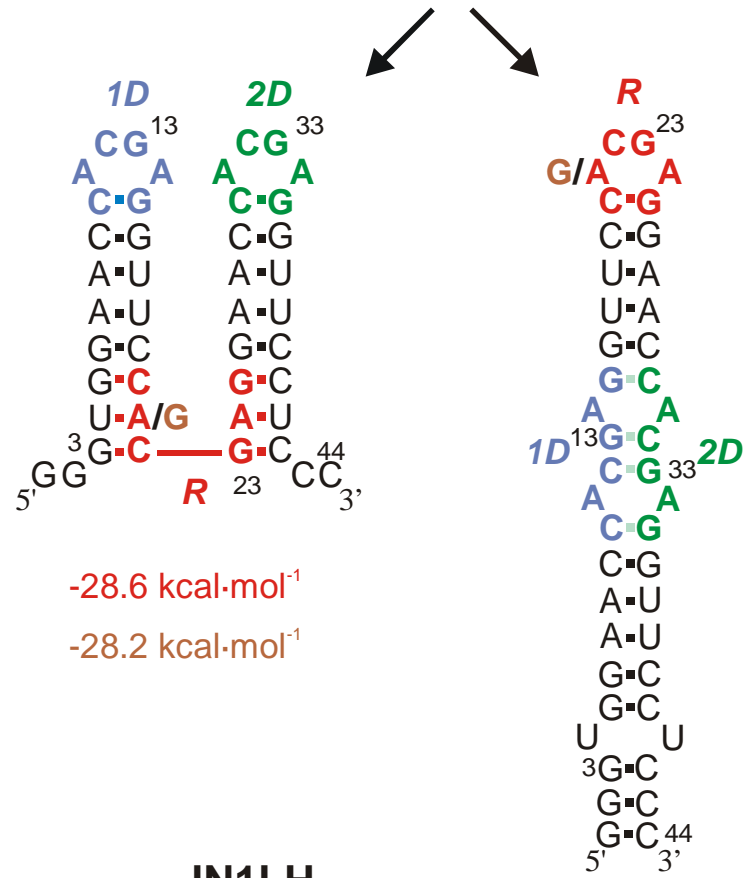
$$P_{ji}(t) = P_j(t) k_{ji} = P_j(t) \exp(-\Delta G_{ji}/2RT) / \Sigma_j$$

$$\Sigma_k = \sum_{k=1, k \neq i}^{m+2} \exp(-\Delta G_{ki}/2RT)$$

The symmetric rule for transition rate parameters is due to Kawasaki (K. Kawasaki, *Diffusion constants near the critical point for time dependent Ising models*. Phys.Rev. **145**:224-230, 1966).



Definition of a ,barrier tree‘



-28.6 kcal·mol⁻¹

-28.2 kcal·mol⁻¹

-28.6 kcal·mol⁻¹

-31.8 kcal·mol⁻¹

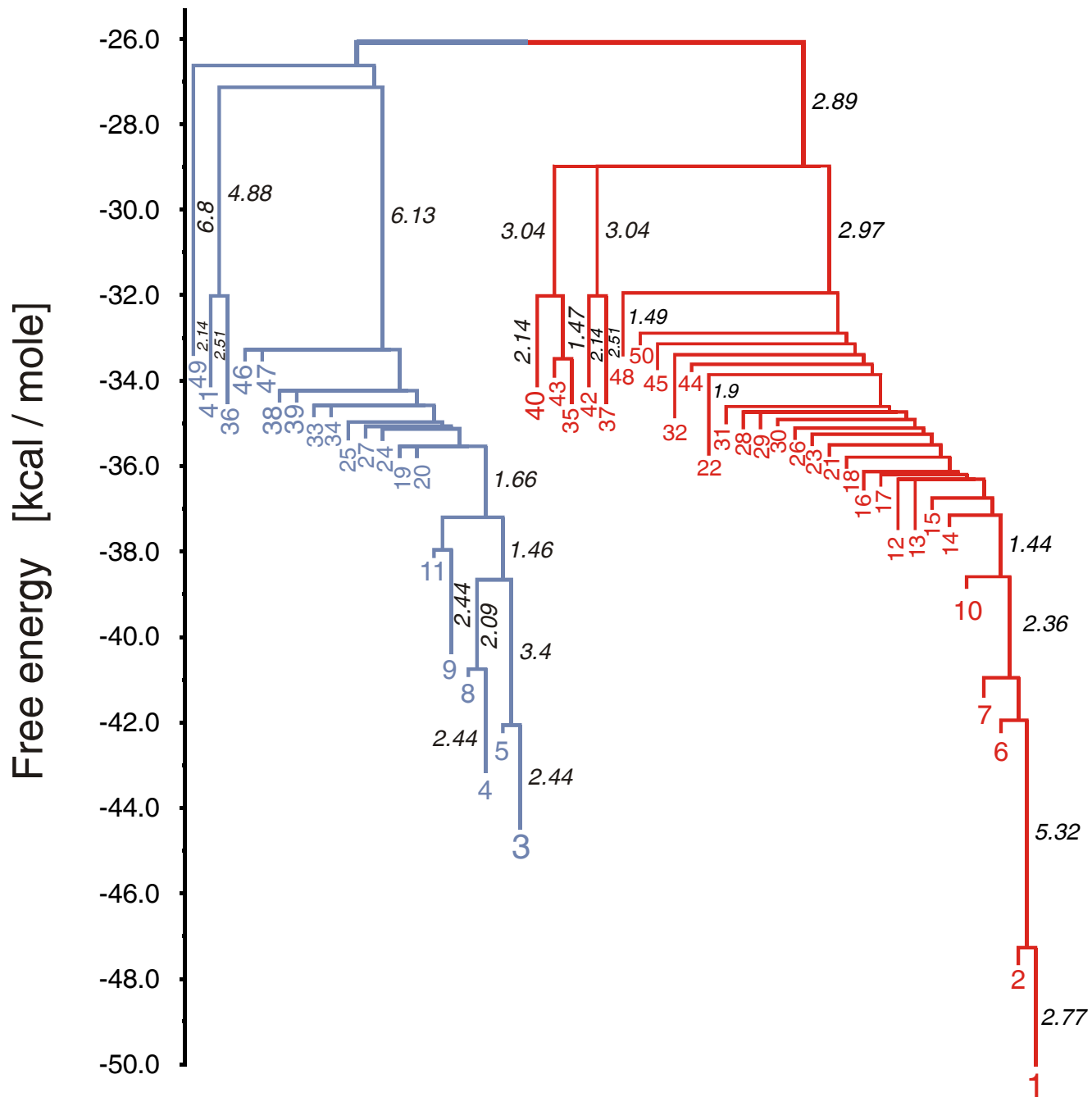
An experimental RNA switch

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)

J1LH barrier tree



- 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 μ M) was incubated with rat liver cytosol (20 mg) or purified recombinant p115 (0.5 μ 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 μ 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/CH₂Cl₂ 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).

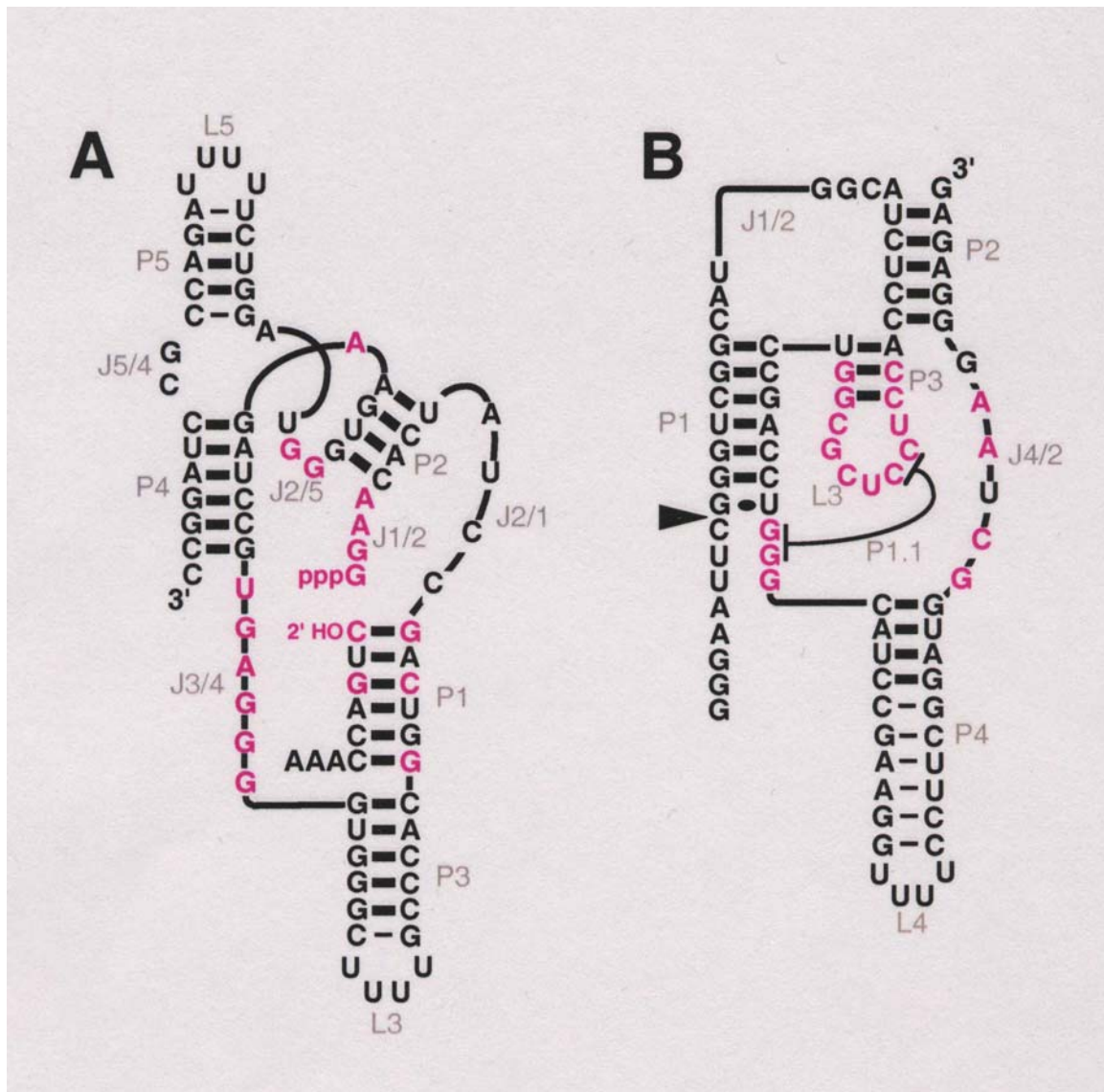
The prototype class III and HDV ribozymes have no more than the 25% sequence identity expected by chance and no fortuitous structural similarities that might favor an intersection of their two neutral networks. Nevertheless, sequences can be designed that simultaneously satisfy the base-pairing requirements

A ribozyme switch

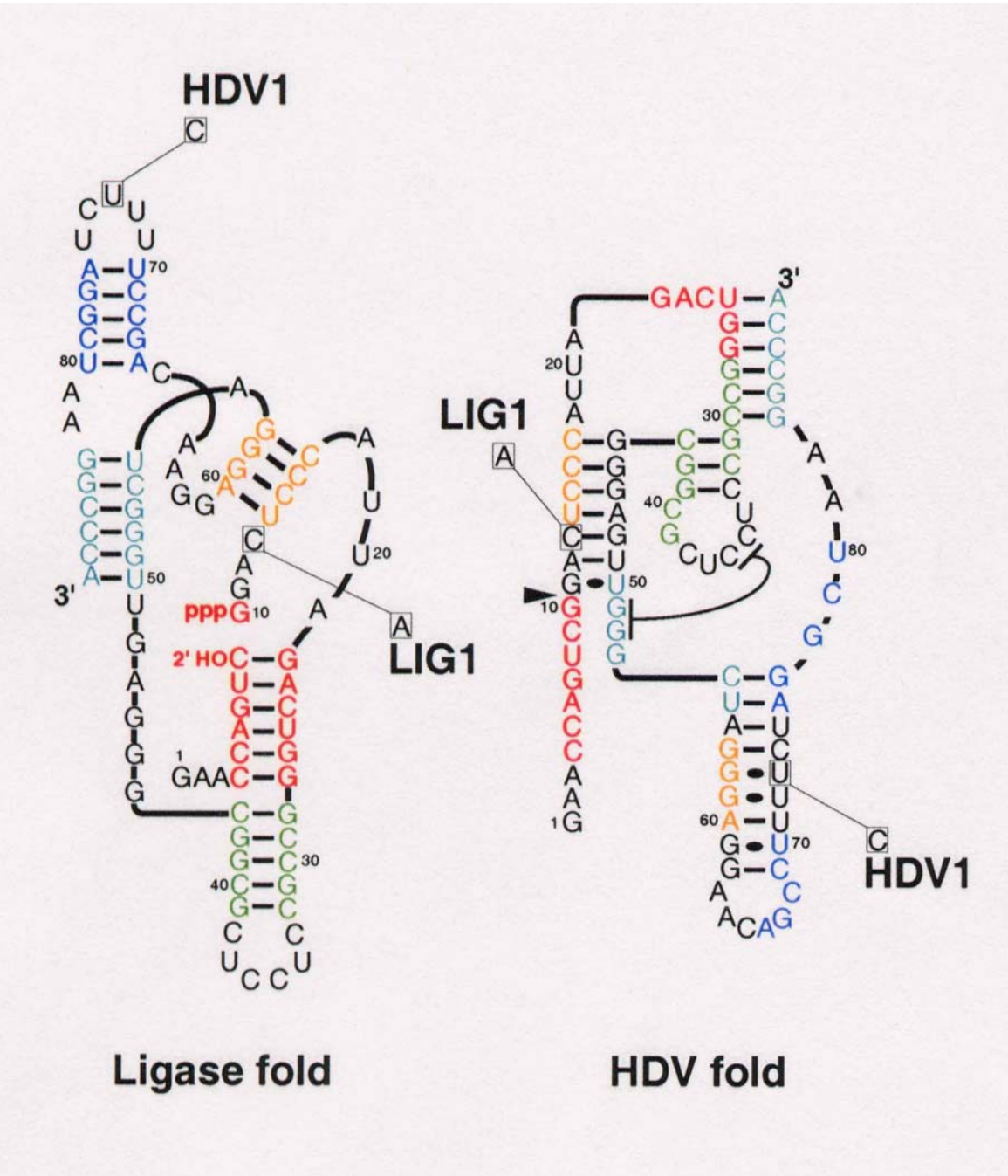
E.A.Schultes, D.B.Bartel, *Science*
289 (2000), 448-452

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

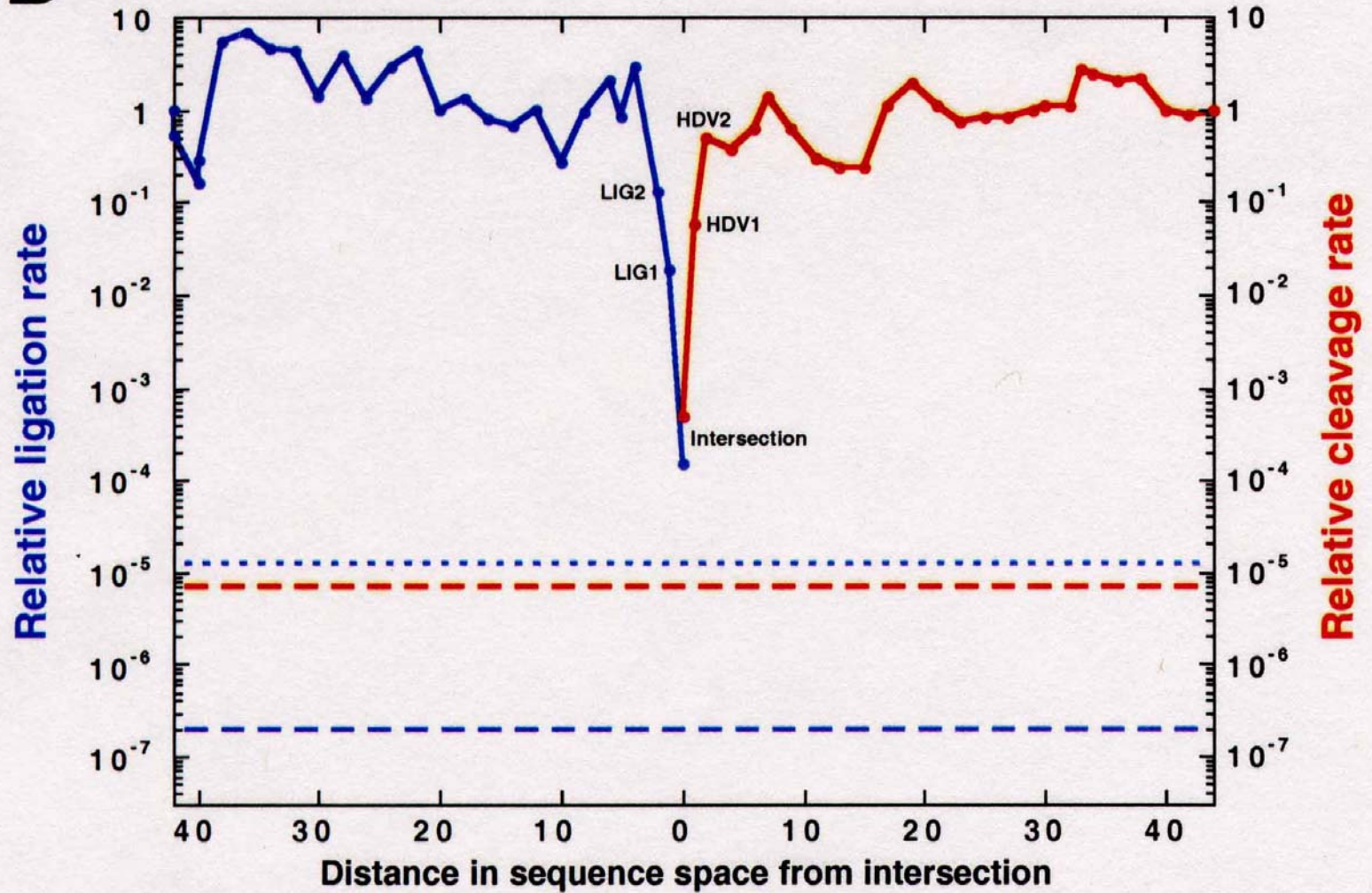


Two ribozymes of chain lengths $n = 88$ nucleotides: An artificial ligase (**A**) and a natural cleavage ribozyme of hepatitis- δ -virus (**B**)

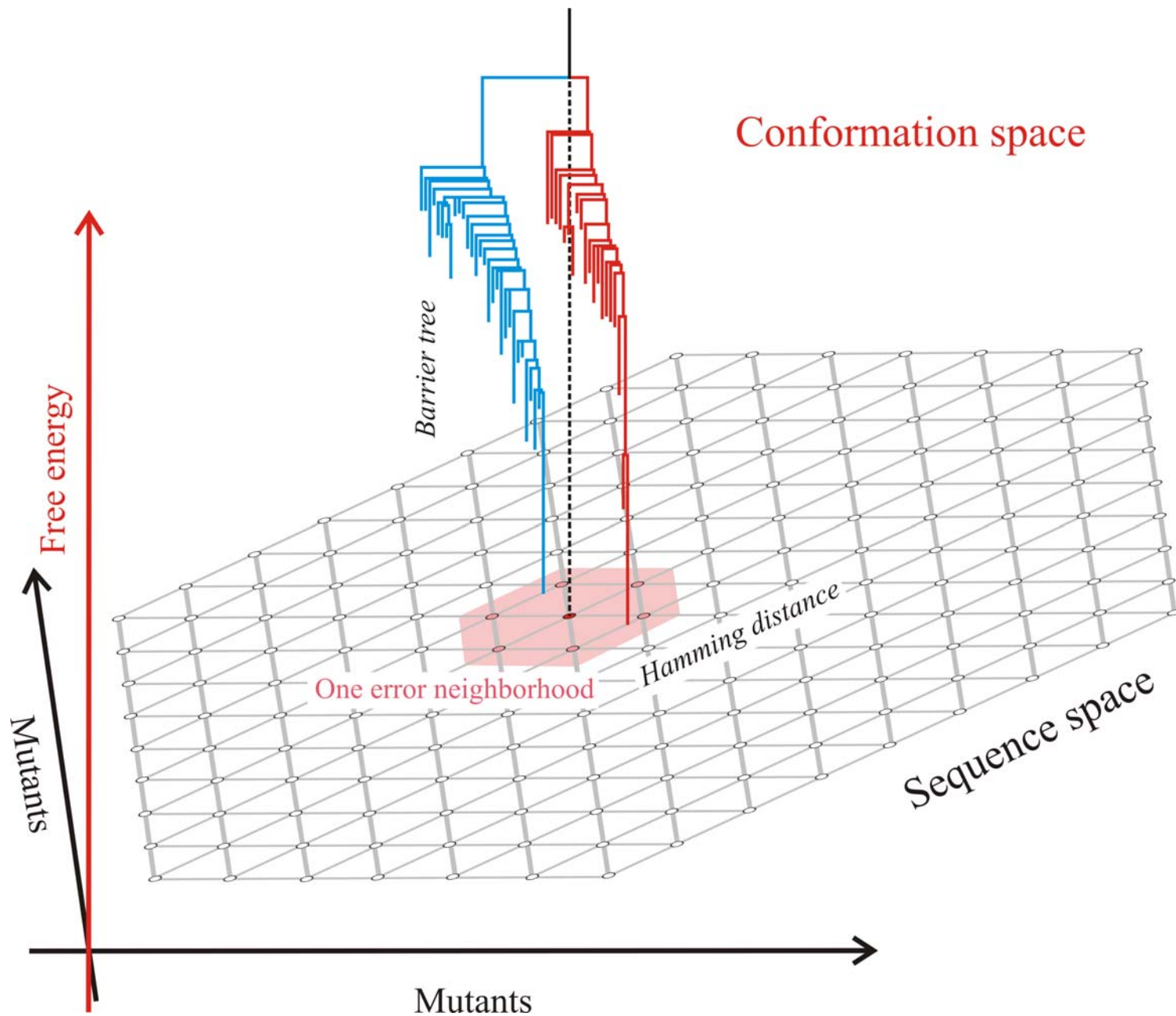


The sequence at the *intersection*:

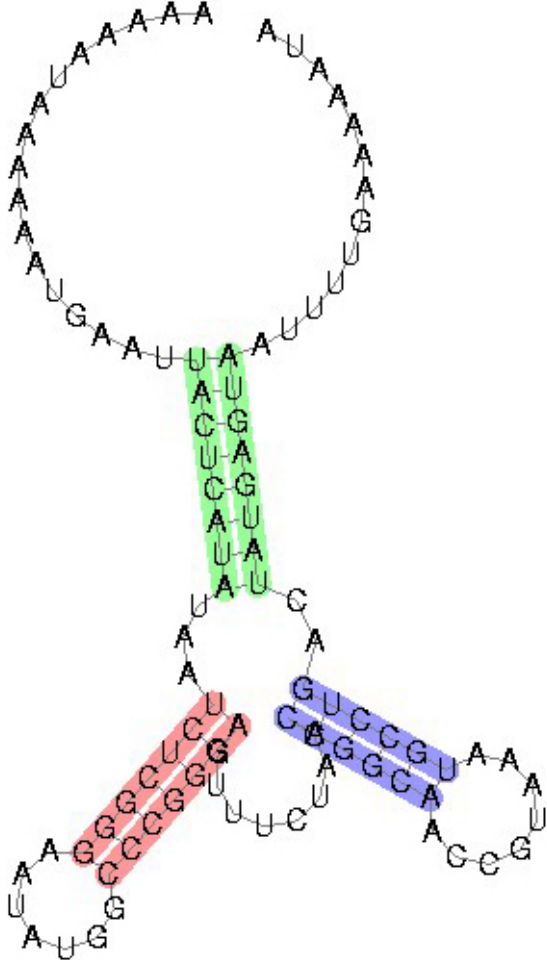
An RNA molecules which is 88 nucleotides long and can form both structures

B

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



A natural metabolic riboswitch



The purine riboswitch

M. Mandal, B. Boese, J.E. Barrick, W.C. Winkler, and R.R. Breaker. 2003. *Molecular Cell*. 11:1419-1420, *Cell* 113:577-586.

AAAAAUAAAAAUGAAUUACUCAUAUAUAUCUCGGGAAUAUGGCCCGGGAGUUUCUAGCAGGCAACCGUAAAUGCCUGACUAUGAGUAAUUUUGAAAAAUA

.....((((((((((...((((((.....))))))))......((((((.....)))))))).)..... -32.10

.....((((((((((...((((((.....))))))))......((((((.....)))))))).)..... -31.80

.....((((((((((...((((((.....))))))))......((((((.....)))))))).)..... -31.80

.....((((((((((...((((((.....))))))))......((((((.....)))))))).)..... -31.80

.....((((((((((...((((((.....))))))))......((((((.....)))))))).)..... -31.00

.....((((((((((...((((((.....))))))))......((((((.....)))))))).)..... -31.00

.....((((((((((...((((((.....))))))))......((((((.....)))))))).)..... -31.00

.....((((((((((...((((((.....))))))))......((((((.....)))))))).)..... -30.70

.....((((((((((...((((((.....))))))))......((((((.....)))))))).)..... -28.60

.....((((((((((...((((((.....))))))))......).....).....).....).....)..... -24.80

.....((((((((((...((((((.....))))))))......).....).....).....).....)..... -24.60

.....((((((((((...((((((.....))))))))......).....).....).....).....)..... -24.60

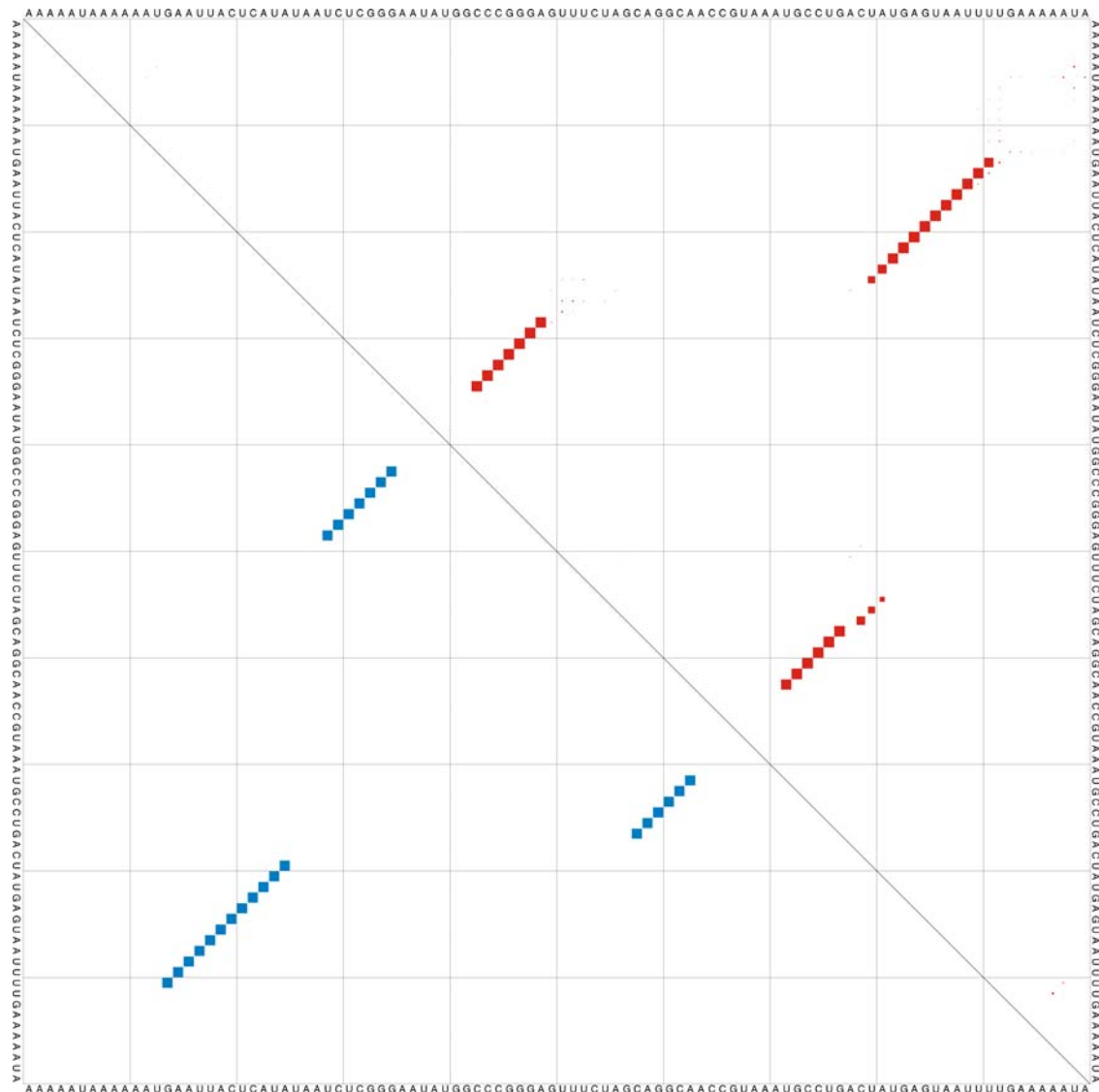
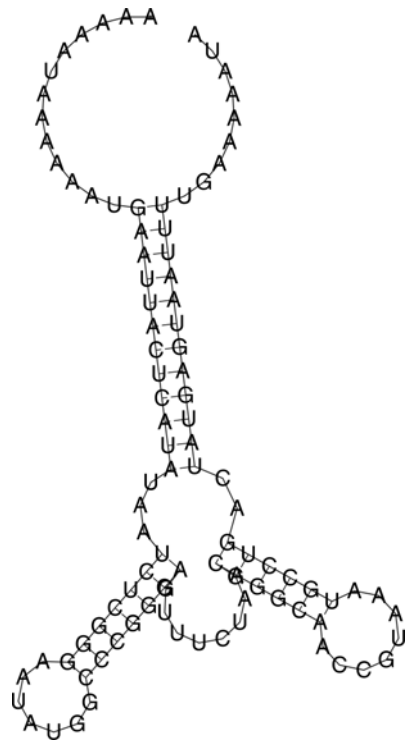
.....((((((((((...((((((.....))))))))......).....).....).....).....)..... -24.60

.....((((((((((...((((((.....))))))))......).....).....).....).....)..... -24.50

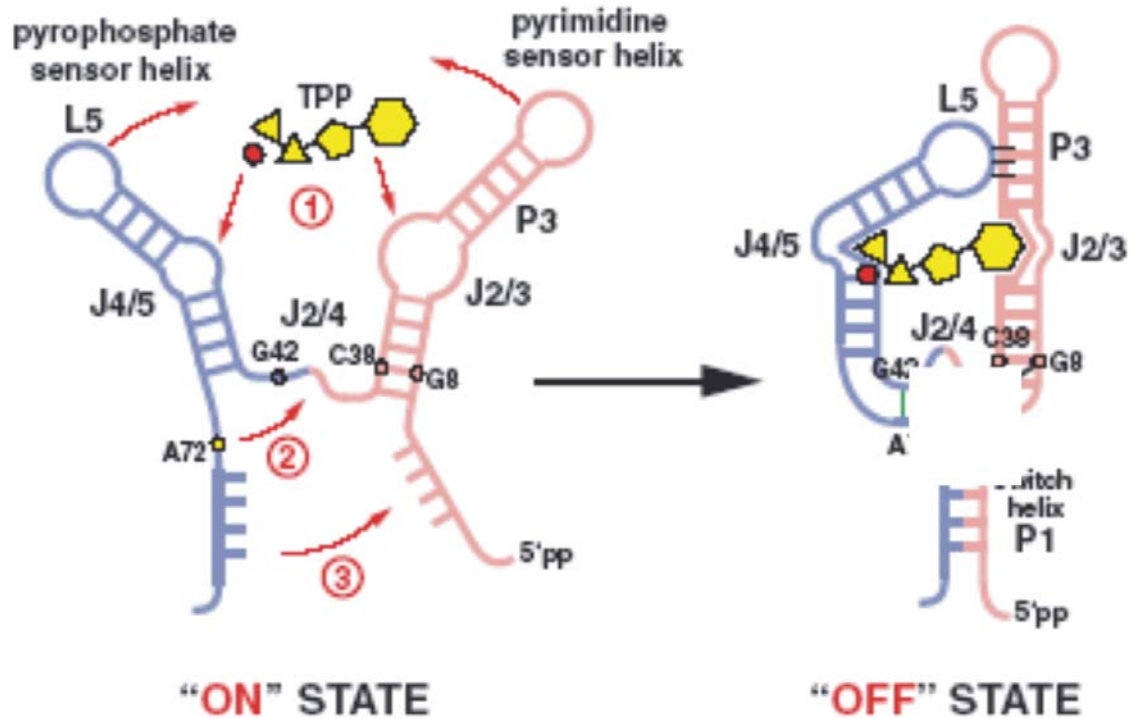
.....((((((((((...((((((.....))))))))......).....).....).....).....)..... -24.50

.....((((((((((...((((((.....))))))))......).....).....).....).....)..... -24.50

The purine riboswitch: *Molecular Cell*. 2003. 11:1419-1420.

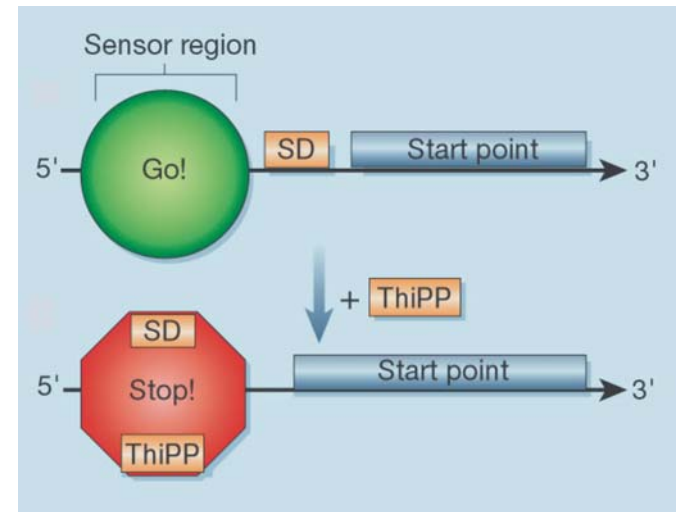


mfe-weight: 0.1459



The thiamine-pyrophosphate riboswitch

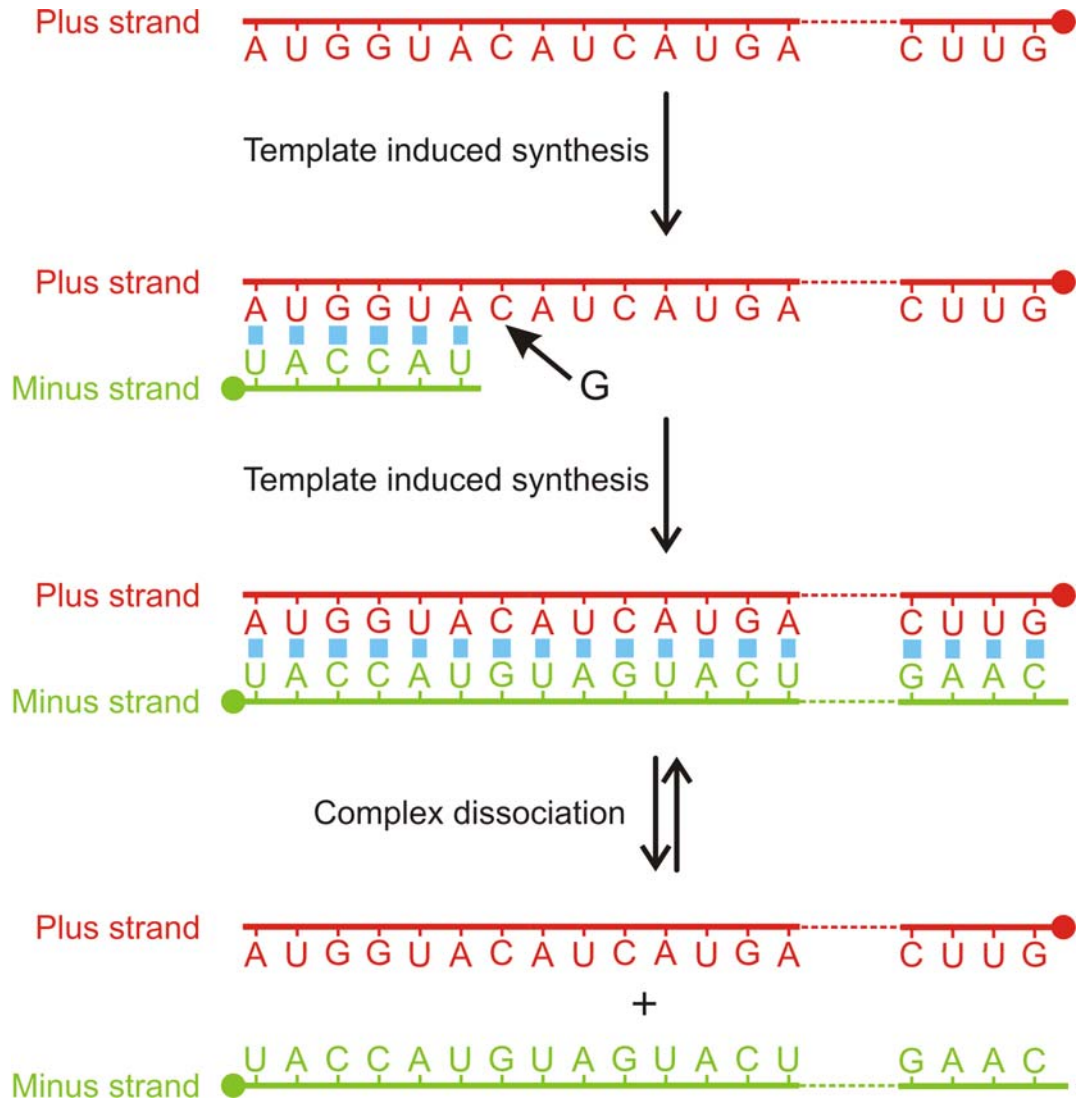
S. Thore, M. Leibundgut, N. Ban.
Science **312**:1208-1211, 2006.



Results from RNA folding kinetics:

- In addition to the minimum free energy structure RNA molecules can exist in one, two or more long-lived metastable structures.
- RNA switches are molecules with two or more long-lived conformations that allow for metabolic control.

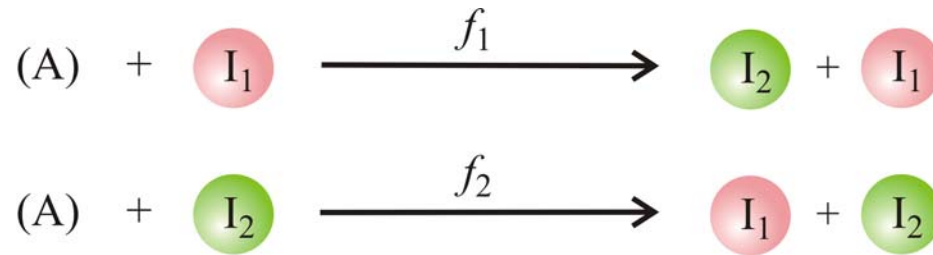
1. Minimum free energy structures of RNA
2. Suboptimal structures of RNA
3. Kinetic folding and RNA switches
- 4. Chemistry of Darwinian evolution**
5. Consequences of neutrality
6. Evolutionary optimization of RNA structure



Complementary replication is the simplest copying mechanism of RNA.

Complementarity is determined by Watson-Crick base pairs:

G≡C and **A=U**



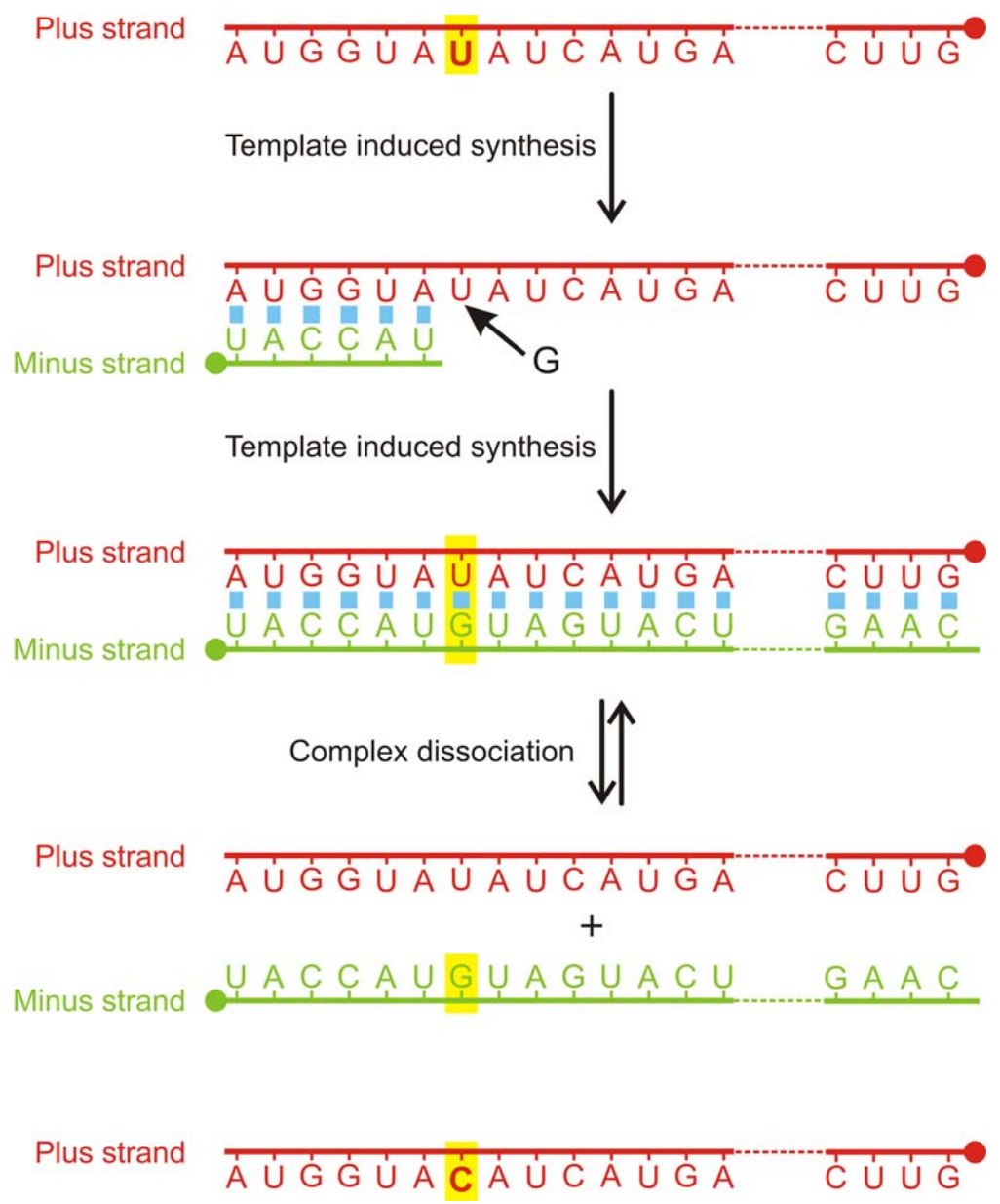
$$\frac{dx_1}{dt} = f_2 x_2 \quad \text{and} \quad \frac{dx_2}{dt} = f_1 x_1$$

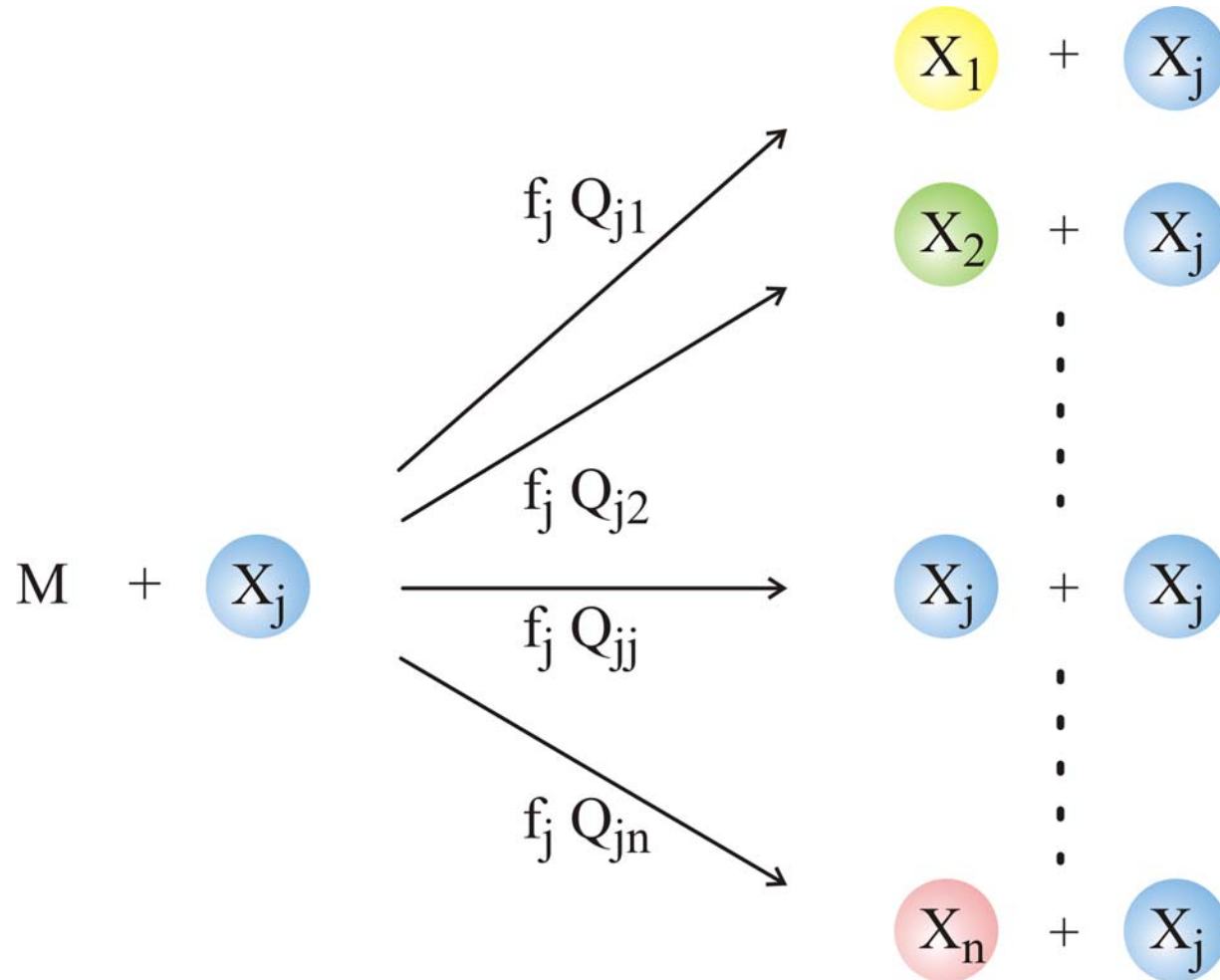
$$x_1 = \sqrt{f_2} \xi_1, \quad x_2 = \sqrt{f_1} \xi_2, \quad \zeta = \xi_1 + \xi_2, \quad \eta = \xi_1 - \xi_2, \quad f = \sqrt{f_1 f_2}$$

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

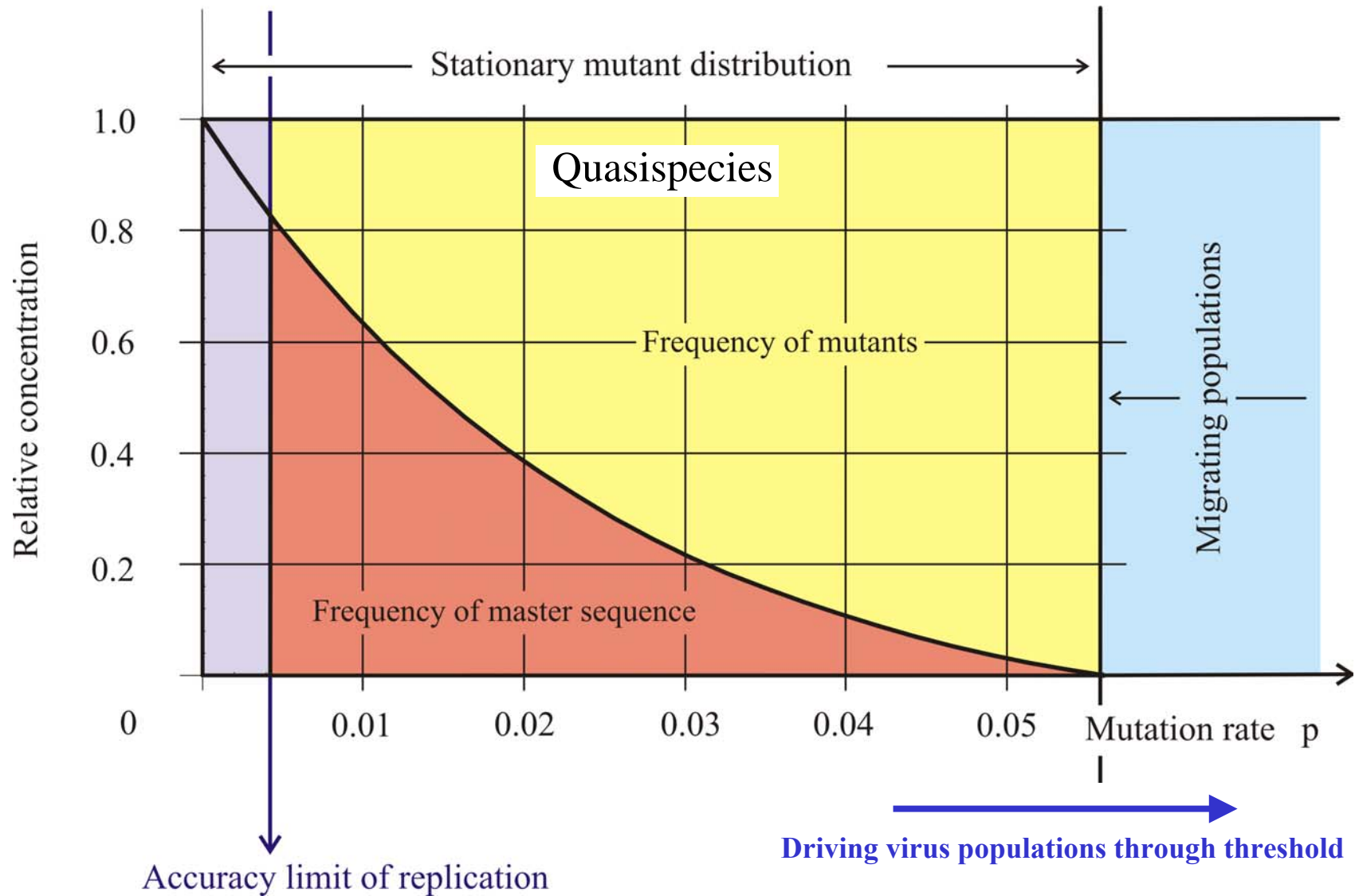
$$\zeta(t) = \zeta(0) e^{ft}$$

Complementary replication as the simplest molecular mechanism of reproduction

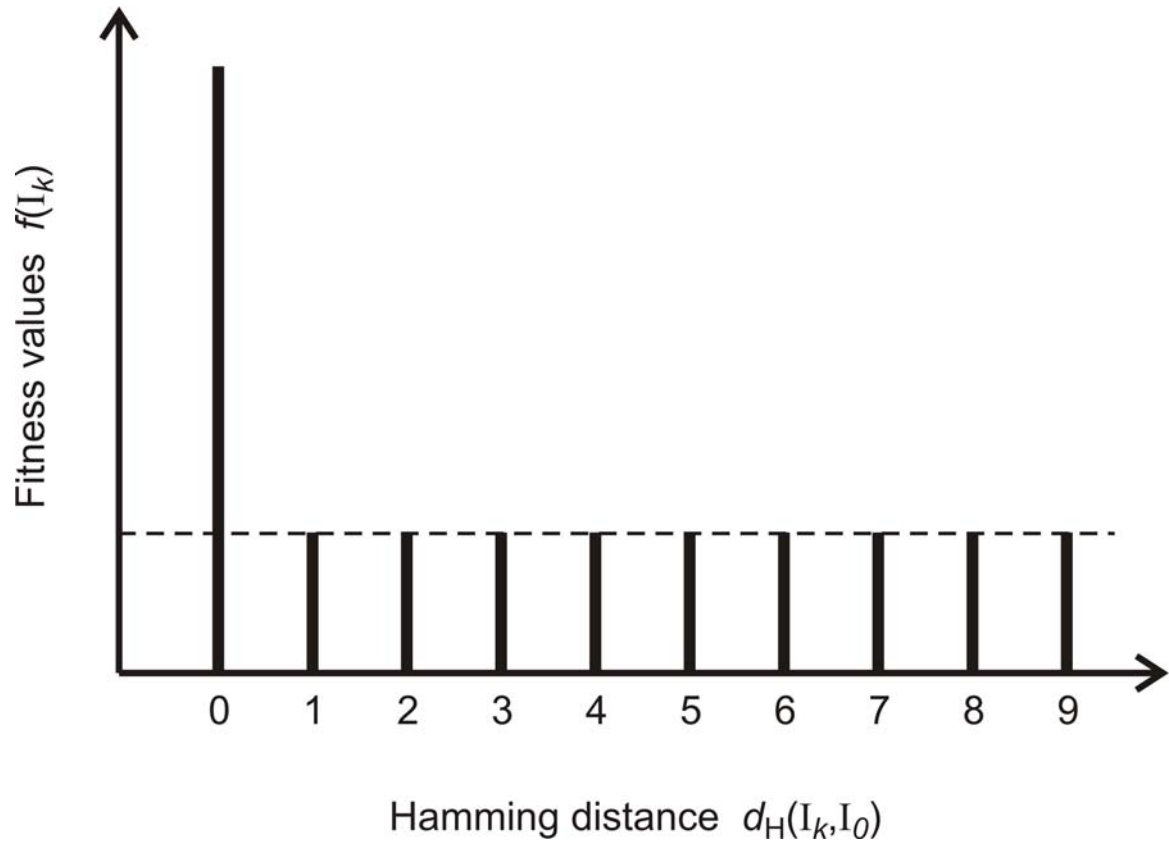




Chemical kinetics of replication and mutation as parallel reactions



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:

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

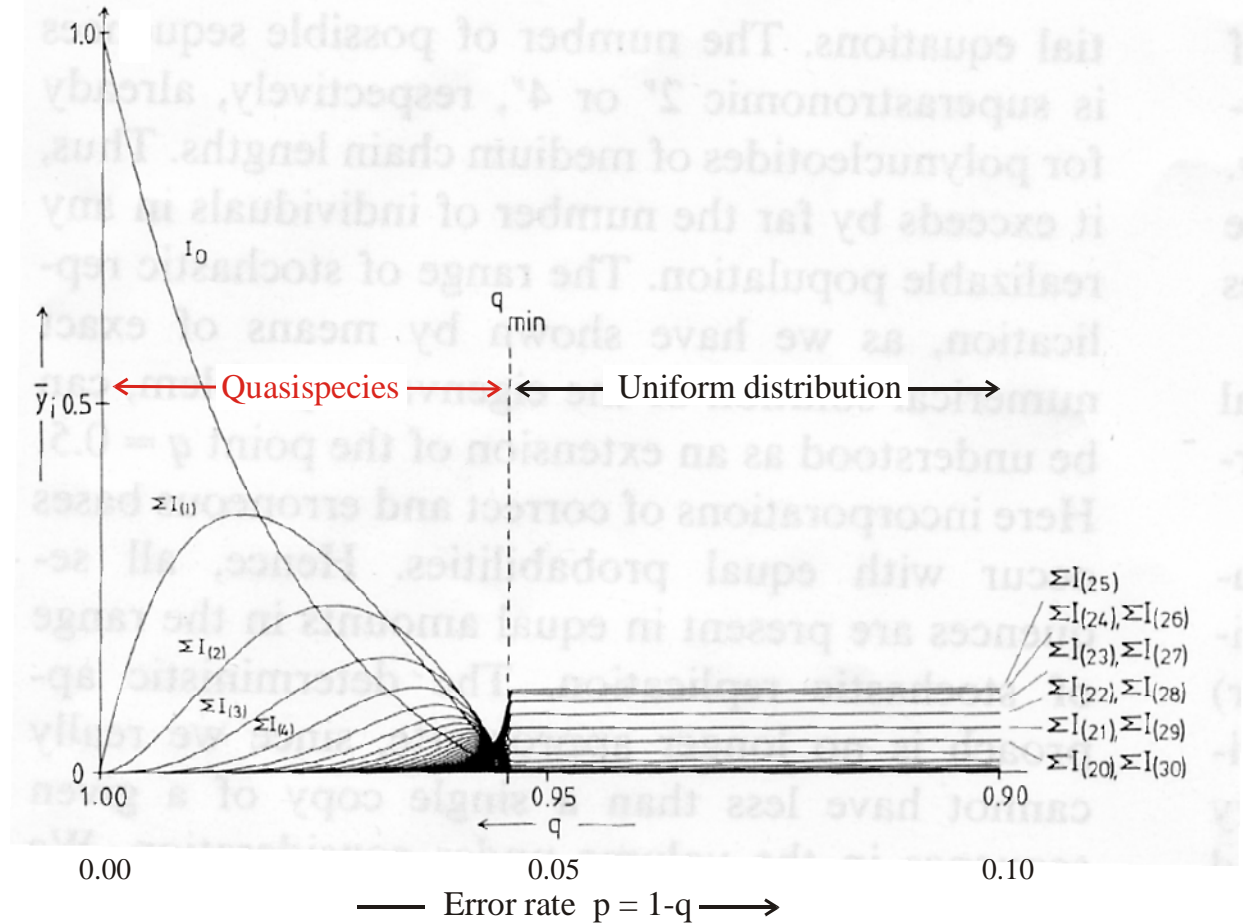
By x_i we denote the population number or concentration of the self-replicating element I_i , i.e., $x_i = [I_i]$. The total population size or total concentration $c = \sum_i x_i$ is kept constant by proper adjustment of the constraint $\phi = \sum_i \sum_j w_{ij} x_j x_i$. Characteristically, this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

(w_{ij}) 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 ($n = 1$).

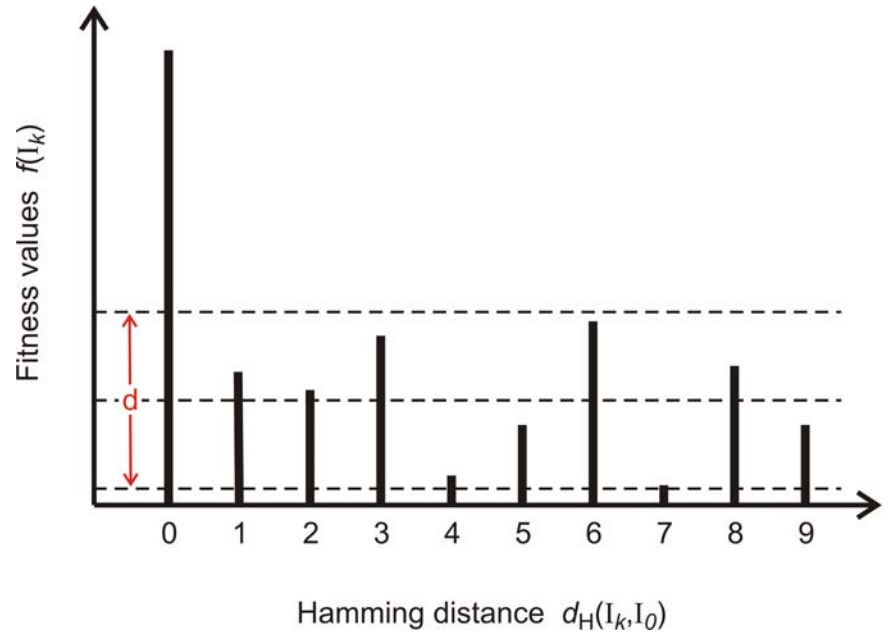
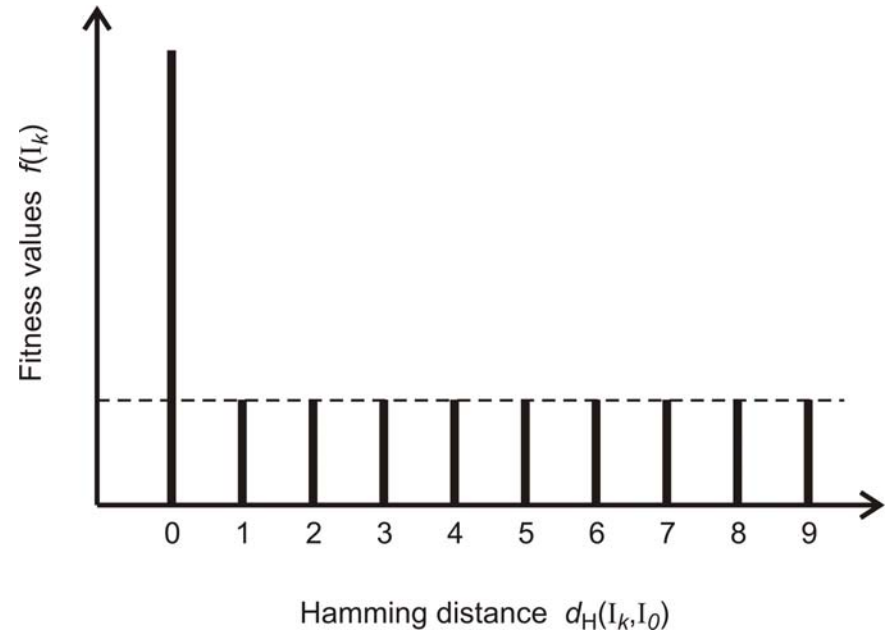
The non-linear differential equation, eq. 1 - the non-linearity is introduced by the definition of ϕ at constant organization - shows a remarkable feature: it leads to selection of a defined ensemble of self-replicating elements above a certain accuracy threshold. This ensemble of a master and its most 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

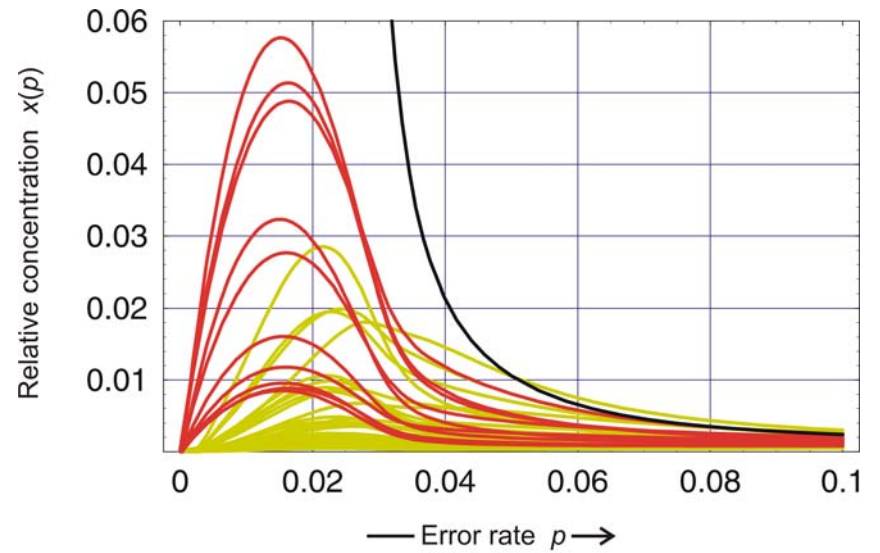
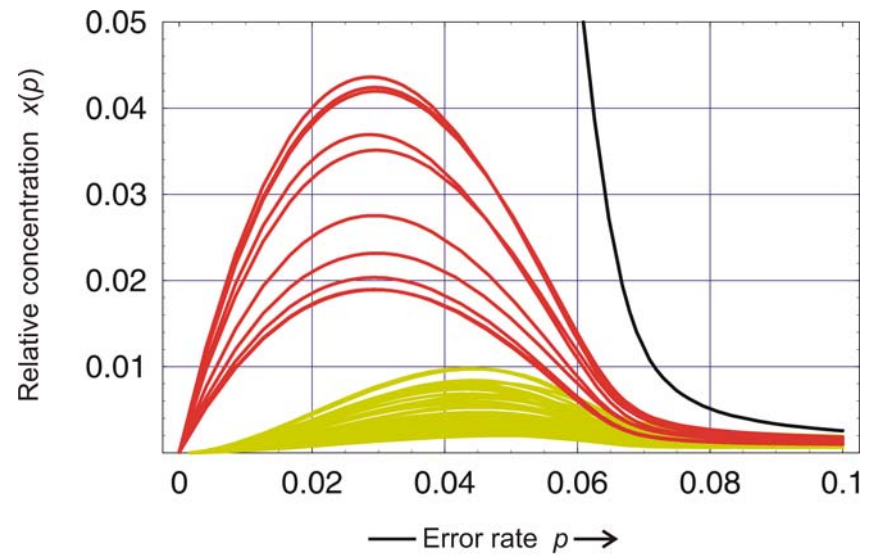
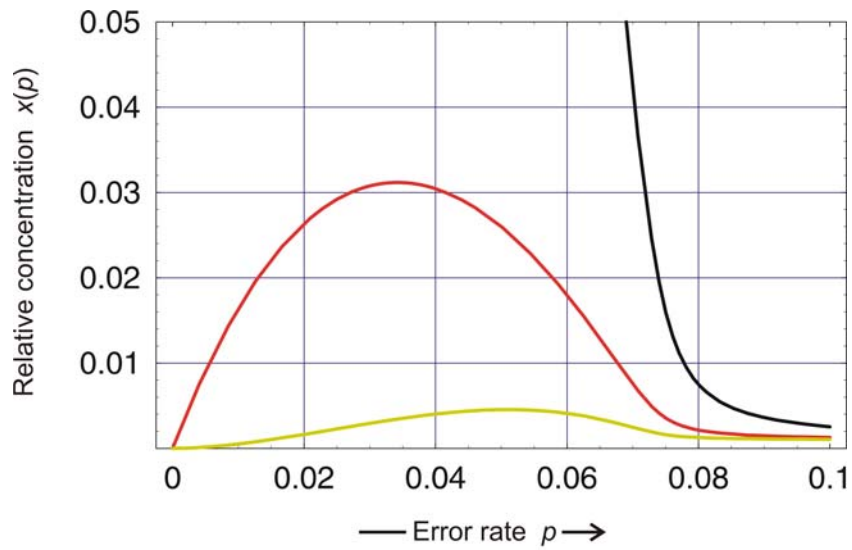
* Dedicated to the late Professor B.L. Jones who was among the first to do rigorous mathematical analysis on the problems described here.
 ** This paper is considered as part II of Model Studies on RNA replication. Part I is by Gassner and Schuster [14].
 † All summations throughout this paper run from 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



Fitness landscapes showing error thresholds



Error threshold: Individual sequences

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

Evolution of RNA molecules based on Q β 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

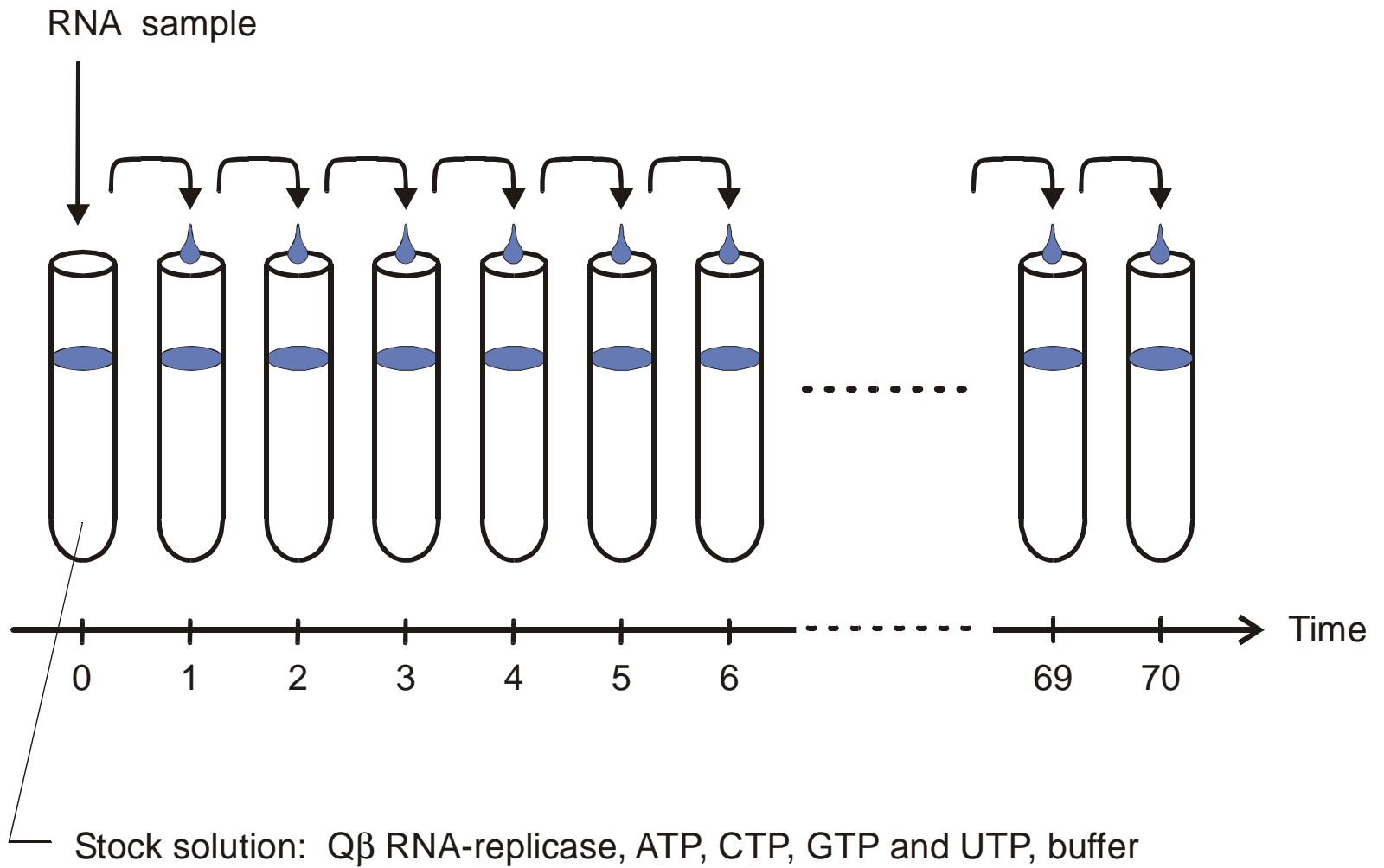
C.K.Biebricher, *Darwinian selection of self-replicating RNA molecules*. Evolutionary Biology **16** (1983), 1-52

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

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

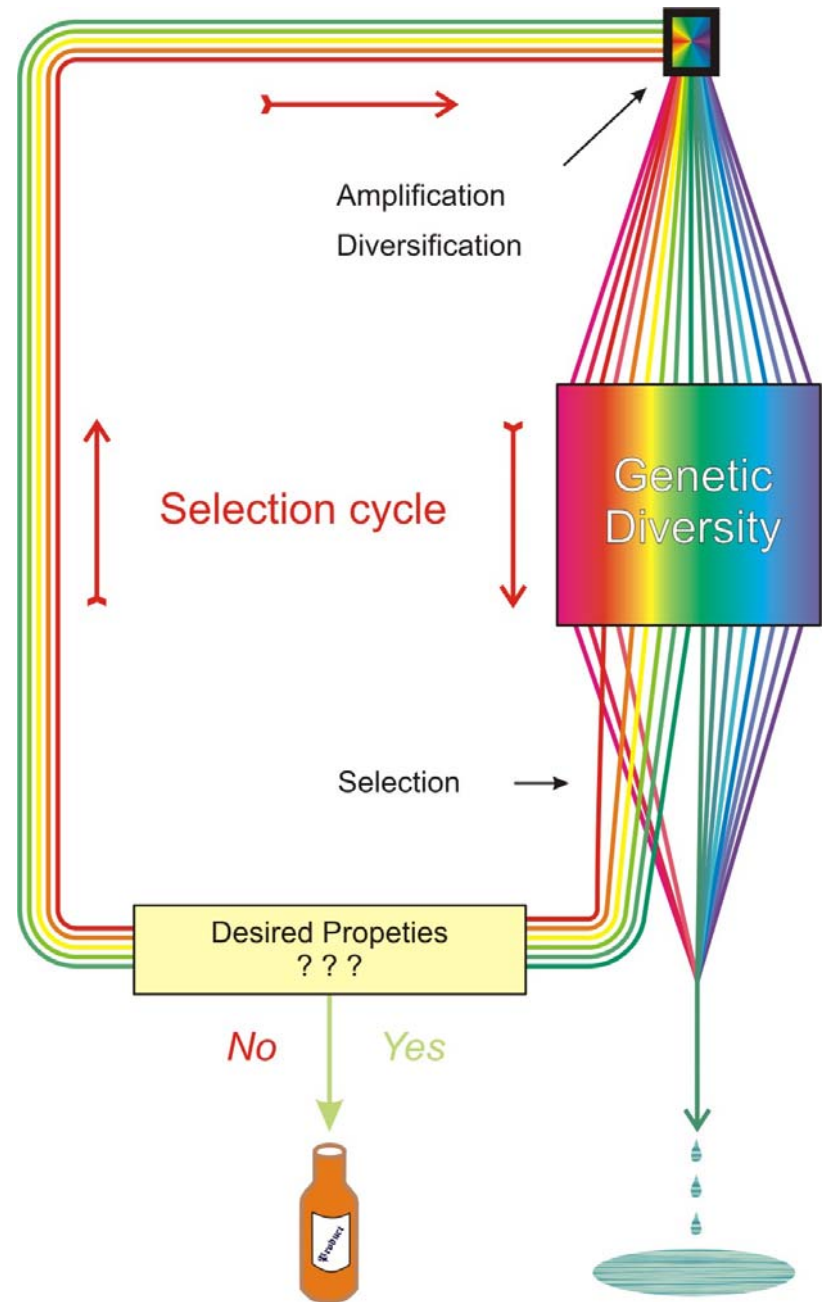
C. Tuerk, L. Gold, *SELEX - Systematic evolution of ligands by exponential enrichment: RNA ligands to bacteriophage T4 DNA polymerase.* Science **249** (1990), 505-510

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

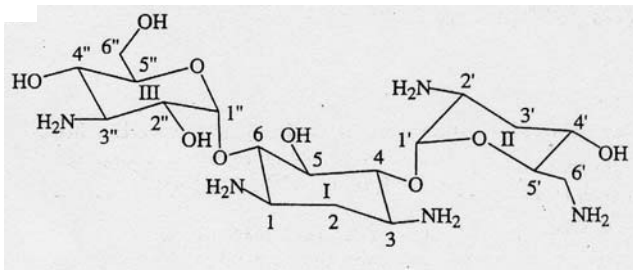
R.D. Jenison, S.C. Gill, A. Pardi, B. Poliski, *High-resolution molecular discrimination by RNA.* Science **263** (1994), 1425-1429

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

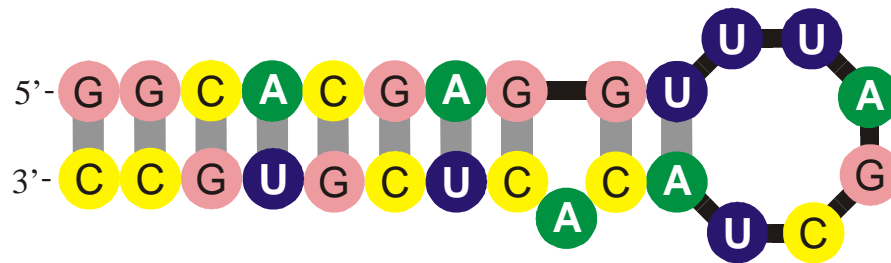
L. Jiang, A. K. Suri, R. Fiala, D. J. Patel, *Saccharide-RNA recognition in an aminoglycoside antibiotic-RNA aptamer complex.* Chemistry & Biology **4** (1997), 35-50



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



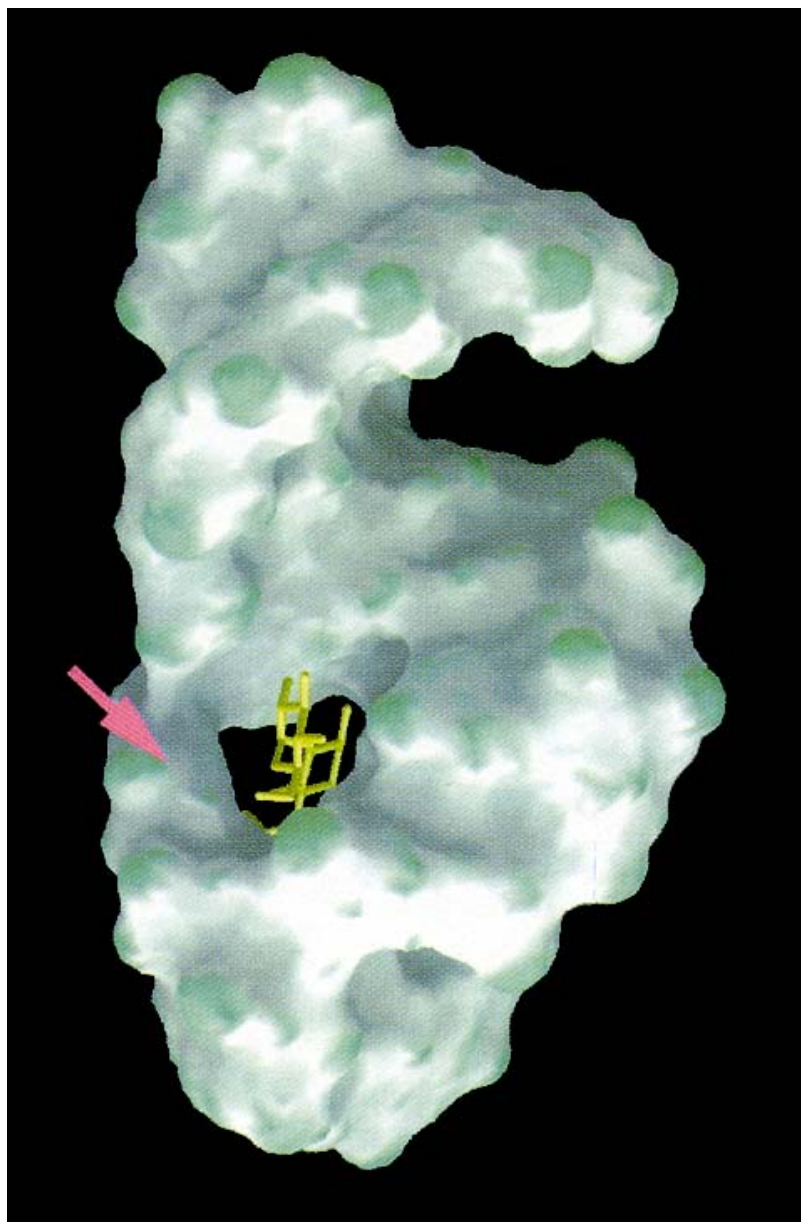
tobramycin



RNA aptamer

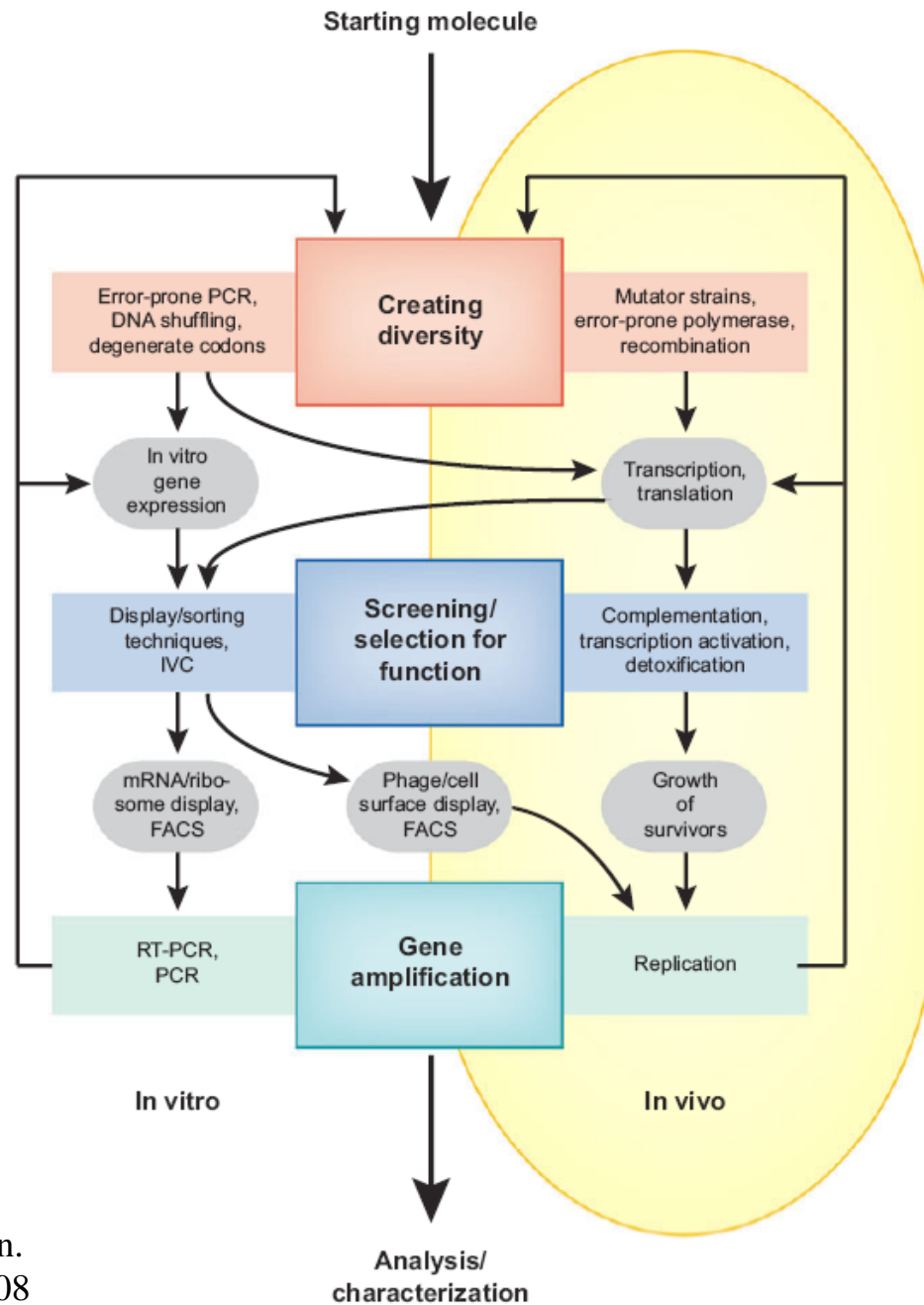
Formation of secondary structure of the tobramycin binding RNA aptamer with $K_D = 9 \text{ nM}$

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)



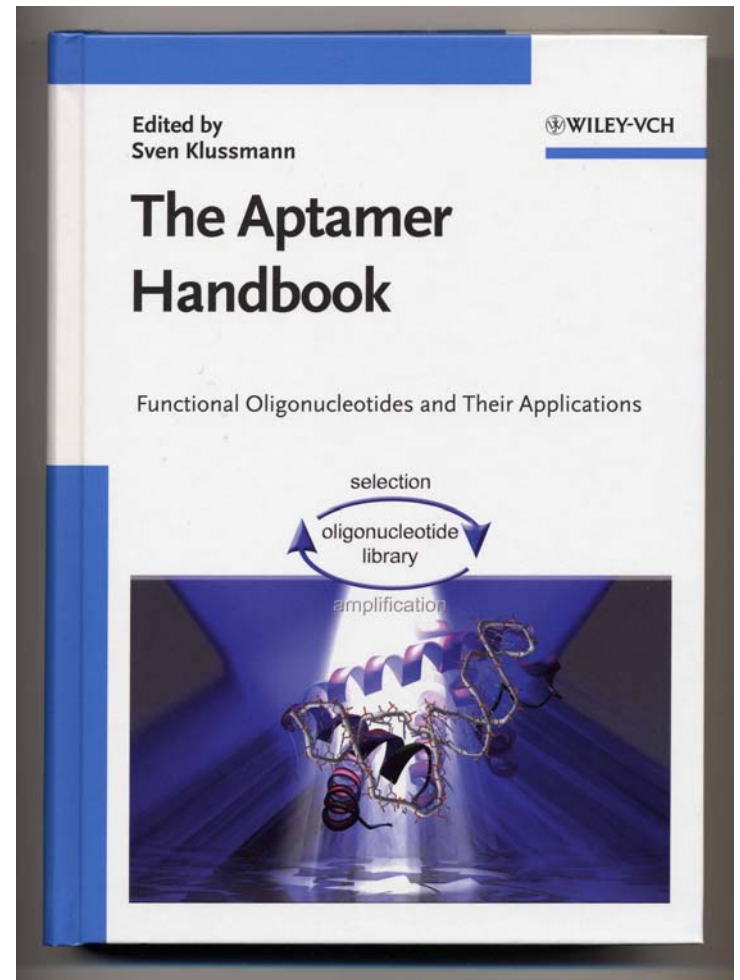
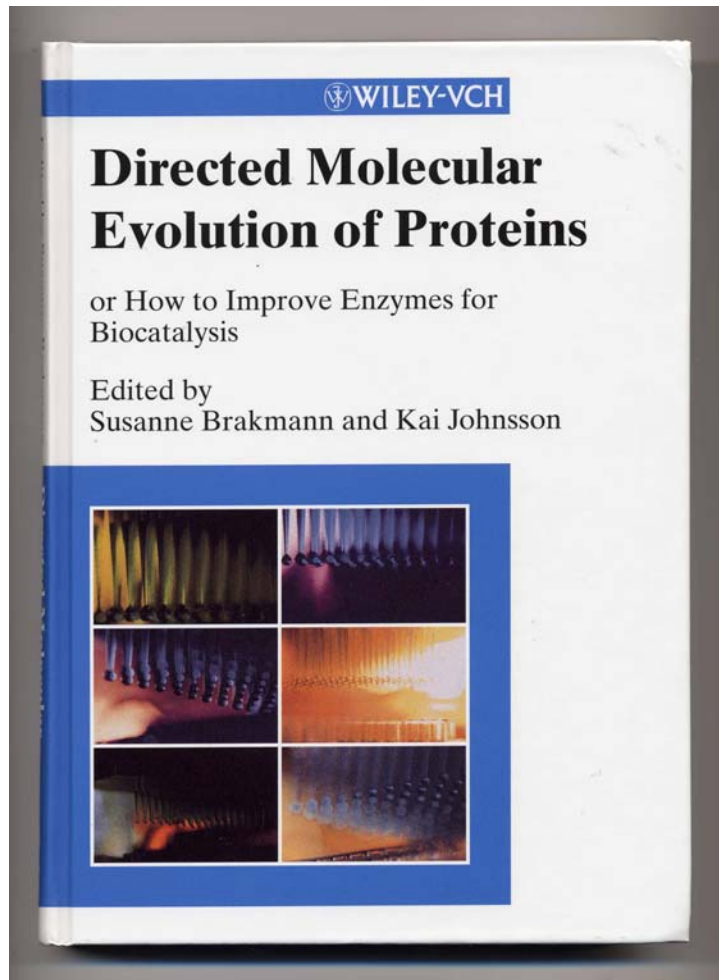
The three-dimensional structure of the
tobramycin aptamer complex

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



Schematic overview of the principal processes, strategies, and techniques of directed evolution. Today, numerous experimental methods are available to perform the fundamental processes of true Darwinian evolution (*central boxes*) in the laboratory, either in vivo within microorganisms or entirely in vitro in the test tube. Arrows indicate possible routes for connecting individual evolutionary steps. Abbreviations: PCR, polymerase chain reaction; RT-PCR, reverse transcription PCR; IVC, in vitro compartmentalization; FACS, fluorescence-activated cell sorting.

Christian Jäckel, Peter Kast, and Donald Hilvert.
 Protein design by directed evolution.
Annu.Rev.Biophys. **37**:153-173, 2008



Application of molecular evolution to problems in biotechnology

Artificial evolution in biotechnology and pharmacology

G.F. Joyce. 2004. Directed evolution of nucleic acid enzymes. *Annu.Rev.Biochem.* **73**:791-836.

C. Jäckel, P. Kast, and D. Hilvert. 2008. Protein design by directed evolution. *Annu.Rev.Biophys.* **37**:153-173.

S.J. Wrenn and P.B. Harbury. 2007. Chemical evolution as a tool for molecular discovery. *Annu.Rev.Biochem.* **76**:331-349.

Results from replication kinetics and 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.
- Novel antiviral strategies were developed from known molecular mechanisms of virus evolution.

1. Minimum free energy structures of RNA
2. Suboptimal structures of RNA
3. Kinetic folding and RNA switches
4. Chemistry of Darwinian evolution
- 5. Consequences of neutrality**
6. Evolutionary optimization of RNA structure

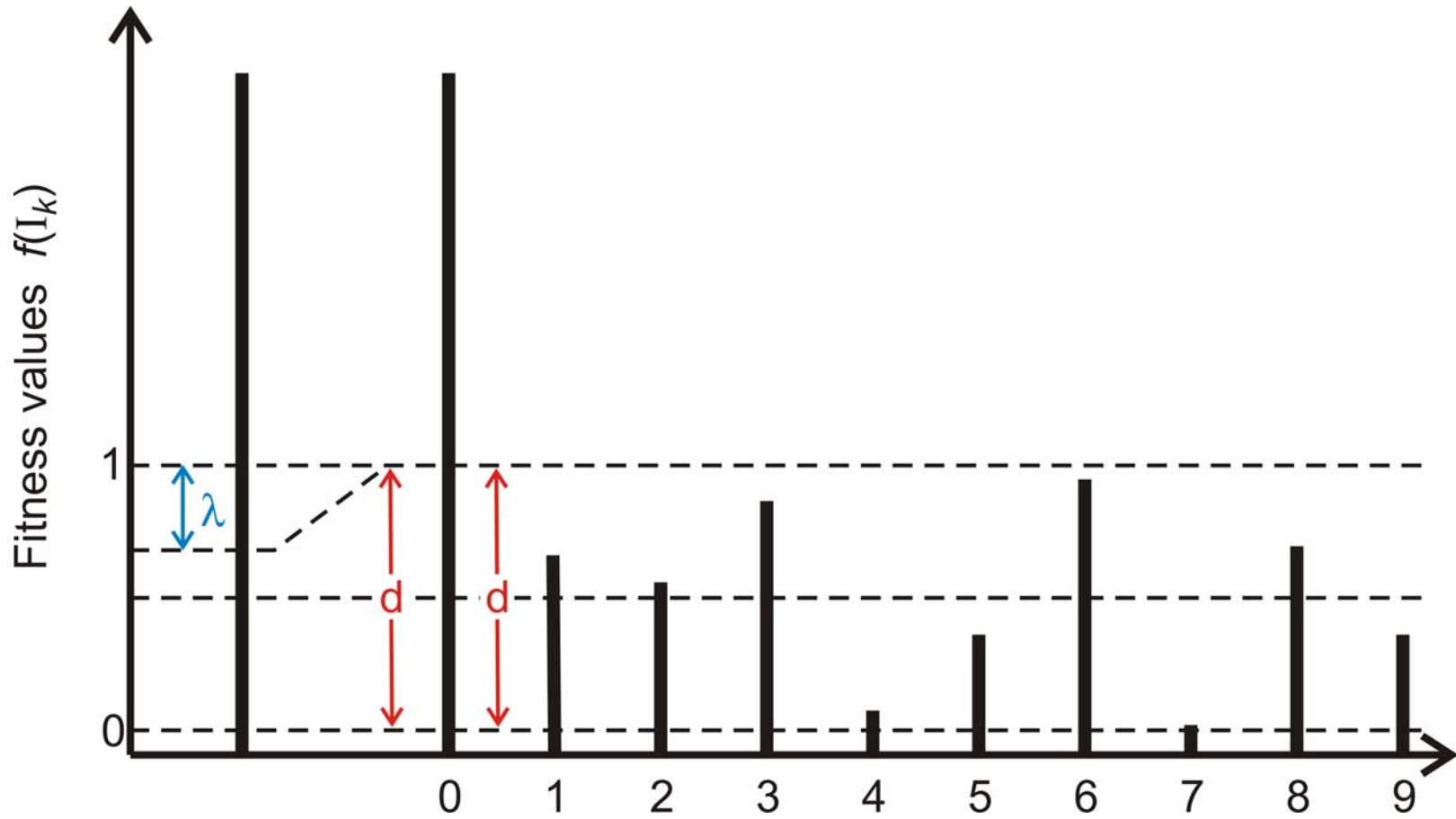
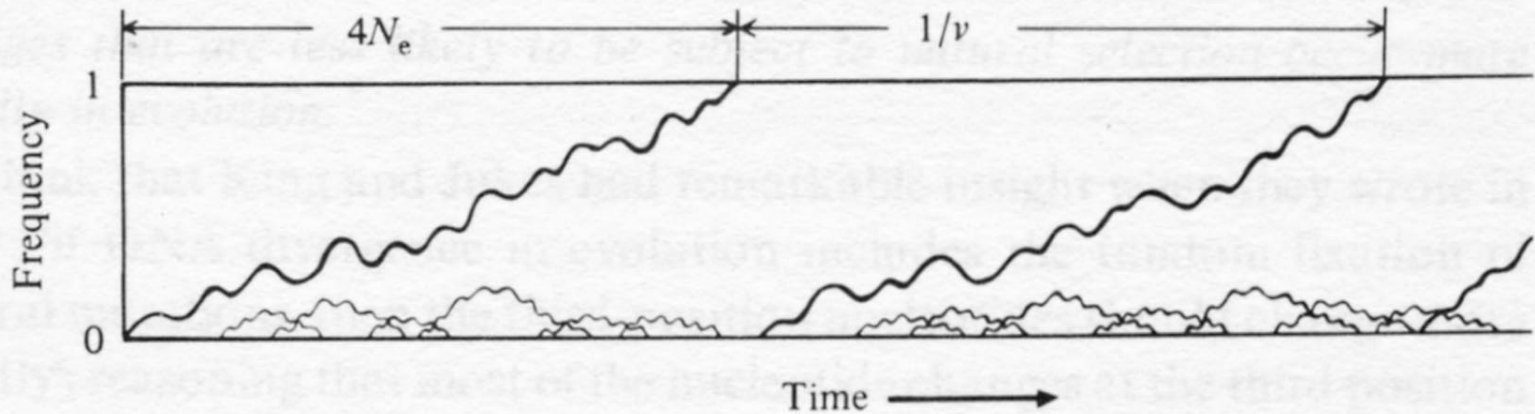


Fig. 3.1. Behavior of mutant genes following their appearance in a finite population. Courses of change in the frequencies of mutants destined to fixation are depicted by thick paths. N_e stands for the effective population size and v is the mutation rate.



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

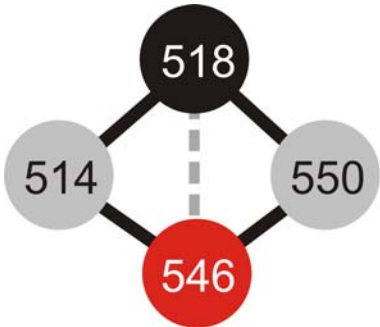


Neutral network

$\lambda = 0.01, s = 367$

$$d_H = 1$$

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



Neutral network

$\lambda = 0.01, s = 877$

$$d_H = 2$$

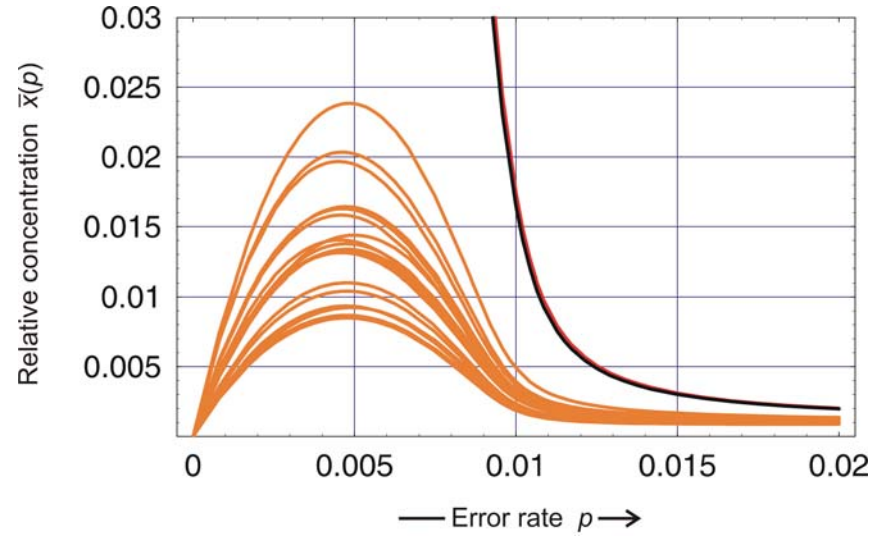
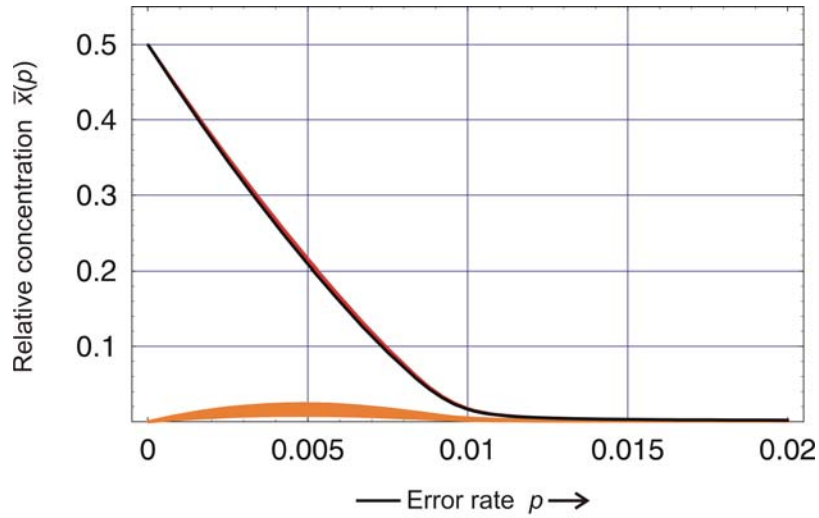
$$\lim_{p \rightarrow 0} x_1(p) = a$$

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

$$d_H = 3$$

random fixation in the sense of
Motoo Kimura

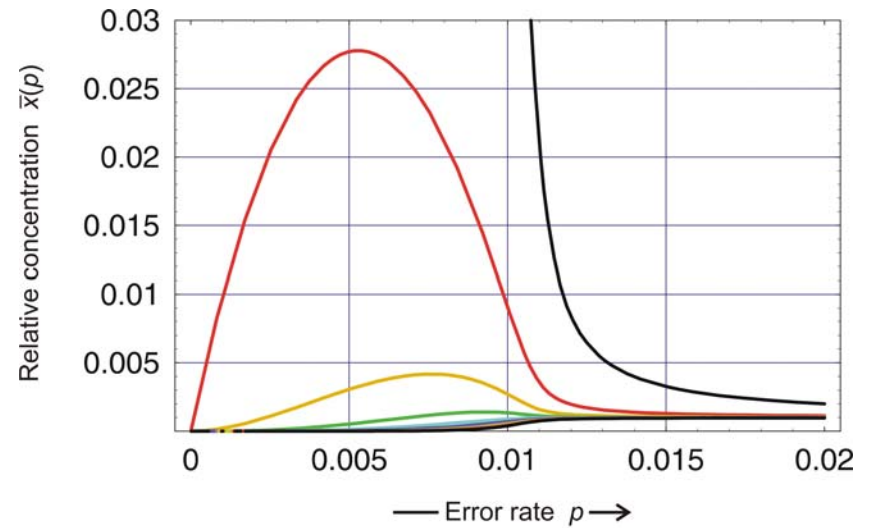
Pairs of genotypes in neutral replication networks

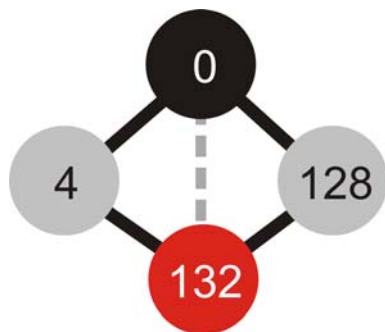


Neutral network
 $\lambda = 0.01, s = 367$

Neutral network: Individual sequences

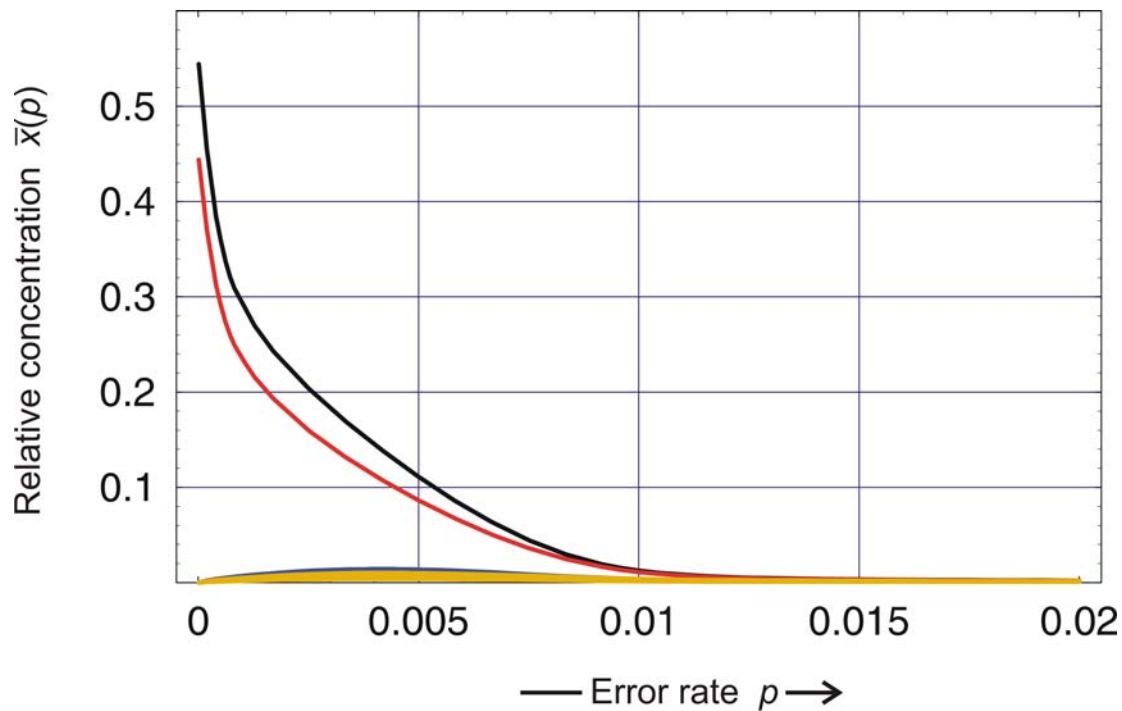
$n = 10, \sigma = 1.1, d = 1.0$





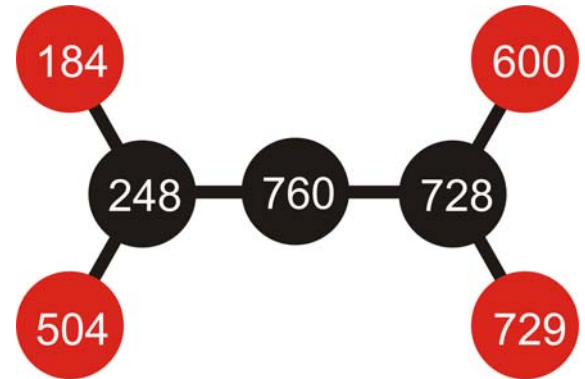
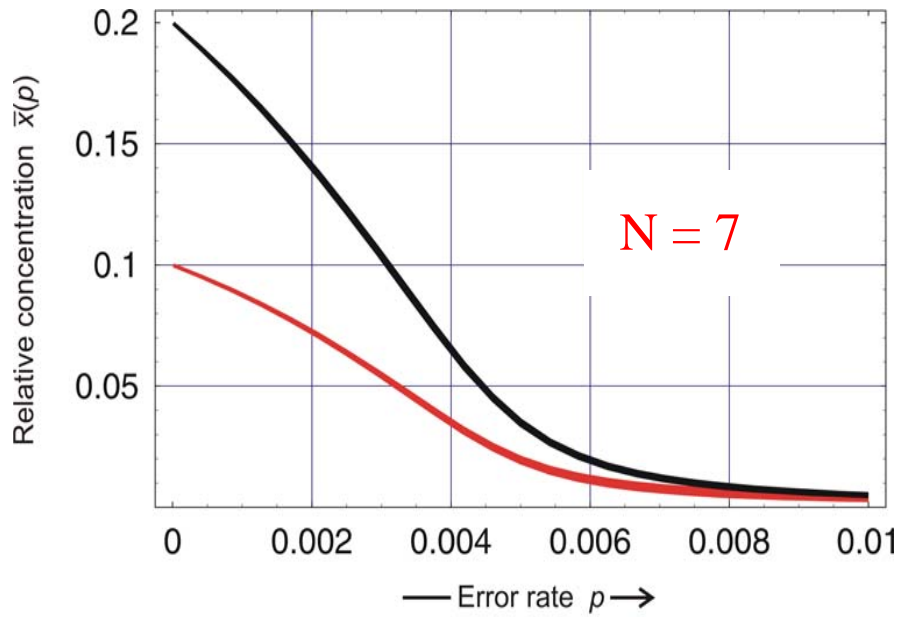
Neutral network

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



Neutral network: Individual sequences

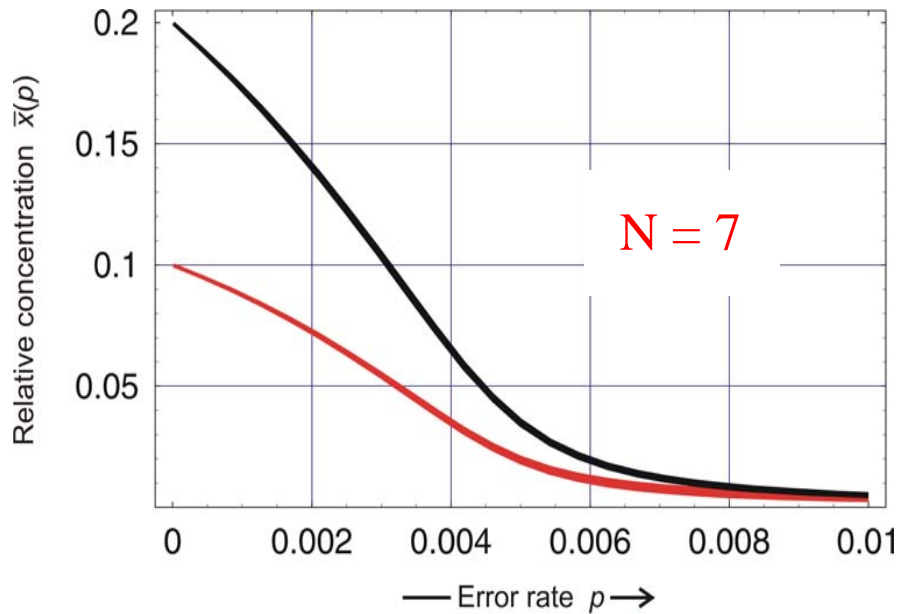
$n = 10$, $\sigma = 1.1$, $d = 1.0$



Neutral network

$\lambda = 0.10, s = 229$

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



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

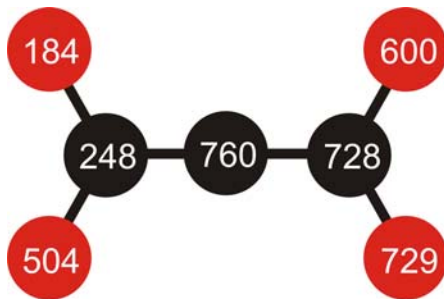
$$\lambda_0 = f + 2\varepsilon,$$

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

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

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

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



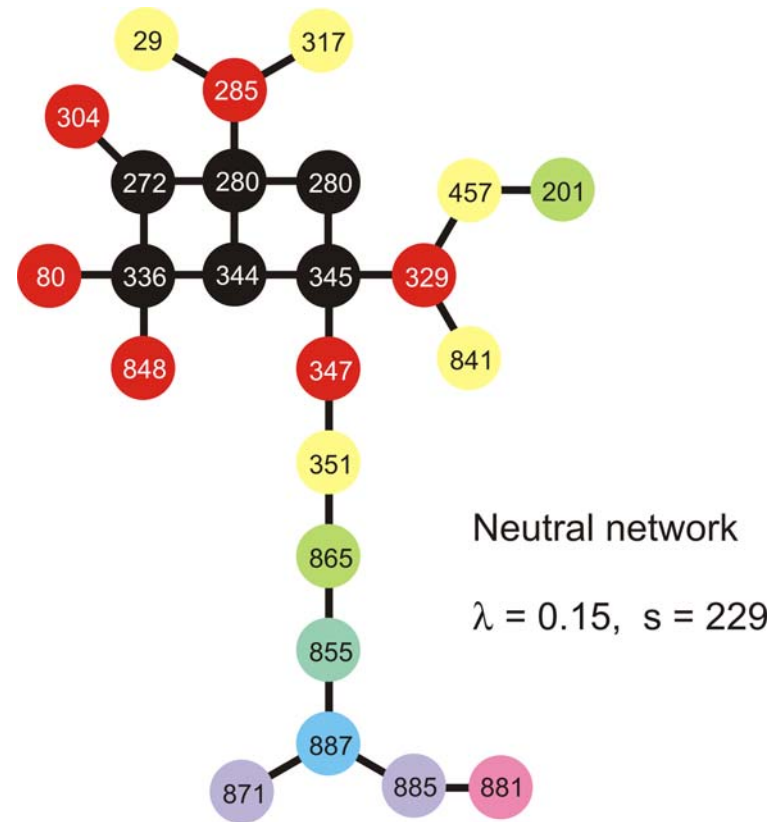
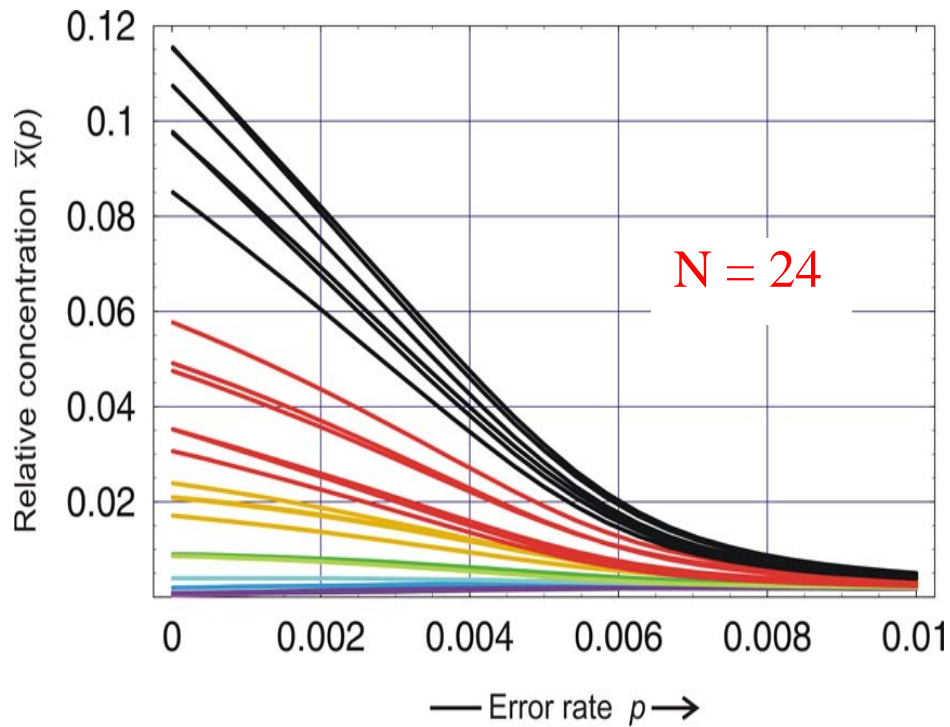
Neutral network

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

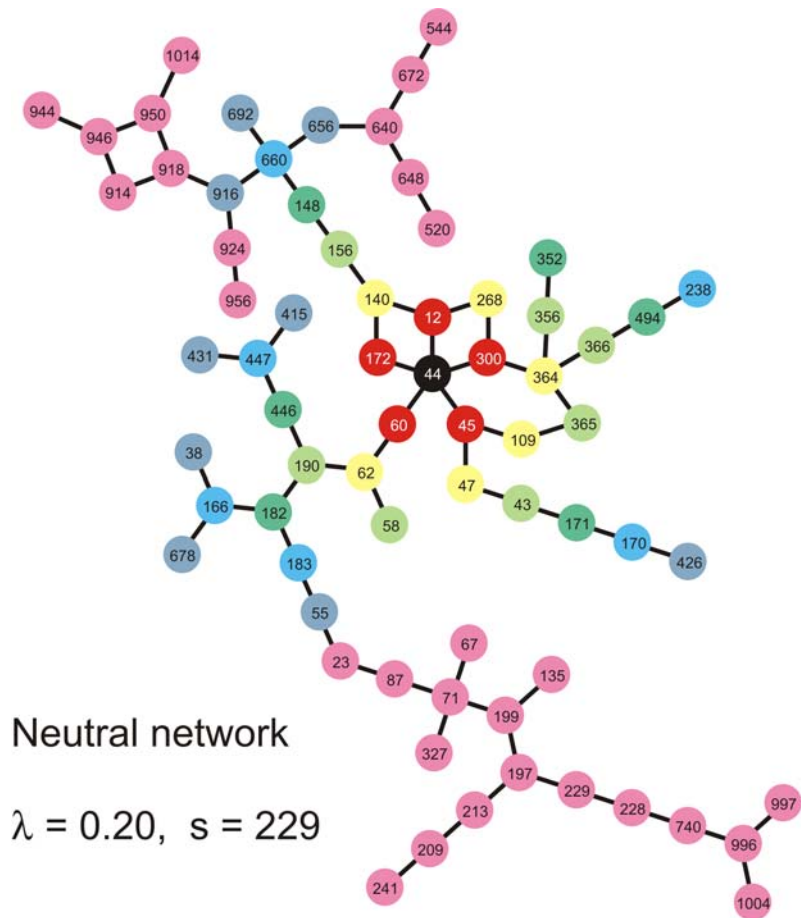
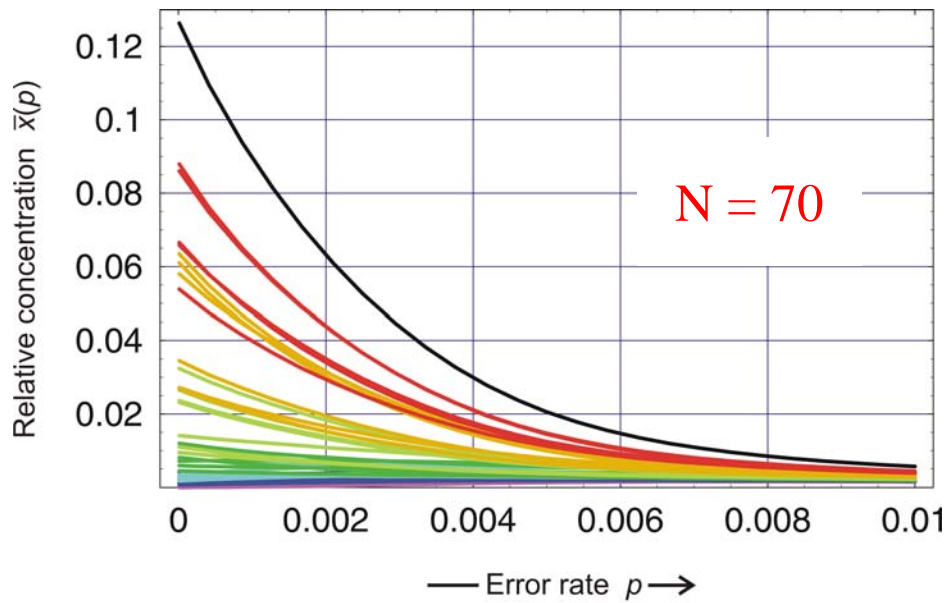
Largest eigenvector of W

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

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



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



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

Results from replication kinetics and RNA neutral networks:

- RNA sequences with Hamming distance $d = 1$ and $d = 2$ form strongly coupled replication ensembles. For $d > 2$ random drift in the sense of Kimura's theory occurs.
- Direct evidence that neutrality is increasing the repertoire of structures and properties in populations.
- Implication for virus replication in infected hosts.

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. Minimum free energy structures of RNA
2. Suboptimal structures of RNA
3. Kinetic folding and RNA switches
4. Chemistry of Darwinian evolution
5. Consequences of neutrality
6. **Evolutionary optimization of RNA structure**

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 μ l using 1 unit of Taq DNA polymerase with each primer at 0.4 μ M, 200 μ M each dATP, dTTP, dCTP, and dGTP, and PCR buffer [10 mM Tris-HCl (pH 8.3), 50 mM KCl, 1.5 mM MgCl₂] 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)⁺ 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)⁺ 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. Torkzadeh, C. Varner, M. Walker, G. Bouffard, and S. Beckstrom-Sternberg (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.

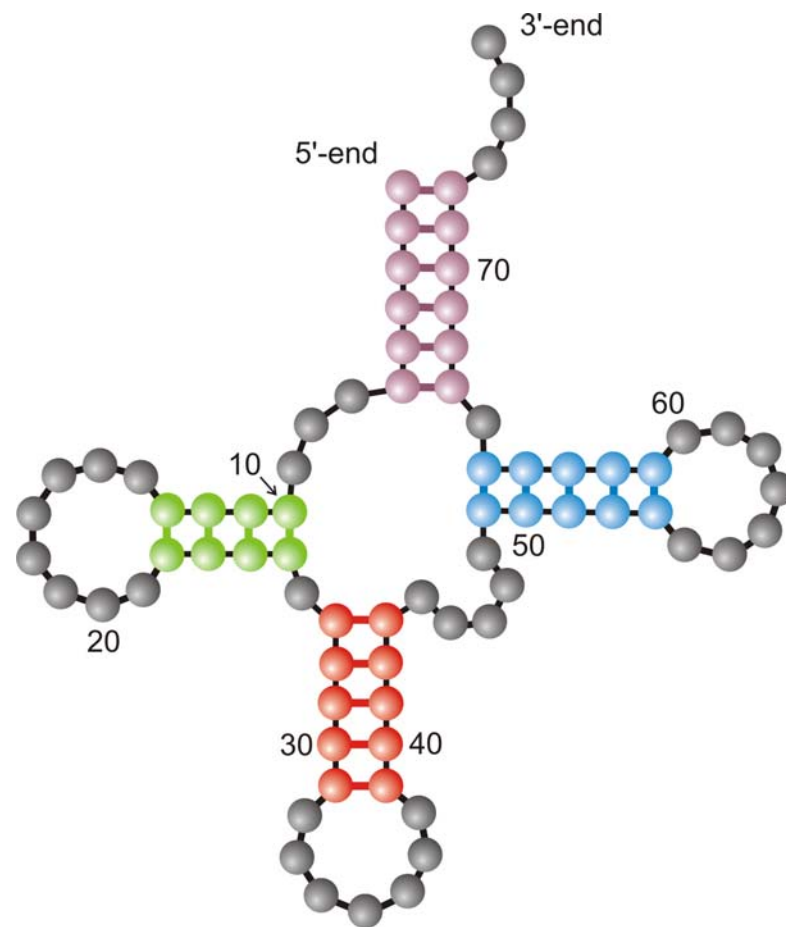
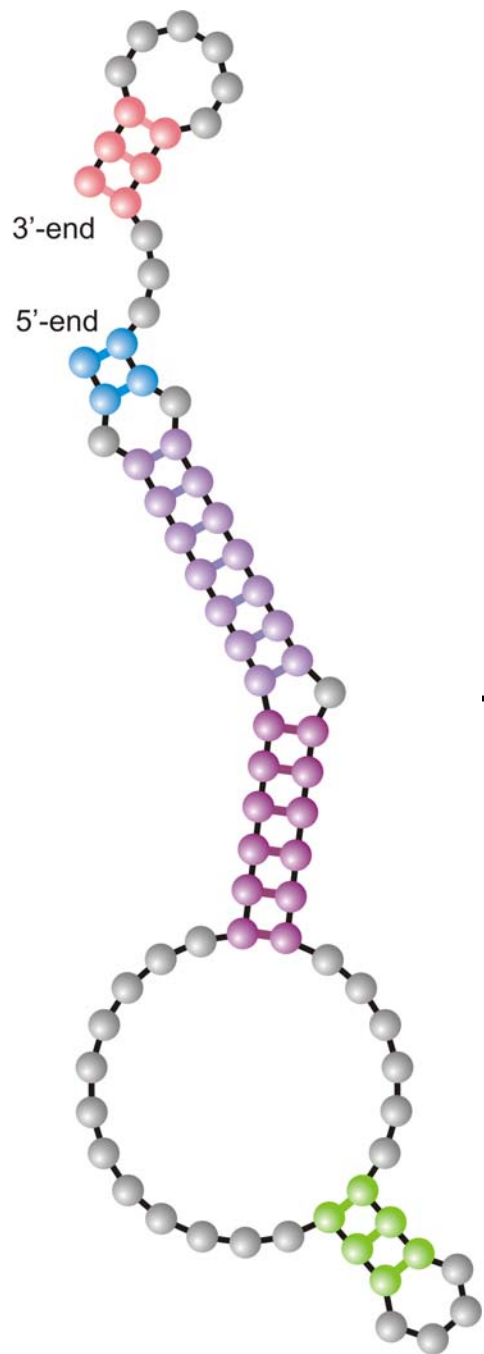
An instance representing in its qualitative features all the simulations we performed is shown in Fig. 1A. Starting with identical sequences folding into a random shape, the simulation was stopped when the population became dominated by the target, here a canonical tRNA shape. The black curve traces 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*

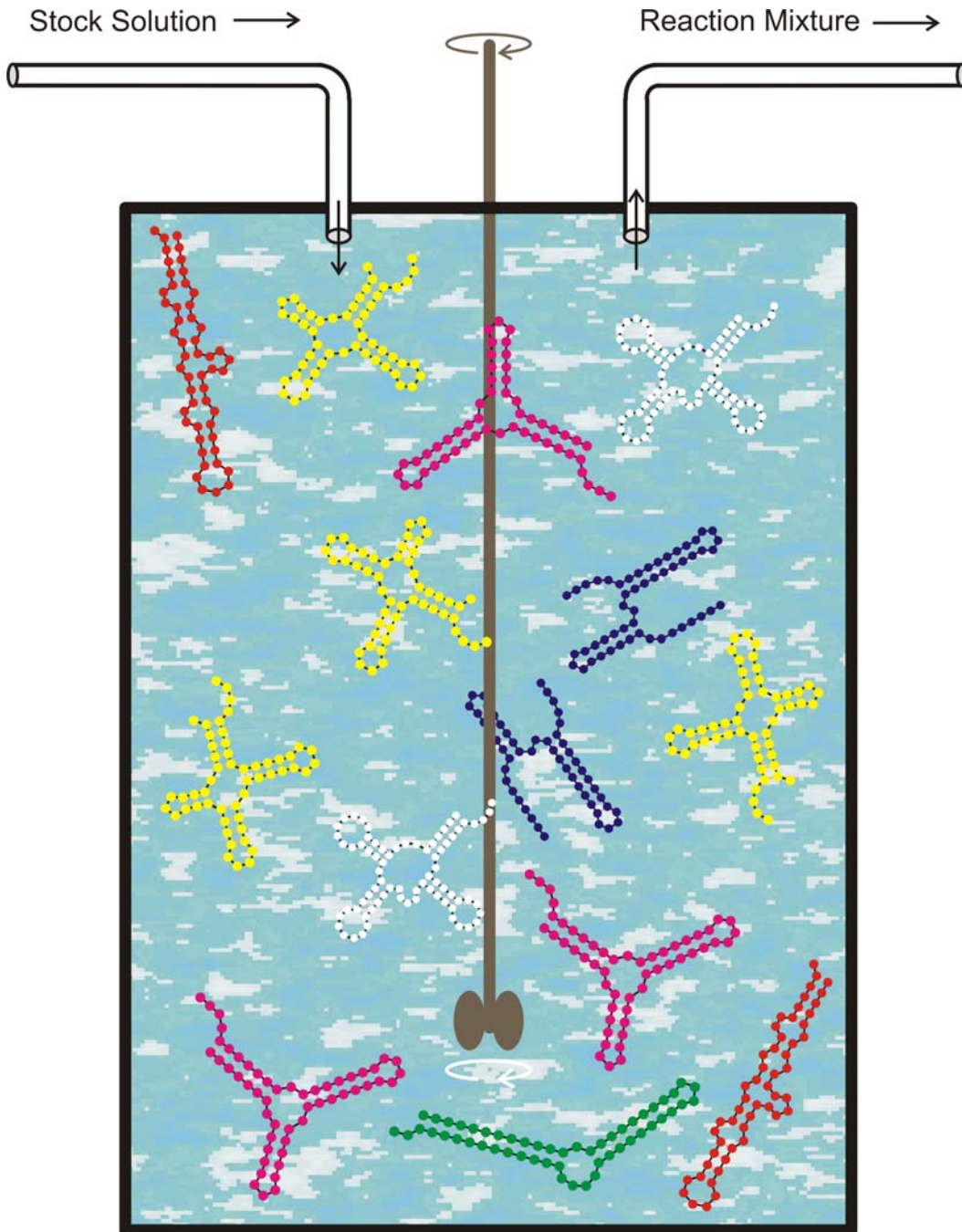
W. Fontana, P. Schuster,
Science **280** (1998), 1451-1455

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,

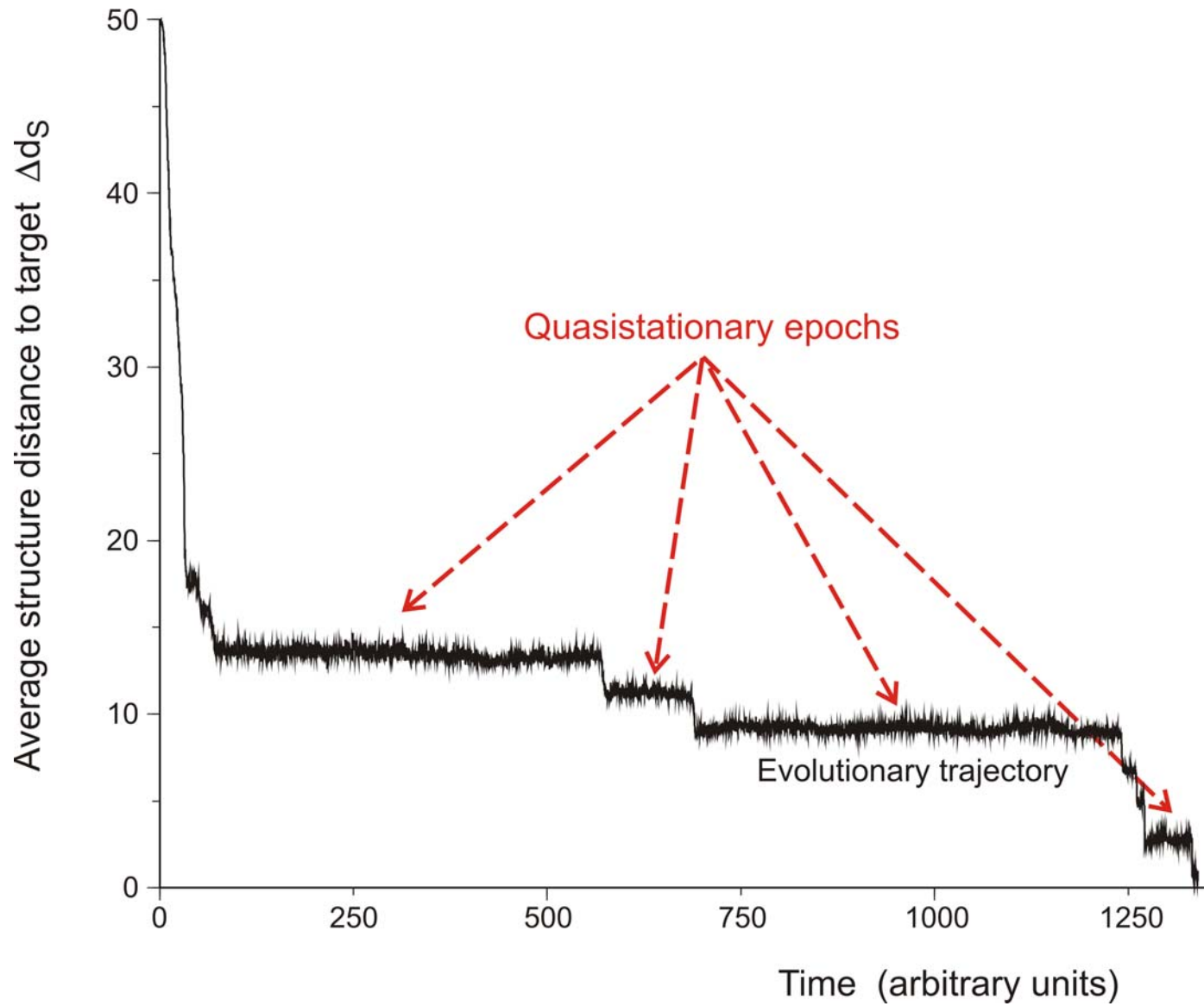
is determined by the flux:

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

Mutation rate:

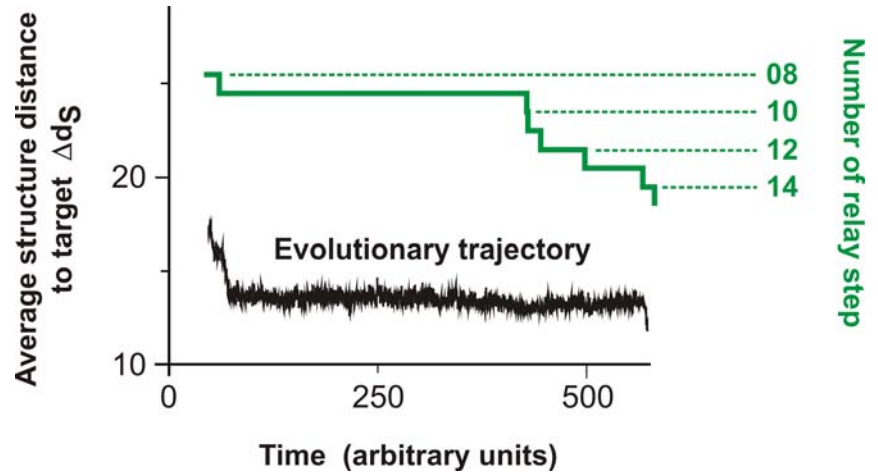
$$p = 0.001 / \text{Nucleotide} \times \text{Replication}$$

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



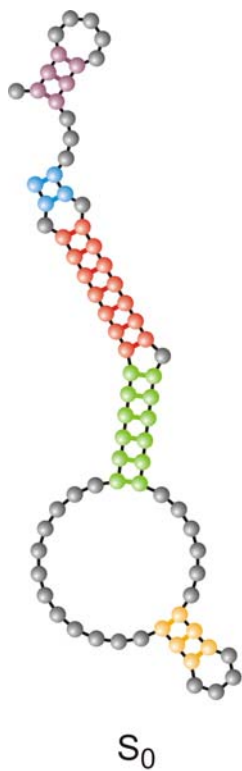
entry GGUAUGGGCGUUGAAUAGGGUUAACCAAUCGGCCAACGAUCUCGUGUGCGCAUUUCAUAUCCCGUACAGAA
 8 .(((((((((((((. ((((.))))))))))(((((.)))))))))
 exit GGUAUGGGCGUUGAAUAAUAGGGUUAACCAAUCGGCCAAACGAUCUCGUGUGCGCAUUUCAUAUCCCAUACAGAA
 entry GGUAUGGGCGUUGAAUAAUAGGGUUAACCAAUCGGCCAAACGAUCUCGUGUGCGCAUUUCAUAUACCAUAACAGAA
 9 .((((((. ((((. ((((.))))))))(((((.)))) .))))
 exit UGGAUGGACGUUGAAUAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACACGUCCCAAG
 entry UGGAUGGACGUUGAAUAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACGCGUCCCAAG
 10 .(((((. . ((((. ((((.))))))))(((((.)))) .))))
 exit UGGAUGGACGUUGAAUAACAAGGUAUCGACCAACAACCAACGAGUAAGUGUGUACGCCCCACACAGCGUCCCAAG

Transition inducing point mutations change the molecular structure

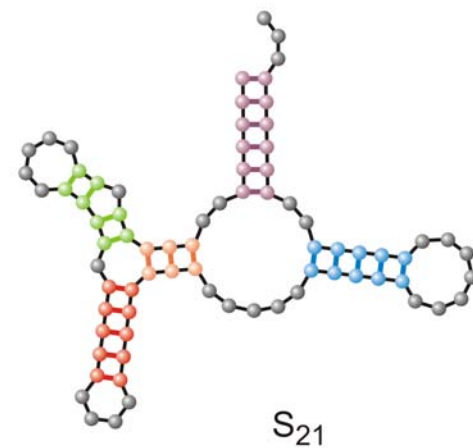
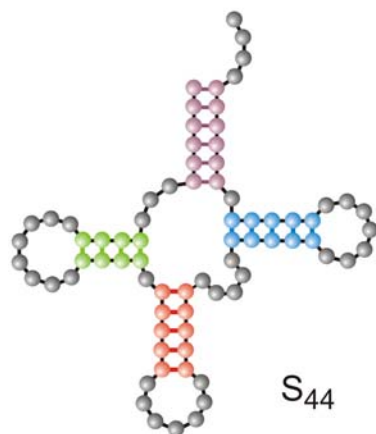
Neutral point mutations leave the molecular structure unchanged

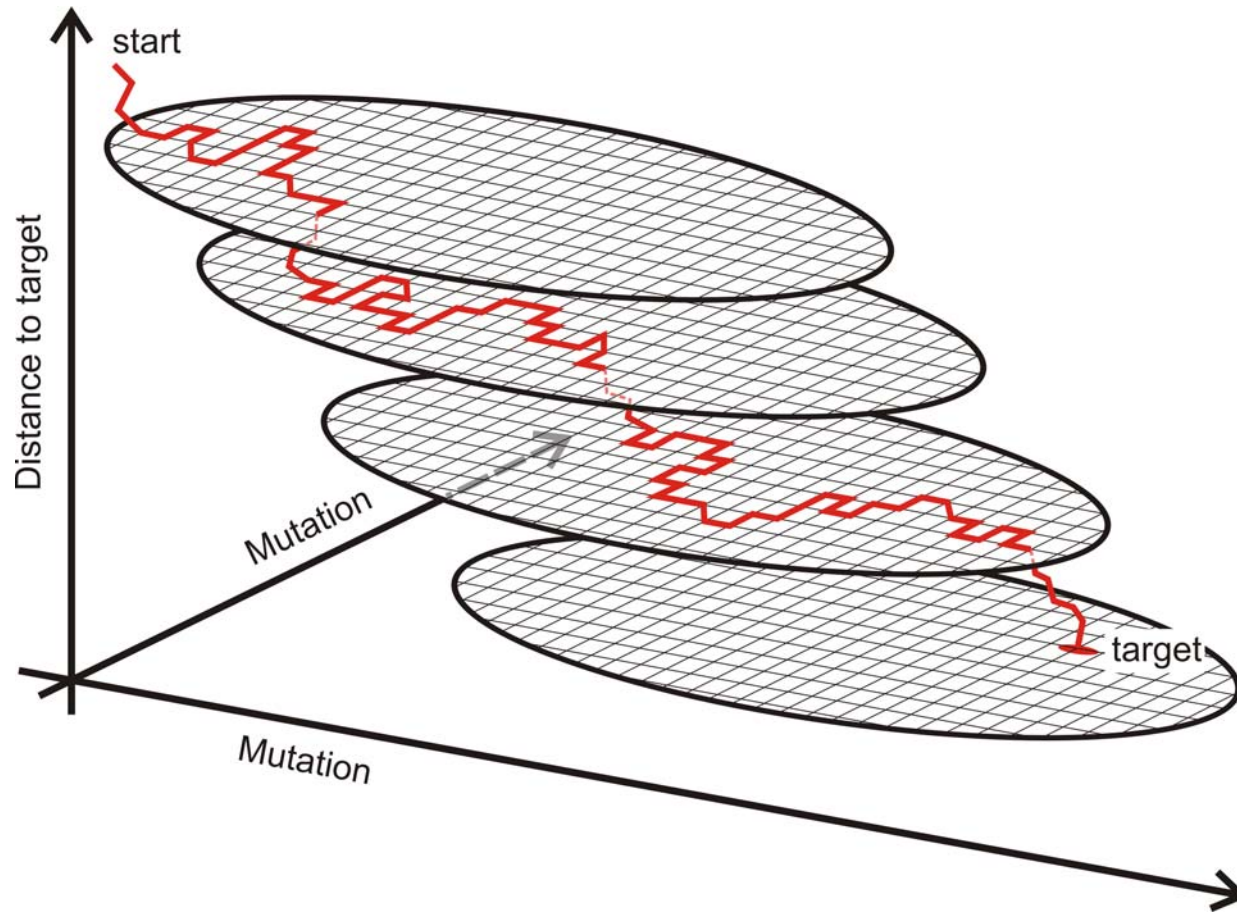
Neutral genotype evolution during phenotypic stasis

Randomly chosen
initial structure



Phenylalanyl-tRNA
as target structure





A sketch of optimization on neutral networks

Results from *in silico* simulation of RNA evolution:

- Evolutionary optimization occurs on two time scales: Fast adaptive phases and random walk on neutral networks.
- Neutral networks are essential for searching sequence space.

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^{1,2}

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

²The Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA

E-mail: pbs@tbi.univie.ac.at

Web-Page for further information:

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

