

LA-UR- 96-1642

CONF. 9605141--2

Title:

EVOLUTIONARY OPTIMIZATION OF BIOPOLYMERS AND  
SEQUENCE STRUCTURE MAPS

RECEIVED

JUN 11 1996

OSTI

Author(s):

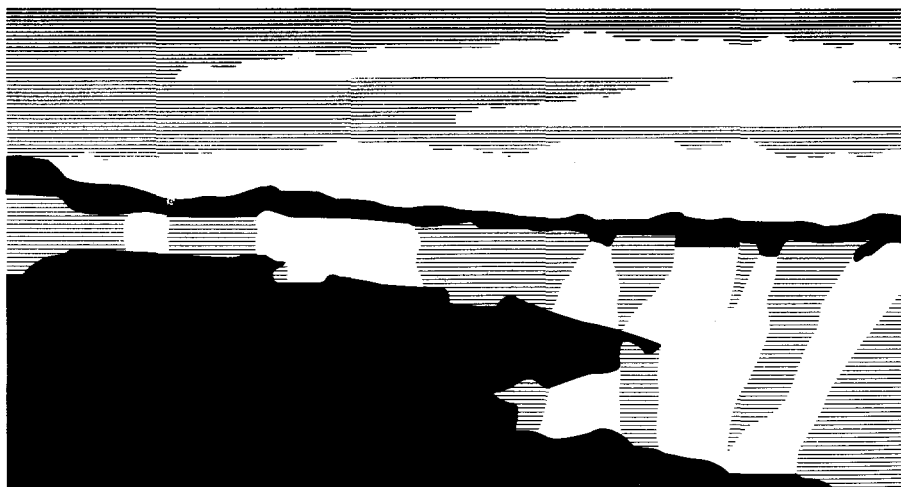
Christian M. Reidys, TSA-DO/SA  
Stephan Kopp, Institut fur Molekulare Biotechnologie  
Peter Schuster, Santa Fe Institute

Submitted to:

Artificial Life V  
Nara, Japan  
May 16-18, 1996

MASTER

**Los Alamos**  
NATIONAL LABORATORY



Los Alamos National Laboratory, an affirmative action/equal opportunity employer, is operated by the University of California for the U.S. Department of Energy under contract W-7405-ENG-36. By acceptance of this article, the publisher recognizes that the U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or to allow others to do so, for U.S. Government purposes. The Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy.

Form No. 836 R5  
ST 2629 10/91

DISTRIBUTION OF THIS DOCUMENT IS UNLIMITED *α*

# Evolutionary Optimization of Biopolymers and Sequence Structure Maps

CHRISTIAN REIDYS<sup>a,b,c</sup>, STEPHAN KOPP<sup>a</sup>, PETER SCHUSTER<sup>a,b</sup>

<sup>a</sup>Institut für Molekulare Biotechnologie  
Beutenbergstraße 11, D-07708 Jena, Germany

<sup>b</sup>Santa Fe Institute  
1399 Hyde Park Rd., Santa Fe, NM 87501, USA

<sup>c</sup>Los Alamos National Laboratory  
Theoretical Division, Los Alamos, New Mexico, USA

E-Mail: duck@santafe.edu

## Abstract

Searching for biopolymers having a predefined function is a core problem of biotechnology, biochemistry and pharmacy. On the level of RNA sequences and their corresponding secondary structures we show that this problem can be analyzed mathematically. The strategy will be to study the properties of the RNA sequence to secondary structure mapping that is essential for the understanding of the search process. We show that to each secondary structure  $s$  there exists a neutral network consisting of all sequences folding into  $s$ . This network can be modeled as a random graph and has the following generic properties: it is dense and has a giant component within the graph of compatible sequences. The neutral network percolates sequence space and any two neutral nets come close in terms of Hamming distance. We investigate the distribution of the orders of neutral nets and show [10] that above a certain threshold the topology of neutral nets allows to find practically all frequent secondary structures.

## 1. Introduction

An RNA structure with a given shape or function is assumed to be formed by many RNA sequences. Their distribution in sequence space is of particular importance for the hardness of the corresponding search problem. The structure of a biopolymer is defined only in the context of some physical conditions. Minimum free energy structures for example fulfill the thermodynamic condition of a molecular ground state, or kinetic structures that are understood as the well defined outcome of a

controlled process of biopolymer formation. In an abstract sense this means that one is interested in a (local) point to point assignment of sequence space and shape space. In general such a mapping will not be one-to-one: many sequences may be mapped into the same structure. The degree of this redundancy will strongly depend on the notion of structure applied. Structure in X-ray crystallography is tantamount to a set of atomic coordinates and at sufficiently high resolution structures are unique in the sense that structures from different sequences will never coincide. Molecular biologists, however, commonly apply another, a coarse-grained notion of structure when, for example, they say intuitively that two proteins have the same structure. An appropriate coarse grained notion of structure apparently is context dependent and thus anything but trivial.

In this paper we are dealing with RNA molecules. Secondary structures are used as appropriate examples for structural coarse-graining. They are sufficiently simple to allow statistical analysis by means of conventional combinatorics [9]. The relation between RNA sequences and secondary structures is understood as a (non invertible) mapping from sequence space into shape space [3, 5, 6]. RNA secondary structures distinguish only paired and unpaired regions irrespective of the particular bases at the individual positions (G, C, A, or U). Therefore many different sequences the so called compatible sequences can meet the base pairing conditions as determined by a given secondary structure. For the biophysical alphabet the number of compatible sequences is readily computed for any given secondary structure  $s$ , with  $u$  unpaired bases and  $p$  base pairs to be  $4^u \cdot 6^p$ . The number of compatible sequences is certainly substantially

## 2.1 Density

We shall discuss in this section the *density property* of random graphs  $\Gamma_n < Q_\alpha^n$ , where  $Q_\alpha^n$  is a generalized hypercube, i.e., the graph formed by all  $n$ -tuples of coordinates  $x_i$  contained in a finite set  $\mathcal{A}$  (of cardinality  $\alpha$ ) where each two tuples are neighbors when they differ in exactly one coordinate.

Let  $H$  be a finite graph. A subgraph  $G < H$  is called *dense* in  $H$  if and only if  $\overline{v[G]} = v[H]$ .

We will establish the existence of a "critical"  $\lambda$ -value,  $\lambda^*$  that has the following property: for  $\lambda < \lambda^*$  a.a.s. (asymptotically almost surely) no random graph  $\Gamma_n$  is dense and for  $\lambda > \lambda^*$  a.a.s. every random graph  $\Gamma_n$  is dense. We will call  $\lambda^*$  the *threshold value for the density property*. For this purpose we consider the random variable

$$\hat{Z}_n(\Gamma_n) := |\{v \in v[Q_\alpha^n] \mid v \notin \overline{v[\Gamma_n]}\}| \quad (1)$$

that is defined on  $\Omega_n$  and counts the number of vertices having no adjacent vertex  $v \in v[\Gamma_n]$ . We first compute the asymptotic distribution of the following sequence of random variables  $(\hat{Z}_n)$  associated to the sequence of probability spaces  $(\Omega_n)$ . For this purpose we make use of the *sieve formula* [2] (p.17) that implies a number of results about the convergence in distribution for a sequence of integer valued random variables  $(\hat{X}_n)$ . **Theorem 1.** Let  $(\hat{X}_i)_{i \in \mathbb{N}}$  be a sequence of non-negative integer valued random variables such that

$$\forall r \in \mathbb{N}: \lim_{n \rightarrow \infty} E[\hat{X}_n]_r = E[\hat{X}]_r$$

and

$$\forall m \in \mathbb{N}: \lim_{r \rightarrow \infty} E[\hat{X}]_r r^m / r! = 0$$

Then we have the following convergence in distribution:  $\hat{X}_n \rightarrow \hat{X}$ .

**Proof.** [2, p.23] ■

**Corollary 1.** Let  $\mu = \mu(n)$  be a bounded, non-negative function on  $\mathbb{N}$  and assume a sequence of non-negative integer valued random variables  $(\hat{X}_i)_{i \in \mathbb{N}}$  to be given. Suppose for an arbitrary natural number  $r$  we have

$$\lim_{n \rightarrow \infty} E[\hat{X}_n]_r - \mu^r = 0.$$

Then the following convergence holds in distribution:

$$d(\hat{X}_n, P_\mu) \rightarrow 0,$$

where  $P_\mu$  is the Poisson measure.

**Lemma 1.** Suppose

$$\mu := \lim_{n \rightarrow \infty} (|Q_\alpha^n| (1 - \lambda)^{\gamma_n + 1}) \in \mathbb{R} \cup \{0\} \cup \{\infty\}$$

In particular we have

$$\lim_{n \rightarrow \infty} \mu_n \{\hat{Z}_n = 0\} = e^{-\mu},$$

and

$$\lim_{n \rightarrow \infty} E[\hat{Z}_n] = \alpha^n (1 - \lambda)^{\gamma_n + 1}.$$

Finally, for  $\mu = \infty$  and  $\ell \in \mathbb{N}$  holds

$$\lim_{n \rightarrow \infty} \mu_n \{\hat{Z}_n \geq \ell\} = 1.$$

The following theorem shows that  $\lambda^* := 1 - \alpha^{-1/\alpha-1}$  is a *threshold value* for the density property of random induced subgraphs as introduced in the basic model. Above  $\lambda^*$  a.a.s. all subgraphs are dense and below  $\lambda^*$  a.a.s. none of them.

**Theorem 2.** Let  $\lambda^* := 1 - \alpha^{-1/\alpha-1}$  then for  $\lambda > \lambda^*$  holds

$$\lim_{n \rightarrow \infty} \mu_n \{\Gamma_n \text{ is dense in } Q_\alpha^n\} = 1$$

and for  $\lambda < \lambda^*$  we have

$$\lim_{n \rightarrow \infty} \mu_n \{\Gamma_n \text{ is dense in } Q_\alpha^n\} = 0.$$

## 2.2 Connectivity and Giant Components

Let  $G$  be a finite graph. Being connected is an equivalence relation on  $v[G]$  and there exist maximal subsets  $V \subset v[G]$  consisting of connected vertices. A *component* of  $G$  is then the induced subgraph  $G' = G[V]$  of such a maximal connected subset of vertices. If  $V = \emptyset$ ,  $G[\emptyset]$  is called a *trivial component*. If  $G$  is disconnected we shall investigate the so called *sequence of components*, i.e., the list of orders of the maximal connected subgraphs of  $G$  into which  $G$  can be decomposed.

Given a graph  $G$ , the *sequence of components* of  $G$  is the ordered tuple  $(|\mathcal{X}_i|)_{1 \leq i \leq |G|}$ , where each  $\mathcal{X}_i$  is a component of  $G$  and  $|\mathcal{X}_i| \geq |\mathcal{X}_{i+1}|$ . We call a component  $\mathcal{X} < G$  a *giant component* if and only if  $|\mathcal{X}| \geq \frac{2}{3}|G|$ .

The key idea in the proof of the connectivity theorem bases on the following observation (formulated as lemma 2 below). A.a.s. each pair of vertices  $v, v' \in v[\Gamma_n]$  with  $d(v, v') = k$ , for fixed natural number  $k$ , is connected by a path in  $\Gamma_n$ .

For this purpose we refer to a certain family of independent paths in  $Q_\alpha^n$ . I.e. for  $v, v' \in v[Q_\alpha^n]$  with  $d(v, v') = k$  we write  $v, v'$  as  $v = (x_1, \dots, x_k, x_{k+1}, \dots, x_n)$  and  $v' = (x'_1, \dots, x'_k, x_{k+1}, \dots, x_n)$ . Then for  $v_1 \in \partial\{v\} \cap B_{1+k}(v')$  we set

$$a_i(v_1) := (x_1, \dots, x_i, x'_1, \dots, x'_i, x_{i+1}, \dots, x_n)$$

## 2.1 Density

We shall discuss in this section the *density property* of random graphs  $\Gamma_n < Q_\alpha^n$ , where  $Q_\alpha^n$  is a generalized hypercube, i.e., the graph formed by all  $n$ -tuples of coordinates  $x$ , contained in a finite set  $\mathcal{A}$  (of cardinality  $\alpha$ ) where each two tuples are neighbors when they differ in exactly one coordinate.

Let  $H$  be a finite graph. A subgraph  $G < H$  is called *dense* in  $H$  if and only if  $\overline{v[G]} = v[H]$ .

We will establish the existence of a "critical"  $\lambda^*$  value,  $\lambda^*$  that has the following property: for  $\lambda < \lambda^*$  a.a.s. (asymptotically almost surely) no random graph  $\Gamma_n$  is dense and for  $\lambda > \lambda^*$  a.a.s. every random graph  $\Gamma_n$  is dense. We will call  $\lambda^*$  the *threshold value for the density property*. For this purpose we consider the random variable

$$\hat{Z}_n(\Gamma_n) := |\{v \in v[Q_\alpha^n] \mid v \notin \overline{v[\Gamma_n]}\}| \quad (1)$$

that is defined on  $\Omega_n$  and counts the number of vertices having no adjacent vertex  $v \in v[\Gamma_n]$ . We first compute the asymptotic distribution of the following sequence of random variables  $(\hat{Z}_n)$  associated to the sequence of probability spaces  $(\Omega_n)$ . For this purpose we make use of the *sieve formula* [2] (p.17) that implies a number of results about the convergence in distribution for a sequence of integer valued random variables  $(\hat{X}_n)$ . **Theorem 1.** Let  $(\hat{X}_i)_{i \in \mathbb{N}}$  be a sequence of non-negative integer valued random variables such that

$$\forall r \in \mathbb{N}: \lim_{n \rightarrow \infty} E[\hat{X}_n]_r = E[\hat{X}]_r$$

and

$$\forall m \in \mathbb{N}: \lim_{r \rightarrow \infty} E[\hat{X}]_r r^m / r! = 0$$

Then we have the following convergence in distribution:  $\hat{X}_n \rightarrow \hat{X}$ .

**Proof.** [2, p.23] ■

**Corollary 1.** Let  $\mu = \mu(n)$  be a bounded, non-negative function on  $\mathbb{N}$  and assume a sequence of non-negative integer valued random variables  $(\hat{X}_i)_{i \in \mathbb{N}}$  to be given. Suppose for an arbitrary natural number  $r$  we have

$$\lim_{n \rightarrow \infty} E[\hat{X}_n]_r - \mu^r = 0.$$

Then the following convergence holds in distribution:

$$d(\hat{X}_n, P_\mu) \rightarrow 0,$$

where  $P_\mu$  is the Poisson measure.

**Lemma 1.** Suppose

$$\mu := \lim_{n \rightarrow \infty} (|Q_\alpha^n| (1 - \lambda)^{\gamma_n + 1}) \in \mathbb{R}_+ \cup \{0\} \cup \{\infty\}$$

exists. Then for  $\mu < \infty$  the random variables  $\hat{Z}_n$  converge in distribution to a Poisson distributed random variable, i.e.,

$$\lim_{n \rightarrow \infty} \mu_n \{\hat{Z}_n = \ell\} = \frac{\mu^\ell}{\ell!} e^{-\mu}. \quad (2)$$

In particular we have

$$\lim_{n \rightarrow \infty} \mu_n \{\hat{Z}_n = 0\} = e^{-\mu},$$

and

$$\lim_{n \rightarrow \infty} E[\hat{Z}_n] = \alpha^n (1 - \lambda)^{\gamma_n + 1}.$$

Finally, for  $\mu = \infty$  and  $\ell \in \mathbb{N}$  holds

$$\lim_{n \rightarrow \infty} \mu_n \{\hat{Z}_n \geq \ell\} = 1.$$

The following theorem shows that  $\lambda^* := 1 - \alpha^{-1/\sqrt{\alpha-1}}$  is a *threshold value* for the density property of random induced subgraphs as introduced in the basic model. Above  $\lambda^*$  a.a.s. all subgraphs are dense and below  $\lambda^*$  a.a.s. none of them.

**Theorem 2.** Let  $\lambda^* := 1 - \alpha^{-1/\sqrt{\alpha-1}}$  then for  $\lambda > \lambda^*$  holds

$$\lim_{n \rightarrow \infty} \mu_n \{\Gamma_n \text{ is dense in } Q_\alpha^n\} = 1$$

and for  $\lambda < \lambda^*$  we have

$$\lim_{n \rightarrow \infty} \mu_n \{\Gamma_n \text{ is dense in } Q_\alpha^n\} = 0.$$

## 2.2 Connectivity and Giant Components

Let  $G$  be a finite graph. Being connected is an equivalence relation on  $v[G]$  and there exist maximal subsets  $V \subset v[G]$  consisting of connected vertices. A *component* of  $G$  is then the induced subgraph  $G' = G[V]$  of such a maximal connected subset of vertices. If  $V = \emptyset$ ,  $G[\emptyset]$  is called a *trivial component*. If  $G$  is disconnected we shall investigate the so called *sequence of components*, i.e., the list of orders of the maximal connected subgraphs of  $G$  into which  $G$  can be decomposed.

Given a graph  $G$ , the *sequence of components* of  $G$  is the ordered tuple  $(|\mathcal{X}_i|)_{1 \leq i \leq |G|}$ , where each  $\mathcal{X}_i$  is a component of  $G$  and  $|\mathcal{X}_i| \geq |\mathcal{X}_{i+1}|$ . We call a component  $\mathcal{X} < G$  a *giant component* if and only if  $|\mathcal{X}| \geq \frac{2}{3}|G|$ . The key idea in the proof of the connectivity theorem bases on the following observation (formulated as lemma 2 below). A.a.s. each pair of vertices  $v, v' \in v[\Gamma_n]$  with  $d(v, v') = k$ , for fixed natural number  $k$ , is connected by a path in  $\Gamma_n$ .

For this purpose we refer to a certain family of independent paths in  $Q_\alpha^n$ . I.e. for  $v, v' \in v[Q_\alpha^n]$  with  $d(v, v') = k$  we write  $v, v'$  as  $v = (x_1, \dots, x_k, x_{k+1}, \dots, x_n)$  and  $v' = (x'_1, \dots, x'_k, x_{k+1}, \dots, x_n)$ . Then for  $v_1 \in \partial\{v\} \cap B_{1+k}(v')$  we set

$$g_j(v_1) := (x_1, \dots, x_j, x'_{j+1}, \dots, x'_k, x_{k+1}, \dots, x_n)$$

$$0 \leq j \leq k \quad \hat{x}_r \neq x_r \quad (3)$$

and inspect  $g_k(v_1) = v_1$ ,  $g_0(v_1) \in B_1(v') \cap B_{1+k}(v)$ . We introduce the random variable  $\hat{Y}_{n,k}^{v,v'}$  that counts the (independent) paths in the random graph  $\Gamma_n$  connecting the vertices  $v, v'$  having distance  $k$ .

**Lemma 2.** Let  $k$  be a natural number,  $Q_\alpha^n$  a generalized hypercube and  $\Gamma_n < Q_\alpha^n$  a random graph with  $\lambda > 1 - \alpha^{-1/\sqrt{\alpha-1}}$ . Then  $\lim_{n \rightarrow \infty} \mu_n\{T\} = 1$  where

$$T := \{\Gamma_n \mid \forall v, v' \in v[Q_\alpha^n], d(v, v') = k : \\ \exists v_1 \in \partial\{v\}, v'_1 \in \partial\{v'\} : \hat{Y}_{n, d(v_1, v'_1)}^{v_1, v'_1} > 0\}.$$

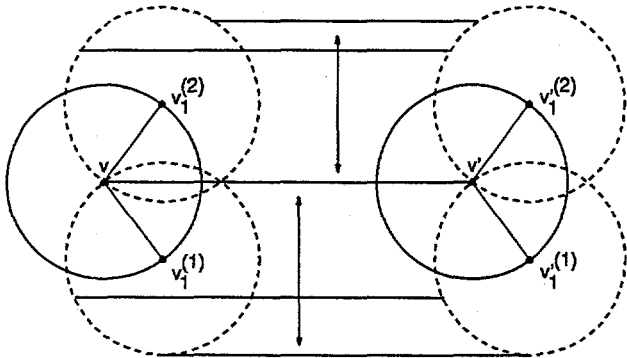


Fig. 1: An illustration for the proof of lemma 2. For given  $v, v' \in v[Q_\alpha^n]$  each pair of vertices  $(v_1^{(i)}, v'_1{}^{(i)})$  leads to "sufficiently many" independent pairwise disjoint paths in  $\Pi(Q_\alpha^n)$ .

Now we are prepared to state the connectivity theorem:

**Theorem 3.** Let  $Q_\alpha^n$  a generalized hypercube and  $\Gamma_n < Q_\alpha^n$  a random induced subgraph. Then

$$\lim_{n \rightarrow \infty} \mu_n\{\Gamma_n \text{ is connected}\} = \begin{cases} 1 & \text{for } \lambda > 1 - \alpha^{-1/\sqrt{\alpha-1}} \\ 0 & \text{for } \lambda < 1 - \alpha^{-1/\sqrt{\alpha-1}} \end{cases} \quad (4)$$

**Remark.** A related result in the special case of the Boolean hypercube can be found in [2]. The corresponding subgraphs  $A_p$  are constructed as follows: We set  $v[A_p] := v[Q_2^n]$  as vertex set and the edge set  $e[A_p]$  is obtained by independent random choices with probability  $p$  in the edge set  $e[Q_2^n]$ . Then the idea of the proof is to establish an *edge boundary* of possible components using an *isoperimetric inequality* due to Harper, Bernstein, and Row [7, 2]. For Boolean hypercubes Ajtai, Komlós and Szemerédi 1982 proved the following related result: for random subgraphs  $A_p$  of  $Q_2^n$  obtained by edge selections, there exists a component of order  $g 2^n$  with constant  $g \in \mathbb{R}_+$  if  $p = c/n$  and  $c > 1$  [1].

Finally, for any positive  $\lambda$  a.a.s. there exists a giant component in random induced subgraphs  $\Gamma_n < Q_\alpha^n$ .

**Theorem 4.** Let  $0 < \lambda \leq 1$ . Then we have

$$\lim_{n \rightarrow \infty} \mu_n\{\Gamma_n \text{ has a giant component}\} = 1.$$

### 3. Sequence Structure Maps Via Random Graphs

#### 3.1 Neutral Networks of Secondary Structures

A pairing rule  $\Pi$  is a symmetric relation in  $\mathcal{A} \times \mathcal{A}$ . Following [15] a secondary structure  $s$  w.r.t.  $\Pi$  is a vertex-labeled graph on  $n$  vertices  $(x_1, \dots, x_n)$  with an adjacency matrix  $A$  fulfilling

- (1)  $a_{i, i+1} = 1$  for  $1 \leq i \leq n-1$ ;
- (2) For each  $i$  there is at most a single  $k \neq i-1, i+1$  such that  $a_{i, k} = 1$  and  $[x_i, x_k] \in \Pi$ ;
- (3) If  $a_{i, j} = a_{k, l} = 1$  and  $i < k < j$  then  $i < l < j$ .

We call an edge  $(x_i, x_k)$ ,  $|i - k| \neq 1$  a *bond* or *base pair*. A vertex  $x_i$  connected only to  $x_{i-1}$  and  $x_{i+1}$  shall be called *unpaired*. The number of base pairs and the number of unpaired bases in a secondary structure  $s$  are  $p(s)$  and  $u(s)$ , respectively. The size of the alphabet is  $\alpha$  and the number of distinct base pairs is given by  $\beta$ .

We proceed by constructing the preimage of a fixed secondary structure as a random induced subgraph of the graph of compatible sequences. Let  $s$  be a secondary structure and

$$\Pi(s) := \{[i, k] \mid a_{i, k} = 1, k \neq i \pm 1\}$$

its *set of contacts*. Then a vertex  $x \in v[Q_\alpha^n]$  is said to be *compatible* to  $s$  if and only if  $\forall [i, j] \in \Pi(s) : [x_i, x_j] \in \Pi$  i.e. the coordinates  $x_i$  and  $x_j$  are in  $\Pi$  for all pairs  $[i, j] \in \Pi(s)$ . We denote the set of all compatible sequences by  $C[s]$ . Finally the graph of compatible sequences w.r.t. the secondary structure  $s$ , is given by

$$C[s] := Q_\alpha^{u(s)} \times Q_\beta^{p(s)}.$$

**Neutral Network:** Let  $\Gamma_u < Q_\alpha^u$  and  $\Gamma_p < Q_\beta^p$  be random subgraphs with underlying parameters  $\lambda_u$  and  $\lambda_p$  as introduced in the basic model. Then we set  $\Gamma_n[s] = \Gamma_u \times \Gamma_p$ ,

$$\mu_{\lambda_u, \lambda_p}(\Gamma_n[s]) := \mu_{\lambda_u, \lambda_u}(\Gamma_u) \times \mu_{\lambda_p, \lambda_p}(\Gamma_p),$$

and  $\mu_{\lambda_u, \lambda_p}$  is a probability measure and  $\Gamma_n[s] < C[s]$ .

We can think of the neutral network,  $\Gamma_n[s]$ , to be obtained by selecting the coordinates  $v_1, v_2$  of the vertex  $(v_1, v_2) \in v[C[s]]$  with the probabilities  $\lambda_u$  and  $\lambda_p$ . This process leads to the vertex set  $V_{\lambda_u, \lambda_p} \subset C[s]$  whose induced subgraph  $C[s][V_{\lambda_u, \lambda_p}]$  is

$$\Gamma_n[s] = \Gamma_u \times \Gamma_p.$$

The theory presented in the previous chapter implies for neutral networks of secondary structures density and connectivity if both factors  $\Gamma_u$  and  $\Gamma_p$  are dense and connected.

### 3.2 Complete Mappings

Once we know how to construct a neutral network  $\Gamma_n[s]$  we order the set of secondary structures  $\mathcal{S}_n$  and define a complete mapping by iterating the construction process of the corresponding neutral network w.r.t. the ordering. Thereby we obtain w.r.t. a given  $\lambda$  parameter a complete sequence to structure mapping.

Let  $C^* : \mathcal{S}_n \rightarrow \mathcal{P}(\mathcal{V}[\mathcal{Q}_\alpha^n])$  and  $r : \mathcal{S}_n \rightarrow \mathbb{N}$  be two mappings such that  $j \leq i \Rightarrow r(s_j) \geq r(s_i)$ . A mapping  $f : \mathcal{Q}_\alpha^n \rightarrow \mathcal{S}_n$  is called  $C^*$ -map if and only if

$$(*) : f(v) = s \Rightarrow v \in C^*[s].$$

A mapping  $f_r : \mathcal{Q}_\alpha^n \rightarrow \mathcal{S}_n$  is called  $C^*$ -random-map if and only if  $f_r$  is given by

$$f_r^{-1}(s_0) := \Gamma_n[s_0]$$

$$f_r^{-1}(s_i) := \Gamma_n[s_i] \setminus \bigcup_{j < i} [\Gamma_n[s_i] \cap \Gamma_n[s_j]].$$

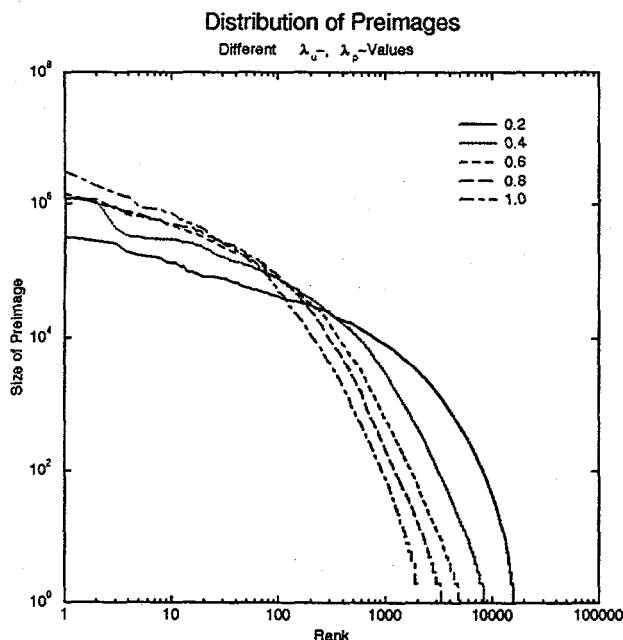


Fig. 2: We report here the logarithm of the sizes of the neutral networks  $f^{-1}(s)$  obtained by a  $C^*$ -random map and  $\lambda$  parameters  $\lambda = 0.2, 0.4, 0.6, 0.8, 1.0$ . The corresponding neutral networks are ordered on the x-axis by the logarithm of their orders. Note that the rank of the secondary structure  $s$  does not necessarily coincide with the size of the corresponding preimage  $|f^{-1}(s)|$ .

In particular any RNA folding map is a  $C^*$ -map setting  $C^*[s] := C[s]$  since the neutral networks are constructed *a priori* in the graph of compatible sequences. We will assume that  $C^* = C$  and then the recursion

reads in the particular case of  $\Gamma_n[s] = C[s]$ :

$$f^{-1}(s_{r_0}) := C[s_0]$$

$$f^{-1}(s_{r_i}) := C[s_i] \setminus \bigcup_{j < i} [C[s_i] \cap C[s_j]]$$

It turns out that the random maps in secondary structures as defined above have a few large preimages and many small ones as reported in figure 2. This observation fits in the computational results about RNA sequence structure maps that exhibit a characteristic rank order function known as a *generalized Zipf's Law*:

$$\psi(i) = a(1 + i/b)^{-c},$$

as shown by extensive numerical calculations [13, 14]. Here  $a$  is a normalization constant,  $b$  is the number of frequent structures and  $c$  describes the power-law decay for rare structures.

### 4. Searching in Shape Space

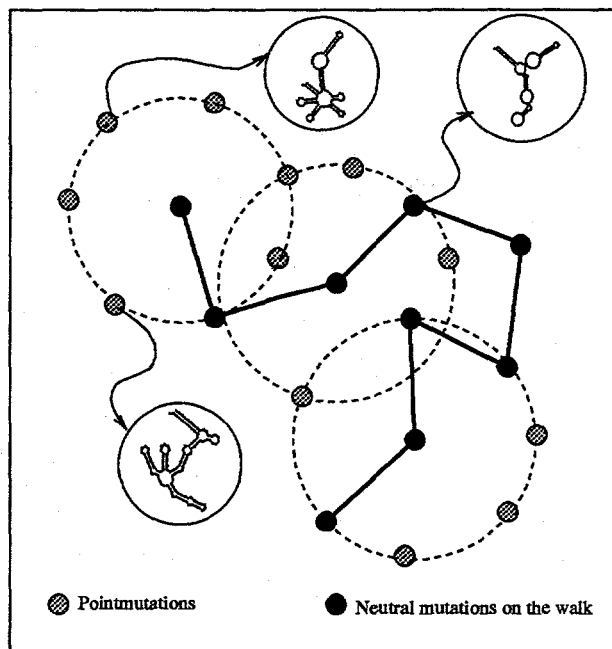


Fig. 3: A random walk on a neutral network. From each step (black) all point mutant sequences (grey) are formed and mapped into their corresponding secondary structures. Thereby, step by step, shape space is searched by w.r.t. one neutral network.

The time evolution of a population of asexually replicating molecules on a flat landscape can be described by a random walk. On more complex fitness landscapes a selective pressure drives the population towards sites of higher fitness. In the simplest case one may think of one or few fit structures. Then the populations searches

"along one neutral network" [4] until sequences of another network corresponding to a fitter structure are found (see fig. 3).

The key question [10] is: to what extent is the shape space searched by a population diffusing [11] on a fixed neutral network? Here the following result [12, 11] gives some insight, implying that *any* two neutral networks that are dense and have a giant component come close in sequence space and therefore allow for transitions from one network to the other [4]. Practically *all* other structures are found forming one error mutants with respect to a fixed neutral net. Explicitly the result reads

**Theorem 5. [Intersection-Theorem]** *Let  $\Pi$  be a non-empty pairing rule on  $\mathcal{A}$  and  $s$  and  $s'$  be arbitrary (non-empty) secondary structures. Then we have w.r.t.  $\Pi$*

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

The topology of neutral networks (following the predictions of random graph theory) plays therefore a crucial role in the optimization process. Only the existence of a giant component in the net and its density guarantee that any other net can be reached. In the tables shown in appendix A, we report the existence of those giant components.

Finally we discuss the so called *search capacity* of our mappings. This means to a fixed  $\lambda$  parameter we compute the size of the preimage of those secondary structures that are found by a random walk of the particular (fixed) neutral network. Passing the critical parameter  $\lambda = 0.5$  for density and connectivity, the shape space is searched effectively by a random walk on the corresponding neutral network. For the small sequence length considered in the computer experiments, it turns out that below  $\lambda = 0.5$  it is difficult to perform a random walk on the network lasting sufficiently many steps.

Fractions of neutral neighbors can be determined numerically by RNA folding based on sequences over an AUGC alphabets [3, 9, 8]. These fractions turn out to be characteristically above the critical values derived from random graph theory, and therefore indicate how well suited nature is for optimization on secondary structures.

## Covering Ability of Neutral Walks

Performed for Chainlength 25, binary Alphabet

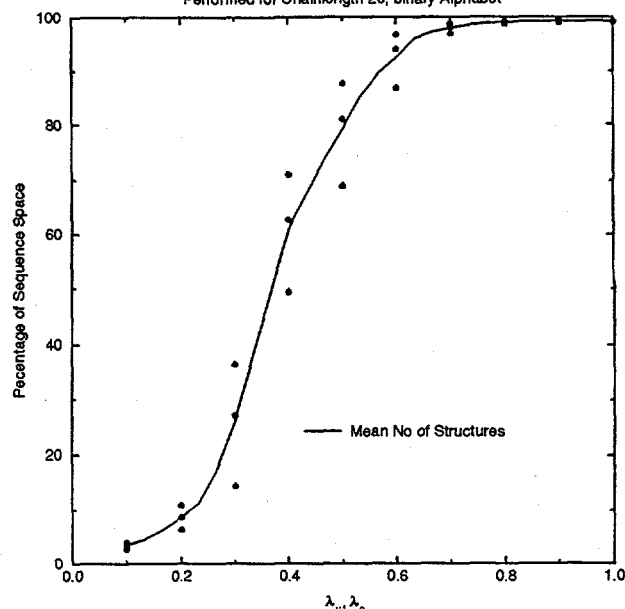


Fig. 4: We compute the percentage of the preimage of those secondary structures that are found by a random walk of one particular (fixed) neutral network. "Found" means that from each sequence realized in the walk we map all sequences of Hamming distance one and thereby obtain a stepwise increasing family of structures realized w.r.t. the random walk on the chosen network.

## Acknowledgments

We want to thank Mr. F. Haubensak for the hints and support he gave us to create the two columns per page.

## References

- [1] Ajtai, Komlós, and Szemerédi. Largest random component of a  $k$ -cube. *Combinatorica*, 2:1 – 7, 1982.
- [2] B. Bollobás. *Random Graphs*. ACADEMIC PRESS, 1985.
- [3] W. Fontana, P. F. Stadler, E. G. Bornberg-Bauer, T. Griesmacher, I. L. Hofacker, M. Tacker, P. Tarrazona, E. D. Weinberger, and P. Schuster. RNA folding and combinatorial landscapes. *Phys. Rev. E*, 47(3):2083 – 2099, March 1993.
- [4] C. V. Forst, C. Reidys and J. Weber. Evolutionary Dynamics and Optimization: Neutral Networks as Model-Landscapes for RNA Secondary-Structure Folding-Landscapes. *Advances in Artificial Life*, 1995.
- [5] W. Gruener, R. Giegerich, D. Strothmann, C. Reidys, J. Weber, I. L. Hofacker, P. F. Stadler, and

- P. Schuster. Analysis of rna sequence structure maps by exhaustive enumeration I. neutral networks. *Monatshefte fuer Chemie*, 127:\*-\*, 1996.
- [6] W. Gruener, R. Giegerich, D. Strothmann, C. Reidys, J. Weber, I. L. Hofacker, P. F. Stadler, and P. Schuster. Analysis of rna sequence structure maps by exhaustive enumeration II. structures of neutral networks and shape space covering. *Monatshefte fuer Chemie*, 127:\*-\*, 1996.
- [7] L. Harper. Minimal numberings and isoperimetric problems on cubes. *Theory of Graphs, International Symposium, Rome*, 1966.
- [8] I. L. Hofacker. *A Statistical Characterisation of the Sequence to Structure Mapping in RNA*. PhD thesis, University of Vienna, 1994.
- [9] I. L. Hofacker, W. Fontana, P. F. Stadler, S. Bonhoeffer, M. Tacker, and P. Schuster. Fast folding and comparison of RNA secondary structures. *Monatshefte f. Chemie*, 125(2):167-188, 1994.
- [10] M. Huynen, P. F. Stadler, and W. Fontana. Evolutionary dynamics of RNA and the neutral theory. *PNAS*, 1994. accepted.
- [11] C. Reidys. *Neutral Networks of RNA Secondary Structures*. PhD thesis, Friedrich Schiller Universität, Jena, May 1995.
- [12] C. Reidys, P. F. Stadler, and P. Schuster. Generic properties of combinatory maps and neutral networks of RNA secondary structures. *Bull. Math. Biol.*, 1995. submitted.
- [13] P. Schuster, W. Fontana, P. F. Stadler, and I. L. Hofacker. From sequences to shapes and back: A case study in RNA secondary structures. *Proc. Roy. Soc. (London) B*, 255:279-284, 1994.
- [14] M. Tacker, W. Fontana, P. Stadler, and P. Schuster. Statistics of RNA melting kinetics. *Eur. Biophys. J.*, 23(1):29 - 38, 1994.
- [15] M. S. Waterman. Secondary structure of single - stranded nucleic acids. *Studies on foundations and combinatorics, Advances in mathematics sup-*

## Appendix A

Rank	$ \Gamma $	Structure	$\mu_u$	$\mu_p$	Sequence of Components
1	1494359	.....(((.....)).)	0.823	0.854	1494359
2	1471814	.(((.....)).).....	0.795	0.858	1471810, 4×1
3	1180817	((.....)).....	0.599	0.487	1180665, 2, 150×1
4	1131781	..((.....)).)...	0.736	0.723	1131775, 6×1
5	895743	....((.....))....((...))	0.726	0.499	829839, 65886, 2, 16×1
6	811332	..((((.....)))..)	0.862	0.871	811332
7	769357	..((((.....)))((.....)).	0.852	0.825	769357
8	764731	..((.....(((.....))).)...	0.880	0.826	764731
9	676878	....((...)).((.....)).	0.608	0.501	676783, 95×1
10	577625	((.....)).(((.....))).	0.740	0.723	577616, 9×1

Tab. 1: Mapping parameters  $\lambda_u = \lambda_p = 0.9$ ,  $\mu_u, \mu_p$  are the percentages of unpaired and paired segments of sequences contained in the neutral net.

Rank	$ \Gamma $	Structure	$\mu_u$	$\mu_p$	Sequence of Components
1	3466927	.....(((.....)).).....	0.567	0.618	3466927
2	1339085	.....(((.....)).).....	0.409	0.306	1337547, 2×3, 30×2, 1472×1
3	718788	....((.....)).)...	0.514	0.576	718747, 2, 39×1
4	650290	....((((.....))).....	0.357	0.305	646820, 6, 5, 14×3, 165×2, 3087×1
5	642094	((.....(((.....))).)...	0.500	0.513	642025, 2, 67×1
6	606699	.(((.....)).).....	0.456	0.553	606464, 6×2, 223×1
7	596554	..((.....)).).....	0.230	0.237	570539, 10, 3×7, 7×6, 20×5, 56×4, 277×3, 1618×2, 21551×1
8	575245	.....((((.....)))..)	0.313	0.336	569428, 6×4, 25×3, 280×2, 5158×1
9	500107	.(((.....)).).....	0.456	0.393	499841, 5×2, 256×1
10	447051	.....(((.....)).).....	0.333	0.000	443329, 4, 15×3, 122×2, 3429×1

Tab. 2: Mapping parameters  $\lambda_u = \lambda_p = 0.6$  (see also caption of tab. 1)



Rank	$ \Gamma $	Structure	$\mu_u$	$\mu_p$	Sequence of Components
1	1604053	.....((.....)).....	0.387	0.481	1604009, $18 \times 2, 8 \times 1$
2	1011867	.....((.....)).....	0.279	0.349	1008832, $4, 6 \times 3, 169 \times 2, 2675 \times 1$
3	649324	..((.....)).....	0.209	0.000	327019, $301123, 2 \times 8, 2 \times 7, 4 \times 6,$ $9 \times 5, 53 \times 4, 215 \times 3, 1299 \times 2,$ $17628 \times 1$
4	511011	.....(((.....)))...	0.244	0.311	506282, $2 \times 6, 5, 8 \times 4, 41 \times 3,$ $274 \times 2, 4009 \times 1$
5	484477	..((...((...)).....))...	0.351	0.555	484319, $3 \times 11, 2 \times 10, 8, 2 \times 7,$ $6, 2 \times 2, 73 \times 1$
6	471135	.....(((.....)))..	0.270	0.219	359724, $107579, 9, 5, 8 \times 4,$ $29 \times 3, 206 \times 2, 3287 \times 1$
7	440221	.....((.....)).....	0.170	0.193	413333, $9, 8, 4 \times 7, 9 \times 6,$ $29 \times 5, 143 \times 4, 411 \times 3, 2060 \times 2,$ $20719 \times 1$
8	375080	....((...((.....))..))...	0.368	0.467	330454, $44442, 7 \times 7, 4 \times 6, 111 \times 1$
9	363837	.....(((.....))).....	0.272	0.181	274624, $85356, 6, 5, 8 \times 4,$ $41 \times 3, 212 \times 2, 3267 \times 1$
10	315083	((...((.....))..)).....	0.248	0.443	313124, $2 \times 7, 3 \times 5, 4, 13 \times 3,$ $97 \times 2, 1693 \times 1$

Tab. 3: Mapping parameters  $\lambda_u = \lambda_p = 0.4$  (see also caption of tab. 1)

## DISCLAIMER

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.