

Statistics of RNA Secondary Structures*

WALTER FONTANA,^{1,2} DANIELLE A. M. KONINGS,³ PETER F. STADLER,⁴ and PETER SCHUSTER^{2,4,5,†}

¹Theoretical Division, Los Alamos National Laboratory, New Mexico, USA 87545; ²Santa Fe Institute, New Mexico, USA 87501; ³Department of Molecular, Cellular and Developmental Biology, University of Colorado at Boulder, Boulder, Colorado, USA 80309; ⁴Institut für Theoretische Chemie der Universität Wien, Währingerstrasse 17, A-1090 Wien, Austria; and ⁵Institut für Molekulare Biotechnologie, Beutenbergstrasse 11, D O-690 Jena, Germany

SYNOPSIS

A statistical reference for RNA secondary structures with minimum free energies is computed by folding large ensembles of random RNA sequences. Four nucleotide alphabets are used: two binary alphabets, AU and GC, the biophysical AUGC and the synthetic GCXK alphabet. RNA secondary structures are made of structural elements, such as stacks, loops, joints, and free ends. Statistical properties of these elements are computed for small RNA molecules of chain lengths up to 100. The results of RNA structure statistics depend strongly on the particular alphabet chosen. The statistical reference is compared with the data derived from natural RNA molecules with similar base frequencies.

Secondary structures are represented as trees. Tree editing provides a quantitative measure for the distance d_t , between two structures. We compute a structure density surface as the conditional probability of two structures having distance t given that their sequences have distance h . This surface indicates that the vast majority of possible minimum free energy secondary structures occur within a fairly small neighborhood of any typical (random) sequence.

Correlation lengths for secondary structures in their tree representations are computed from probability densities. They are appropriate measures for the complexity of the sequence–structure relation. The correlation length also provides a quantitative estimate for the mean sensitivity of structures to point mutations. © 1993 John Wiley & Sons, Inc.

INTRODUCTION

A great variety of natural biopolymers were studied extensively by sequence analysis, the x-ray diffraction, and spectroscopic techniques. Their molecular structures are well known by now. In some cases the studies were extended to closely related variants differing in one or a few positions from the wild-type sequences. Despite the availability of very detailed information on many individual biomolecules, very little, if anything, is known about the statistics of

structural features of biopolymers. General information on the sensitivity of structures against changes in the sequences is also very rare. In the case of RNA, direct interest in the properties of molecules with random sequences is steadily growing: random RNA sequences are readily available by the current synthetic techniques and they are frequently used in applied molecular evolution.^{1–3} A straightforward strategy for gaining knowledge about molecules with random sequences is to compute and analyze the structures of statistical ensembles of biopolymers. The computation of three-dimensional structures is, however, highly time-consuming, and still encounters substantial theoretical and algorithmic difficulties.

The folding of RNA sequences into three-dimensional structures is decomposed into two steps:

Biopolymers, Vol. 33, 1389–1404 (1993)

© 1993 John Wiley & Sons, Inc.

CCC 0006-3525/93/091389-16

* Dedicated to Professor Manfred Eigen.

† To whom correspondence should be addressed.

1. folding of the string of bases into a secondary structure by base-pair formation, and
2. formation of the spatial structure by folding the secondary structure into a three-dimensional object.

The biophysical rationale for considering RNA secondary structures as a first crude approach to real RNA structures is based on several facts: RNA secondary structure formation covers the major part of the folding energy, RNA secondary structures can be used successfully to interpret RNA function and reactivity, and RNA secondary structures are conserved in evolution.

Secondary structures of RNA molecules, represented as the list of Watson-Crick (WC)-like base pairs, are much easier to predict than three-dimensional structures. These structures are mainly determined by the conventional base-pairing rules of RNA: ($G \equiv C$, $A \equiv U$, and $G \text{---} U$). In this paper we also consider the non-natural xanthine-2,6-diaminopyrimidine pair, $X \equiv K$, which was recently incorporated enzymatically into synthetic RNA and DNA molecules.⁴ Base-pairing and base-pair stacking energies are generally larger than those of other interactions involved in the formation of spatial structures and arrange the RNA structure in helical segments interrupted by various types of loop regions. In the most narrow sense the term "secondary structure" is used to refer to the major list of WC-like base pairs that constitute an unknotted structure that may be drawn as a planar graph. Thus, loop-loop interactions are not included in this definition. All the remaining WC-like base pairs (e.g., those occurring in pseudoknots) and other types of interactions are then referred to as "tertiary interactions" within this concept. Although this classification is in a sense arbitrary, in particular with reference to pseudoknot structures,^{5,6} it reflects the partitioning of the folding process into several stages, the first one being the nucleation of double-stranded helical regions. It is therefore meaningful to consider the folding into a secondary structure (even in this narrow sense) as a first step toward the formation of the full three-dimensional structure. The computational approach used here for the analysis of gross features of RNA secondary structures is based on this concept of "unknotted" secondary structures.

The predominantly used thermodynamic folding algorithm for RNA secondary structures was originally conceived by Zuker, Stiegler, and Sankoff,^{7,8} and is based on an application of dynamic programming to the RNA problem.⁹ It was primarily de-

signed to compute the minimum free energy structure, but derivative algorithms allow to obtain suboptimal foldings as well.¹⁰⁻¹² Alternatively, one may consider suboptimal foldings with the corresponding Boltzmann weights and compute partition functions for RNA secondary structures directly.¹³ The empirical parameters used in the folding algorithm were obtained some years ago.¹⁴ (For a more recent updated version, see Ref. 15.) Here we shall be dealing exclusively with minimum free energy secondary structures computed by a variant of the Zuker algorithm. The computer code was originally designed for fast folding as part of a simulation package for molecular evolution.^{16,17} The present version of the software package uses the updated version of the empirical parameter set,¹⁵ and includes a statistics program as well as tree editing routines. For the XK base pair we use the GC parameter set, which seems to come closest to the base pairing strength of the synthetic base pair.⁴

In this paper we are concerned with the statistical properties of RNA secondary structures and their dependence on various base alphabets: the natural four-letter alphabet (AUGC), a non-natural GCXK alphabet, and the two-letter alphabets (AU and GC). The GCXK alphabet provides an interesting case since it contains two base pairs of approximately equal strength, and it is free of complementarity violations ($G \text{---} U$ in the natural alphabet). Comparison of structures derived from different alphabets allows us to separate effects of different origin and helps us to understand the complex superposition of contributions from the size of the alphabet, from the pairing rules, and from the strength of base pairs in the natural RNAs.

The statistics of structural elements for small RNA molecules (with chain length $\nu < 100$) is presented in the next section. These data are considered as a statistical reference to be compared with data from natural sequences in the third section. In the section after that we present a quantitative measure for the distance between RNA secondary structures. Sequence-structure relations are viewed as *combinatorial maps* (CMs) from sequence space into a shape space. The notion of *structure density surfaces* (SDS) is introduced in the fifth section. The SDS casts statistical aspects of the relation between sequences and structures into a condensed form, and provides a tool to derive and calculate global properties for classes of RNA molecules. Autocorrelation functions and correlation lengths of structures finally characterize the sequence-structure relation by a single function or a single number, respectively. They provide a useful measure for the sensitivity of

RNA structures against point mutations (penultimate section).

STATISTICS OF ELEMENTS OF RNA SECONDARY STRUCTURES

The folding algorithm is a procedure that converts an RNA primary sequence, say

$$I_k = \{ \text{AUGCGUUGGACGUGCAGUCCAGUCAG} \\ \dots \text{AAACGC} \}$$

into a secondary structure $S_k = \mathcal{S}(I_k)$ where $\mathcal{S}(\cdot)$ stands for the folding algorithm, which computes a

unique structure for every sequence I_k . An example is shown in Figure 1. Many widely different sequences, however, may fold into the same secondary structure. This fact makes the *reverse folding problem*—the problem to determine all sequences which fold into a given secondary structure—a particularly hard task.

A secondary structure is viewed conventionally as a combination of *structure elements*, which fall into seven classes:

1. *stacks* (S), which are double-helical regions consisting of stacked base pairs;
2. *hairpin loops* (H), representing stretches of unpaired bases that close terminal stacks;

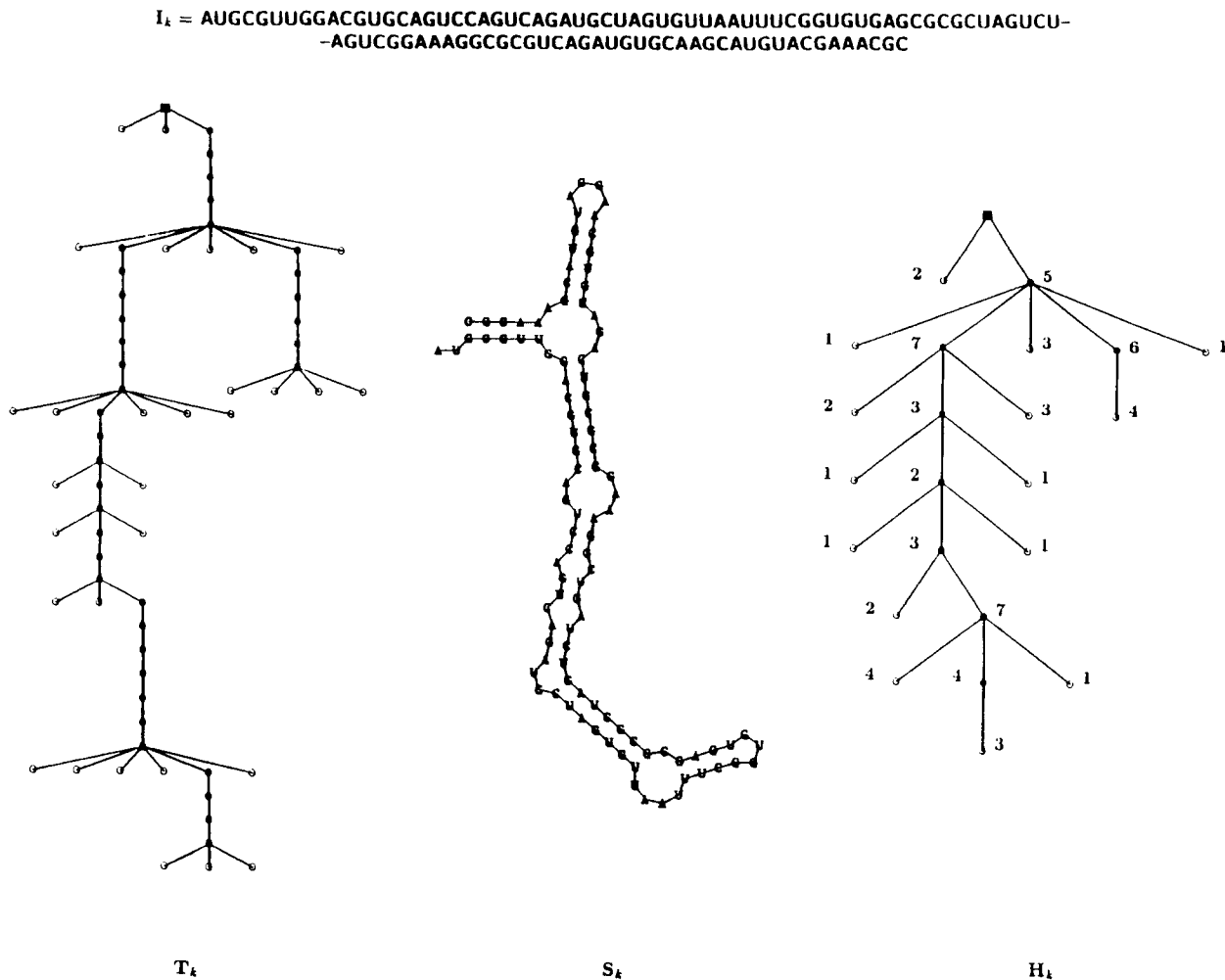


Figure 1. An example for folding an RNA sequence I_k into a secondary structure S_k and its conversion into a (full) tree T_k . In this tree representation, single-stranded bases are shown as open circles (\circ) and base pairs as a full circles (\bullet), respectively. A root (\blacksquare), not corresponding to a physical unit of the RNA, is added. The full tree T_k is transformed into a *homeomorphically irreducible tree* (HIT) H_k by assigning a weight w to every node of the HIT, which counts the number of nodes contracted into a single one.

3. *bulges* (B), which connect two stacks by an unpaired stretch;
4. *internal loops* (I), joining two stacks with two single-stranded stretches;
5. *multiloops* (M), consisting of several single-stranded stretches that connect more than two stacks;
6. *joints* (J), which are stretches of unpaired bases joining freely movable substructures; and
7. *free ends* (E).

Stacks stabilize secondary structures whereas loops have a destabilizing effect whose magnitude depends on the size of the loop. Isolated single base pairs are considered as stacks as well. The degree of a loop is the number of stacks connected to it. It is often useful to lump loops of all degrees together into one class and to consider, for example, the total number of loops

$$n_L = n_H + n_B + n_I + n_M \quad (1)$$

which must be identical to the number of stacks, $n_L = n_S$. Nucleotides in joints and free ends are often termed *external* bases.

We generated a statistical reference for RNA secondary structures on sequences of length $20 \leq \nu \leq 100$ built from various alphabets of size κ ($\kappa = 2$: AU, GC; and $\kappa = 4$: AUGC, GCXK). Equal probability of digits was assumed and hence the base frequencies of random sequences lie around the most probable distributions, (0.25, 0.25, 0.25, 0.25) or (0.5, 0.5), respectively. In this section we report on the average occurrence of secondary structure elements.

Base Pairs. The mean number of base pairs (n_{BP}) increases linearly with the chain length ν for sufficiently long sequences (Figure 2). Deviations at small chain lengths ($\nu < 50$) are found with AU, AUGC, and GCXK sequences. The influence of the alphabet is readily interpreted by considering the *stickiness* P of the sequences, which is understood as the probability that two arbitrarily chosen bases can form a base pair. For the uniform base compositions one has $P_{AU} = P_{GC} = 0.5$, $P_{AUGC} = 0.375$, and $P_{GCXK} = 0.25$. As expected, and as seen in Figure 2, the pure GC sequences are leading with respect in the number of base pairs since they have the highest possible stickiness and form the strongest base pairs. Destabilization by loops is more readily compensated by GC pairs than by AU pairs, and hence AU sequences form fewer base pairs on the average than GC sequences. Sequences derived from four-letter

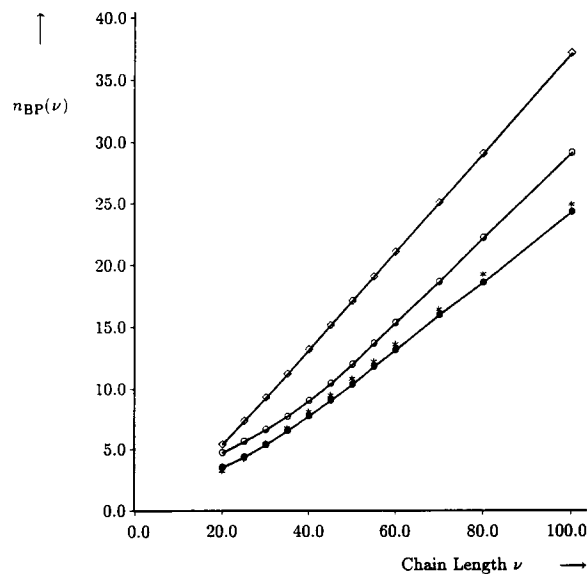


Figure 2. The mean number of base pairs (n_{BP}) as a function of the chain length ν . Values are shown for binary GC sequences (\diamond), for binary AU sequences (\circ), for four-letter GCXK sequences with GC parameters ($*$), and for natural AUGC sequences (\bullet). Depending on the chain length ν , the mean values are computed from samples of 50,000 ($\nu = 20$) to 500,000 sequences ($\nu = 100$) that were sampled by the technique described in the section on comparison with natural sequences. Unstable structures, i.e., structures with non-negative free energies ($f \geq 0$) are not considered for structure statistics.

alphabets are less sticky and form still fewer base pairs. That AUGC and GCXK sequences have almost the same mean number of base pairs is fortuitous: the former are more sticky, the latter form stronger base pairs and the two effects cancel by accident.

Loops and Stacks. The mean number of loops n_L per structure, which is identical to the mean number of stacks n_S , also scales linearly with chain length ν (Figure 3). Weaker base pairing implies that closing of loops is more difficult and hence the structures derived from AUGC or AU sequences have fewer loops than their GCXK or GC counterparts. Sequences with lower stickiness values have on the average more loops than stickier sequences. The effect of base pair-strength is apparently stronger than that of stickiness.

Loop Degree. The (branching) degree of a loop is the number of stacks that are attached to the loop. The average degree of loops (Figure 4) is in the range $1 < n_{LD} < 2$. It converges to a constant value with

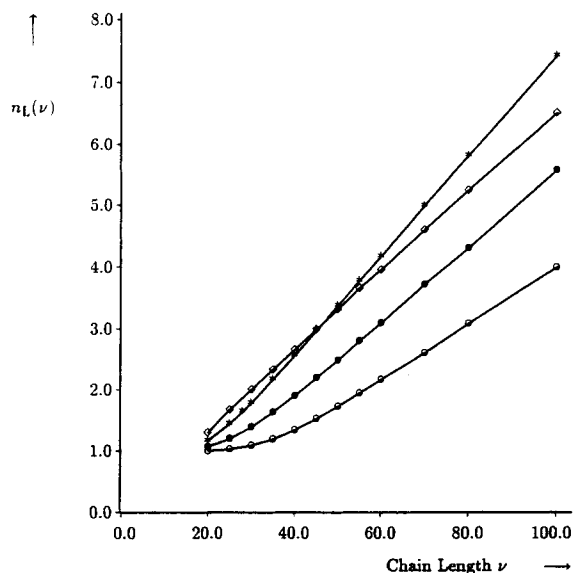


Figure 3. The mean number of loops (n_L), which is identical to the mean number of stacks (n_S), as a function of the chain length ν . Values are shown for binary GC sequences (\diamond), for binary AU sequences (\circ), for four-letter GCXK sequences with GC parameters ($*$), and for natural AUGC sequences (\bullet). Computations were performed as described in Figure 2.

increasing chain length ν . Structures derived from sequences with strong base pairs (GC, GCXK) have more higher order branches than those obtained from AUGC and AU sequences.

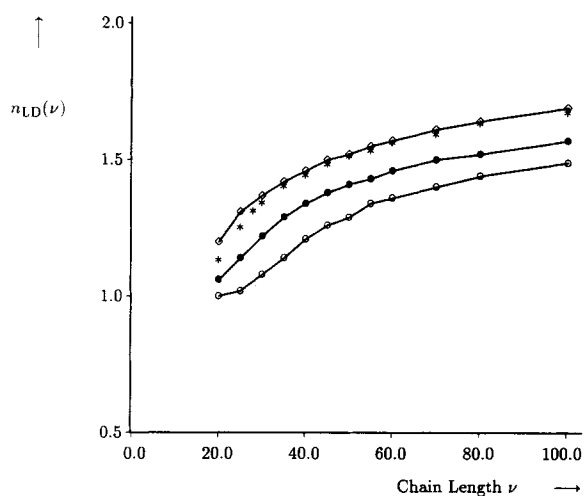


Figure 4. The mean degree of loops (n_{LD}) as a function of the chain length ν . Values are shown for binary GC sequences (\diamond), for binary AU sequences (\circ), for four-letter GCXK sequences with GC parameters ($*$), and for natural AUGC sequences (\bullet). Computations were performed as described in Figure 2.

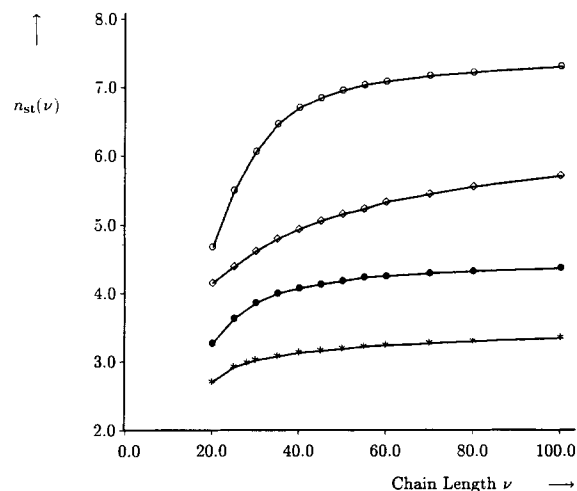


Figure 5. The mean number of base pairs in one stack (n_{st}) as a function of the chain length ν . Values are shown for binary GC sequences (\diamond), for binary AU sequences (\circ), for four-letter GCXK sequences with GC parameters ($*$), and for natural AUGC sequences (\bullet). Computations were performed as described in Figure 2.

Stack Sizes and Loop Sizes. Mean stack sizes (n_{st}) and mean loop sizes (n_{lp}) converge to almost constant values at fairly small chain lengths ($\nu \approx 50$) as shown in Figures 5 and 6. Apparently the convergence of n_{st} and n_{lp} in structures of GC sequences is slower than in structures from the other three alphabets. Stickiness is important for stack sizes:

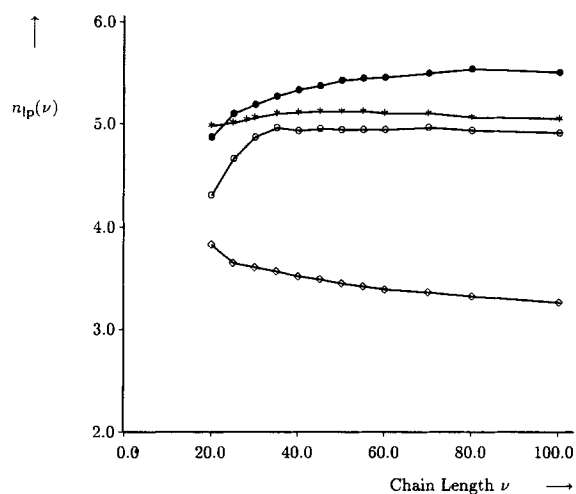


Figure 6. The mean number of bases in one loop (n_{lp}) as a function of the chain length ν . Values are shown for binary GC sequences (\diamond), for binary AU sequences (\circ), for four-letter GCXK sequences with GC parameters ($*$), and for natural AUGC sequences (\bullet). Computations were performed as described in Figure 2.

on the average the two-letter sequences form longer stacks than AUGC sequences. The mean stack lengths of GCXK sequences are shortest. Weak base pairing makes nucleation of stacks more difficult. Stacks are, therefore, on average longer in the AU case than in the GC case. Large loops are favored by both low stickiness and weak base pairs since a weak closing pair is more likely to be destabilized by a loop than a strong one. Thus, on average, GC-sequences form the smallest loops.

Components. Components are substructures that are connected by joints. The mean number of components n_C shows a characteristic lag phase before it starts off to increase with increasing chain length ν . This lag phase reflects the fact that a certain minimum chain length is required in order to form local substructures. The lag phase is more pronounced in structures built from alphabets with weaker base pairs (AUGC, AU). The increase of n_C with ν is much stronger in the case of the four-letter alphabets. The data shown in Figure 7 suggest that this increase is roughly linear. In order to be able to study large ensembles of longer sequences, the folding algorithm was adapted to a parallel computer.¹⁸ These computations have shown, however, that the chain length dependence of the number of components is not yet linear at chain lengths around $\nu \approx 500$. GC sequences are special in a way since there the mean

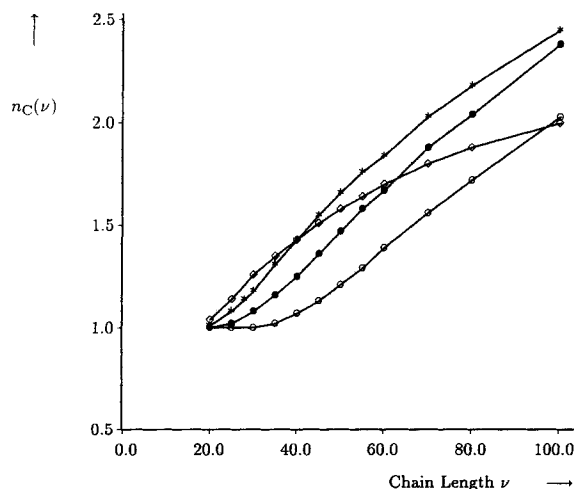


Figure 7. The mean number of components (n_C) connected by $n_J = n_C - 1$ joints as a function of the chain length ν . Values are shown for binary GC sequences (\diamond), for binary AU sequences (\circ), for four-letter GCXK sequences with GC parameters ($*$), and for natural AUGC sequences (\bullet). Computations were performed as described in Figure 2.

number of components converges to a value in the range $2.75 \leq n_C \leq 3$.

External Digits. The mean number of external bases (n_{ext}) shows a rather complex increase with the chain length ν (Figure 8). It contains two contributions, the number of terminal free bases and the number of free bases in joints. The number of joints as we saw before shows a nonlinear increase with the chain length in the range in question (Figure 7). The strength of base pairs has a strong influence on n_{ext} : weak pairing implies larger mean numbers of external digits. Interestingly, n_{ext} is almost constant for GC sequences where the number of components converges to a constant value for longer sequences ($\nu > 300$) too.¹⁸

So far we focused only on mean values. Resolved distributions for stack sizes, loop sizes, and loop degrees will be presented in the next section, where data from the statistical reference will be compared with those obtained from natural sources with similar base compositions.

COMPARISON WITH NATURAL SEQUENCES

In order to compare data from the statistical reference for RNA secondary structures with those from natural RNA sequences, samples with base distributions as close as possible to the uniform distribution, (n_A, n_U, n_G, n_C) $\approx \nu(0.25, 0.25, 0.25, 0.25)$, are required. In addition, we expect best agreement when the biological function of the RNA demands as little secondary structure as possible. A sample that meets both requirements consists of 12 full mature mRNA molecules, i.e., with the introns removed, coding for β -globin molecules* from different animals. Within this sample the chain lengths vary from 534 to 627, the sequence that deviates most strongly from the uniform distribution comes from *Xenopus laevis*: (n_A, n_U, n_G, n_C) $\approx \nu(0.29, 0.26, 0.21, 0.24)$. The 12 sequences were folded and the five structures with lowest free energies were considered for statistical analysis (the sample thus contains 60 structures). The five structures span a energy band of about 1–2% of the absolute free energy of the optimal structures. Some mean values are shown in Table I together with the statistical data for the

* Genbank names: gothbbaa, hsbgl1, hsdgl1, hsggl4, lebglob, mushbbmaj, ptggglog, rabhbba, rabbbb3, ratglbr, xebbeta, and xlbgl1r.

AUGC, the AU, and the GC alphabet. The mean number of base pairs per 100 nucleotides of the β -globin mRNA sequences is 31. It is slightly larger than the corresponding quantity derived from the statistical reference at approximately the same chain length, $\nu = 500$ ¹⁸: 29.0 base pairs. The β -globin mRNAs are slightly more stable than the statistical reference in the sense that they have lower free energies in their minimum free energy structures. The stability argument has been used previously¹⁹ in sequence and structure comparisons of viroids and virusoids that, as commonly assumed, have been selected for stability by evolution. There the thermodynamically calculated structures from natural sequences form indeed more base pairs than the corresponding sample of random sequences, for example, 35 base pairs per 100 nucleotides in the CSV viroid RNA as compared to 29.5 in the random standard. The mean loop size (Table I) is somewhat smaller in the β -globin mRNA sample than in the reference. The mean branching degrees of loops and the mean stack sizes of the two samples agree well.

In addition to β -globin mRNAs, two samples of natural RNA molecules from other sources were considered as well: 14 eubacterial 16s rRNAs[†] and 8 mitochondrial 16s rRNAs.[‡] These two samples fulfill the two conditions for a meaningful comparison with the statistical reference less well: the ribosomal RNAs show substantially larger deviation from the uniform base distribution than the β -globin mRNAs and their secondary structures are certainly important for biological function. For statistical analysis we again choose the minimum free energy structure together with the four most stable sub-optimal folding patterns (which were again in the 1–2% energy range, as with the β -globin mRNAs). As far as the mean values are concerned, the eubacterial rRNA sample fits the reference data as well as the mRNA sample. There are, however, significant differences with the mitochondrial rRNAs: the preference for larger loops is remarkable, and it is reflected by the smaller number of base pairs. Despite rather small sample sizes, we also computed statistical data for phylogenetically derived secondary structures of rRNAs.²⁰ In order to be able to compare with computed minimum free energy structures, only Watson–Crick and GU base pairs were considered. Differences between computed and

phylogenetic structures are quite moderate. Deviations from the statistical reference are smaller and larger depending on the quantity considered, but they point in the same directions as observed with the computed mitochondrial rRNA sample: fewer base pairs, shorter stacks, and larger loops.

In order to make the comparison with the experimental data more precise, we computed probability densities for stack sizes, loop sizes, and loop branching degrees. The results are shown in Figures 9–11. For stack size frequencies, the agreement between the β -globin mRNAs and the statistical reference from AUGC sequences is very good. The data computed from the other two samples from rRNAs fit the curve for the reference not nearly as well.

The distribution of loop sizes for natural sequences are compared with the statistical reference (AUGC, $\nu = 500$) in Figure 10. In essence, the results are the same as with the probability distribution for stack sizes: the data from mRNAs of β -globins fit the curves computed for the statistical reference better than the points obtained from the rRNAs. In detail, however, the agreement between the β -globin mRNAs and the random RNA sequences is not as good as for the stack sizes. The data for natural sequences show much larger deviations. Natural sequences have significantly more loops of size 1 and less loops of size 3 than the reference. The samples

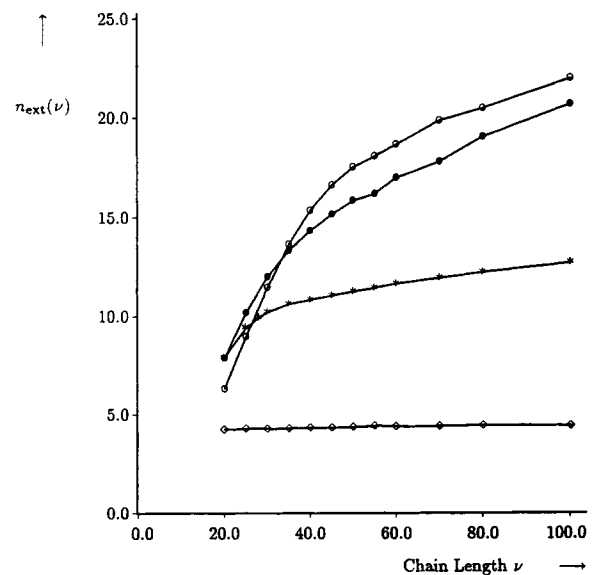


Figure 8. The mean number of unpaired external bases (n_{ext}) as a function of the chain length ν . Values are shown for binary GC sequences (\diamond), for binary AU sequences (\circ), for four-letter GCXK sequences with GC parameters ($*$), and for natural AUGC sequences (\bullet). Computations were performed as described in Figure 2.

[†] Genbank names: anlmttgrg, bovmt, ceumtfvla, frgmtcr12s, gotmttgrg, hummtgc, hyrmtfvla, mmumtfvla, musmt, odomtfvla, palmtcg, ratmtrgpd, trgmttgrg, and xelmtrrza.

[‡] Genbank names: bacgrrrnb, deirgda, hclrgda, mpocpcg, m27040, prirrgda, and stmrrnb.

Table I Comparison of Mean Number of Base Pairs, Mean Stack Size, Mean Loop Size, and Mean Branching Degree of Loops Between the Statistical Reference and Samples of Natural RNA Sequences

Source	\bar{n}_{BP}/ν	\bar{n}_{st}	\bar{n}_{lp}	\bar{n}_{LD}
Statistical reference of RNA sequences ^a				
AUGC	0.290	4.57	5.42	1.82
AU	0.354	7.66	4.69	1.78
GC	0.403	6.46	2.98	1.92
Natural RNA sequences, minimum free energy structures ^b				
β -Globin mRNAs	0.31	4.49	4.42	1.89
Mitochondrial rRNAs	0.26	4.44	6.53	1.74
Eubacterial rRNAs	0.33	4.59	4.62	1.92
Natural RNA sequences, phylogenetic structures ^c				
Mitochondrial rRNAs	0.24	3.76	6.00	1.90
Eubacterial rRNAs	0.28	4.35	5.81	1.93
All rRNAs	0.27	4.06	5.88	1.93

^a A sample of about 50,000 structures from sequences of chain length $\nu = 500$ with different base-pairing alphabets computed according to the section on landscapes, combinatory maps, and density surfaces.¹⁸

^b Sample sizes: 12 β -globin mRNAs (see asterisked footnote in text), 14 mitochondrial 16S rRNAs (see daggered footnote in text), and 8 eubacterial 16S rRNAs (see double-daggered footnote in text). The sample contains the minimum energy structure together with the four suboptimal foldings of lowest free energies.

^c Sample sizes: 2 eubacterial and 2 mitochondrial phylogenetic structures.¹⁹ The row "all" refers to these 4 structures and, in addition, to 2 archebacterial, 2 eukaryotic, and 1 chloroplast phylogenetic rRNA structures.

of natural sequences are inevitably small and hence the deviations at higher loop sizes need not be statistically significant. Apparently, loop size is a more sensitive parameter than stack size.

Frequency counts for the branching degree of loops (Figure 11) show excellent agreement between β -globin mRNAs and random sequences. Substantial deviations are observed only with the two samples derived from mitochondrial and eubacterial 16S rRNAs.

Phylogenetically derived structures of RNA molecules are considered biologically more significant than minimum free energy structures. For the purpose of comparison we also analyzed a sample of phylogenetically derived 16S rRNA structures²⁰ and included the data for the stack size distribution in Figure 9. We observe much larger scatter from the statistical reference than for the data points from

minimum free energy structures. One feature seems to be of special interest: the phylogenetic 16S rRNA structures show a significant local peak for a stack size of 7 base pairs which does neither exist in the minimum free energy samples nor in the phylogenetic data for 23S rRNAs (not shown here). This is a clear hint that there is an unusual favoring of stacks with 7 base pairs in the 16S rRNAs, and accordingly these stacks may have a specific function. Loop size and loop degree distributions of phylogenetic structures also show larger scatter than the minimum free energy samples of the natural sequences.

The results obtained in the comparison of secondary structures of natural sequences with the statistical reference suggest the proposition that this procedure as a tool for testing biological function. If a sufficiently large sample of natural sequences shows detailed agreement with the statistical reference of the same base composition, then it is very unlikely that biological function imposes a severe constraint on structures. Since the same folding algorithm is used in the computation of the reference and the sample to be compared, the empirical pa-

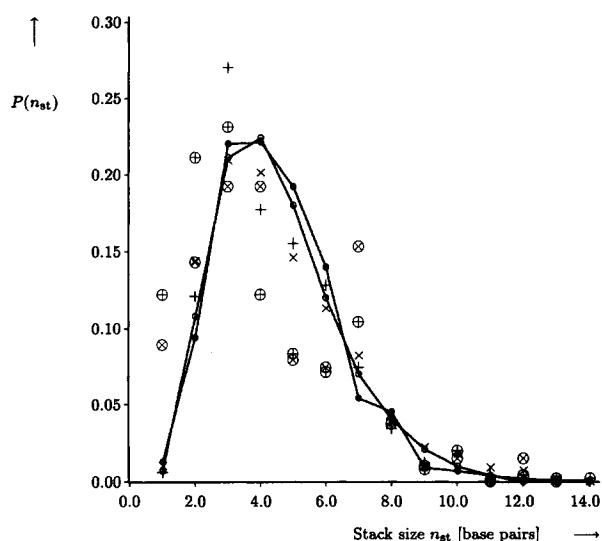


Figure 9. Probability densities $P(n_{st})$ of stack sizes in natural RNA sequences compared with the statistical reference of chain lengths $\nu = 500$. One curve shows the data computed from minimum free energy structures for mRNAs of β -globins (\bullet), and individual points are given for eubacterial 16S rRNAs (\times) and for mitochondrial 16S rRNAs ($+$). Further points are given for the phylogenetically derived secondary structures¹ of the eubacterial (\otimes) and the mitochondrial (\oplus) rRNAs. The second curve refers to the statistical reference built from the AUGC alphabet with chain lengths $\nu = 500$ ¹⁸ (\circ).

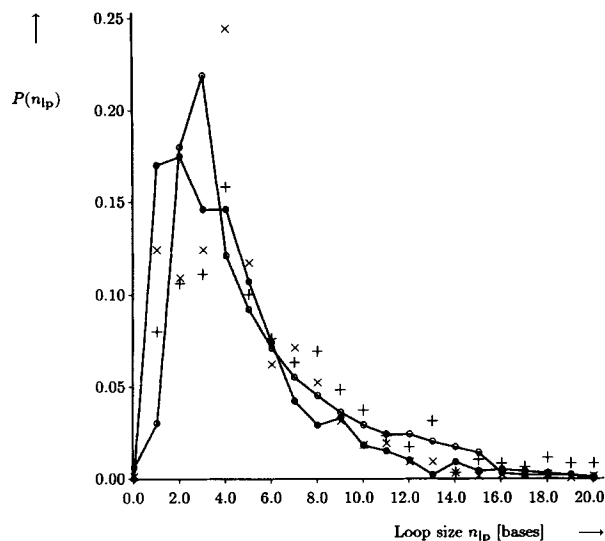


Figure 10. Probability densities $P(n_{lp})$ of loop sizes in natural RNA sequences compared with the statistical reference of chain lengths $\nu = 500$. One curve shows the data computed from minimum free energy structures for mRNAs of β -globins (\bullet), and individual points are given for eubacterial 16s rRNAs (\times) and for mitochondrial 16S rRNAs ($+$). The second curve refers to the statistical reference built from the AUGC alphabet with chain lengths $\nu = 500$ ¹⁸ (\circ).

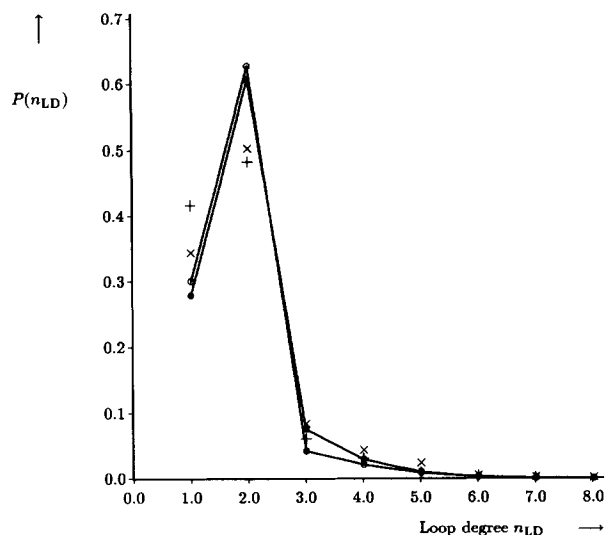


Figure 11. Probability densities $P(n_{LD})$ of degrees of branching in the loops of natural RNA sequences compared with the statistical reference of chain lengths $\nu = 500$. One curve shows the data computed from minimum free energy structures for mRNAs of β -globins (\bullet), and individual points are given for eubacterial 16s rRNAs (\times) and for mitochondrial 16s rRNAs ($+$). The second curve refers to the statistical reference built from the AUGC RNA sequences built from the AUGC alphabet with chain lengths $\nu = 500$ ¹⁸ (\circ).

rameters and the particular folding routine chosen have only minor influence on the result. The comparison is done on the level of ensembles and hence individual errors are fairly unimportant.

DISTANCES BETWEEN STRUCTURES

RNA secondary structures can be represented as trees.^{8,21,22} A secondary structure S_k is converted one to one into a tree T_k by assigning an internal node to each base pair and a leaf node to each unpaired digit (Figure 1). The conversion starts with a root that does not correspond to a physical unit of the RNA molecule. It is introduced to prevent the formation of a tree forest for RNA structures with ex-

Table II Interconversion of Secondary Structures and Trees

Accessibility of bases and base pairs⁸:

The bases are labeled with their position i in the sequence. Consider a base pair $i - j$ with $i < j$ in the secondary structure graph. A base r is said to be *accessible* from the base pair $i - j$, if there is no base pair $x - y$ such that $i < x < r < y < j$. A base pair $r - s$ is said to be accessible if r and s are accessible from the base pair $i - j$.

A secondary structure graph (SSG) is converted into a tree graph (TG) by the following procedure:

1. Assign to each unpaired base a leaf node, and to each base pair (two SSG nodes) one internal TG node.
2. The parent of each TG node is the internal node corresponding to that accessible base pair $k - l$ ($k < l$) that contains the smallest k . The parent of nodes corresponding to external bases and external base pairs is a virtual node representing the root of the tree.
3. Connect each node to its parent.

A TG is converted into a SSG by the following procedure:

1. Replace each internal TG node by a connected pair of SSG nodes (base pair). The TG edges are inherited by the left SSG node of the pair.
2. For any TG node with a left sibling replace the edge to its parent with an edge to its left sibling.
3. An SSG node that is base paired has three edges (two from the backbone, and one from the pairing); otherwise it has two. 5' and 3' ends have one less. Complete the SSG by inserting all missing edges at a node as connections to the corresponding parent in the TG, proceeding from deep to shallow levels.

ternal elements. For details of the interconversion of secondary structures and trees, see Table II.

As shown in Figure 1, the trees T_k can be rewritten as HITs, which will be denoted by H_k . The apparently simpler tree structure of the HIT is compensated by the assignment of weights (w) to the internal nodes and leaves. A weight reflects the number of nodes or leaves in the full tree T_k , which are lumped into a single node or leaf in the HIT representation. The transformation from the full tree to the HIT retains complete information on the structure. Secondary structure, full tree as well as HIT, are equivalent:

$$I_k \Rightarrow S_k \Leftrightarrow T_k \Leftrightarrow H_k \quad (2)$$

Tree editing induces a distance in the space of trees, and hence also in the space of RNA secondary structures. A tree is transformed into another tree by a series of editing operations with predefined costs.²¹⁻²⁴ The distance between two trees,

$$d(T_j, T_k) = d_t(j, k) \quad (3)$$

is the smallest sum of the costs along an editing path. It can be shown that this tree distance forms a metric in the space of trees.²⁴ The parameters used

in our tree editing are summarized in Table III. The editing operations preserve the relative traversal order of the tree nodes. Tree editing can therefore be viewed as a generalization of sequence alignment. In fact, for trees that consist solely of leaves, tree editing becomes the standard sequence alignment.

Using the parameters of Table III for editing operations on weighted trees, distances between HITs can be computed as well:

$$d(H_j, H_k) = d_t^{(\text{HIT})}(j, k) \quad (4)$$

Tree distances between full trees and HITs fulfill the general relation

$$d_t \leq d_t^{(\text{HIT})} \quad (5)$$

A proof of the inequality (5) will be given elsewhere.

An alternative graphical method for the comparison of RNA secondary structures²⁵⁻²⁷ encodes secondary structures as linear strings with balanced parentheses representing the base pairs, and some other symbol coding for unpaired positions. Distances are computed by direct end-to-end alignment of the strings.

Tree representations in full resolution make it often difficult to focus on the major structural fea-

Table III Cost Table for Tree Edit Operations in the Full Tree Representation and in the HIT Representation and the Edit Cost Parameters Used in this Paper^a

Edit Operation	Symbolic Notation	Edit Parameters	
		Symbol	Cost
Full tree edit operations			
Single base insertion	$0 \rightarrow \circ$	δ_{\circ}	1
Base pair insertion	$0 \rightarrow \bullet$	δ_{\bullet}	2
Single base deletion	$\circ \rightarrow 0$	δ_{\circ}	1
Base pair deletion	$\bullet \rightarrow 0$	δ_{\bullet}	2
Relabel		ρ	0
Substitution	$\bullet \leftrightarrow \circ$	σ	1
HIT edit operations			
Insertion of single bases	$0 \rightarrow \circ w$	$w \cdot \delta_{\circ}$	w
Insertion of base pairs	$0 \rightarrow \bullet w$	$w \cdot \delta_{\bullet}$	$2w$
Deletion of single bases	$w \circ \rightarrow 0$	$w \cdot \delta_{\circ}$	w
Deletion of base pairs	$w \bullet \rightarrow 0$	$w \cdot \delta_{\bullet}$	$2w$
Relabel of single bases	$w \circ \leftrightarrow \circ v$	$ w - v \cdot \delta_{\circ}$	$ w - v $
Relabel of base pairs	$w \bullet \leftrightarrow \bullet v$	$ w - v \cdot \delta_{\bullet}$	$2 w - v $
Substitution	$w \bullet \leftrightarrow \circ v$	$\sigma \min\{w, v\} + w - v \cdot \begin{cases} \delta_{\bullet} & \text{if } w > v \\ \delta_{\circ} & \text{if } w < v \end{cases}$	

^a Weights in the HIT representation are denoted by w and v .

tures of RNA molecules since they are often overloaded with irrelevant details. Coarse-grained tree representations were invented previously to solve this problem.^{21,22} They are based on plausible ad hoc assumptions. The HITs introduced here are particularly well suited for unambiguous coarse-graining. The weights w allow straightforward definition of a resolution parameter: all nodes up to some predefined value, $0 < w \leq w_r$, are omitted, and the secondary structure corresponding to the coarse-grained tree thus contains only structural features above a certain size.

LANDSCAPES, COMBINATORY MAPS, AND DENSITY SURFACES

The notion of *landscapes* was introduced into biology by the pioneering paper of Sewall Wright.²⁸ He visualized evolution as an optimization process on a *fitness landscape*. Since fitness is an exceedingly complex quantity, and since the occurrence of unconstrained optimization is at least debatable, Wright's metaphor turned out to be no more than a heuristic principle. The landscape concept in biology saw a recent revival in model studies of evolution and coevolution,²⁹⁻³¹ in theoretical immunology,^{32,33} and in studies based on RNA structures and properties.^{16,17,34}

Comparison of RNA sequences is standard in molecular evolution. The distance between two sequences of the same length ν is given by the Hamming metric $d_h(k, l)$,^{35,36} which counts the number of different digits in the two end-to-end aligned sequences I_k and I_l . The set of all sequences of a given length ν is of combinatorial complexity: an alphabet with κ letters yields κ^ν different sequences. Embedding this set in Euclidean space such that pairs of sequences with Hamming distance $d_h = 1$ are closest neighbors yields the so-called *sequence space*^{35,37} for which d_h is a metric.³⁶ The sequence space of binary sequences ($\kappa = 2$) is a hypercube of dimension ν . In the case of four-letter alphabets ($\kappa = 4$) the sequence space is more complex and can be understood as a union of two hypercubes with the appropriate connections. An earlier variant of sequence space is the *protein space* introduced by John Maynard-Smith.³⁸

A landscape is obtained by assigning a scalar quantity, for example, a fitness value, a free energy, or a rate constant, to every point in sequence space. Thus it is a mapping of the sequence space into the real numbers:

$$\mathcal{L}: (X, d_h) \Rightarrow (\mathbb{R}^1) \quad (6)$$

In order to point to the combinatorial complexity of the support, we refer to such an object as a *combinatory landscape* (CL).

Secondary structures—or trees—are also elements from a space of combinatorial complexity with the tree distance $d_t(k, l)$ as a metric. The notion of a *shape space*, originally conceived for antibody-antigen recognition in theoretical immunology,^{39,40} appears to be a useful concept for RNA secondary structures as well. The process of folding assigns to an element I_k of the sequence space a tree representing a secondary structure S_k . RNA folding, therefore, is a mapping Φ from a sequence space X into a shape space Y :

$$\Phi: (X, d_h) \Rightarrow (Y, d_t) \quad (7)$$

The elements of both spaces are discrete structures of combinatorial complexity. Both spaces are endowed with appropriate metrics. We refer to a mapping of this kind as a *combinatory map* (CM).

As a generic tool for the study of CMs we propose the use of a two-dimensional probability density $P(.,h)$. A tree distance density surface, for example, expresses the joint probability of two sequences of length ν , I_i and I_j , to have Hamming distance $d_h(i, j) = h$, as well as folding into secondary structures with tree distance $d_t(i, j) = t$. The surface $P(t, h)$ is biased along the Hamming distance axis with the distribution

$$p(h) = \binom{\nu}{h} (\kappa - 1)^{h} \kappa^{\nu-h} \quad (8)$$

since there are $\kappa^\nu \cdot p(h)$ sequences at Hamming distance h from a given sequence. In order to compensate for the bias, we consider the conditional probability of finding a tree distance t between two structures whose sequences have a Hamming distance h :

$$P(t|h) = \frac{P(t, h)}{p(h)} \approx \frac{n(t, h)}{\sum_{t=0}^{t_{\max}} n(t, h)} \quad (9)$$

As indicated in Eq. (9), the conditional probability $P(t|h)$ is computed by sampling tree distances of pairs of sequences with Hamming distance h . By $n(t, h)$ we denote the number of pairs of sequences in the sample that have Hamming distance h and tree distance t ; t_{\max} is used for the maximal tree distance. The total sample size is $N = \sum_{t=0}^{t_{\max}} \sum_{h=0}^{\nu} n(t, h)$,

and from the definition of conditional probabilities

it follows that $\sum_{t=0}^{t_{\max}} P(t|h) = 1$ for each h .

We use a sampling technique that directly computes the conditional probabilities. It may be characterized as uniform sample statistics and proceeds as follows:

1. we choose a reference sequence at random;
2. we sample exactly l sequences at each Hamming distance $h = 1, 2, \dots, \nu$ from the reference sequence;
3. we fold these sequences into secondary structures and compute tree distances relative to the structure of the reference sequence;
4. we bin them in (t, h) boxes by counting numbers of instances; and
5. continue with item 1 until convergence, or some desired accuracy, has been achieved after, say, r reference points.

The sampling procedure directly corrects for the bias of the binomial distribution. By definition we have

$$\sum_{t=0}^{t_{\max}} n(t, h) = r \cdot l = \text{const} \quad (10)$$

The sample size now is $N = r \cdot l \cdot \nu$. For binary sequences a slight modification of this procedure is required since we have only a single sequence at Hamming distance ν .

The tree distance probability surfaces represent the average distribution of tree distances from a reference sequence. In the examples shown in Figure 12 only 1000 reference sequences and 10^6 data points were already sufficient to yield a roughly constant surface. Convergence of this sampling technique is remarkably fast. Both surfaces from Figure 12 show an overall shape and, superimposed upon it, rich and bizarre-looking details. The overall shape of density surfaces for binary sequences (e.g., the GC surface in Figure 12) is horseshoe-like, and exhibits mirror-plane symmetry. This symmetry is only approximate. In the binary case a sequence at distance ν from the reference has a nucleotide arrangement that is clearly compatible with the secondary structure of the reference. As a matter of fact, such a sequence indeed folds into the reference structure or into a quite related one. Since for binary alphabets there is only one sequence at Hamming distance ν from a reference, this symmetry shows up. In contrast, AUGC sequences at distance ν from a reference sequence need no longer be compatible with the ref-

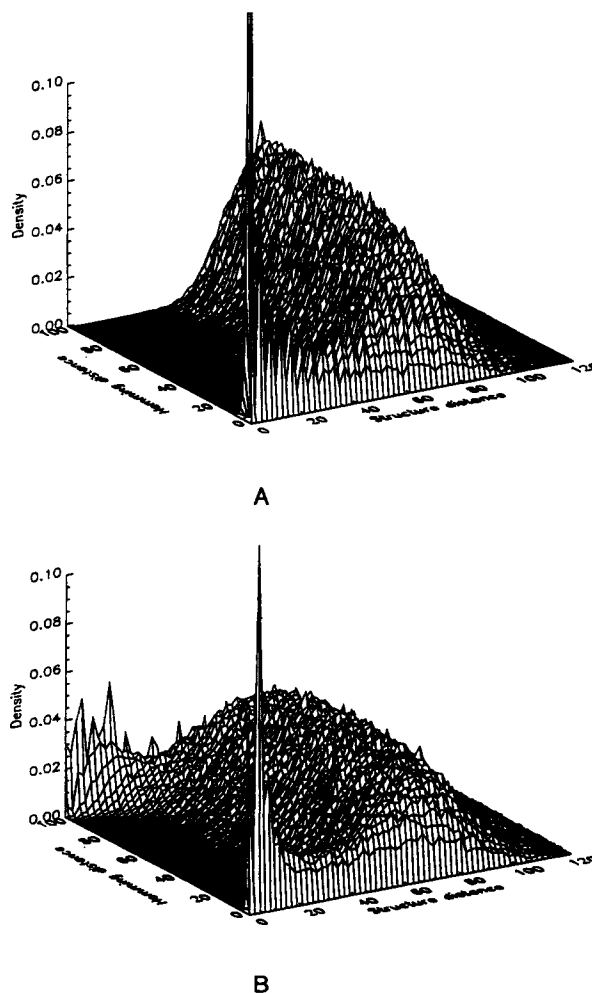


Figure 12. Probability density surfaces for tree distances of secondary structures of AUGC (A) and GC sequences (B) of chain length $\nu = 100$. In this computation $r = 1000$ reference points and a sample of $l = 10$ sequences were used, which amounts to a total sample size of 10^6 .

erence structure. The overall shape of the tree density surface of AUGC sequences, therefore, looks like one half of a horseshoe. The first half of the density surface ($0 \leq h \leq \nu/2$) closely resembles the corresponding part of the surface of binary sequences. The second half ($\nu/2 \leq h \leq \nu$), however, is like a mountain ridge parallel to the Hamming axis, and indicates that the distribution of tree distances is more or less independent of the Hamming distance h .

The overall shape of the tree density surfaces does not depend on the usage of full trees. Essentially the same features are observed with density surfaces whose structure distances are computed from HITs (Figure 13). Moreover, the shape of the density sur-

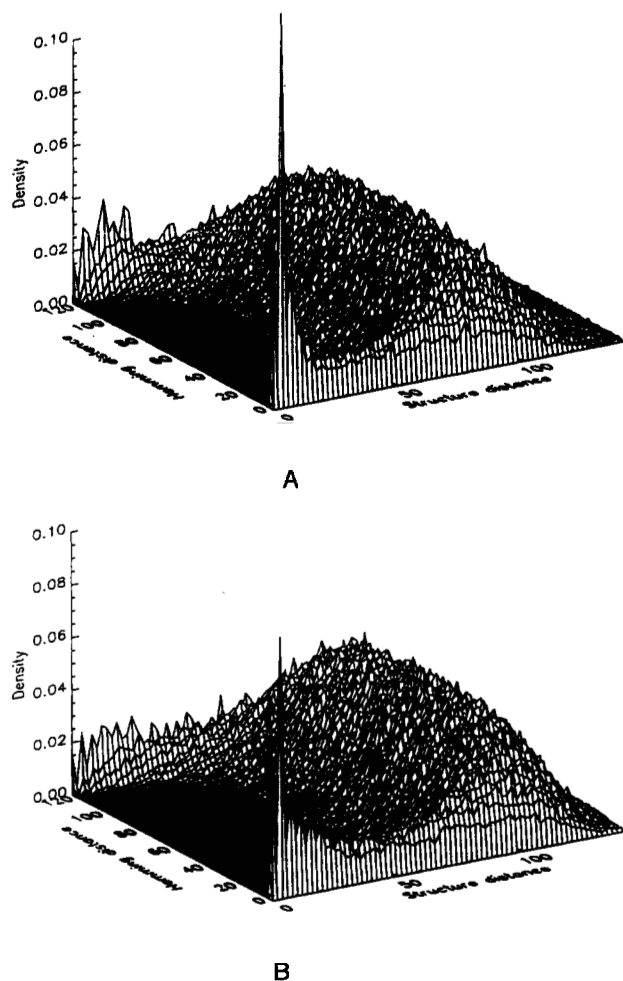


Figure 13. Probability density surfaces for tree distances of full trees (A) and HITs (B) from structures of GC sequences of chain length $\nu = 120$.

face is not altered when the structure distance measure is based on the above-mentioned linear encoding of secondary structures.⁴¹ The overall shape of distance density surfaces is a robust statistical feature of RNA secondary structure folding. It is, therefore, appropriate to use the generic term *structure density surface* (SDS) for this object.

First we interpret the shape of the first half of the SDS ($0 \leq h \leq \nu/2$), which is common to two-letter and four-letter alphabets. At very small Hamming distances from the reference sequence ($d_h = h = 1, 2, 3$) the most probable structures are similar to that of the reference, and the probability density $P(t|h)$ has its maximum at small tree distances. With increasing Hamming distance this maximum shifts toward larger tree distances, and the probability to find the same or a closely related structure is already very small at intermediate Hamming dis-

tances ($4 \leq h \leq h_{cr}$). For Hamming distances $h \geq h_{cr}$ the probability densities are very similar and, apart from finer details, independent of h . The critical Hamming distance lies at the (more or less sharp) turn of the horseshoe and represents the change from local to global features of the SDS. Thus, the value of h_{cr} defines a neighborhood size around an average sequence within which secondary structures can be effectively randomized. An extension of the interpretation of SDS to the second half of the sequence space ($\nu/2 \leq h \leq \nu$) is straightforward by the arguments on structure compatible sequences given in the last but one paragraph.

Shape Space Covering

Since the above-mentioned aspect of density surfaces does neither depend on the chain length ν , nor on the base-pairing alphabet, nor on the particular measure for structure distances, we turn it into a conjecture on RNA shape space covering: almost all typical RNA secondary structures—these are structures obtained as minimum free energy structures on random sequences—occur already within a small neighborhood of any random sequence. This conjecture has been strengthened by further independent techniques on which we will report elsewhere.

Probability density surfaces can be computed for any quantity for which a measure of distance exists. For example, a free energy density surface may be constructed and computed by substituting the absolute value of the difference in free energies,

$$d_f(i, j) = |f(S_i) - f(S_j)| = f \quad (11)$$

for the tree distance $d_t(i, j)$. All other scalar properties of RNA molecules can be studied analogously.

AUTOCORRELATION FUNCTIONS AND CORRELATION LENGTHS

Landscapes and CMs in general can be characterized statistically by autocorrelation functions,^{37,42} which can be expressed in terms of mean square distances:

$$\rho(h) = 1 - \frac{\langle d^2(h) \rangle}{\langle d^2 \rangle} \quad (12)$$

$\langle d^2 \rangle$ is the mean square distance sampled over the entire sequence space, and the mean square distance conditioned on sequences with Hamming distance

h is denoted by

$$\langle d^2(h) \rangle := \langle d^2(i, k | d_h(i, k) = h) \rangle \quad (13)$$

The autocorrelation function of tree distances, for example, can be computed from the density surface $P(t|h)$. The conditional mean square distance is simply the expectation value of t^2 computed for a given Hamming distance h

$$\langle d_i^2(h) \rangle = \sum_{t=0}^{t_{\max}} t^2 P(t|h) \approx \frac{\sum_{t=0}^{t_{\max}} t^2 n(t, h)}{\sum_{t=0}^{t_{\max}} n(t, h)} \quad (14)$$

Recalling that the mean square distance on the entire sequence space can be expressed as a weighted sum of the conditional mean square distances $\langle d_i^2(h) \rangle$, we find

$$\begin{aligned} \langle d_i^2 \rangle &= \sum_{h=0}^{\nu} \langle d_i^2(h) \rangle \cdot p(h) \\ &\approx \sum_{h=0}^{\nu} \frac{\sum_{t=0}^{t_{\max}} t^2 n(t, h) p(h)}{\sum_{t=0}^{t_{\max}} n(t, h)} \end{aligned} \quad (15)$$

For uniform sampling, $\sum_{t=0}^{t_{\max}} n(t, h) = l \cdot r = \text{const}$, the autocorrelation function can now be written in terms of elements of the sampling array $n(t, h)$,

$$\rho_i(h) = 1 - \frac{\sum_{t=0}^{t_{\max}} t^2 n(t, h)}{\sum_{h=0}^{\nu} \sum_{t=0}^{t_{\max}} t^2 n(t, h) p(h)} \quad (16)$$

Equation (16) is directly applicable during the computation of the density surface and hence convergence of the autocorrelation function can be monitored at runtime.

Computation of autocorrelation functions from density surfaces by means of Eq. (12) represents an alternative to the random walk technique.^{34,42} Since the more distant mutant classes are treated with higher numerical accuracy by uniform sampling, this method has an advantage over random walks. In practice, convergence of the computed autocorrelation functions is faster in the uniform sampling case.

Autocorrelation functions of tree distances $\rho_i(h)$ are approximated by an exponential fit in order to calculate a *correlation length* (l_i) for secondary structures in sequence space: $\ln \rho_i(l_i) = -1$. As we

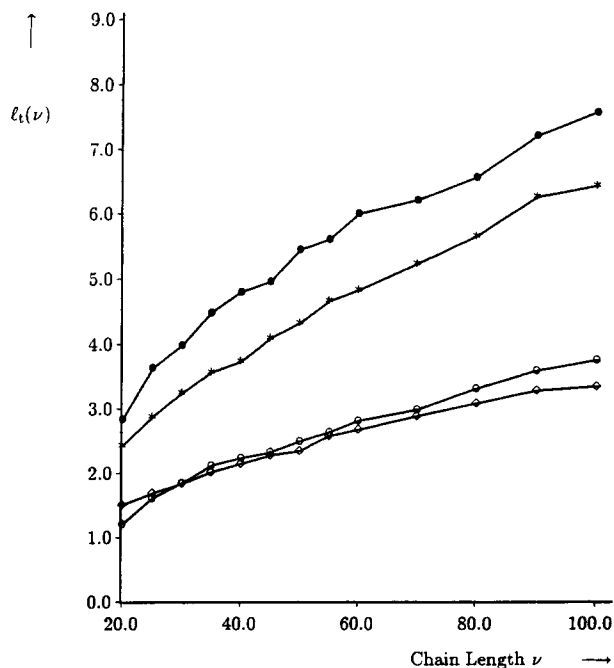


Figure 14. Correlation lengths of tree distances (l_i) of RNA molecules in their most stable secondary structures as functions of the chain length ν . Values are shown for binary GC sequences (\diamond), for binary AU sequences (\circ), for four-letter GCXK sequences with GC parameters ($*$), and for natural AUGC sequences (\bullet). Correlation lengths are calculated from $[\ln \rho_i(h), h]$ plots by means of a least rms deviation fit.

conclude from Figure 14, the correlation length increases roughly linearly with the chain length ν . It depends strongly on the base-pairing alphabet. The correlation lengths of structures from binary sequences, AU or GC, are very similar and substantially shorter than those derived from GCXK sequences and AUGC sequences. There seems to be a trend toward longer correlation lengths with weaker base pairing: for sequences of chain lengths $\nu \geq 30$ the correlation lengths for AU are slightly longer than those for GC, and the l_i values for AUGC are always longer than the values for GCXK.

The correlation length represents a global characterization of a CM by a single value. The shorter the correlation length, the more complicated the sequence-structure relation, and optimization of properties related to structure occurs on a very rugged landscape. The correlation length provides a measure of the mean stability of structures against point mutations in their sequences. If the correlation length is long, a sequence will tolerate on average many point mutations without changing its minimum free energy structure. Apparently, natural se-

quences are least sensitive to point mutations. This sensitivity can be changed by variation of the relative GC content of sequences: the structures of sequences with predominantly one WC pair have lower stability against point mutations. An observation of extensive structural variations in 28s rRNAs from different vertebrates may be related to this lack in resilience against point mutation: GC-rich phylogenetically variable regions show large-scale structural variations as detected by electron microscopy.⁴³

CONCLUDING REMARKS

A novel approach for investigating biopolymer structures that does not focus on the properties derived from single sequences has been proposed and carried out for RNA secondary structures. A statistical reference is computed for structures with given chain lengths and base-pairing alphabets. It provides information on the statistics of structural elements. These statistics is complemented by SDS, which describe local and global features of sequence–structure relations in condensed form. Correlation lengths of secondary structures in sequence space provide a measure for the predictability of structural changes as sequences are modified by point mutations.

Although a particular empirical parameter set had to be used in the computation, most of the features discussed here are robust with respect to corrections in individual parameters. What does depend on the choice of parameters are, of course, the predictions of structures for single sequences as well as some numerical values of the statistical reference. For example, the frequencies of structural elements are sensitive to the parameters determining the stability of these elements. The results of the comparison with natural sequences (whose structures were computed with the same folding routine as the statistical reference), the shapes of the structure density distributions, or the correlation lengths and their dependence on chain lengths and base-pairing alphabets will be largely independent of parameters and other details of the folding algorithms. This has actually been tested in one particular case: consideration of especially stable tetraloops did not significantly change correlation lengths.

Our approach yields three major results:

1. Comparison of sequence data at the level of structures computed for sufficiently large samples with the statistical reference is pro-

posed as a tool for the detection of structural features determined by biological function.

2. Most typical structures are found in close neighborhoods of any random sequence. To find sequences that fold into predetermined structures, only regions of a small size in sequence space have to be searched, and any random sequence is an equally valid starting point. The size of this region can be read off directly from the SDS.
3. A comparison of correlation lengths for secondary structures in sequence space shows that structures derived from GC- or AU-rich sequences are much more sensitive against mutation than those from AUGC sequences with uniform base distributions.

Financial support for the work reported here was provided by the Austrian Fonds zur Förderung der wissenschaftlichen Forschung (projects P 6864, P 8526, and S 5305), by the Jubiläumsfonds der Österreichischen Nationalbank (project no. 3819), by the Austrian Bundesministerium für Wissenschaft und Forschung (GZ 30.330/2-23/90), by the German Volkswagen-Stiftung, by the John D. and Catherine T. MacArthur Foundation, by the National Science Foundation (PHY-8714918), and by the U.S. Department of Energy (ER-FG05-88ER25054). DAMK thanks the W. M. Keck Foundation for its generous support of RNA science on the Boulder campus. An IBM 6000 RISC workstation was generously supplied by the EDV-Zentrum der Universität Wien within the EASI project of IBM. The Institut für Molekulare Biotechnologie (Jena, Germany) is sponsored by core funding from the Thüringische Ministerium für Wissenschaft and Kunst (Erfurt) and the Bundesministerium für Forschung und Technologie (Bonn). We are grateful to Dr. R. R. Gutell for providing the ribosomal RNA structure files, and to Dipl. Ings. Erich Bauer, Manfred Tacker, and Ivo L. Hofacker for performing some computations referred to here. Useful hints in stimulating discussions given by Professors Doyne Farmer, Paulien Hogeweg, Stuart Kauffman, Ted Maden, John McCaskill, Alan Perelson, and Detlef Riesner are gratefully acknowledged. Dr. Michael Ramek provided TEX MACROS⁴⁴ for drawing diagrams.

REFERENCES

1. Horwitz, M. S. Z., Dube, D. K. & Loeb, L. A. (1989) *Genome* **31**, 112–117.
2. Tuerk, C. & Gold, L. (1990) *Science* **249**, 505–510.
3. Ellington, A. D. & Szostak, J. W. (1990) *Nature* **346**, 818–822.
4. Piccirilli, J. A., Krauch, T., Moroney, S. E. & Benner, S. A. (1990) *Nature* **343**, 33–37.

5. Pleij, C. W. A. (1990) *Trends Biochem. Sci.* **15**, 143–147.
6. Westhof, E. & Jaeger, L. (1992) *Current Opinion Struct. Biol.* **2**, 327–333.
7. Zuker, M. & Stiegler, P. (1981) *Nucleic Acids Res.* **9**, 133–148.
8. Zuker, M. & Sankoff, D. (1984) *Bull. Math. Biol.* **46**, 591–621.
9. Waterman, M. S. & Smith, T. F. (1978) *Math. Biosci.* **42**, 257–266.
10. Williams, A. L. & Tinoco, I. (1986) *Nucleic Acids Res.* **14**, 299–315.
11. Zuker, M. (1989) *Science* **244**, 48–52.
12. Jaeger, J. A., Turner, D. H. & Zuker, M. (1990) *Methods Enzymol.* **183**, 281–306.
13. McCaskill, J. S. (1990) *Biopolymers* **29**, 1105–1119.
14. Freier, S. M., Kierzek, R., Jaeger, J. A., Sugimoto, N., Caruthers, M. H., Neilson, T. & Turner, D. H. (1986) *Proc. Natl. Acad. Sci. USA* **83**, 9373–9377.
15. Jaeger, J. A., Turner, D. H. & Zuker, M. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 7706–7710.
16. Fontana, W. & Schuster, P. (1987) *Biophys. Chem.* **26**, 123–147.
17. Fontana, W., Schnabl, W. & Schuster, P. (1989) *Phys. Rev. A* **40**, 3301–3321.
18. Hofacker, I. L., Fontana, W., Konings, D. A. M. & Schuster, P. (1992) in preparation.
19. Steger, G., Hofmann, H., Förtsch, J., Gross, H. J., Randles, J. W., Sänger, H. L. & Riesner, D. (1984) *J. Biomol. Struct. Dynam.* **2**, 543–571.
20. Gutell, R. R. (1992) *The Origin and Evolution of Prokaryotic and Eukaryotic Cells*, Hartman, H. & Matsuno, K., Eds., World Scientific, Singapore, in press.
21. Shapiro, B. A. (1988) *CABIOS* **4**, 387–397.
22. Shapiro, B. A. & Zhang, K. (1990) *CABIOS* **6**, 309–318.
23. Tai, K. (1979) *J. ACM* **26**, 422–433.
24. Sankoff, D. & Kruskal, J. B., Eds. (1983) *Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparison*, Addison Wesley, London.
25. Hogeweg, P. & Hesper, B. (1984) *Nucleic Acid Res.* **12**, 67–74.
26. Konings, D. A. M. (1989) *Pattern Analysis of RNA Secondary Structure*. Proefschrift, Rijksuniversiteit te Utrecht.
27. Konings, D. A. M. & Hogeweg, P. (1989) *J. Mol. Biol.* **207**, 597–614.
28. Wright, S. (1932) in *Proceedings of the Sixth International Congress on Genetics*, Vol. 1, pp. 356–366.
29. Kauffman, S. & Levin, S. (1987) *J. Theor. Biol.* **128**, 11–45.
30. Kauffman, S. A. (1989) in *Complex Systems, SFI Studies in the Sciences of Complexity*, Stein, D., Ed., Addison Wesley, Redwood City, CA, pp. 527–618.
31. Kauffman, S. A. & Johnsen, S. (1991) *J. Theor. Biol.* **149**, 467–505.
32. Macken, C. A. & Perelson, A. S. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 6191–6195.
33. Weisbuch, G. (1990) *J. Theor. Biol.* **143**, 507–522.
34. Fontana, W., Griesmacher, T., Schnabl, W., Stadler, P. F. & Schuster, P. (1991) *Mh. Chem.* **122**, 795–819.
35. Hamming, R. W. (1986) *Coding and Information Theory*, 2nd ed., Prentice-Hall, Englewood Cliffs, NJ, pp. 44–47.
36. Avis, D. (1981) *Can. J. Math.* **33**, 795–802.
37. Eigen, M., McCaskill, J. & Schuster, P. (1989) *Adv. Chem. Phys.* **75**, 149–263.
38. Maynard-Smith, J. (1970) *Nature* **225**, 563–564.
39. Perelson, A. S. & Oster, G. F. (1979) *J. Theor. Biol.* **81**, 645–670.
40. Segel, L. A. & Perelson, A. S. (1988) in *Theoretical Immunology*, Part Two, Perelson, A. S., Ed., Addison Wesley, Redwood City, CA, pp. 321–343.
41. Huynen, M., Konings, D. A. M. & Hogeweg, P. (1992) The effect of multiple coding on the evolutionary properties of RNA secondary structure. Preprint.
42. Weinberger, E. (1990) *Biol. Cybern.* **63**, 325–336.
43. Wakeman, J. A. & Maden, B. E. H. (1989) *Biochem. J.* **258**, 49–56.
44. Ramek, M. (1990) in Clark, M., Ed., *T_EX: Applications, Uses, Methods*. Proceedings of the T_EX 88 Conference. Ellis Horwood, Chichester, UK, pp. 227–258.

Received February 13, 1992

Accepted January 27, 1993