

Simulating an exploration of RNA conformation space with an appropriate parallel-updating strategy

Ariel Fernández

*Department of Biochemistry and Molecular Biology, The Medical School, Miami, Florida 33101-6129
and The Frick Laboratory, Princeton University, Princeton, New Jersey 08544*

(Received 5 March 1993)

A thermodynamic ensemble of stable RNA structures should emerge in the long-time limit as a result of an expeditious exploration of conformation space. The design of a simulation of this exploration giving consistent results in the long-time limit has been hindered by two main factors: (1) the need to incorporate the kinetic or activation-energy barriers for structure conversion and (2) the possibility of competing folding pathways. In this work, we implement a parallel kinetically controlled simulation that encompasses both aspects and ultimately yields all significant contributors to the thermodynamic ensemble. *We provide evidence supporting the conjecture that thermodynamic representativity may be reproduced by introducing an appropriate learning or updating strategy in a kinetically controlled parallel simulation.* The results are specialized to the illustrative cases of a transfer RNA and a midvariant-1 $Q\beta$ RNA for which the dominant native and non-native structures have been independently established. In both cases the thermodynamic ensemble is reproduced.

PACS number(s): 87.15.He, 36.20.Ey

The simulation of competing folding pathways for an RNA molecule treated *in vitro* calls for a parallel computation that can handle simultaneously different folding alternatives [1,2]. In implementing a parallel strategy we ultimately aim at reproducing the thermodynamic ensemble of secondary structures that should result from an expeditious exploration of conformation space.

In this context, two different levels of parallelism may be envisioned: (a) Massive parallelism [3,4], aimed at revealing cooperative effects in structure formation. The processing units consist of single putative Watson-Crick base pairs. (b) Multiprocessed floating-point Monte Carlo simulations of sequential folding [1,2], in which each processor deals at a given time with a single folding pathway and the processing units are topologically admissible (nonknotted) secondary structures with a time-dependent activation.

We shall focus on the second level of parallelism to implement a simulation in which the final destination processing units constitute the most significant contributions to the thermodynamic ensemble. In rigorous terms, this implies that the level of activation of the destination structure should be very close to its Boltzmann weight. This result would not be unexpected were it not for the fact that the activation of each processing unit is not controlled by the net stabilities of structures connected to it but is governed by the activation-energy barriers of structure conversion [1,5]. On the other hand, this kinetic control in the folding process cannot be overlooked since relevant experimental time scales invariably constrain the search in conformation space [5,6].

Within these premises, the central problem is to construct a learning or connection-updating strategy that makes kinetic control compatible in the long-time limit with thermodynamic representativity. That is, the final destination structures should be the significant contribu-

tors to the ensemble of most stable structures. An updating of connectivities that converges relatively fast to a winner-takes-all strategy [1,3,4] would not be adequate since it would lead to a dominant kinetically arrested structure, as shown in our first illustrative example. Instead, we show that an adequate strategy consists in strengthening a connectivity whenever it provides maximum stimulation to the activation of a given structure at a given stage of the simulation.

The basic tenets of a parallel kinetically controlled simulation will be introduced, the updating strategy will be implemented, and the results will be illustrated with the computation of all significant folding pathways for two RNA molecules of general interest: (a) The so-called midvariant-1 RNA (MDV-1RNA), a natural replication template for the enzyme $Q\beta$ -replicase, whose active structure is known [5,7]; and (b) a transfer RNA (tRNA) from *E. coli*, whose optimal active folding in the form of a cloverleaf base-pairing pattern has been firmly established to be a folding motif in most tRNA's [8,9].

The search for pre-mRNA structures starts with the very first refolding event during the synthesis of the molecule [1,2,5,10]. Given the relatively short biological times scales involved (approximately 15 s for the synthesis by sequential incorporation of nucleotides of a fragment 220 nucleotides long [5]), the exploration of configuration space concomitant with chain growth becomes heavily time constrained. Thus any algorithm which attempts to predict biologically relevant structures must incorporate the restrictions which bias exploration of configuration space. Earlier we proposed a Monte Carlo simulation which handles refolding events *concurrent* with sequential polymerization events, in an attempt to account for biological time constraints. The simulation mimics a Markov process such that if at a given stage a refolding event has a larger transition rate than a polymerization event,

the former is chosen, whereas if the reverse holds, the chain grows by incorporation of one nucleotide.

For the sake of completion, we shall first sketch the general tenets of the simulation. The Markov process is comprised of three different kinds of kinetically governed elementary events: (I) intrachain partial helix formation, (II) intrachain helix decay, and (III) chain growth by incorporation of a single nucleotide, with a fixed rate of phosphodiester linkage of 50 s^{-1} . The transition time for each event is a Poisson random variable. Eighteen simulations at most may be run in parallel, with each floating-point VAX 8650 computer processor dealing with a single competing folding pathway. The branching of folding pathways itself is the result of flexibility in the choice of structures to be formed at each stage of polymerization.

We shall start by describing the computation of the optimal pathway and then determine how to compute the branching leading to competing pathways. The interactive units are RNA secondary structures accessible for specific lengths of the chain, from $N=1$ to, say, $N \sim 10^3$. Only structural elements such as loop-stem systems, internal loops, and bulge loops are allowed, as discussed presently [1,2,5]. Given two arbitrarily chosen structures i and j , the rate for the interconversion $i \rightarrow j$ is denoted $k(i \rightarrow j)$ or, alternatively, k_{ij} , and is defined as

$$k_{ij}^{-1} = k_{d(i)}^{-1} + k_{f(j)}^{-1}, \quad (1)$$

where $k_{d(i)}^{-1}$ denotes the time span for dismantling the minimal portion of structure i which must be refolded to yield structure j and $k_{f(j)}^{-1}$ is the time span of formation of structure j starting from a partially dismantled structure i . The only substructures whose disruption or formation we allow are single or multiple stem-loop systems. Thus two structures whose interconversion requires two transformations (dismantling and refolding) of the type indicated *might* be connected and they are disconnected if their interconversion requires the occurrence of more than two such events.

The inverse mean time for intrachain helix dismantling (an elementary event of type II) may be obtained from the kinetics for helix decay [1,6]:

$$k_{d(i)} = t^{-1} = fn \exp[G_h/RT], \quad (2)$$

where f is the kinetic constant for a single base-pair formation (estimated at 10^6 s^{-1} , cf. [1,6]), n is the number of base pairs in the helix, and G_h is the (negative) free-energy contribution of the set of base-pairs in the helix. Thus the essentially enthalpic term $-G_h$ should be regarded as the activation energy for helix disruption. If an admissible helix formation (an elementary event of type I) happens to be the event favored, the inverse of the mean time for the transition will be given by

$$k_{f(j)} = t^{-1} = fn \exp(-\Delta G_{\text{loop}}/RT), \quad (3)$$

where ΔG_{loop} is the change in free energy due to the closure of the loop concurrent with helix formation. Loop closure is the rate-limiting event; therefore this quantity is essentially the activation energy of helix formation, corresponding to a loss in conformational entropy. Since

water is a relatively good solvent for RNA, excluded-volume effects might be significant [11]. Therefore we shall treat the negative entropic change due to loop formation as a convex function of the number of unpaired bases in the loop. In order to allow for the possibility of large loops we shall adopt an accurate function [11,12] to incorporate excluded-volume effects: $\Delta S_{\text{loop}} = -2.1 - [23/12 \ln u]$, where u is the number of unpaired bases in the loop. This substantially increases the time span of the simulation as well as the order of the algorithm.

We can see from general Eqs. (1)–(3) that the kinetically controlled simulation is dependent on a compilation of thermodynamic parameters. The computation of the rates of refolding events generated rely on a compilation of free-energy increments for the formation of helical regions [13] and loops [14]. In particular, the enthalpic parameters carry an uncertainty of at least 2% [13]. On the other hand, the essentially entropic contributions for loop closure are particularly unreliable since they have been derived for small synthetic oligonucleotides. The parameters have been corrected to encompass excluded-volume effects which are particularly conspicuous for small loops (containing less than six nucleotides).

In view of this situation, the robustness of the simulations has been tested making use of parallelism, attributing random variations in the thermodynamic parameters to the different processors according to a Gaussian distribution representing a margin of 8.8% uncertainty. The simulations were specialized to a specific RNA species, the phenylalanine transfer RNA from the organism *E. coli*: *E. coli* phe tRNA [8,9]. Transfer RNA's, with their well-established cloverleaf base-pairing pattern, provide an optimal testing ground for our computations. Thus any weaknesses resulting from uncertainties in the thermodynamic parameters should reflect themselves in the prediction of the dominant structure. All 18 processors used in this test reproduced 94% of the consensus secondary structure for *E. coli* phe tRNA. The results of the parallel simulation are shown explicitly below. It appears that errors do not accumulate and propagate in sequential folding algorithms the same way they do in free-energy minimization algorithms. Thus an error in the estimation of a thermodynamic parameter has only a local effect in sequential folding, involving the rate of retention of an upstream structure, whereas in free-energy minimization algorithms, each structure competes on a global level against all other folding alternatives.

So, far, we have sketched the basic elements needed to determine the optimal pathway. Along this pathway, if structure i_o is chosen by the processor at time t , the next structure chosen would be j_o , the structure which realizes the maximum

$$\max_j k(i_o \rightarrow j) = k(i_o \rightarrow j_o). \quad (4)$$

We are, at this point, in a position to implement a *parallel* extension of the kinetically controlled simulation [1–4]. Competing pathways result from perturbations of the optimal pathway due to occasional base-pair disruptions. However, in dealing with such base-pair disruptions of the transient structures, we shall not be able to distinguish or specify the agent causing the perturbation:

The parallel algorithm only reveals the specific stages of sequential folding in which a perturbation is required to ultimately generate the structures with the highest statistical weights.

To model the emergence of competing pathways, we need to introduce a connectivity matrix, as follows:

$$\mathbf{W} = [w_{ij}]_{ij}, \quad (5)$$

$$w_{ij} = k(j \rightarrow i) / \left[\sum_n k(j \rightarrow n) \right]. \quad (6)$$

In Eq. (6), i , j , and n denote transient structures such that i and any of the n 's are accessible from j . The connectivity matrix represents in a compact fashion the links between structures and determines the architecture of the simulation, since the choice of two different structures

linked to a specific one causes a branching of the folding pathway.

Two essential features in the architecture of the algorithm restrict and direct the computation and make the problem of branching of folding pathways tractable:

(1) The system consists of a set of *hierarchically layered processing units* (structures). Each layer corresponds to a fixed length N of the chain. Thus, if two structures i and j satisfy $N(i) = N(j)$, then they belong to the same layer $L_{N(i)}$. Each layer connects via excitatory links with the layer immediately above and receives an input from the layer immediately below.

(2) The output function $f_{N(i)}$ for layer $L_{N(i)}$ allows only certain units to produce an output signal. A unit i with state of activation a_i produces an output $O_i = f_{N(i)}(a_i)$ according to the following scheme:

$$f_{N(i)}(a_i) = \begin{cases} a_i & \text{if } |a_i - a_{i0}| < h(N(i)), \text{ where } h(N(i)) \sim \exp\{-\beta\Delta[N(i)]^{1/4}\} \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

The output function has been chosen in this way to incorporate structure fluctuations in the form of base pairings and base-pair disruptions. Such fluctuations determine the branching of the optimal folding pathway. Thus the constant $\Delta \approx 1.81k_B T$ (k_B is the Boltzmann constant) from Eq. (7) is the scaling factor for the minimal activation energy of a refolding event [1]. The minimal activation energy between mutually accessible foldings for a chain of length N is $E_a \sim \Delta N^{1/4}$. Thus branching of a folding pathway becomes increasingly rarer as we approach higher layers.

We may define the input of structure i at times t as

$$a_i(t+1) = \sum_j w_{ij} O_j(t), \quad 0 \leq a_i \leq 1. \quad (8)$$

Upon examination of Eqs. (5), (6), and (8), we may conclude that the input a_i is to be interpreted as the probability or statistical weight of structure i .

Our updating strategy consists in systematically reinforcing or weakening connections at each time t according to a prescription defined by a linear map $\Omega = \Omega(t)$:

$$\mathbf{W} = \mathbf{W}(0) \xrightarrow{\Omega(0)} \mathbf{W}(1) \xrightarrow{\Omega(1)} \mathbf{W}(2) \xrightarrow{\Omega(2)} \cdots \xrightarrow{\Omega(t)} \mathbf{W}(t+1) \xrightarrow{\Omega(t)} \cdots. \quad (9)$$

This construction allows us to modify Eq. (8) by progressively updating it:

$$a_i(t+1) = \sum_j w_{ij}(t) O_j(t), \quad (10)$$

$$\Omega(t-1) \cdots \Omega(1) \Omega(0) \mathbf{W} = [\omega_{ij}(t)].$$

In order to understand how kinetic control could be made compatible with thermodynamics, two updating maps $\Omega_I(t)$ and $\Omega_{II}(t)$ will be defined according to two different learning strategies. The map $\Omega_I(t)$ is defined by

$$[\Omega_I(t) \mathbf{W}(t)]_{ij} = w_{ij}(t+1) = \frac{w_{ij}(t) + \mu_I O_j(t) w_{ij}(t) \delta_{L(j)i}}{1 + \mu_I O_j(t) w_{L(j)j}(t)}, \quad (11)$$

where $\delta_{L(j)i}$ is the Kronecker delta function, the subindex $L(j)$ satisfies

$$\max_i w_{ij}(t) O_j(t) = w_{L(j)j}(t) O_j(t), \quad (12)$$

and the denominator on the right-hand side of Eq. (11) is the normalization factor ensuring that

$$\sum_i w_{ij}(t) = 1 \text{ for all } t. \quad (13)$$

The updating map $\Omega_I(t)$ is parametrically dependent on the positive number μ_I which determines the rate of convergence of the process to a winner-takes-all process, as revealed by iterating the result of Eq. (11). Thus, given a structure j , the connection with the fastest formed immediate destination structure $L(j)$, which results from a refolding of j , is strengthened whereas all other connections $j \rightarrow i$, $i \neq L(j)$, are weakened. As inferred from Eqs. (11)–(13), this updating strategy ends up yielding the optimal kinetically controlled pathway.

The map $\Omega_{II}(t)$ is defined by

$$[\Omega_{II}(t) \mathbf{W}(t)]_{ij} = w_{ij}(t+1) = \frac{w_{ij}(t) + \mu_{II} O_j(t) w_{ij}(t) \delta_{J(i)j}}{1 + \mu_{II} O_{J(i)}(t) w_{iJ(i)}(t)}, \quad (14)$$

where $J(i)$ satisfies

$$\max_j w_{ij}(t) O_j(t) = w_{iJ(i)}(t) O_{J(i)}(t). \quad (15)$$

The updating map $\Omega_{II}(t)$ reinforces connections differently from $\Omega_I(t)$. Given a destination state i , the connection with the structure $j = J(i)$ that converts into i

with the fastest transition rate is reinforced, whereas all other connections to destination state i are weakened. Again, this updating map is parametrically dependent on a positive number μ_{II} . At this point, a conjecture will be framed which is supported by our results: *For appropriate ranges in the learning parameter μ_{II} , the updating strategy II makes the kinetically controlled parallel simulation thermodynamically compatible in the long-time limit.* The evidence for this conjectured will be presented now.

In order to quantitatively compare the two strategies, we shall first indicate how to represent the results of the simulation in the long-time limit. Since the level of description entails only base-pair patterns, we shall represent the outcome by means of a matrix of base-pair probabilities which results as we overlap the final activation states $a_i(t \rightarrow \infty)$'s. Thus, if x and $y=1,2,\dots,N$ (N =length of chain; $x \neq y$) represent two bases along the RNA chain, the complete long-time information is contained in the $N \times N$ matrix $[p_{xy}]$ defined by

$$p_{xy} = \sum_i a_i(t \rightarrow \infty) \Pi_{xy}(i), \quad (16)$$

where $\Pi_{xy}(i)=1$ if bases x and y are engaged in a standard Watson-Crick base pairing in structure i , and 0 otherwise.

We may visualize this matrix as a shade matrix, where a fraction of the xy entry is shaded according to the probability p_{xy} .

Since the base-pair probability matrix is symmetrical, we shall conveniently condense the information in the upper right or lower left triangle in Fig. 1.

The results of the parallel simulations comprised of 10^6 Monte Carlo (MC) steps for the species MDV-1RNA ($N=220$) [5,7] are presented in Fig. 1. The choice of parameters is $\mu_I=4.4$ and $\mu_{II}=2.0$. Larger values of the

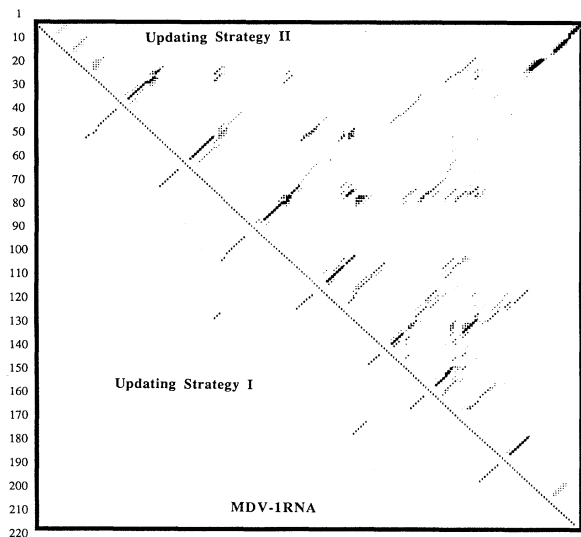


FIG. 1. The 220×220 base-pair probability matrix $[p_{xy}]$ for MDV-1 $Q\beta$ RNA generated in the long-time limit by a parallel kinetically controlled Monte Carlo simulation. The upper right triangle corresponds to updating strategy II and the lower left triangle to updating strategy I.

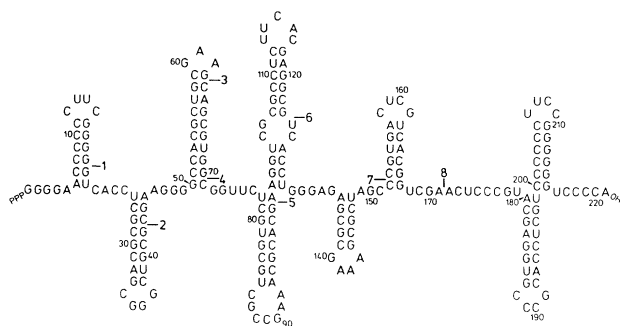


FIG. 2. The most probable kinetically arrested structure of MDV-1 $Q\beta$ RNA. This structure is biologically active and regulates replication of this RNA species. It is also one of the two main contributors to the thermodynamic ensemble for MDV-1 $Q\beta$ RNA. The base pairs $A-U$, $G-C$ (A =adenine, U =uracil, G =guanosine, C =cytosine) correspond to standard Watson-Crick complementarity.

parameters do not lead to any appreciable difference in the statistical weights of the dominant structures. Although a detailed investigation of the ranges of convergence as a function of the length of the chain would be desirable, it is beyond the scope of this discussion. The lower left triangle corresponds to learning strategy I and the upper right triangle to learning strategy II. Two significant contributions are apparent from direct examination of the upper right triangle: (a) The optimal structure resulting from a single kinetically controlled simulation [5], represented in Fig. 2; and (b) the most stable structure, schematically represented in Fig. 3, re-

Global Minimum (Schematic)

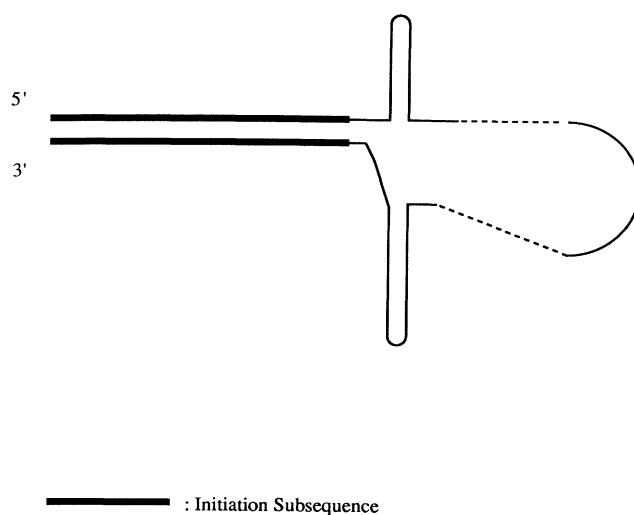


FIG. 3. Scheme of the most stable structure for MDV-1 $Q\beta$ RNA. The two 21-nucleotide-long extremities featuring a high level of Watson-Crick complementarity (see Fig. 2) are now bound to each other.

sulting from the binding of the highly complementary extremities of the molecule (see Fig. 2) [7] and otherwise identical to the previous structure. The near-diagonal base pairings in Fig. 1 correspond to the various hairpins of the structures depicted in Fig. 2, while the base pairings at the upper right corner correspond to the two extremities bound together as depicted in Fig. 3.

Most importantly, the structure represented in Fig. 2 has been shown to be biologically significant and experimental evidence reveals it is the active structure required to regulate the replication of the molecule [5,7]. On the other hand, the most stable structure, depicted in Fig. 3, is biologically inert since the initiation of RNA replication is blocked [7]. Both structures represent the dominant contributors to the thermodynamic ensemble [7]. We can see that only strategy II is thermodynamically compatible. When strategy I is used, a rapid convergence to the optimal folding pathway results only in the kinetically most favored structure [5] shown in Fig. 2.

The phe tRNA from *E. coli* ($N=76$), on the other hand, provides a good additional illustration since its stable cloverleaf secondary structure has been independently established [8,9]. The results of both strategies are again depicted by the 76×76 shade matrix shown in Fig. 4 for $\mu_I \geq 2.24$ and $\mu_{II} \geq 1.88$. The stable cloverleaf structure with hairpin base pairing near the diagonal and near end-to-end binding is the dominant contributor to the matrix. As expected, the convergence range for this species is much broader relative to the longer MDV-1RNA. For the whole parametric range indicated, both strategies yield essentially a constant statistical weight in the long-time limit corresponding to 10^6 MC steps. Had we dealt only with the tRNA example alone,

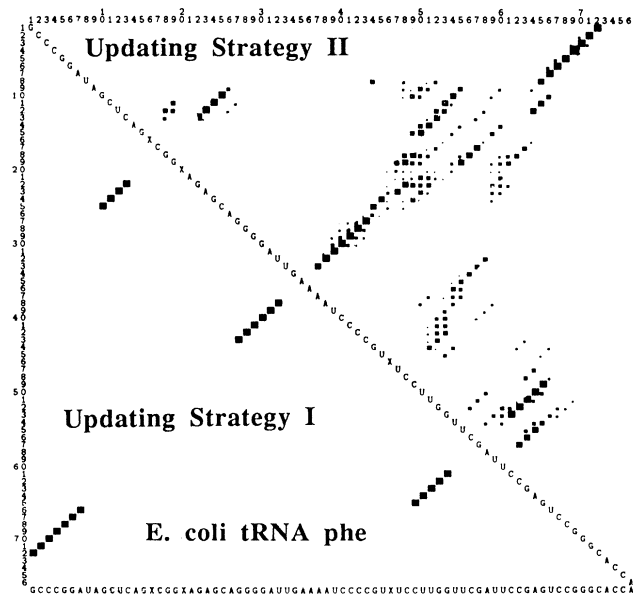


FIG. 4. The 76×76 base-pairing probability matrix for the *E. coli* phe tRNA in the long-time limit.

we would have not been able to properly distinguish between the two updating strategies since the thermodynamic ensemble for the tRNA is basically dominated by a single stable structure [8,9].

This work was supported by the Camille and Henry Dreyfus Foundation.

-
- [1] A. Fernández, *J. Theor. Biol.* **157**, 487 (1992).
 [2] A. Fernández, *Phys. Rev. Lett.* **64**, 2328 (1990).
 [3] D. Rummelhart and J. L. McClelland, in *Parallel Distributed Processing, Vol. 1: Foundations*, edited by D. Rummelhart, J. L. McClelland, and the PDP Research Group (MIT Press, Boston, 1988), Chap. 2, pp. 45–76.
 [4] D. Rummelhart and J. L. McClelland, in *Parallel Distributed Processing, Vol. 1: Foundations* (Ref. [3]), Chap. 5, pp. 151–193.
 [5] A. Fernández, *Eur. J. Biochem.* **182**, 161 (1989).
 [6] V. V. Anshelevich, V. A. Vologodskii, A. V. Lukashin, and M. D. Frank-Kamenetskii, *Biopolymers* **23**, 39 (1984).
 [7] (a) D. R. Mills, C. Dobkin, and F. R. Kramer, *Cell* **15**, 541 (1978); (b) A. Fernández, *Physica A* **176**, 499 (1991).
 [8] A. Rich and U. L. RajBhandary, *Annu. Rev. Biochem.* **45**, 805 (1976).
 [9] P. Schimmel, in *Nucleic Acids and Molecular Biology*, edited by F. Eckstein and D. M. Lilley (Springer-Verlag, Berlin, 1990), Vol. 4, pp. 274–287.
 [10] S. E. LaFlamme, F. R. Kramer, and D. R. Mills, *Nucleic Acids Res.* **13**, 8425 (1985).
 [11] A. Fernández, *Phys. Rev. A* **44**, R7910 (1991).
 [12] H. S. Chan and K. A. Dill, *J. Chem. Phys.* **90**, 492 (1989).
 [13] D. H. Turner, N. Sugimoto, and S. M. Freier, *Annu. Rev. Biophys. Chem.* **17**, 167 (1988).
 [14] D. R. Groebe and O. C. Uhlenbeck, *Nucleic Acids Res.* **16**, 11 725 (1988).