# An Algorithm for Computing Nucleic Acid Base-Pairing Probabilities Including Pseudoknots

ROBERT M. DIRKS,<sup>1</sup> NILES A. PIERCE<sup>2</sup>

<sup>1</sup>Department of Chemistry, California Institute of Technology, Pasadena, California 91125 <sup>2</sup>Departments of Applied & Computational Mathematics and Bioengineering, California Institute of Technology, Mail Code 114-96, Pasadena, California 91125

> Received 21 January 2004; Accepted 19 March 2004 DOI 10.1002/jcc.20057 Published online in Wiley InterScience (www.interscience.wiley.com).

**Abstract:** Given a nucleic acid sequence, a recent algorithm allows the calculation of the partition function over secondary structure space including a class of physically relevant pseudoknots. Here, we present a method for computing base-pairing probabilities starting from the output of this partition function algorithm. The approach relies on the calculation of recursion probabilities that are computed by backtracking through the partition function algorithm, applying a particular transformation at each step. This transformation is applicable to any partition function algorithm that follows the same basic dynamic programming paradigm. Base-pairing probabilities are useful for analyzing the equilibrium ensemble properties of natural and engineered nucleic acids, as demonstrated for a human telomerase RNA and a synthetic DNA nanostructure.

© 2004 Wiley Periodicals, Inc. J Comput Chem 25: 1295-1304, 2004

Key words: DNA; RNA; base-pairing probabilities; partition function; pseudoknots

## Introduction

Thermodynamic models based on nucleic acid secondary structure and nearest-neighbor identities<sup>1–5</sup> underly dynamic programming algorithms for predicting the minimum energy secondary structure<sup>6–10</sup> and calculating the partition function over secondary structure space.<sup>10–12</sup> In their original forms, these algorithms exclude the possibility of pseudoknots, a biologically relevant class of secondary structures<sup>13</sup> that also arises in DNA nanotechnology applications.<sup>14,15</sup> Pseudoknots result when two base pairs  $i \cdot j$  and  $d \cdot e$ , with i < d, fail to satisfy the nesting property i < d < e < j (see, e.g., Fig. 1). Recent extensions of the structure prediction<sup>16–18</sup> and partition function<sup>18</sup> algorithms allow the inclusion of certain pseudoknots.

For an ensemble of secondary structures  $s \in \Omega$ , the partition function

$$Q = \sum_{s \in \Omega} e^{-G_s/RT}$$

may be used to compute the probability

$$p(s^*) = \frac{1}{Q} e^{-G_{s^*}/RT} \tag{1}$$

that secondary structure  $s^*$  is sampled at thermodynamic equilibrium. The ensemble equilibrium can also be characterized by the matrix of base-pairing probabilities with entries  $p_{i,j}$  corresponding to the probability that base *i* is paired with base *j* in  $\Omega$ .

McCaskill's original article<sup>11</sup> defines elegant dynamic programs to compute the partition function and base-pairing probabilities over the ensemble of unpseudoknotted secondary structures. The partition function algorithm builds up recursively from short subsequences to the full strand, and then the pair probabilities are computed by working backwards to short subsequences using intermediate results from the partition function calculation. In the absence of pseudoknots, the partition function algorithm is sufficiently succinct that McCaskill is able to determine the form of the pair probability backtrack algorithm simply by considering the few possible forms of enclosing secondary structure for any given base pair. Although this approach is simple and efficient, it is not easily

Contract/grant sponsor: NSF graduate research fellowship (R.M.D.).

Correspondence to: Niles A. Pierce; e-mail: niles@caltech.edu

Contract/grant sponsor: Defense Advanced Research Projects Agency (DARPA) and Air Force Research Laboratory under F30602-010200561 (N.A.P.).

Contract/grant sponsor: Ralph M. Parsons Foundation (N.A.P.).

Contract/grant sponsor: Charles Lee Powell Foundation (N.A.P.).

Pseudoknot



**Figure 1.** Secondary structures of competing pseudoknot and hairpin constructs in human telomerase RNA. The wild-type sequence is shown. For the two-point mutant implicated in dyskeratosis congenita, GC is replaced by AG in the shaded boxes, disrupting two base pairs in the pseudoknot construct. For the experimental studies of the hairpin structure,<sup>20</sup> the 18 nucleotides at the 3' end are excluded to prevent formation of the pseudoknot.

generalizable to algorithmic extensions, such as the inclusion of pseudoknots. Here, we describe a general method for mechanically transforming the new pseudoknot partition function algorithm<sup>18</sup> to compute recursion probabilities, which can be used in turn to compute base-pairing probabilities. The transformation approach is generalizable to any future partition function extensions that follow the same dynamic programming paradigm.

Base-pairing probabilities assist in the analysis of biologically relevant pseudoknots. Here, we examine human telomerase RNA, which exists at equilibrium in both hairpin and pseudoknotted forms.<sup>19</sup> A two-point mutation, implicated in the disease dyskeratosis congenita, alters the thermodynamic balance between these competing structures.<sup>20</sup> This shift in equilibrium is clearly identifiable when the base-pairing probabilities for the two sequences are compared. Base-pairing probabilities that permit pseudoknots are also useful in analyzing synthetic DNA nanostructures.<sup>14,15</sup>

# Algorithm

For clarity, we begin by considering the class of secondary structures excluding pseudoknots and then address the additional complexity that arises when pseudoknots are introduced.

### Partition Function Recursions

For a strand of length N, the partition function may be computed over all unpseudoknotted secondary structures in  $O(N^4)$  using the algorithm<sup>10,11</sup> summarized in Figure 2 (see ref. 18 for a detailed description in the same notation). [The complexity may be reduced to  $O(N^3)$  by exploiting the formulation of the nearest-neighbor energy model for long interior loops.<sup>18,21</sup>] Partition function recursions are nonredundant in the sense that every secondary structure in the ensemble  $\Omega$  is visited exactly once using a unique sequence of recursions. The algorithm computes the partition function  $Q_{i,j}$  for each subsequence [i, j] ignoring all bases exterior to [i, j], starting from subsequences of length l = 1 and building up incrementally to l = N. The recursions that define  $Q_{i,j}$  rely on additional restricted partition functions  $Q_{i,j}^{b}$  and  $Q_{i,j}^{m}$ .  $Q_{i,j}^{b}$  represents the partition function for subsequence [i, j] given that i and j are base paired and  $Q_{i,i}^m$  is used to calculate multiloop contributions. At the end of the recursive process, the full partition function Q is given by  $Q_{1,N}$  and the values of  $Q_{i,j}$ ,  $Q_{i,j}^b$ ,  $Q_{i,j}^m$  are stored in matrices for  $1 \le i, j \le N$ . These intermediate results will play a critical role in the new algorithm described below.

#### **Recursion Probabilities**

Following the execution of the partition function calculation, a second algorithm can be implemented to calculate probability matrices, P,  $P^b$ ,  $P^m$ , corresponding to the Q,  $Q^b$ ,  $Q^m$  matrices. The values stored in these *P*-type matrices will be termed *recursion probabilities*.

Recursion probabilities can be intuitively described as follows. Consider sampling the ensemble of secondary structures  $s \in \Omega$ where the probability of selecting structure  $s^*$  is given by the Boltzmann probability (1). For each secondary structure  $s^*$ , the contribution to Q is computed by a unique recursion sequence involving specific  $Q_{i,j}$ ,  $Q_{i,j}^b$ , and  $Q_{i,j}^m$  intermediates. Associating these intermediates with structure  $s^*$ , the recursion probability  $P_{i,j}$ ,  $P_{i,j}^b$  or  $P_{i,j}^m$  corresponds to the probability that the sampled structure  $s^*$  requires the use of the corresponding intermediate  $Q_{i,j}$ ,  $Q_{i,j}^b$  or  $Q_{i,j}^m$  to calculate the partition function contribution.

Recent work by Ding and Lawrence<sup>22</sup> exploits quantities related to recursion probabilities to statistically sample the distribution of unpseudoknotted secondary structures for a given sequence. Here, we develop a general approach for computing *P*-type ma-

```
 \begin{array}{ll} \mbox{Initialize} \left(Q,Q^b,Q^m\right) // O(N^2) \mbox{ space} \\ \mbox{Set all values to 0 except } Q_{i,i-1} = 1 \\ \mbox{for } i = 1, N \\ \mbox{for } i = 1, N - l + 1 \\ j = i + l - 1 \\ // Q^b \mbox{ recursion} \\ Q_{i,j}^b = \exp\{-G_{i,j}^{\rm hairpin}/RT\} \\ \mbox{for } e = d + 4, j - 1 \\ Q_{i,j}^b + e \exp\{-G_{i,d,e,j}^{\rm interior}/RT\} Q_{d,e}^b \\ Q_{i,j}^b + = Q_{i+1,d-1}^m Q_{d,e}^b \exp\{-[\alpha_1 + 2\alpha_2 + \alpha_3(j - e - 1)]/RT\} \\ // Q, Q^m \mbox{ recursions} \\ Q_{i,j} = 1 \ // \mbox{environs} \\ Q_{i,j} = 1 \ // \mbox{loop over all possible rightmost pairs } d \cdot e \\ \mbox{for } e = d + 4, j \\ \mbox{for } e = d + 4, j \\ \mbox{ quark is } - 2 \ (i - 1) \ (i - 1)
```

**Figure 2.**  $O(N^4)$  partition function algorithm that excludes pseudoknots.



**Figure 3.** Recursion diagram corresponding to recursive update (2), depicting the addition to  $Q_{i,j}$  of partition function contributions for those structures with rightmost base pair  $d \cdot e$ . See ref. 18 for a thorough description of the partition function algorithm (with or without pseudoknots) in terms of recursion diagrams.

trices given a set of *Q*-type matrices and corresponding partition function recursions.

An algorithm for computing recursion probabilities can be formulated in a mechanical way starting from a set of partition function recursions. The crux of this formulation is the repeated application of a single transformation to the partition function code. In particular, updates of the form

$$Q_{i,j} + = Q_{i,d-1} Q_{d,e}^{b}$$
(2)

(equivalent to the recursion diagram of Fig. 3) are converted to the following series of statements

$$\Delta p = P_{i,j} \underbrace{Q_{i,d-1}Q_{d,c}^{b}/Q_{i,j}}_{P_{i,d-1} + = \Delta p}$$

$$P_{d,e}^{b} + = \Delta p$$
(3)

Specifically, the right-hand side (RHS) of each recursive update is divided by the left-hand side (LHS), and the *P* term corresponding to the new denominator is multiplied by this quotient. The resulting probabilities, temporarily stored as  $\Delta p$ , are subsequently added to every *P*-type value corresponding to the *Q*-type terms on the RHS of the original statement (2).

To understand this transformation, recall that  $Q_{i,j}$ ,  $Q_{i,j}^{b}$  and  $Q_{i,j}^{m}$  are partition functions for structural subclasses of the full sequence. In recursive updates such as (2), the ratio of the RHS to the fully computed LHS corresponds to the probability that a structure drawn from an equilibrium ensemble defined by the LHS partition function is in the subensemble defined by the RHS partition function. As an example, transformation (3) states that for any *i*, *d*, *e*, *j*, the structures represented by  $Q_{i,d-1}$  and  $Q_{d,e}^{b}$ . Consequently, once the probability  $P_{i,j}$  is determined, it can be used to augment  $P_{i,d-1}$  and  $P_{d,e}^{b}$  because the frequencies of the corresponding substructures within the  $Q_{i,j}$  ensemble can be derived from  $Q_{i,d-1}$  and  $Q_{d,e}^{b}$ . By backtracking through the partition function algorithm and transforming all recursive updates analagously to (3), probabilities can be calculated for each recursion.

Starting from the partition function algorithm of Figure 2, the recursion probability algorithm is obtained by performing three modifications: (1) the two outermost loops are altered so that the algorithm starts with the full strand of length l = N and decrements down to subsequences of length l = 1; (2) all recursive updates are transformed as for (3) above; (3) the order of the recursion blocks  $(Q^b, [Q, Q^m])$  is reversed  $([P, P^m], P^b)$ . This last modification is necessary because the recursion order in the partition function algorithm ensures that if one quantity (e.g.,  $Q_{i,j})$  recurses to another quantity of the same length (e.g.,  $Q_{i,j}^b$ ) then the "lower level" quantity (i.e.,  $Q_{i,j}^b$ ) is calculated first. The reverse ordering is needed for the recursion probability algorithm, because  $P_{i,j}^b$  cannot be used until it has been fully computed in the  $P_{i,j}$  loop.

The pseudocode in Figure 4 details the outcome of these transformations for the unpseudoknotted case. This modified algorithm reverses the flow of the partition function calculation and incrementally determines all recursion probabilities (frequencies of families of structures), based on the probabilities of all superstructures that directly contain them. Once recursion probabilities are computed for all *i* and *j*, the base-pairing probability  $p_{i,j}$  is simply  $P_{i,j}^b$ , because  $Q_{i,j}^b$  is associated with every structure *s* in which  $i \cdot j$  appears, and  $i \cdot j$  is associated with exactly one  $Q_{i,j}^b$ . By starting from a more complicated  $O(N^3)$  partition function algorithm, <sup>18,21</sup> the computational complexity of the recursion probability algorithm can also be reduced to  $O(N^3)$  as described in the Appendix.

#### Pseudoknots

The procedure outlined above for obtaining recursion probability algorithms is equally applicable to a new partition function algorithm that includes pseudoknots (see the pseudocode of Fig. 21 in ref. 18). For the unpseudoknotted algorithm, all base pairs stem

```
\begin{array}{l} \mbox{Compute } Q, Q^b, Q^m \mbox{ using } O(N^4) \mbox{ partition function algorithm } \\ \mbox{Initialize } (P, P^b, P^m) \space \\ \mbox{Set all } P\mbox{-type values to } 0 \\ P_{1,N} = 1 \space \space \\ \mbox{Set all } P\mbox{-type values to } 0 \\ P_{1,N} = 1 \space \space \space \\ \mbox{Set all } P\mbox{-type values to } 0 \\ \space \space \space \space \\ \mbox{Set all } P\mbox{-type values to } 0 \\ \space \\ \space \
```

**Figure 4.**  $O(N^4)$  recursion probability algorithm that excludes pseudoknots. For simplicity, we omit details such as checking for updates with zero in the denominator (in which case the numerator will also evaluate to zero and the expression should be skipped).

from  $Q^b$  recursions, so the values stored in  $P^b$  are precisely the desired probabilities (i.e.,  $p_{i,j} = P_{i,j}^b$ ). For the pseudoknotted case,  $P_{i,j}^b$  only gives the probability that *i* and *j* form a nested pair. The full base-pairing probability must also take into consideration those base pairs that are nonnested and lead to pseudoknotted structures (termed gap-spanning pairs in ref. 18). For these gap-spanning pairs, there is no single recursion probability that represents the contribution to  $p_{i,j}$ . However, this contribution may be succinctly represented in terms of Q-type and P-type matrices for the full pseudoknotted algorithm.

A new set of quantities,  $P_{i,j}^{bg}$ , will be used to store the base pairing probabilities of  $i \cdot j$  gap-spanning pairs in pseudoknots. The most pertinent recursion probability,  $P_{i,d,e,j}^{g}$ , stores the probability of a gap structure with outer gap-spanning pair  $i \cdot j$  and inner gap-spanning pair  $d \cdot e$  corresponding to the partition function recursion  $Q_{i,d,e,j}^{g}$  (see Fig. 19 in ref. 18). Due to the structure of the  $Q_{i,d,e,j}^{g}$  recursion, the sum of  $P_{i,d,e,j}^{g}$  over all values of d, e precisely gives the probability of an outer pair  $i \cdot j$ 

$$P_{i,j}^{bg} + = \sum_{i < d < e < j} P_{i,d,e,j}^{g}.$$
 (4)

However, the sum of  $P_{i,d,e,j}^g$  over all values of *i*, *j* does *not* give the probability of an inner pair  $d \cdot e$ , because the same inner pair may be present in multiple recursions required to define the same secondary structure. To correctly determine the probabilities of inner gap-spanning pairs, only the portion of  $P^g$  that corresponds to calling  $Q^g$  directly from  $Q^{g^l}$  should be included

$$P_{df}^{bg} + = \sum_{i \le d \le e \le f \le j} P_{i,e,f,j}^{gl} Q_{i,d,f,j}^{g} Q_{d+1,e}^{z} \exp(-\beta_2 / RT) / Q_{i,e,f,j}^{gl}.$$
 (5)

Here,  $Q^{gl}$  and  $Q^z$  are partition function recursions used to define the interior structure of a pseudoknot, and  $\beta_2$  is a pseudoknot energy parameter (see Figs. 18 and 12 in ref. 18). Allowing pseudoknots, the total probability of a base pair  $i \cdot j$  is then

$$p_{i,j} = P_{i,j}^b + P_{i,j}^{bg}.$$

Pseudocode detailing the algorithm for computing recursion probabilities in the pseudoknotted case is provided in Figure 5, where the calculation of  $P_{i,j}^{bg}$  using (4) and (5) has been embedded at little additional cost. [Note that (4) and (5) use different indices for  $P^{bg}$  to maintain consistency with the pseudocode.] In the Appendix, we describe how to reduce the complexity of the pseudoknotted algorithm from  $O(N^6)$  to  $O(N^5)$ .

# Methods

The standard energy model<sup>4</sup> and pseudoknot extension<sup>18</sup> are implemented as described previously,<sup>18</sup> including dangle energies and penalties for helices not terminated by a  $G \cdot C$  pair. These terms do not change the structure of the recursions described in the pseudocode and are omitted for clarity. Coaxial stacking contributions are not included in the physical model, as it is unclear how

to treat different stackings associated with the same secondary structure in the context of the partition function. To maintain consistency with a recent design study,<sup>23</sup> dangle energies are treated analogously to the d2 option in the Vienna package.<sup>10</sup> Following this approach, dangle energies are included even if two helices are separated by one or zero bases, providing some compensation for the neglect of coaxial stacking bonuses.

# Applications

The recursion probability algorithm provides a simple, general method for calculating the frequency of various substructures in the ensemble of states for a given nucleic acid. Base-pairing probabilities derived from the recursion probabilities are particularly useful for analyzing secondary structure via dot plot analyses.<sup>11</sup> A traditional dot plot depicts the probabilities of forming all possible  $i \cdot j$  base pairs. In the case of pseudoknots, the dot plot can be decomposed into two dot plots—one for nested pairs and one for nonnested gap-spanning pairs.

To see the utility of this decomposition, calculations were run on wild-type and mutant sequences of a pseudoknot construct derived from human telomerase RNA.<sup>20</sup> Experimental evidence suggests that this pseudoknot exists in equilibrium with an alternative, hairpin structure, and that this equilibrium functions as a biological switch.<sup>19</sup> A two-point mutant, found in a small percentage of people, shifts the equilibrium towards the hairpin structure, leading to a disease known as dyskeratosis congenita.<sup>19</sup> Feigon and coworkers<sup>20</sup> examine this shift in equilibrium for segments of the wild-type and mutant sequences described in Figure 1, revealing that the pseudoknot conformation dominates the hairpin for the wild-type sequence ( $\sim 95\%$  to  $\sim 5\%$ ) but competes roughly equally in the mutant sequence ( $\sim 50\%$  to  $\sim 50\%$ ). Using preliminary pseudoknot parameters,18 energies were computed for both the wild-type sequence and the two-point mutant on the pseudoknotted and hairpin structures. The predicted energies in Table 1 match reasonably well with experimental values.<sup>20</sup> For the wild-type sequence, the disparity between the pseudoknot and hairpin energies suggests an equilibrium that favors the more stable pseudoknot. In contrast, the energies for the double mutant sequence suggest a more balanced equilibrium. Figures 6 and 7 illustrate that the hairpin conformation has a significant impact on the pair probabilities for the mutant, but not for the wild-type sequence.

Base-pairing probabilities can also be used to construct metrics for evaluating nucleic acid designs. The secondary structure *s* may be described by a symmetric  $N \times N$  matrix *S* with entries  $S_{i,j} =$ 1 if *s* contains base pair  $i \cdot j$  and  $S_{i,j} = 0$  otherwise. We augment this matrix by an additional column with entries  $S_{i,N+1} = 1$  if base *i* is unpaired and  $S_{i,N+1} = 0$  otherwise. Hence, each row sum is one. For a given sequence of length *N*, the metric<sup>23</sup>

$$n(s^*) = N - \sum_{\substack{1 \le i \le N \\ 1 \le j \le N+1}} p_{i,j} S^*_{i,j}$$

represents the average number of nucleotides that differ from the target secondary structure  $s^*$  at thermodynamic equilibrium. This

Compute  $Q, Q^b, Q^m, Q^p, Q^z, Q^g, Q^{gl}, Q^{gr}, Q^{gls}, Q^{grs}$ // Set all  $Q^x$ -type values to 0 for  $O(N^5)$  version Set all P-type values to 0 Set an *P*-type values to 0  $P_{1,N} = 1$  //probability of "recursing" to the entire strand is 1 for l = N, 1 //decrement // Initialize  $Q^x = Q^{x1}, Q^{x1} = Q^{x2}, Q^{x2} = 0$  for  $O(N^5)$  version // Initialize  $P^x = P^{x1}, P^{x1} = P^{x2}, P^{x2} = 0$  for  $O(N^5)$  version for i = 1, N - l + 1j = i + l - 1//P,  $P^m$ ,  $P^z$  recursions for d = i, j-4 // all possible rightmost pairs  $d \cdot e$ for e = d+4, j $\Delta p = P_{i,j} Q_{i,d-1} \, Q^{\mathsf{b}}_{d,e} / Q_{i,j}$  $\begin{array}{l} & \sum_{p_{i,d-1} \neq i, d-1} \varphi_{d,e} \neq i, j \\ P_{i,d-1}^{i} = \Delta p \\ P_{d,e}^{b} + = \Delta p \\ P_{d,e}^{b} + = P_{i,j}^{m} \exp\{-[\alpha_{2} + \alpha_{3}(d-i) + \alpha_{3}(j-e)]/RT\} Q_{d,e}^{b} / Q_{i,j}^{m} \end{array}$  $\begin{array}{l} \sum_{i,j=1}^{n} \sum_{i,j=1}^{m} \sum_{i=1}^{m} \sum_{i=1}^{m} \sum_{j=1}^{m} \sum_{i=1}^{m} \sum_{j=1}^{m} \sum_{i=1}^{m} \sum_{j=1}^{m} \sum_{i=1}^{m} \sum_{j=1}^{m} \sum_{i=1}^{m} \sum_{j=1}^{m} \sum_{i=1}^{m} \sum_{j=1}^{m} \sum_{j=1}$  $P_{d,e}^{b} += \Delta p$  $\begin{aligned} & \Delta p = P_{i,j}^z Q_{i,d-1}^z Q_{d,e}^b \exp\{-[\beta_2 + \beta_3(j-e)]/RT\}/Q_{i,j}^z \\ & P_{i,d-1}^z + \Delta p \\ & P_{d,e}^b + \Delta p \end{aligned}$ for d = i, j-8 // all possible rightmost pseudoknots filling [d, e]for e = d+8, j $\Delta p = P_{i,j} Q_{i,d-1} Q_{d,e}^p \exp\{-\beta_1 / RT\} / Q_{i,j}$  $P_{i,d-1} += \Delta p$  $P_{i}^{p} += \Delta p$  $\begin{array}{l} P_{d,e}^{p_{d-1}} = \Delta p \\ P_{d,e}^{p_{d-1}} + = \Delta p \\ P_{d,e}^{p_{d-1}} + P_{i,j}^{m} \exp\{-[\beta_{1}^{m} + 2\alpha_{2} + \alpha_{3}(d-i+j-e)]/RT\}Q_{d,e}^{p}/Q_{i,j}^{m} \\ \Delta p = P_{i,j}^{m}Q_{i,d-1}^{m}Q_{d,e}^{p}\exp\{-[\beta_{1}^{m} + 2\alpha_{2} + \alpha_{3}(j-e)]/RT\}/Q_{i,j}^{m} \end{array}$  $P_{i,d-1}^m += \Delta p$   $P_{i,d-1}^m += \Delta p$  $p_{d,e}^{p} + = \Delta p$  $\Delta p^{a,e} = P_{i,j}^{z} Q_{i,d-1}^{z} Q_{d,e}^{p} \exp\{-[\beta_{1}^{p} + 2\beta_{2} + \beta_{3}(j-e)]/RT\}/Q_{i,j}^{z}$  $P_{i,d-1}^{z} += \Delta p$   $P_{d,e}^{p} += \Delta p$   $//P^{p} \text{ recursion}$ for d = i+2, j-4for  $e = \max(d+2, i+5), j-3$ for  $f = e^{+1}, j^{-2}$   $\Delta p = P^p_{i,j} Q^{gr}_{d,d-1,e,f} Q^{gr}_{d,e-1,f+1,j} / Q^p_{i,j}$  $P_{i,d-1,e,f}^{gl} + = \Delta p$   $P_{d,e-1,f+1,j}^{gr} + = \Delta p$   $// P_{d,e-1,f+1,j}^{gr} + = \Delta p$ for d = i+1, j-4for e = d+3, j-1for f = e, j-1  $\Delta p = P_{i,d,e,j}^{gr} Q_{i,d,f,j}^{gl} Q_{e,f-1}^{z} / Q_{i,d,e,j}^{gr}$  $\begin{array}{c}P_{i,d,f,j}^{gl} += \Delta p\\P_{e,f-1}^{z} += \Delta p\end{array}$  $// P^{gl}$  recursion for d = i + 1, j - 5for f = d + 4, j - 1for f = d+4, j-1for e = d, f-3  $\Delta p = P_{i,e,f,j}^{gl} Q_{i,d,f,j}^{g} Q_{d+1,e}^{z} \exp(-\beta_2/RT)/Q_{i,e,f,j}^{gl}$   $P_j^{gd} Q_{i,f,j}^{gl} = \Delta p$   $P_{d+1,e}^{dg} + \Delta p$   $P_{d,f}^{bg} = \Delta p //P^{bg}$  inner gap-spanning base-pairing prob  $P_j^{dgrs}$  recursion for d = i+1, j-10for e = d+4, j-6 $\begin{array}{l} \int e = a + 4, j - 5 \\ \text{for } f = e + 1, j - 5 \\ \Delta p = P_{i,d,e,j}^{grs} Q_{i,d,e,f}^{g} Q_{f+1,j}^{m} \exp(-\alpha_2/RT) / Q_{i,d,e,j}^{grs} \\ P_{i,d,e,f}^{g} + = \Delta p \\ P_{f+1,j}^{m} + = \Delta p \end{array}$ 

 $// P^{gls}$  recursion for c = i+5, j-6for d = c+1, j-5for e = d + 4, j - 1  $\Delta p = P_{i,d,e,j}^{gl_s} \exp(-\alpha_2/RT)Q_{i,c-1}^m Q_{c,d,e,j}^g/Q_{i,d,e,j}^{gl_s}$   $P_{i,d,e,j}^m = \Delta p$   $P_{c,d,e,j}^{gl_s} + = \Delta p$  $P_{c,d,e,j}^{j,c-1} += \Delta p$ //  $P^{g}$  recursion for d = i+2, j-6 // simple interior loops for e = d + 4, j - 2for c = i+1, d-1for f = e+1, j-1 $\begin{array}{l} \text{IO} f = e^{-1}, j-1\\ P_{c,d,e,f}^g + = P_{i,d,e,j}^g \exp(-G_{i,c,f,f}^{\text{interior}}/RT)Q_{e,d,e,f}^g/Q_{i,d,e,j}^g\\ // \text{ for } O(N^5) \text{ version, replace provious five lines with:}\\ // \text{ call fastiloops} N5(i, j, l, Q^g, Q^x, Q^{x2}, P^g, P^x, P^{x2}) \end{array}$ for d = i + 6, j - 5 $\begin{array}{l} \text{for } e = d + 4, j - 1 \\ \Delta p = P^g_{i,d,e,j} Q^m_{i+1,d-1} \exp\{-[\alpha_1 + 2\alpha_2 + \alpha_3(j - e - 1)]/RT\}/Q^g_{i,d,e,j} \\ P^m_{i+1,d-1} + = \Delta p \end{array}$ for d = i+1, j-10for e = d+4, j-6 $\Delta p = P_{i,d,e,j}^{g} \exp\{-[\alpha_1 + 2\alpha_2 + \alpha_3(d-i-1)]/RT\}Q_{e+1,j-1}^m/Q_{i,d,e,j}^g$ 
$$\begin{split} & \Delta p = i_{i,d,e,j} \exp[-|\alpha_1 + 2\alpha_2 + \alpha_3(a - i - i)]/Rt \int \mathcal{Q}_{e+1,j-1}/\mathcal{Q}_i \\ & P_{e+1,j-1}^m + \Delta p \\ & \text{for } d = i + 6, j - 10 \\ & \text{for } e = d + 4, j - 6 \\ & \Delta p = P_{i,d,e,j}^g \mathcal{Q}_{i+1,d-1}^m \exp\{-[\alpha_1 + 2\alpha_2]/Rt\} \mathcal{Q}_{e+1,j-1}^m / \mathcal{Q}_{i,d,e,j}^g \\ & P_{i+1,d-1}^m + \Delta p \\ & P_{e+1,j-1}^m + \Delta p \\ & \text{for } d = i + 7, j - 6 \\ & \text{for } d = i + 7, j - 6 \end{split}$$
for e = d + 4, j - 2for f = e+1, j-1 $\Delta p = P_{i,d,e,j}^{g} Q_{i+1,d,e,f}^{gls} \exp\{-[\alpha_1 + \alpha_2 + \alpha_3(j-f-1)]/RT\}/Q_{i,d,e,j}^{g}$  $P_{i+1,d,e,f}^{gls} \stackrel{i,d,e,j \neq i_{i+1}}{\models} \Delta p$ for d = i+2, j-11for e = d + 4, j - $\begin{array}{l} \text{for } e = a + \alpha_{3,j} - i \\ \text{for } c = i + 1, d - 1 \\ \Delta p = P_{i,d,e,j}^{g} \exp\{-[\alpha_{1} + \alpha_{2} + \alpha_{3}(c - i - 1)]/RT\}Q_{c,d,e,j-1}^{grs} / Q_{i,d,e,j}^{g} \\ P_{c,d,e,j-1}^{grs} + \pm \Delta p \\ \text{for } d = i + 7, j - 11 \\ \end{array}$  $\begin{array}{l} \sum_{i=1}^{n} c_{i} = a + 4, j - 7 \\ \text{for } c_{i} = i + 6, d - 1 \\ \sum_{i,d,e,j} \Delta p = P_{i,d,e,j}^{g} Q_{i+1,c-1}^{a} Q_{c,d,e,j-1}^{grs} \exp\{-[\alpha_{1} + \alpha_{2}]/RT\}/Q_{i,d,e,j}^{g} \\ P_{i+1,c-1}^{m} + = \Delta p \\ P_{c,d,e-i-1}^{grs} + = \Delta p \end{array}$ for e = d + 4, j - 4 $P_{c,d,e,j-1}^{q+1} += \Delta p$ //  $P^{bg}$  outer gap-spanning base-pairing prob for d = i+1, j-5for e = d + 4, j - 1 $P_{i,j}^{bg} + = P_{i,d,e,j}^{g}$ //  $P^{b}$  recursion for d = i+1, j-5 // all possible rightmost pairs  $d \cdot e$ for e = d + 4, j - 1For e = d + 4, j - 1  $P_{d,e}^{b} + = P_{i,j}^{b} \exp\left(-G_{i,d,e,j}^{\text{interior}}/RT\right)Q_{d,e}^{b}/Q_{i,j}^{b}$   $\Delta p = P_{i,j}^{b}Q_{i+1,d-1}^{m}$   $P_{i+1,d-1}^{m} + \Delta p$   $P_{d,e}^{b} + \Delta p$ for d = i+1, j-9 // all possible rightmost pseudoknots filling [d, e]for a = d+2, i-1for e = d+8, j-1 $G^{\text{recursion}} = \alpha_1 + \beta_1^m + 3\alpha_2 + \alpha_3(j - e - 1)$  $\begin{array}{c} {}_{d,e} := r_{i,j} \exp\{-[G^{\text{recursion}} + \alpha_3(d-i-1)]/RT\}Q_d^p \\ \Delta p = P_{i,j}^b Q_{i+1,d-1}^m Q_{d,e}^p \exp\{-G^{\text{recursion}}/RT\}/Q_{i,j}^b \\ P_{i+1,d-1}^m + = \Delta p \\ P_{d,e}^p + = \Delta p \end{array}$  $P_{d,e}^{p} + = P_{i,j}^{b} \exp\{-[G^{\text{recursion}} + \alpha_{3}(d-i-1)]/RT\}Q_{d,e}^{p}/Q_{i,j}^{b}$ 

**Figure 5.**  $O(N^6)$  recursion probability algorithm that includes a class of pseudoknots. Modifications required to produce an  $O(N^5)$  version of the algorithm are noted in comments. See the Appendix for details.

Table 1. Energ	gy Comparisons	for Human Te	lomerase RNA	Constructs.
----------------	----------------	--------------	--------------	-------------

RNA		Energies (kcal/mol)		
	Conformation	$\Delta G_{\mathrm{exp}}$	$\Delta G_{ m calc}$	
Wild-type	Pseudoknot Hairpin	-17.8 -9.8ª	-18.5 <sup>b</sup> -11.5 <sup>c</sup>	
Mutant	Pseudoknot Hairpin	$-11.2 -10.5^{a}$	-11.3 <sup>b</sup> -11.5 <sup>c</sup>	

<sup>a</sup>Experiments were performed on partial sequences that excluded the 18 nucleotides on the 3' end to prevent the formation of pseudoknots.<sup>20</sup> This truncation does not affect the corresponding  $\Delta G_{calc}$ .

<sup>b</sup>A related pseudoknot structure that is otherwise identical but omits the three consecutive  $A \cdot U$  pairs in the stem with the bulge loop is predicted to be 0.5 kcal/mol more stable.

<sup>c</sup>The secondary structure energy calculation ignores the four consecutive noncanonical base pairs that are observed to close the interior loop in the hairpin stem.<sup>20</sup>

is a less stringent metric than  $p(s^*)$ , the probability that the sequence exactly adopts structure  $s^*$ ; even if  $p(s^*)$  is not close to unity,  $n(s^*)$  can still be small if the equilibrium ensemble is dominated by structures that differ only slightly from  $s^*$ .

It is illustrative to compare the two metrics on a real design problem involving pseudoknots. For example, Winfree et al.<sup>14</sup> designed and constructed DNA double-crossover molecules<sup>24</sup> that interact to form a two-dimensional lattice with a pseudoknotted unit cell. These sequence designs were performed using sequence symmetry minimization<sup>25</sup> to ensure that incorrectly paired subsequences of length six would always contain at least one mismatch and most incorrectly paired subsequences of length five would also contain a mismatch.<sup>14</sup> Lacking DNA pseudoknot parameters, we examine an RNA analog of their sequence for the portion of the pseudoknotted unit cell depicted in Figure 8a. The probability of adopting the target structure is  $p(s^*) = 0.1$  and the average number of incorrect nucleotides is  $n(s^*) = 4.0$ . The low value of



**Figure 6.** Dot plots for wild-type human telomerase RNA. (a) Pseudoknot (bottom left) and hairpin (top right) constructs. For (b) and (c), large dots indicate a  $p_{i,j} \ge 0.5$  and small dots indicate  $0.5 > p_{i,j} \ge 0.05$ . (b) Base-pairing probabilities including pseudoknots (bottom left) and excluding pseudoknots (top right). (c) A decomposition of the full base-pairing probabilities into gap-spanning pairs (bottom left) and nested pairs (top right). Note that there are no nested pairs with significant probability, indicating that pseudoknot conformations are dominating the equilibrium.



**Figure 7.** Dot plots for double mutant human telomerase RNA. The plots are analogous to those of Figure 6. The key difference is observed in (c), where the hairpin stem appears as both gap-spanning pairs and nested pairs, indicating the increased significance of hairpin conformations.

 $p(s^*)$  might possibly indicate a cause for concern, but for a structure with 90 nucleotides and helices of length eight, the average number of incorrect nucleotides is relatively small. Hence, it is not surprising that the sequence behaves well experimentally, demonstrating the correct base-pairing topology despite slight predicted variations at the ends of helices. The dot plot in Figure 8b illustrates the similarity between the average structure and the desired target.

Interestingly, design methods described in previous work<sup>23</sup> can be used, in conjunction with the pseudoknot partition function algorithm, to find sequences that achieve  $p(s^*) = 0.98$  and  $n(s^*) \ll 1$ . It is unclear whether these sequences would provide any experimental benefit for this system (even given a perfect energy model), because the difference between  $n(s^*) \approx 4$  and  $n(s^*) \ll 1$ may be lost in experimental noise. By contrast, if a sequence produced  $p(s^*) = 0.1$  with  $n(s^*) \gg 4$ , then the equilibrium ensemble could include important structures differing significantly from the target structure.

# Conclusions

A general transformation rule extends nucleic acid partition function algorithms to calculate recursion probabilities, which in turn, can be used to compute base-pairing probabilities. We use this approach to derive an algorithm for computing base-pairing probabilities starting from a partition function algorithm that includes a class of pseudoknots. The same strategy will apply to future partition function extensions that follow the same dynamic programming paradigm.

To demonstrate the utility of base-pairing probabilities, calculations were performed on a pseudoknot/hairpin construct thought to represent an important biological switch. In agreement with experimental evidence, the computational results indicate that the pseudoknot dominates the hairpin for the wild-type sequence, but not for the double mutant. Base-pairing probabilities were also used to examine the ensemble properties of a synthetic nucleic acid sequence designed to assemble into a pseudoknotted double-crossover molecule. The average number of incorrect nucleotides was found to be small, suggesting that the relatively low computed probability of adopting the



Figure 8. Computational examination of a pseudoknotted DNA nanostructure. (a) Secondary structure for a double-crossover molecule that forms a portion of the unit cell in a two-dimensional lattice.<sup>14</sup> For our computational study, we join the blue and orange strands (arrows denote 3') into a single strand using auxiliary nucleotides (green) to facilitate the use of the single-stranded partition function algorithm.<sup>18</sup> In the absence of DNA pseudoknot parameters, we consider the RNA analog 5'-CCAACUCCUAGCGAUUUUUCGCUAGGUUUACCA-GAUCCACAAGCCGACGUUACA-UUUU-GGAUCUGGUAAG-UUGGUGUAACGUCGGCUUGU-3', where the interior hyphens denote the boundaries of the auxiliary linker segment. (b) Dot plot analysis of the designed sequence. The bottom left depicts the basepairs in the target structure, and the upper right depicts the base-pairing probabilities. Large dots indicate a  $p_{i,j} \ge 0.5$  and small dots indicate  $0.5 > p_{i,i} \ge 0.05$ . The circles indicate the major differences between the target structure and the calculated pair probabilities.

target secondary structure should not significantly affect the experimental performance of the molecule.

#### Software Download

The algorithms described in this article are available for download at http://www.acm.caltech.edu/~niles as part of the NUPACK software suite.

# Acknowledgments

We wish to thank C. Ueda for discussions on human telomerase RNA and E. Winfree for discussions on the DNA lattice.

## **Appendix: Reducing Computational Complexity**

The algorithms presented in the main text provide an inefficient treatment of interior loops. By exploiting the form of the interior loop potential function, the computational complexity of the partition function algorithms excluding and including pseudoknots can be reduced by a factor of N, where N is the sequence length.<sup>18,21</sup> A detailed description of the "fastiloops" treatment is provided in ref. 18 and the corresponding Supplementary Material. The "fastiloops" modification detracts from the simplicity of the presentation because the necessary recursions do not conform to the same structure as the other terms in the algorithm. Here, we describe the extension of this approach to recursion probability algorithms.

In the unpseudoknotted case, pseudocode for an  $O(N^3)$  partition function algorithm is provided in Figure 11 of ref. 18, which employs the "fastiloops" function of Supplementary Material Figure S2. To this point, we have assumed that all Q-type values are accessible at the end of the partition function calculation. For the "fastiloops" methods, the values  $Q^x$ ,  $Q^{x1}$  and  $Q^{x2}$  are computed on the fly and discarded to save memory. Hence, for the recursion probability algorithm, it is necessary to recompute the  $Q^{x}$ -type terms at the same time that the corresponding  $P^{x}$ -type terms are calculated. An  $O(N^3)$  recursion probability algorithm that excludes pseudoknots is described in Figure A1, which references the function "fastiloopsN3" of Figure A2. If pseudoknots are included, the computational complexity of the recursion probability algorithm in Figure 5 is reduced to  $O(N^5)$  using "fastiloopsN5" described in Figure A3. A few aspects of the "fastiloopsN3" and "fastiloopsN5" routines deserve mention. It is advisable to review the relevant sections of ref. 18 and the corresponding Supplementary Material before proceeding.

An interior loop with closing pair  $i \cdot j$  and interior pair  $d \cdot e$  has energy  $G_{i,d,e,j}^{\text{interior}}$ , sides of lengths

$$L_1 \equiv (d - i - 1), \quad L_2 \equiv (j - e - 1),$$
 (6)

and size  $L_1 + L_2$ . Loops with both  $L_1 \ge 4$  and  $L_2 \ge 4$  are termed "extensible" and their contributions to the partition function algorithm are calculated using  $Q^x$ . Furthermore,  $Q^x$  also contains information about "possible extensible loops" for which the definitions of  $L_1$ ,  $L_2$  are the same but *i* and *j* are not required to base-pair.

The partition function algorithm examines subsequences of length l = j - i + 1, starting with l = 1 and ending with l = N.  $Q^x$  is efficiently calculated using the extension identity [see eq. (15) of ref. 18],

$$Q_{i-1,s+2}^{x} = \Gamma(s+2) \Big|_{L_{1}+L_{2}=s+2}^{L_{1}=4} + \Gamma(s+2) \Big|_{L_{1}+L_{2}=s+2}^{L_{2}=4} + \left[ Q_{i,s}^{x} \exp\{-\left[\gamma_{1}(s+2) - \gamma_{1}(s)\right]/RT\} \right]$$
(7)

Compute  $Q, Q^b, Q^m, Q^s, Q^{ms}$  using  $O(N^3)$  partition function algorithm Initialize  $(Q^x, Q^{x1}, Q^{x2}, P, P^b, P^m, P^s, P^{ms}, P^x, P^{x1}, P^{x2}) // O(N^2)$  space Set all  $Q^x$ -type and P-type values to 0  $P_{1,N} = 1 //$ probability of "recursing" to the entire strand is 1 for l = N, 1 //dccrement subsequence length Initialize  $Q^x = Q^{x1}, Q^{x1} = Q^{x2}, Q^{x2} = 0$ Initialize  $P^x = P^{x1}, P^{x1} = P^{x2}, P^{x2} = 0$ for i = 1, N - l + 1 j = i + l - 1  $//P, P^m$  recursions for d = i, j - 4  $\Delta p = P_{i,j}Q_{i,d-1}Q_{d,j}^s/Q_{i,j}$   $P_{i,d-1}^s + \Delta p$   $P_{d,j}^{s,j} + \Delta p$   $P_{d,j}^{m,j} + D_{p}$   $P_{d,j}^{m,j} + \Delta p$   $P_{d,j-1}^{m,j} + \Delta p$   $P_{d,j-1}^{m,j} + \Delta p$  $P_{d,j-1}^{m,j} + \Delta p$ 

**Figure A1.**  $O(N^3)$  recursion probability algorithm that excludes pseudoknots. The algorithm proceeds from longer subsequences to shorter ones, so in contrast to the analogous partition function algorithm (see Fig. 11 of ref. 18),  $Q^{x1}$  and  $Q^{x2}$  refer to subsequences whose lengths are shorter (by 1 and 2, respectively) than the current subsequence of length *l*.

which relates  $Q_{i,s}^x$  (for subsequences of length l) to  $Q_{i-1,s+2}^x$  (for subsequences of length l + 2). The first line "seeds"  $Q^x$  with cases at an extension border ( $L_1 = 4$  or  $L_2 = 4$ ) for subsequent extension to longer subsequences. For conciseness, we have introduced the definition

$$\Gamma(s) \equiv \exp\{-[\gamma_1(s) + \gamma_2(|L_1 - L_2|) + \gamma_3(e, d, e+1, d-1)]/RT\}Q_{d,e}^b$$

where *d* and *e* are defined implicitly in terms of  $L_1$  and  $L_2$ . For implementation purposes, the second line of (7) is calculated during the *l*, *i* loop and temporarily stored in  $Q_{i-1,s+2}^{x^2}$ . The first line of (7) is added to this contribution in the l + 2, i - 1 loop. As a result of this two step procedure, we adopt the convention that  $L_1$  and  $L_2$  are defined with respect to the loop index in which they are calculated (i.e., *l*, *i* for the second line and l + 2, i - 1 for the first line). This convention facilitates the comparison of the extension identity with pseudocode.

The recursion probability algorithm examines subsequences of length l starting with l = N and ending with l = 1. To recompute  $Q^x$  in this context, we use the contraction identity

$$Q_{i+1,s-2}^{x} = \sum_{\substack{L_{1} \geq 4, L_{2} \geq 4\\L_{1}+L_{2}=s-2}} \Gamma(s-2) + \sum_{\substack{L_{1} \geq 4, L_{2} \geq 4\\L_{1}+L_{2}=s-2}} \Gamma(s-2) + \left[ (Q_{i,s}^{x} - \Gamma(s) \Big|_{L_{1}+L_{2}=s} - \Gamma(s) \Big|_{L_{1}+L_{2}=s} \right] + \left[ (Q_{i,s}^{x} - \Gamma(s) \Big|_{L_{1}+L_{2}=s} - \Gamma(s) \Big|_{L_{1}+L_{2}=s} \right] \exp\{-\left[\gamma_{1}(s-2) - \gamma_{1}(s)\right]/RT\}\right]$$
(8)

which relates  $Q_{i,s}^x$  (for subsequences of length l) to  $Q_{i+1,s-2}^x$  (for subsequences of length l-2). The first line "seeds"  $Q^x$  with cases that are both extensible ( $L_1 \ge 4$  and  $L_2 \ge 4$ ) and at an end of the strand (i = 1 or j = N). For implementation purposes, the second line of (8) is calculated during the l, i loop and temporarily stored in  $Q_{i+1,s-2}^{x2}$ . The first line of (8) is added to this contribution in the l-2, i + 1 loop. We retain the convention that  $L_1$  and  $L_2$  are defined with respect to the loop index in which they are calculated (i.e., l, i for the second line and l-2, i + 1 for the first line).

Derivation of the algorithm to compute  $P^x$  requires careful consideration. The quantities  $Q^x$  and  $Q^{x2}$  contain incomplete partition function information for "possible extensible loops," but they do not represent subsequence partition functions in the manner of other Q-type matrices. In a normal recursion relation,

```
function fastiloopsN3(i, j, l, Q^b, Q^x, Q^{x2}, P^b, P^x, P^{x2})
  //Add small non-contractible interior loop terms to P^{\dot{b}} as special cases
 for L_1 = 0, 3
d = i + L_1 + 1
         for L_2 = 0, min(3, j-d-5)
      For D_2 = 0, matrix, j = 1,

e = j - L_2 - 1

P_{d,e}^b += P_{i,j}^b \exp\{-G_{i,d,e,j}^{internal}/RT\} Q_{d,e}^b/Q_{i,j}^b

/Add bulge loops and large asymmetric loops as special cases
 for L_1 = 0, 3 //Cases L_1 = 0, 1, 2, 3, L_2 \ge 4
d = i + L_1 + 1
         for L_2 = 4, j - d - 5
 e = j - L_2 - 1
P_{d,e}^b += P_{i,j}^b \exp\{-G_{i,d,e,j}^{\text{internal}}/RT\} Q_{d,e}^b/Q_{i,j}^b
for L_2 = 0, 3//\text{Cases } L_1 \ge 4, L_2 = 0, 1, 2, 3
          e = j - L_2 - 1
        \begin{array}{l} e = j - D_2 - 1\\ \text{for } L_1 = 4, e - i - 5\\ d = i + L_1 + 1\\ P_{d,e}^b + = P_{i,j}^b \exp\{-G_{i,d,e,j}^{\text{internal}}/RT\} Q_{d,e}^b/Q_{i,j}^b\\ / \operatorname{Seed} Q^x \text{ with contractible cases} \end{array}
             Also add cases that are at an end with L_1 \ge 4, L_2 \ge 4
  if (i = 1 \text{ or } j = N) and l \ge 15
         for d = i + 5, j - 9
L_1 = d - i - 1
                 for e = d + 4, j - 5
    L_2 = j - e - 1
s = L_1 + L_2
G^{\text{partial}} = \gamma_1(s) + \gamma_2(|L_1 - L_2|) + \gamma_3(e, d, e+1, d-1)
Q_{i,s}^x + \exp\{-G^{\text{partial}}/RT\} Q_{d,e}^b
//Use Q^x to finish calculation of P^x.
() Use Q_{i} to infinite distribution of i

if (sequence permits i \cdot j base pair)

for s = 8, l - 7

P_{i,s}^{x} += P_{i,j}^{b} Q_{i,s}^{x} \exp\{-\gamma_{3}(i, j, i+1, j-1)/RT\}/Q_{i,j}^{b}

//Calculate P^{b} contribution using Q^{x} and P^{x}
 if (l \ge 15) // smallest subsequence not added to P^b as special case
         L_1=4 // explicitly add in terms for L_1=4, L_2\geq 4 d=i+L_1+1
          for L_2 = 4, j - d - 5
               \begin{aligned} & \sum_{l=1}^{T} \sum_{i=1}^{L_{2}} \sum_{j=1}^{L_{2}} \sum_{j=1}^{L_{2}} \sum_{i=1}^{L_{2}} \sum_{j=1}^{L_{2}} \sum_{j=1}^{L_{
                                                                                                                                                                                                                                                                                                                 (*)
        \begin{array}{l} \sum_{i,s} P_{i,s} = p \\ P_{d,e}^{d} = \Delta p \\ P_{i,s}^{d} = \sum p \\ P_{i,s}^{d} = \sum p \\ P_{i,s}^{d} = \exp\{-G^{\text{partial}}/RT\} Q_{d,e}^{b} / / \text{ Remove border cases} \\ Q_{i,s}^{2} = \exp\{-G^{\text{partial}}/RT\} Q_{d,e}^{b} / / \text{ Remove border cases} \\ L_{2} = 4 \\ / / \text{ explicitly add in terms for } L_{1} \geq 5, L_{2} = 4 \end{array}
              = j - L_2 - 1
         for L_1 = 5, e - i - 5
d = i + L_1 + 1
                 Insert (*)
   // Store partial values for Q^{x2} and P^{x2}
\begin{array}{l} \int \sigma(s) = 10, l-7 \\ Q_{i+1,s-2}^{x} = Q_{i,s}^{x} \exp\{-[\gamma_{1}(s-2) - \gamma_{1}(s)]/RT\} \\ P_{i+1,s-2}^{x} = P_{i,s}^{x} \end{array}
```

**Figure A2.** Pseudocode for computing interior loop contributions to  $P^b$  in  $O(N^3)$  as an alternative to the  $O(N^4)$  interior loop recursion of Figure 4.

function fastiloopsN5 $(i, j, l, Q^g, Q^x, Q^{x2}, P^g, P^x, P^{x2})$ for d = i + 1, j - 5for e = d + 4, j - 1//Add small non-contractible interior loop terms to  $P^g$  as special cases for  $L_1 = 0$ , min(3, d - i - 2) $c = i + L_1 + 1$  $c = i + L_1 + 1$ for  $L_2 = 0$ , min(3, j - e - 2)  $f = j - L_2 - 1$  $P_{c,d,e,f}^g + = P_{i,d,e,j}^g \exp\{-G_{i,c,f,j}^{internal}/RT\} Q_{c,d,e,f}^g / Q_{i,d,e,j}^g$ //Add bulge loops and large asymmetric loops as special cases for  $L_1 = 0$ , min(3, d - i - 2) //Cases  $L_1 = 0, 1, 2, 3, L_2 \ge 4$  $c = i + L_1 + 1$ for  $L_2 = 4, j - e - 2$ f = j - L<sub>2</sub> - 1  $F_{e,d,e,f}^g + = P_{e,d,e,j}^g \exp\{-G_{i,c,l,j}^{\text{internal}}/RT\} Q_{e,d,e,f}^g/Q_{i,d,e,j}^g$ for  $L_2 = 0, \min(3, j - e - 2)//\text{Cases } L_1 \ge 4, L_2 = 0, 1, 2, 3$  $f = j - L_2 - 1$ for  $L_1 = 4, d - i - 2$  $\begin{array}{l} \text{ for } L_1 = i_1 L_1 + 1 \\ c = i + L_1 + 1 \\ P_{c,d,e,f}^g = P_{i,d,e,f}^g \exp\{-G_{i,c,f,j}^{\text{internal}}/RT\} Q_{c,d,e,f}^g/Q_{i,d,e,j}^g \\ // \operatorname{Seed} Q^x \text{ with contractible cases} \end{array}$ // Seed Q with contractine cases // Also add cases that are at an end with  $L_1 \ge 4, L_2 \ge 4$ if (i = 1 or j = N) and  $l \ge 17$ for d = i + 6, j - 10for e = d + 4, j - 6for c = i + 5, d - 1 $L_1 = c - i - 1$  $\begin{array}{l} L_1 = c - i - 1 \\ {\rm for} \ f = e + 1, j - 5 \\ L_2 = j - f - 1 \\ s = L_1 + L_2 \\ G^{\rm partial} = \gamma_1(s) + \gamma_2(|L_1 - L_2|) + \gamma_3(f, c, f + 1, c - 1) \\ Q^x_{1,d,c,s} + exp\{-G^{\rm partial}/RT\} Q^g_{c,d,c,f} \\ / \operatorname{Use} Q^x \ to \ fnish \ calculation \ of \ P^x \\ / \operatorname{Use} Q^x \ to \ fnish \ calculation \ of \ P^x \\ (sequence \ nermit \ i, \ base \ nermit) \end{array}$ // Use Q to much tack that of the form of Pif (sequence permits i, j base pair) for d = i + 1, j - 5for e = d + 4, j - 1for s = 8, l - 9 $P_{i,d,e,s}^x += P_{i,d,e,j}^g Q_{i,d,e,s}^x \exp\{-\gamma_3(i, j, i+1, j-1)/RT\}/Q_{i,d,e,j}^g$ // Calculate  $P^g$  contribution using  $Q^x$  and  $P^x$ if  $(l \ge 17)$ for d = i + 6, j - 10 $\begin{array}{l} \text{ or } d = i+6, j-10 \\ \text{ for } e = d+4, j-6 \\ L_1 = 4 \ // \ \text{ explicitly add in terms for } L_1 = 4, L_2 \geq 4 \\ c = i+L_1+1 \\ \text{ for } L_2 = 4, j-e-2 \\ f = j-L_2-1 \\ s = L_1+L_2 \\ \text{ G}^{\text{partial}} = \gamma_1(s) + \gamma_2(|L_1-L_2|) + \gamma_3(f,c,f+1,c-1) \\ \Delta p = P_{i,d,e,s}^{i,d} \exp\{-G^{\text{partial}}/RT\} \ \mathcal{Q}_{c,d,e,f}^{g}/\mathcal{Q}_{i,d,e,s}^{i,d} \\ P_{j-s-e}^{g} + = \Delta p \end{array}$  $\begin{array}{c} \overset{i,u,e,s}{\underset{c,d,e,s}{\overset{g}{\underset{c}}}{\underset{c}}{\underset{c}}{\underset{c}}{\underset{c}}{\underset{c}}{\underset$  $\begin{array}{l} r_{i,d,e,s} = -\Delta p \ / \ \text{Remove border cases} \\ Q_{i,d,e,s}^{x} = \exp\{-G^{\text{partial}}/RT\} \ Q_{c,d,e,f}^{g} \ / \ \text{Remove border cases} \\ L_{2} = 4 \ / \ \text{explicitly add in terms for } L_{1} \geq 5, L_{2} = 4 \end{array}$  $f = j - L_2 - 1$ for  $L_1 = 5, d - i - 2$  $c = i + L_1 + 1$ Insert (\*) / Store partial values for  $Q^{x2}$  and  $P^{x2}$ for s = 10, l - 9 $\begin{array}{l} Q_{i+1,d,e,s-2}^{x_2} = Q_{i,d,e,s}^x \, \exp\{-[\gamma_1(s-2) - \gamma_1(s)]/RT\} \\ P_{i+1,d,e,s-2}^{x_2} = P_{i,d,e,s}^x \end{array}$ 

**Figure A3.** Pseudocode for computing interior loop contributions to  $P^g$  in  $O(N^5)$  as an alternative to the  $O(N^6)$  interior loop recursion of Figure 5.

Q-type matrices on the right-hand side are subsequence partition functions describing a local structural motif that contributes to the larger subsequence partition function on the left-hand side.  $Q_{i,s}^x$ contains information about possible extensible loops that may not actually exist (if *i* and *j* are not complementary). The extension identity (7) passes this potentially useful information on to  $Q_{i-1,s+2}^{x2}$ . Consider, for example, a chain of  $Q^x$  values related by the extension identity in a case where no complementary  $i \cdot j$  base pair is encountered while incrementing *l* by 2 until an end of the strand is reached. In this scenario, the values of  $Q^x$  computed in this chain should not contribute to the corresponding recursion probabilities  $P^x$  because the values of  $Q^x$  are not identified with any secondary structure in the equilibrium ensemble. Hence, the calculation of  $P^x$  requires information about which  $Q^x$  quantities ultimately contribute to secondary structures in the ensemble. As a result, the extension identity (7) cannot simply be transformed using the standard recursion probability approach, which assumes that both sides of the equation represent subsequence partition functions that are assured of contributing to the equilibrium ensemble. This realization suggests computing  $P_{i,s}^x$  by adding the probabilities of all internal loops that rely on  $Q_{i,s}^x$  to incorporate information in the partition function.

To calculate  $P_{i,s}^x$  (for a fixed l), note that  $Q_{i,s}^x$  will be invoked for all interior loops (i', d, e, j') with interior pair  $d \cdot e$  and closing pair  $i' \cdot j'$  such that

$$i - i' = j' - j \ge 0, \quad L_1 \ge 4, \quad L_2 \ge 4, \quad L_1 + L_2 = s, \quad (9)$$

where  $L_1$ ,  $L_2$  and *s* are defined with respect to *i* and *j*. Hence, a particular loop (i', d, e, j') is identified with a set of  $Q_{i,s}^x$  terms that are related by the extension identity (7). Alternatively, a particular  $Q_{i,s}^x$  term is identified with all of the interior loops (i', d, e, j') to which it ultimately contributes via the extension identity. Consequently, from the notion of recursion probabilities introduced earlier,  $P_{i,s}^x$  (for a fixed *l*) should be the sum of the probabilities of all interior loops (i', d, e, j') that satisfy the properties (9). For the case where  $i - 1 \le N - j$  (the case  $i - 1 \ge N - j$  yields analogous results), it follows that

$$P_{i,s}^{x} = \sum_{i'=1}^{i} \sum_{\substack{L_{1} \ge 4, L_{2} \ge 4\\L_{1}+L_{2}=s}} p(i', d, e, j'),$$
(10)

where p(i', d, e, j') is the probability of the (i', d, e, j') interior loop in the equilibrium ensemble of secondary structures. Because  $P_{i+1,s-2}^{x2}$  is defined similarly, with *l* and *s* decremented by 2, it follows that

$$P_{i+1,s-2}^{x2} = \sum_{i'=1}^{i+1} \sum_{\substack{L_1 \ge 5, L_2 \ge 5\\ L_1 + L_2 = s}} p(i', d, e, j'),$$
(11)

where  $L_1$  and  $L_2$  are temporarily defined with respect to *i* and *j* to retain the size constraint  $L_1 + L_2 = s$ . Comparing (10) and (11), we then identify the relationship

$$P_{i+1,s-2}^{x2} = \sum_{\substack{L_1 \ge 5, L_2 \ge 5\\L_1 + L_2 = s}} p(i', d, e, j')|_{i'=i+1,j'=j-1} + \left[ P_{i,s}^x - \sum_{i'=1}^i p(i', d, e, j')|_{L_1 = 4, L_2 \ge 4}^{L_1 = 4, L_2 \ge 4} - \sum_{i'=1}^i p(i', d, e, j')|_{L_1 \ge 5, L_2 = 4}^{L_1 \ge 5, L_2 = 4} \right]$$

where  $L_1$  and  $L_2$  continue to be defined with respect to *i* and *j*. Finally, we shift the indices in the first line so that  $L_1$  and  $L_2$  are defined with respect to i + 1 and j - 1

$$P_{i+1,s-2}^{x^{2}} = \sum_{\substack{L_{1} \ge 4, L_{2} \ge 4\\ L_{1}+L_{2}=s-2}} p(i', d, e, j')|_{i'=i+1,j'=j-1} + \left[ P_{i,s}^{x} - \sum_{i'=1}^{i} p(i', d, e, j')|_{L_{1}=4,L_{2}=s}^{L_{1}=4,L_{2}\geq 4} - \sum_{i'=1}^{i} p(i', d, e, j')|_{L_{1}=5,L_{2}=4}^{L_{1}\geq 5,L_{2}=4} \right].$$
(12)

This identity relates  $P_{i,s}^{x}$  (for subsequences of length l) to  $P_{i+1,s-2}^{x}$  (for subsequences of length l-2). For implementation purposes, the second line is calculated during the l, i loop and temporarily stored in  $Q_{i+1,s-2}^{x2}$ . Each of the sums of form  $\sum_{i'=1}^{i}$  operates on a single term, which is a subset of the terms in the definition of  $P_{i,s}^{x}$  (10). Hence, the sums of form  $\sum_{i'=1}^{i}$  in (12) may be evaluated implicitly as  $P_{i,s}^{x}$  times a quotient with  $Q_{i,s}^{x}$  in the denominator and the corresponding subset of  $Q_{i,s}^{x}$  in the numerator. The first line is added to this contribution in the l-2, i+1 loop. There, the summation corresponds to exactly those loops treated by  $Q_{i+1,s-2}^{x}$  in the case where i+1 and j-1 base pair. As usual,  $L_1$  and  $L_2$  are defined with respect to the loop index in which they are calculated (i.e., l, i for the second line and l-2, i+1 for the first line).

# References

- 1. Tinoco, I., Jr.; Uhlenbec, O.; Levine, M. Nature 1971, 230, 362.
- 2. Turner, D. H.; Sugimoto, N.; Freier, S. Annu Rev Biophys Biophys Chem 1988, 17, 167.
- 3. SantaLucia, J., Jr. Proc Natl Acad Sci USA 1998, 95, 1460.

- 4. Mathews, D.; Sabina, J.; Zuker, M.; Turner, D. J Mol Biol 1999, 288, 911.
- 5. Zuker, M. Curr Opin Struct Biol 2000, 10, 303.
- Waterman, M. In Studies in Foundations and Combinatorics: Advances in Mathematics Supplemental Studies; Academic Press: New York, 1978, 1, 167.
- 7. Waterman, M.; Smith, T. Math Biosci 1978, 42, 257.
- Nussinov, R.; Pieczenik, J.; Griggs, J.; Kleitman, D. SIAM J Appl Math 1978, 35, 68.
- 9. Zuker, M.; Stiegler, P. Nucleic Acids Res 1981, 9, 133.
- Hofacker, I.; Fontana, W.; Stadler, P.; Bonhoeffer, L.; Tacker, M.; Schuster, P. Chem Monthly 1994, 125, 167.
- 11. McCaskill, J. Biopolymers 1990, 29, 1105.
- Bonhoeffer, S.; McCaskill, J.; Stadler, P.; Schuster, P. Eur Biophys J 1993, 22, 13.
- van Batenburg, F.; Gultyaev, A.; Pleij, C.; Ng, J. Nucleic Acids Res 2000, 28, 201.
- 14. Winfree, E.; Liu, F.; Wenzler, L.; Seeman, N. C. Nature 1998, 394, 539.
- Yan, H.; LaBean, T.; Feng, L.; Reif, J. Proc Natl Acad Sci USA 2003, 100, 8103.
- 16. Rivas, E.; Eddy, S. J Mol Biol 1999, 285, 2053.
- 17. Akutsu, T. Discrete Appl Math 2000, 104, 45.
- 18. Dirks, R.; Pierce, N. A. J Comput Chem 2003, 24, 1664.
- Comolli, L.; Smirnov, I.; Xu, L.; Blackburn, E.; James, T. Proc Natl Acad Sci USA 2002, 99, 16998.
- Theimer, C.; Finger, L.; Trantirek, L.; Feigon, J. Proc Natl Acad Sci USA 2003, 100, 449.
- 21. Lyngso, R.; Zuker, M.; Pedersen, C. Bioinformatics 1999, 15, 440.
- 22. Ding, Y.; Lawrence, C. Nucleic Acids Res 2003, 31, 7280.
- 23. Dirks, R.; Lin, M.; Winfree, E.; Pierce, N. A. Nucleic Acids Res 2004, 32, 1392.
- 24. Fu, T.-J.; Seeman, N. C. Biochemistry 1993, 32, 3211.
- 25. Seeman, N. C. J Theor Biol 1982, 99, 237.