

Supplementary Material: A Partition Function Algorithm for Nucleic Acid Secondary Structure Including Pseudoknots

ROBERT M. DIRKS¹, NILES A. PIERCE²

¹*Department of Chemistry, California Institute of Technology, Pasadena, CA 91125*

²*Department of Applied & Computational Mathematics and Department of Bioengineering,
California Institute of Technology, Pasadena, CA 91125*

Received 16 December 2002; Accepted 5 March 2003; J Comput Chem 24: 1664–1677, 2003

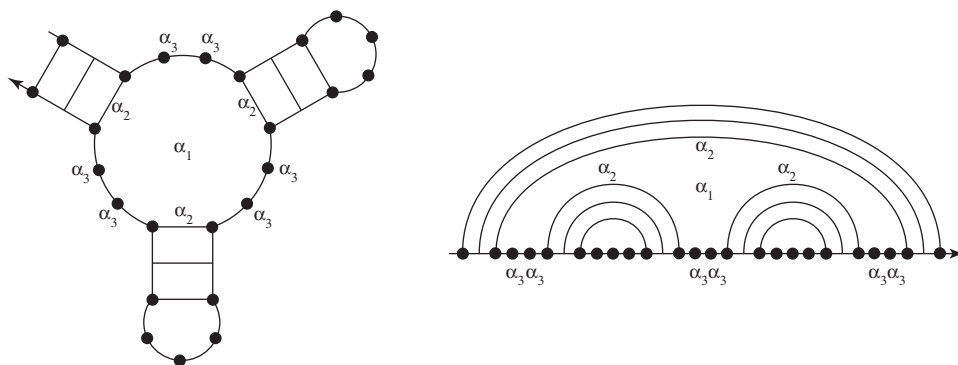


Figure S1. Multiloop energy expression. For this example there are three base pairs defining the multiloop ($B = 3$) and six unpaired bases in the multiloop ($U = 6$) so $G^{multi} = \alpha_1 + 3\alpha_2 + 6\alpha_3$.

```

function fastiloops( $i, j, l, Q^b, Q^x, Q^{x2}$ )
// $Q^x$  recursion for  $O(N^3)$  internal loop contributions to  $Q^b$ 
if ( $l \geq 15$ ) //smallest subsequence not added to  $Q^b$  as special case
 $L_1 = 4$  //explicitly add in terms for  $L_1 = 4, L_2 \geq 4$ 
 $d = i + L_1 + 1$ 
for  $L_2 = 4, j - d - 5$ 
 $s = L_1 + L_2$ 
 $e = j - L_2 - 1$ 
 $G^{\text{partial}} = \gamma_1(s) + \gamma_2(|L_1 - L_2|) + \gamma_3(e, d, e+1, d-1)$ 
 $Q_{i,s}^x += \exp\{-G^{\text{partial}}/RT\} Q_{d,e}^b$ 
 $L_2 = 4$  //explicitly add in terms for  $L_1 \geq 5, L_2 = 4$ 
 $e = j - L_2 - 1$ 
for  $L_1 = 5, e - i - 5$ 
 $s = L_1 + L_2$ 
 $d = i + L_1 + 1$ 
 $G^{\text{partial}} = \gamma_1(s) + \gamma_2(|L_1 - L_2|) + \gamma_3(e, d, e+1, d-1)$ 
 $Q_{i,s}^x += \exp\{-G^{\text{partial}}/RT\} Q_{d,e}^b$ 
//Next convert  $Q^x$  into interior loop energies
for  $s = 8, l - 7$ 
if (sequence permits  $i \cdot j$  base pair)
 $Q_{i,j}^b += Q_{i,s}^x \exp\{-\gamma_3(i, j, i+1, j-1)/RT\}$ 
//Extend loops from  $s$  to  $s+2$  for future calculation
//of  $Q_{i-1, j+1}^b$  with subsequence length  $l+2$ 
if ( $i \neq 1$  &  $j \neq N$ )
for  $s = 8, l - 7$ 
 $Q_{i-1, s+2}^{x2} = Q_{i,s}^x \exp\{-[\gamma_1(s+2) - \gamma_1(s)]/RT\}$ 
//Add small inextensible interior loop terms to  $Q^b$  as special cases
for  $L_1 = 0, 3$ 
 $d = i + L_1 + 1$ 
for  $L_2 = 0, \min(3, j - d - 5)$ 
 $e = j - L_2 - 1$ 
 $Q_{i,j}^b += \exp\{-G_{i,d,e,j}^{\text{internal}}/RT\} Q_{d,e}^b$ 
//Add bulge loops and large asymmetric loops as special cases
for  $L_1 = 0, 3$  //Cases  $L_1 = 0, 1, 2, 3, L_2 \geq 4$ 
 $d = i + L_1 + 1$ 
for  $L_2 = 4, j - d - 5$ 
 $e = j - L_2 - 1$ 
 $Q_{i,j}^b += \exp(-G_{i,d,e,j}^{\text{internal}}/RT) Q_{d,e}^b$ 
for  $L_2 = 0, 3$  //Cases  $L_1 \geq 4, L_2 = 0, 1, 2, 3$ 
 $e = j - L_2 - 1$ 
for  $L_1 = 4, e - i - 5$ 
 $d = i + L_1 + 1$ 
 $Q_{i,j}^b += \exp\{-G_{i,d,e,j}^{\text{internal}}/RT\} Q_{d,e}^b$ 

```

Figure S2. Pseudocode for computing interior loop contributions to Q^b in $O(N^3)$ as an alternative to the $O(N^4)$ interior loop recursion of Figure 8. Here, N is the length of the strand and $l = j - i + 1$ is the length of the substrand under consideration at any given point during the recursive process. A schematic representation of “fastiloops” is provided in Supplementary Material Figure S3. The smallest “possible extensible loop” is the case $L_1 = L_2 = 4$ with size $s = 8$. Therefore, the smallest subsequence for which Q^x can be employed is $l = 15$ (adding the four bases for i, d, e, j and a minimum hairpin of three bases between $d \cdot e$). For a given i and j , $Q_{i,s}^x$ already contains the contributions to $Q_{i,j}^b$ for all extensible loops of size s except for the two cases when either $L_1 = 4$ or $L_2 = 4$ (which cannot be obtained by extending smaller loops that use a different energy expression). Enriching $Q_{i,s}^x$ with these two new possible extensible loops, we then convert $Q_{i,s}^x$ into contributions to $Q_{i,j}^b$ by introducing the term for closing these loops with pair $i \cdot j$. $Q_{i,s}^x$ is then extended to provide future values of $Q_{i-1, s+2}^x$. All other interior loop contributions (cases with either $L_1 \leq 3$ or $L_2 \leq 3$) are then added directly to $Q_{i,j}^b$ using the special energy expressions of the standard model implied by $G_{i,d,e,j}^{\text{internal}}$. Note that the subsequence length l is fixed inside each call to the function “fastiloops”. Hence, specifying i implies $j = i + l - 1$. For subsequences of length l , we use $Q_{i,s}^x$ (j implied) to compute $Q_{i-1, s+2}^x$ ($j+1$ implied) which will later be used to compute contributions to $Q_{i-1, j+1}^b$ for subsequences of length $l+2$. Thus, for a given value of l , the values of $Q_{i,s}^x$ need only be stored for all legal values of i and s until l has been incremented 3 times, at which point it can be discarded. This is accomplished by using $Q_{i,s}^{x1}$ and $Q_{i,s}^{x2}$ to store future contributions for subsequences of length $l+1$ and $l+2$.

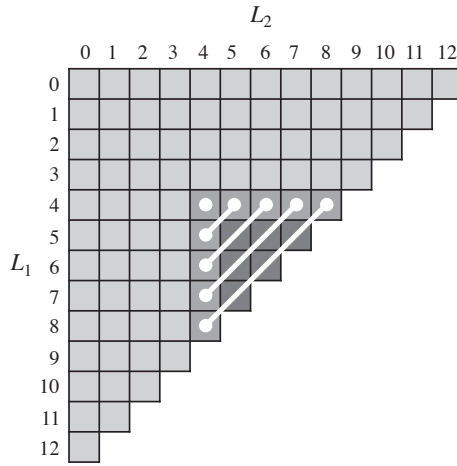


Figure S3. Schematic for interpreting the pseudocode of Supplementary Material Figure S2, which computes the interior loop contributions to the partition function in $O(N^3)$. For a given i (with j implied by the current subsequence length l), the grid illustrates the method of computing contributions to $Q_{i,j}^b$ for all interior loops with sides of length L_1 and L_2 . For the depicted case, the maximum interior loop size is $s = L_1 + L_2 = 12$. Adding seven bases to account for the closing bases i, d, e, j and the smallest allowed hairpin of three bases between $d-e$, the subsequence under consideration is therefore of length $l = 19$. The $O(N)$ pale gray cases with either $L_1 \leq 3$ or $L_2 \leq 3$ use special energy functions and are added explicitly to $Q_{i,j}^b$. The medium gray and dark gray cases are the possible extensible loop contributions that are computed using $Q_{i,s}^x$. Each diagonal line spans the terms that will be stored in $Q_{i,s}^x$ for a particular value of s . The medium gray terms are the new possible extensible loops with either $L_1 = 4$ or $L_2 = 4$. The dark gray terms were previously incorporated into $Q_{i,s}^x$ by extending smaller possible extensible loops.

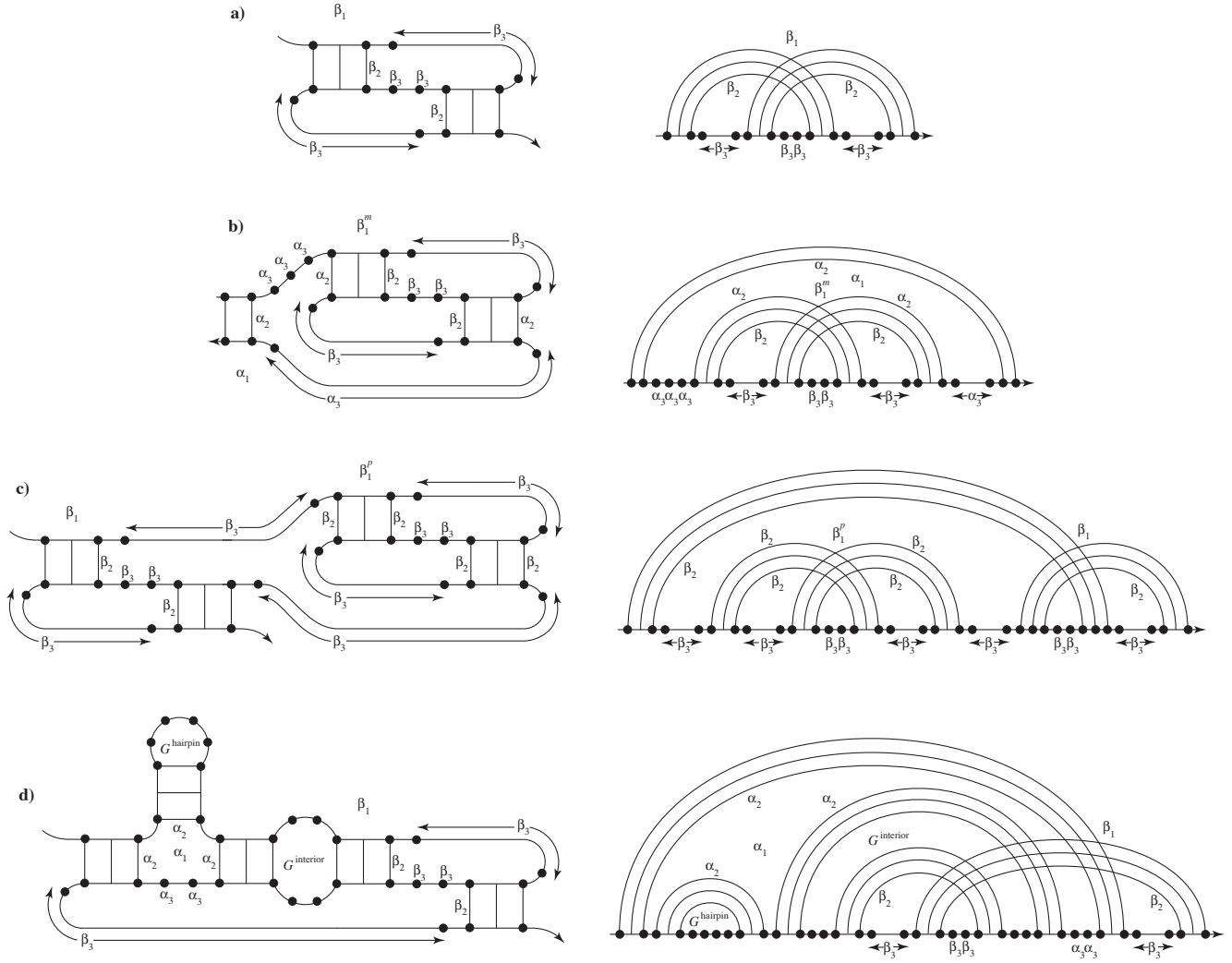


Figure S4. Illustration of the pseudoknot energy expression. a) An external pseudoknot. Bases external to the pseudoknot have no associated energy. The penalty for forming an external pseudoknot is β_1 . Base pairs that border the pseudoknot interior receive penalty β_2 while unpaired bases on the pseudoknot interior receive penalty β_3 . The energies associated with the stacked base pairs are described using the standard model. b) A pseudoknot inside a multiloop. The penalty for forming a pseudoknot inside a multiloop is β_1^m . The treatment of β_2 and β_3 inside the pseudoknot remains the same. In addition, the standard penalty for formation of a multiloop is α_1 , the two pseudoknot base pairs that border the multiloop are given the standard multiloop penalty α_2 , and unpaired bases that are inside the multiloop are given penalty α_3 . c) A pseudoknot within a pseudoknot. The energy treatment for the exterior pseudoknot remains the same. The penalty for forming a pseudoknot inside a pseudoknot is β_1^p . Base pairs from the interior pseudoknot that border the exterior pseudoknot receive penalty β_2 . Otherwise, the energetic treatment of the interior pseudoknot is the same as for an exterior pseudoknot. The partition function recursions allow an arbitrary number of levels of pseudoknots within pseudoknots. d) Pseudoknot with a hairpin and an interior loop inside a spanning region of the pseudoknot. The multiloop that forms at the base of the hairpin is treated using the standard multiloop potential. In general, the spanning region of a pseudoknot may contain interior loops, hairpins, multiloops or additional pseudoknots.

```

Initialize  $(Q, Q^b, Q^m, Q^p, Q^z)$  //  $O(N^2)$  space
Initialize  $(Q^g)$  //  $O(N^4)$  space
Set all values to 0 except  $Q_{i,i-1} = Q_{i,i-1}^z = 1$ 
for  $l = 1, N$ 
  for  $i = 1, N-l+1$ 
     $j = i+l-1$ 
    //  $Q^b$  recursion
     $Q_{i,j}^b = \exp(-G_{i,j}^{\text{hairpin}}/RT)$ 
    for  $d = i+1, j-5$  // all possible rightmost pairs  $d \cdot e$ 
      for  $e = d+4, j-1$ 
         $Q_{i,j}^b += \exp(-G_{i,d,e,j}^{\text{interior}}/RT) Q_{d,e}^b$ 
         $Q_{i,j}^b += Q_{i+1,d-1}^m Q_{d,e}^b \exp\{-[\alpha_1 + 2\alpha_2 + \alpha_3(j-e-1)]/RT\}$ 
      for  $d = i+1, j-9$  // all possible rightmost pseudoknots filling  $[d, e]$ 
        for  $e = d+8, j-1$ 
           $G^{\text{recursion}} = \alpha_1 + \beta_1^m + 3\alpha_2 + \alpha_3(j-e-1)$ 
           $Q_{i,j}^b += \exp\{-[G^{\text{recursion}} + \alpha_3(d-i-1)]/RT\} Q_{d,e}^p$ 
           $Q_{i,j}^b += Q_{i+1,d-1}^m Q_{d,e}^p \exp\{-G^{\text{recursion}}/RT\}$ 
        //  $Q^g$  recursion
        for  $d = i+1, j-5$  // set inner pair  $d \cdot e$ 
          for  $e = d+4, j-1$ 
             $Q_{i,d,e,j}^g += \exp(-G_{i,d,e,j}^{\text{interior}}/RT)$ 
          for  $d = i+2, j-6$  // set inner pair  $d \cdot e$ 
            for  $e = d+4, j-2$ 
              for  $c = i+1, d-1$  // recursion on middle pair  $c \cdot f$ 
                for  $f = e+1, j-1$ 
                   $Q_{i,d,e,j}^g += \exp(-G_{i,c,f,j}^{\text{interior}}/RT) Q_{c,d,e,f}^g$ 
            for  $d = i+6, j-5$  // set inner pair  $d \cdot e$ 
              for  $e = d+4, j-1$ 
                 $Q_{i,d,e,j}^g += Q_{i+1,d-1}^m \exp\{-[\alpha_1 + 2\alpha_2 + \alpha_3(j-e-1)]/RT\}$ 
            for  $d = i+1, j-10$  // set inner pair  $d \cdot e$ 
              for  $e = d+4, j-6$ 
                 $Q_{i,d,e,j}^g += \exp\{-[\alpha_1 + 2\alpha_2 + \alpha_3(d-i-1)]/RT\} Q_{e+1,j-1}^m$ 
            for  $d = i+6, j-10$  // set inner pair  $d \cdot e$ 
              for  $e = d+4, j-6$ 
                 $Q_{i,d,e,j}^g += Q_{i+1,d-1}^m \exp\{-[\alpha_1 + 2\alpha_2]/RT\} Q_{e+1,j-1}^m$ 
            for  $d = i+7, j-6$  // set inner pair  $d \cdot e$ 
              for  $e = d+4, j-2$ 
                for  $c = i+6, d-1$  // recursion on middle pair  $c \cdot f$ 
                  for  $f = e+1, j-1$ 
                     $G^{\text{recursion}} = \alpha_1 + 2\alpha_2 + \alpha_3(j-f-1)$ 
                     $Q_{i,d,e,j}^g += Q_{i+1,c-1}^m Q_{c,d,e,f}^g \exp\{-G^{\text{recursion}}/RT\}$ 
            for  $d = i+2, j-11$  // set inner pair  $d \cdot e$ 
              for  $e = d+4, j-7$ 
                for  $c = i+1, d-1$  // recursion on middle pair  $c \cdot f$ 
                  for  $f = e+1, j-6$ 
                     $G^{\text{recursion}} = \alpha_1 + 2\alpha_2 + \alpha_3(c-i-1)$ 
                     $Q_{i,d,e,j}^g += \exp\{-G^{\text{recursion}}/RT\} Q_{c,d,e,f}^g Q_{f+1,j-1}^m$ 
            for  $d = i+7, j-11$  // set inner pair  $d \cdot e$ 
              for  $e = d+4, j-7$ 
                for  $c = i+6, d-1$  // recursion on middle pair  $c \cdot f$ 
                  for  $f = e+1, j-6$ 
                     $Q_{i,d,e,j}^g += Q_{i+1,c-1}^m \exp\{-[\alpha_1 + 2\alpha_2]/RT\} Q_{c,d,e,f}^g Q_{f+1,j-1}^m$ 
          //  $Q^p$  recursion
          for  $a = i+1, j-7$  // place points from left to right
            for  $b = a+1, j-6$ 
              for  $c = b+1, j-5$ 
                for  $d = \max(c+1, a+4), j-3$ 
                  for  $e = d+1, j-2$ 
                    for  $f = \max(e+1, c+4), j-1$ 
                       $Q_{i,j}^p += Q_{a+1,b-1}^z Q_{c+1,d-1}^z Q_{e+1,f-1}^z$ 
                       $\cdot Q_{i,a,d,e}^g Q_{b,c,f,j}^g \exp\{-2\beta_2/RT\}$ 
        //  $Q, Q^m, Q^z$  recursions
         $Q_{i,j} = 1$  // empty recursion
         $Q_{i,j}^z = \exp(-[\beta_3(j-i+1)]/RT)$ 
        for  $d = i, j-4$  // all possible rightmost pairs  $d \cdot e$ 
          for  $e = d+4, j$ 
             $Q_{i,j} += Q_{i,d-1} Q_{d,e}^b$ 
             $Q_{i,j}^m += \exp\{-[\alpha_2 + \alpha_3(d-i) + \alpha_3(j-e)]/RT\} Q_{d,e}^b$ 
             $Q_{i,j}^m += Q_{i,d-1}^m Q_{d,e}^b \exp\{-[\alpha_2 + \alpha_3(j-e)]/RT\}$ 
             $Q_{i,j}^z += Q_{i,d-1}^z Q_{d,e}^b \exp\{-[\beta_2 + \beta_3(j-e)]/RT\}$ 
        for  $d = i, j-8$  // all possible rightmost pseudoknots filling  $[d, e]$ 
          for  $e = d+8, j$ 
             $Q_{i,j} += Q_{i,d-1} Q_{d,e}^p \exp\{-\beta_1/RT\}$ 
             $Q_{i,j}^m += \exp\{-[\beta_1^m + 2\alpha_2 + \alpha_3(d-i) + \alpha_3(j-e)]/RT\} Q_{d,e}^p$ 
             $Q_{i,j}^m += Q_{i,d-1}^m Q_{d,e}^p \exp\{-[\beta_1^m + 2\alpha_2 + \alpha_3(j-e)]/RT\}$ 
             $Q_{i,j}^z += Q_{i,d-1}^z Q_{d,e}^p \exp\{-[\beta_1^p + 2\beta_2 + \beta_3(j-e)]/RT\}$ 
// Partition function is  $Q_{1,N}$ 

```

Figure S5. Pseudocode implementation of an $O(N^8)$ dynamic programming partition function algorithm for nucleic acids with pseudoknots. Here, N is the length of the strand and $l = j - i + 1$ is the length of the substrand under consideration at any given point during the recursive process. The recursions are described schematically in Figures 12-16. In this pseudocode, care has been taken to define programming loop bounds so as to consider only valid secondary structures. The standard model requires three or more unpaired bases in a hairpin, so a b -curve must satisfy $j - i \geq 4$. Looking at Figure 15, imposing this requirement on all base pairs implies that a p -curve must satisfy $j - i \geq 8$. In the interior of a pseudoknot, the steric constraints on hairpins sometimes lead to two conflicting requirements that are incorporated using a “max” function to define the bounds for d and f . The disallowed structures have infinite energies according to the physical model so it is computationally efficient to exclude them from consideration.

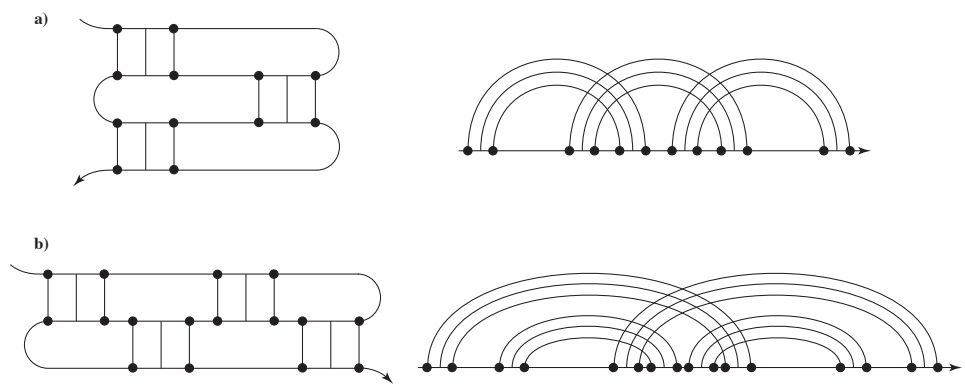


Figure S6. Examples of pseudoknots that are excluded from the partition function recursions. Neither structure can be decomposed into two spanning regions as required by Figure 15. The structure prediction recursions of Rivas and Eddy²⁵ include both structures while the the structure prediction recursions of Akutsu²⁴ include the latter.

```

function fastiloops( $i, j, l, Q^g, Q^x, Q^{x2}$ )
// $Q^x$  recursion for  $O(N^5)$  internal loop contributions to  $Q^g$ 
if ( $l \geq 17$ ) //smallest subsequence not added to  $Q^g$  as special case
  for  $d = i+6, j-10$ 
    for  $e = d+4, j-6$ 
       $L_1 = 4$  //explicitly add in terms for  $L_1 = 4, L_2 \geq 4$ 
       $c = i + L_1 + 1$ 
      for  $L_2 = 4, j-e-2$ 
         $s = L_1 + L_2$ 
         $f = j - L_2 - 1$ 
         $G^{\text{partial}} = \gamma_1(s) + \gamma_2(|L_1 - L_2|) + \gamma_3(f, c, f+1, c-1)$ 
         $Q_{i,d,e,s}^x += \exp\{-G^{\text{partial}}/RT\} Q_{c,d,e,f}^g$ 
      if ( $d \geq i+7$ )
         $L_2 = 4$  //explicitly add in terms for  $L_1 \geq 5, L_2 = 4$ 
         $f = j - L_2 - 1$ 
        for  $L_1 = 5, d-i-2$ 
           $s = L_1 + L_2$ 
           $c = i + L_1 + 1$ 
           $G^{\text{partial}} = \gamma_1(s) + \gamma_2(|L_1 - L_2|) + \gamma_3(f, c, f+1, c-1)$ 
           $Q_{i,d,e,L_1+L_2}^x += \exp\{-G^{\text{partial}}/RT\} Q_{c,d,e,f}^g$ 
    for  $d = i+1, j-5$ 
      for  $e = d+4, j-1$ 
        //Convert  $Q^x$  into interior loop energies
        if ( $l \geq 17$  & sequence permits  $i \cdot j$  base pair)
          for  $\text{size} = 8, l-9$ 
             $Q_{i,d,e,j}^g += Q_{i,d,e,s}^x \exp\{-\gamma_3(i, j, i+1, j-1)/RT\}$ 
          //Extend loops for future use
          if ( $i \neq 1$  &  $j \neq N$ )
            for  $s = 8, l-9$ 
               $Q_{i-1,d,e,s+2}^{x2} = Q_{i,d,e,s}^x \exp\{-[\gamma_1(s+2) - \gamma_1(s)]/RT\}$ 
            //Add small inextensible interior loops to  $Q^g$  as special cases
            for  $L_1 = 0, \min(3, d-i-2)$ 
               $c = i + L_1 + 1$ 
              for  $L_2 = 0, \min(3, j-e-2)$ 
                 $f = j - L_2 - 1$ 
                 $Q_{i,d,e,j}^g += \exp\{-G_{i,c,f,j}^{\text{interior}}/RT\} Q_{c,d,e,f}^g$ 
            //Add bulge loops and large asymmetric loops as special cases
            for  $L_1 = 0, \min(3, d-i-2)$  //Cases  $L_1 = 0, 1, 2, 3, L_2 \geq 4$ 
               $c = i + L_1 + 1$ 
              for  $L_2 = 4, j-e-2$ 
                 $f = j - L_2 - 1$ 
                 $Q_{i,d,e,j}^g += \exp\{-G_{i,c,f,j}^{\text{interior}}/RT\} Q_{c,d,e,f}^g$ 
            for  $L_2 = 0, \min(3, j-e-2)$  //Cases  $L_1 \geq 4, L_2 = 0, 1, 2, 3$ 
               $f = j - L_2 - 1$ 
              for  $L_1 = 4, d-i-2$ 
                 $c = i + L_1 + 1$ 
                 $Q_{i,d,e,j}^g += \exp\{-G_{i,c,f,j}^{\text{interior}}/RT\} Q_{c,d,e,f}^g$ 

```

Figure S7. Pseudocode for computing interior loop contributions to Q^g in $O(N^5)$ as an alternative to the $O(N^6)$ interior loop recursion of Figure 19. Here, N is the length of the strand and $l = j - i + 1$ is the length of the substrand under consideration at any given point during the recursive process. The smallest “possible extensible loop” is the case $L_1 = L_2 = 4$ with size $s = 8$. Therefore, the smallest subsequence for which Q^x can be employed is $l = 17$ (adding the four closing bases i, c, f, j , the additional spanning pair $d \cdot e$, and a minimum hairpin loop of three bases). For given values of i, d and e , $Q_{i,d,e,s}^x$ already contains the contributions to $Q_{i,d,e,j}^g$ for all extensible loops of size s except for the two cases when either $L_1 = 4$ or $L_2 = 4$ (which cannot be obtained by extending smaller loops that use a different energy expression). Enriching $Q_{i,d,e,s}^x$ with these two new possible extensible loops, we then convert $Q_{i,d,e,s}^x$ into contributions to $Q_{i,d,e,j}^g$ by introducing the term for closing these loops with pair $i \cdot j$. $Q_{i,d,e,s}^x$ is then extended to provide future values of $Q_{i-1,d,e,s+2}^x$. All other interior loop contributions (cases with either $L_1 \leq 3$ or $L_2 \leq 3$) are then added directly to $Q_{i,j}^b$ using the special energy expressions of the standard model implied by $G_{i,d,e,j}^{\text{interior}}$. Note that the subsequence length l is fixed inside each call to the function “fastiloops”. Hence, specifying i implies $j = i + l - 1$. For subsequences of length l , we use $Q_{i,d,e,s}^x$ (j implied) to compute $Q_{i-1,d,e,s+2}^x$ ($j+1$ implied) which will later be used to compute contributions to $Q_{i-1,d,e,j+1}^g$ for subsequences of length $l+2$. Thus, for a given value of l , the values of $Q_{i,d,e,s}^x$ need only be stored for all legal values of i, d, e , and s until l has been incremented 3 times, at which point it can be discarded. This is accomplished by using $Q_{i,d,e,s}^{x1}$ and $Q_{i,d,e,s}^{x2}$ to store future contributions for subsequences of length $l+1$ and $l+2$.