

DNA SECONDARY STRUCTURE PREDICTION WITH WORD COMBINATIONS

DANIELLE DEES AND LAUR SLAYBAUGH

1. INTRODUCTION

Zuker's original work on RNA secondary structures provides an algorithm to predict the secondary structure of a single RNA strand which will output a negative number if this strand folds into a secondary structure and a positive number if this strand remains unfolded. This algorithm is quite useful to learn about a single sequence of RNA nucleotides; however, many of strands that are currently used in DNA computing are made up of smaller strands that are called words. The science of choosing these words has become quite complex, with many rules governing the order of nucleotides in order to obtain words which can be combined into a strand that will not fold onto itself. The algorithm that we present here will take in $2n$ words, where a word is defined as a sequence of RNA nucleotides of length l . These words will be labeled as $w(1), \overline{w(1)}, w(2), \overline{w(2)}, \dots, w(n), \overline{w(n)}$. The resultant strand will be the $\sum_{i=1}^n w(i)$ or $\overline{w(i)}$. Our algorithm outputs a negative number if one of the 2^n strands that result from the concatenation of these words will fold into a secondary structure or a positive number if each possible strand will not form a secondary structure. Although the same result could be obtained by running each of the 2^n structures through Dr. Zuker's algorithm, which has a runtime of $O(n^4)$, this new algorithm will solve this problem in polynomial time according to n .

2. BASIC ALGORITHM

A secondary structure formed from a strand $s \in S$ where S is the set of all strands that can be created, will have $1 \leq i \leq |s|$ where S_i is the base at position $\#(i)$ in word $w(\text{word}\#(i))$ or $\overline{w(\text{word}\#(i))}$. The position $\#(x)$ and word $\#(x)$ functions determine the word number and position within the word based on the word length and parameter x . It is assumed that no pseudo-knots exist in any of the secondary structures and that each base can only bond to one other base. Each pseudo-knot free structure can be viewed as a collection of stacked pairs and loops with exterior base pairs closing each loop. The types of loops are hairpin, internal, and multibranch, with a bulge loop being a particular instance of an internal loop. RNA secondary structure prediction is really just determining the most stable configuration of the given nucleotides. The optimal structure will have the lowest free energy. Therefore, we can use recurrence relations to determine the free energy over all 2^n structures by taking advantage of dynamic programming and modifying Zuker's recurrence relations for single strand folding. The modifications that we have made to the original algorithm is that we have added additional parameters into all the equations that specify whether S_i is in $w(\text{word}\#(i))$ or $\overline{w(\text{word}\#(i))}$ based on the parameter $b_i \in \{T, B, E\}$ in which case T would represent $w(\text{word}\#(i))$, B would represent $\overline{w(\text{word}\#(i))}$, and E is the case where it does not matter to which word S_i belongs. The word specification is only T or B if the recurrence calls functions with a parameter that would be in the same word as a second parameter. Specifying the word makes sure that the optimal structure is a concatenation of whole words and not parts of words.

- The energy of the optimal structure from S_1 to S_j where j is the length of $\forall s \in S$ is:

$$(1) \quad W'_S(j) = W'_S(E, j)$$

$$(2) \quad W'_S(b_j, j) = \min \begin{cases} W'_S(x, j-1) & x = b_j \text{ if } j-1 \bmod l \neq 0 \\ & x = E \text{ if } j-1 \bmod l = 0 \\ \min_{\substack{1 \leq i \leq j-1 \\ i-1 \bmod l \neq 0 \\ b_i \in \{B, T\}}} V'_S(b_i, b_j, i, j) + W'_S(b_i, i-1) \\ \min_{\substack{1 \leq i \leq j-1 \\ i-1 \bmod l = 0}} V'_S(E, b_j, i, j) + W'_S(E, i-l) \end{cases}$$

- The energy of a loop which is closed with the base pair S_i, S_j is:

$$(3) \quad V'_S(b_i, b_j, i, j) = \min \begin{cases} eH'_S(b_i, b_j, i, j) \\ eS'_S(b_i, b_j, i, j) + V'_S(x, y, i+1, j-1) \\ VBI'_S(b_i, b_j, i, j) \\ VM'_S(b_i, b_j, i, j) \end{cases}$$

where x which corresponds to b_{i+1} and y which corresponds to b_{j-1} are determined by the following pseudocode:

```

1   x = E;
2   y = E;
3   if(word#(i) == word#(i+1)) then
4   x = $b_i$;
5   if(word#(j) == word#(j-1)) then
6   y = $b_j$;
7   if(word#(i+1) == word#(j-1)) then
8   if(x == E and y == E) then
9   x = b_1;
10  y = b_1;
11  else if(x == E) then
12  x = y;
13  else
14  y = x;

```

- The energy of an internal loop which is closed by the pair S_i, S_j is:

$$(4) \quad VBI'_S(b_i, b_j, i, j) = \min \begin{cases} +\infty \text{ for } j < i + 4 \\ \min_{i < i' < j' < j} eL(b_i, b_j, x, y, i, j, i', j') + V'_S(x, y, i', j') \end{cases}$$

where x which corresponds to $b_{i'}$ and y which corresponds to $b_{j'}$ are determined by the following pseudocode:

```

1   x = E;
2   y = E;
3   if(word#(i) == word#(i')) then
4   x = $b_i$;
5   if(word#(j) == word#(j')) then
6   y = $b_j$;
7   if(word#(i') == word#(j')) then
8   if(x == E and y == E) then
9   x = b_1;
10  y = b_1;
11  else if(x == E) then
12  x = y;
13  else
14  y = x;

```

- The energy of a multibranch loop closed by the base pair S_i, S_j is:

$$(5) \quad VM'_S(b_i, b_j, i, j) = \min_{i+1 < h \leq j-1} WM(w, y, i+1, h-1) + WM(x, z, h, j-1) + a$$

$$(6) \quad WM_S(b_i, b_j, i, j) = \min \begin{cases} V'_S(b_i, b_j, i, j) + b \\ \min_{i < h \leq j} WM'_S(b_i, y, i, h-1) + WM'_S(x, b_j, h, j) \end{cases}$$

where w which corresponds to b_{i+1} , x which corresponds to b_h , y which corresponds to b_{h-1} , and z which corresponds to b_{j-1} are determined by the following pseudocode:

```

1   w = E;
2   x = E;
3   y = E;
4   z = E;
5   if(word#(i+1) == word#(i)) then
6   w = $b_i$;
7   else if(word#(i+1) == word#(h-1)) then
8   w = $b_1$;
9   if(word#(h-1) == word#(i+1)) then
10  x = w;
11 else if(word#(h-1) == word#(h)) then
12 x = $b_2$;
13 if(word#(h) == word#(h-1)) then
14 y = x;
15 else if(word#(h) == word#(j-1)) then
16 y = $b_3$;
17 if(word#(j-1) == word#(h)) then
18 z = y;
19 if(word#(j-1) == word#(j)) then
20 if(z == $b_{num}$) then
21 $b_{num}$ = $b_j$;
22 else
23 z = $b_j$;

```

This algorithm is dependent on eH'_S , eS'_S , and eL'_S which are the free energy equations for a hairpin loop, a stacked pair, and an internal loop respectively. The details of these calculations are discussed in greater detail later in this paper.

3. BACKGROUND

Let $w(i)$ and $\overline{w(i)}$ be the two possible words used at location i . Let S be the set of all strands that created by $\sum_{i=1}^n w(i)$ or $\overline{w(i)}$.

Additionally, we have a table *stack* that returns the energy of stacking the pairs $i, i+1$ on $j, j-1$ and functions that return the penalty for various structures.

4. ES

eS originally takes 2 values, i and j . $eS(i, j)$ in turn calls $stack(i, i+1, j-1, j)$. When we call eS we actually want $\min_{s \in S} eS(s, i, j)$. However, we do not need to test every possible $s \in S$ because the nucleotides relevant to eS are only $i, i+1, j-1$, and j . Thus, what we actually do is create a set of strings

$R : \{r_1, \dots, r_n | r_q[0] = s_i, r_q[1] = s_{i+1}, r_q[2] = s_{j-1}, r_q[3] = s_j | s \in S \text{ and } 1 \leq q \leq n\}$

Thus, we actually call $\min_{r \in R} stack(r_0, r_1, r_2, r_3)$.

5. EH

eH originally takes i and j as well. This returns the value $stack(i, i + 1, j - 1, j) + hPenalty(j - i - 1)$. When we call eH we actually want $\min_{s \in S} eH(s, i, j)$ Which simplifies to $\min_{r \in R} \{stack(r_0, r_1, r_2, r_3)\} + hPenalty(j - i - 1)$.

6. EL

eL originally takes i, i', j' , and j , where i and j close exterior end of the bulge or loop and i' and j' close the interior end of the bulge or loop. eL returns the value $stack(i, i + 1, j - 1, j) + stack(i', i' + 1, j' - 1, j')$. For our case then, we first define a new set

$RL : \{rl_1, \dots, rl_n \mid rl_q[0] = s_i, rl_q[1] = s_{i+1}, rl_q[2] = s_{i'}, \dots, rl_q[7] = s_j \mid s \in S \text{ and } 1 \leq q \leq n\}$
 we want eL to return $\min_{rl \in RL} \{stack(rl_0, rl_1, rl_6, rl_7) + stack(rl_2, rl_3, rl_4, rl_5)\}$

7. EM

eM is the difficult one, however, the values for this function are calculated using the tables $VM(i, j)$ and $WM(i, j)$, so we never actually call eM , thus, it doesn't matter.