



PreRkTAG: Prediction of RNA Knotted Structures Using Tree Adjoining Grammars

Hesamedin Torabi Dashti ¹, Fatemeh Zare-Mirakabad ², Nima Aghaeepour ³, Hayedeh Ahra-
bian ⁴, Abbas Nowzari-Dalini ^{4,*}

¹ Department of Mathematics, University of Wisconsin, Madison, USA

² Department of Computer Science, Faculty of Mathematics and Computer Science, Amirkabir University of Technology, Tehran, IR Iran

³ Department of Computer Science, University of British Columbia, Vancouver, Canada

⁴ Department of Computer Science, School of Mathematics, Statistics, and Computer Science, University of Tehran, Tehran, IR Iran

ARTICLE INFO

Article type:
Research Article

Article history:
Received: 09 Jun 2011
Revised: 11 Aug 2011
Accepted: 15 Oct 2012

Keywords:
Genetic Algorithms
RNA Secondary Structure
Pseudoknot
Tree Adjoining Grammars

ABSTRACT

Background: RNA molecules play many important regulatory, catalytic and structural roles in the cell, and RNA secondary structure prediction with pseudoknots is one of the most important problems in biology. An RNA pseudoknot is an element of the RNA secondary structure in which bases of a single-stranded loop pair with complementary bases outside the loop. Modeling these nested structures (pseudoknots) causes numerous computational difficulties and so it has been generally neglected in RNA structure prediction algorithms.

Objectives: In this study, we present a new heuristic algorithm for the Prediction of RNA Knotted structures using Tree Adjoining Grammars (named PreRkTAG).

Materials and Methods: For a given RNA sequence, PreRkTAG uses a genetic algorithm on tree adjoining grammars to propose a structure with minimum thermodynamic energy. The genetic algorithm employs a subclass of tree adjoining grammars as individuals by which the secondary structure of RNAs are modeled. Upon the tree adjoining grammars, new crossover and mutation operations were designed. The fitness function is defined according to the RNA thermodynamic energy function, which causes the algorithm convergence to be a stable structure.

Results: The applicability of our algorithm is demonstrated by comparing its results with three well-known RNA secondary structure prediction algorithms that support crossed structures.

Conclusions: We performed our comparison on a set of RNA sequences from the RNaseP database, where the outcomes show efficiency and practicality of the proposed algorithm.

Published by Kowsar Corp, 2013. cc 3.0.

► Implication for health policy/practice/research/medical education:

This article is recommended for those who researches in RNA area.

► Please cite this paper as:

Torabi Dashti H, Zare-Mirakabad F, Aghaeepour N, Ahra-bian H, Nowzari-Dalini A. PreRkTAG: Prediction of RNA Knotted Structures Using Tree Adjoining Grammars. *Iran J Biotech.* 2013; **11**(1): 3-13. DOI: 10.5812/ijb.9213

* Corresponding author: Abbas Nowzari-Dalini, Department of Computer Science, School of Mathematics, Statistics, and Computer Science, University of Tehran, Tehran, IR Iran. Tel: + 98-2166412478, Fax: + 98-2166412478, E-mail: nowzari@ut.ac.ir

DOI: 10.5812/ijb.9213

Copyright © 2013, National Institute of Genetic Engineering and Biotechnology; Published by Kowsar Corp.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/3.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

is converging to local minimums, which is highly possible especially with a wide range of feasible solutions. In PreRkTAG we attempt to tackle this problem by using novel evolutionary functions. The evaluation of PreRkTAG against celebrated energy-based prediction algorithms, namely, the STAR, Hotknots, and Pknots, shows high sensitivity and high specificity of PreRkTAG in the most cases. The evaluation was performed on 24 RNA sequences chosen from the RNaseP database. Our algorithm correctly identifies more than 80% of base pairs.

3. Materials and Methods

The employed data structure for capturing RNA crossed structures was introduced in preceding research (14) and in order to synchronize the notifications the first part of this section is focused on ESL-TAG. Consequently, the details of PreRkTAG algorithm are described including the employed inner functions.

3.1. RNA Secondary Structure Modeling

Generally, an ESL-TAG is a 5-tuple grammar $G = (N, T, S, I, A)$ where N and T are finite sets of non-terminal and terminal symbols, respectively. S is the start symbol and I is finite set of initial trees (center trees). Finally, the finite set A is a set of adjunct trees (auxiliary trees). The members of the set N are divided into two subsets; *Active* and *Inactive*. During transition process of the grammar G , the active nodes can be replaced by a member of the set A , where inactive nodes do not. Based on the division, when a tree does not have any active node, it would be a *matured tree*. By adjoining two trees, s and t , a tree t' would be obtained, and $t \rightarrow s t'$ represents the adjuncted transition. The transitive closure of \rightarrow is \rightarrow^* and a tree t' is derived from t , where $t \rightarrow^* t'$ for some $t \in (A \cup I)$. A set of trees of an ESL-TAG G is defined as $T(G) = \{t \mid s \rightarrow^* t, s \in I \text{ where } t \text{ is matured tree}\}$. For modeling the secondary structure of an RNA with pseudoknot, the tree adjoining grammar $G = (N, T, S, I, A)$ is defined as follows. The set N is a finite set of non-terminals $\{X, Y, Z\}$, where active non-terminals are marked with asterisk $\{X^*, Y^*, Z^*\}$. The set T is finite set of terminals $\{a, u, c, g\}$ as abbreviation of nucleotide acids. Moreover, S is a start symbol and the set of initial trees I is composed of just one tree as depicted in Figure 2 (this tree is tree of type 1 and denoted by center tree). The last parameter A , is a finite set of adjunct trees, such as given in Figure 2. This set is composed of four types. As depicted in the Figure 2, type 2 and type 3 represent base pairs. Moreover via preorder traversal of trees of these types, the type 2 represents crossed base pairs denoted by T2u and T2d, and the type 3 represents nested base pairs denoted by T3L and T3R. In addition, type 4 is introduced to generate a single base without any pairing. Trees of type 5 are used to combine substructures by using trees of type T5Ld and T5Rd, and inserting a tree into the current one by using trees of

type T5Lu and T5Ru. As mentioned, the key idea of this modeling is from Uemura et al. (16). Figure 3 depicts the process of constructing the crossed base pairs "agacuu" which is not generated by any context free grammar. Uemura et al. (16) described a computational obstacle of ESL-TAG trees via an example that is shown in Figure 4. In this example, it is shown that two dissimilar sequences can be captured by one ESL-TAG structure.

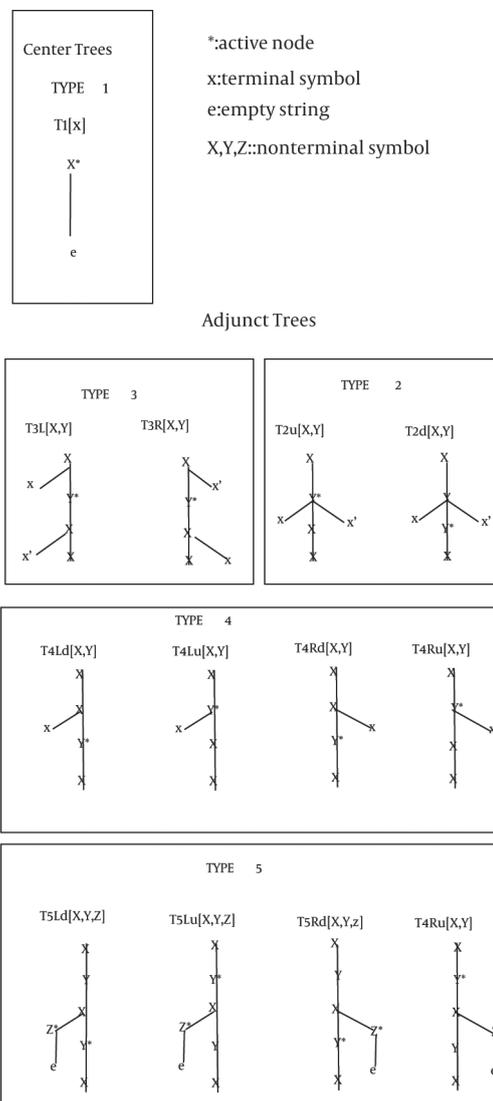


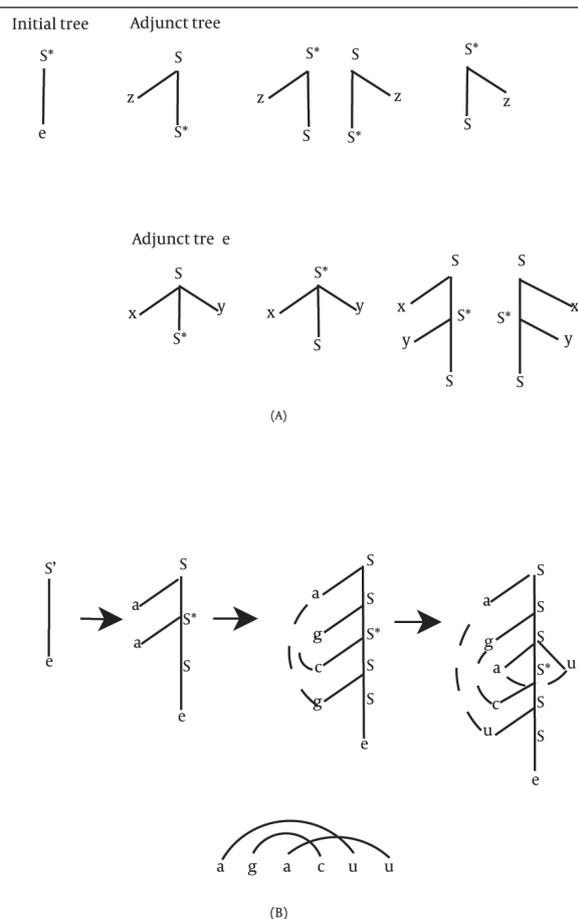
Figure 2. Five Types of Adjunct Trees for Modeling the Secondary Structure of RNA by ESL-TAG (16)

3.2. Genetic Algorithm

Genetic Algorithm (GA) is kind of Evolutionary Algorithms (EA) which via broad wide searching inside of probable solution space, analyzes a population of solutions instead of a single candidate. Goldberg (24) proved the GA convergence theory. Thereafter applicability of

the algorithm was extended in wide range of optimization problems.

Figure 3. The Process of Constructing the Crossed Base Pairs “Agacuu”

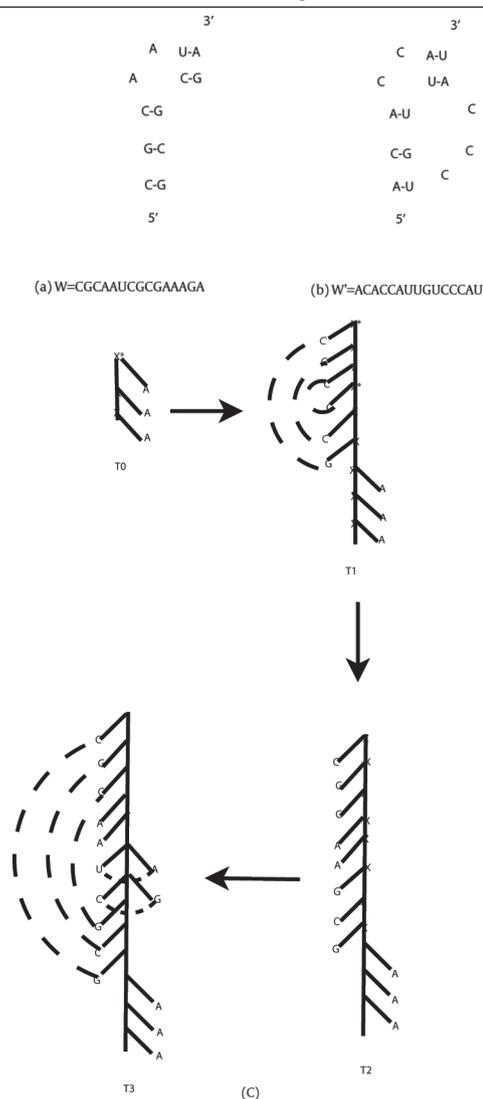


a) Initialization of a sample ESL-TAG grammar; b) Capturing a sequence “agacuu” using ESL-TAG trees in part (a) (22).

The principle of GA is competition of the solutions where a scoring function (fitness function) measures their optimality. GA methodology is composed of two evolutionary processes; the first one is applied on to a solution and by randomly modifying the solution’s properties creates a new solution (Mutation). Another one works on two different randomly selected solutions and generates two new solutions by combining predetermined features of the randomly selected solutions (Crossover). These two evolutionary functions are applied during the progress of the algorithm, where both these functions are derived from biological concepts (24). In this optimization method, the evolutionary functions are supposed to inhibit solutions to converge to local minimums. Therefore, designing efficient and practical algorithms for these functions play significant role as they converge to desirable results. Since these functions apply on the domain of feasible solutions, consideration of the data structure which represents

the solutions is a critical step in designing the evolutionary functions. Novel evolutionary functions are designed according to tree adjoining grammars and equip the PreRkTAG algorithm towards global convergence. In advance to illustrating the evolutionary functions, a short description of feasible population and data structure is presented as follows.

Figure 4. Derivation RNA Structures Using ESL-TAG



a) An RNA sequence and its secondary structures. b) ESL-TAG representation of RNA structure given in (a). c) Another sequence with the same ESL-TAG representation (16)

3.3. Generating Population

Population of a GA is a set of feasible solutions and a solution is feasible if the length of its corresponding structure is the same as the length of the given sequence. Individuals of the population in PreRkTAG are generated

randomly by using ESL-TAG. As described in the previous section, an ESA-TAG starts with an initial state and could be extended easily to cover as many nodes as required. Hence, using ESL-TAG we are able to generate a tree with any given numbers of residues. To construct each element of population, we randomly select a transitive closure operator and apply it on initial trees. In this way a matured tree can be constructed in length of a given RNA sequence. Hence, a tree shows a probable structure for RNA sequences with the same length. Generating the population only depends on the length of the given RNA sequence and the rest of the algorithm adjusts these structures to nucleotide arrangement of the given RNA sequence. It is noted that outcome of the algorithm is an optimal structure for the given RNA sequence, which is a feasible structure for all RNA sequences with the same length but is not the optimal structure for them. Generating the population in this way leads to a wide range of feasible solutions where the fitness function will narrow them down to the optimal solution of the given RNA sequence. *Figure 5* shows an example of a constructed structure for a sequence of length 5.

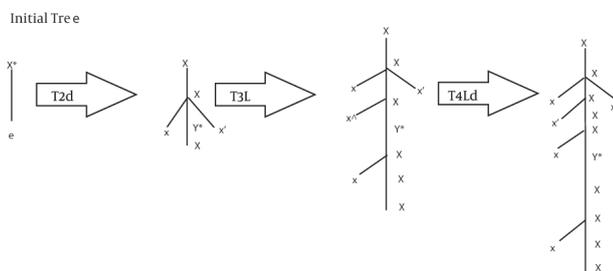


Figure 5. The Generated Sample Tree for a Sequence of Length 5

3.4. Evolutionary Functions

We present a detailed description of the evolutionary functions. These functions perform blindly in terms of considering the context of the given RNA. The mutation function takes a structure from the population and based on its length calculates the number of terminal nodes that might be mutated. Subsequently the mutated terminal nodes are selected by uniform distribution. Terminal nodes inside of adjunct trees have three different types, thus this function acts in three different ways. These treatments are discussed as follows:

- 1) If the selected node is a single node, then the mutation function selects another single node uniformly at random and pairs them up. The novel structure has a new pair which its ESL-TAG transitive closure is unknown. Hence, the new structure must be parsed by ESL-TAG parser to obtain the corresponding transitive closure.
- 2) If the selected node is a member of nested pair, then the tree reconstruction is not needed and just their ad-

junct tree is replaced with two adjunct trees with single terminals.

- 3) If the selected node is a member of crossed pair, then breaking up this pair is equal to replacing the corresponding adjunct tree by two adjunct trees; each of them has a terminal node in opposite side of the tree.

Algorithm 1 shows the mutation function. This algorithm includes two sub-functions; the "Random (L)" returns a random value L' ($0 < L' < L$), and the "Length (tree)" returns the number of terminal nodes of the given tree. In addition, as illustrated (*Algorithm 1*), in some cases new obtained structures must be parsed again to reconstruct the ESL-TAG tree, so in this algorithm the "Reconstruct" function is employed to find ESL-TAG tree of new structures. The next evolutionary function is Crossover that crosses two selected solutions and merges complementary parts. The crossover function selects two solutions via normal distribution from the population. Thereafter this function extracts two sub-trees with the same length (number of terminal nodes) and random starting points from each of the selected trees and the crossover function swaps the sub-trees. Removing a sub-tree from a tree might come with breaking pairs whenever a pair contains a node inside and another node outside of the sub-tree (this pair is called an exteriorly pair). These nodes will remain single when the swapping process is done. Constructing the trees blindly, with respect to context of the given RNA sequence, makes the crossover function much easier and helps to avoid local minimums. The crossover algorithm is presented in *Algorithm 2*. Two introduced inner functions in the Mutation function (*Algorithm 1*) are used in this algorithm as well. In addition, two other functions are applied; "Move Down (node, length)" is a function for traversing the given tree from its node to attain length number of terminal nodes, and "Swap (Sub Tree0, Sub Tree1)" is a function for swapping sub-trees which results in two new trees. *Figure 6* demonstrates a sample of the crossover process.

3.5. Fitness Function

To measure the optimality of a predicted RNA structure, as it has been recently proposed by most methods, thermodynamics rules (17, 18, 25) are applied as the fitness function. In this case, the fitness function estimates the discrepancy of simulated structures by measuring energy of structures, where better structure has lower energy. Upon thermodynamics rules, constructing a base pair needs a specific amount of energy, where in some cases the base pair releases energy (negative value) and in other cases it gains energy (positive energy). In the thermodynamics rules, many parameters effect on the amount of energy and in a simplified model, this amount depends on nucleotide acids that compose a pair and also their neighbor base pairs. For example when 'AA' in 5' - 3' and 'UU' in 3' - 5', the orientation $\begin{pmatrix} AA \\ UU \end{pmatrix}$

makes two consecutive pairs, and the energy of corresponding structure will be reduced by -0.9. One example of assigning energy to two consecutive pairs is borrowed from (5) that is shown in Table 1. In this table, every row and every column shows a pair and the entry (i, j) shows amount of energy of appearing j 'th pair after i 'th pair. Since these are thermodynamics energies and also were assessed experimentally, there is no need to have a symmetric table. We also used the same energy values in our fitness function. For a given sequence S , with no base pair, the energy value is equal to zero. Assume $e(i, j)$ denotes the energy value of base pair between i th and j th nucleotides, thus $\sum_{i,j} e(i, j)$ shows the energy value of sequence S . The energy of RNA sequences obtained by the thermodynamic rules implies that number of single bases (without pair), coaxial, and stacked base pairs have impressive effects on the energy value of the corresponding sequence (17, 18). Our fitness function uses the current version of Turner group thermodynamic rules (5).

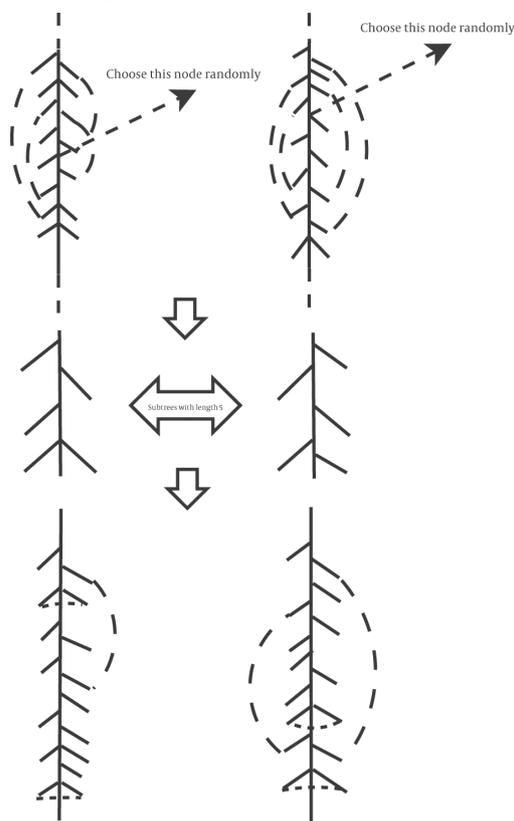


Figure 6. Graphical Representation of the Crossover Process

Within the i th iteration of the genetic algorithm, ten percent of highest score of the current population (that is called P_i) are directly copied into the next population (P_{i+1}). Then remaining 90 percent of P_{i+1} are produced by applying the roulette wheel method on the union of P_i and evaluated individuals (Emembers) of P_i . The evaluated individuals (Emembers) are new members generated by performing the evolution functions on the members of P_i .

Function Mutation (t : Tree)

```

Begin
  MutationNumber = Random (Length( $t$ ))
  For  $i=1$  to MutationNumber Do Begin
     $node_0$  = A random node from  $t$ ;
    If  $node_0$  is a bases of a pair Then Begin
       $node_1$  = pair of  $node_0$ ;
      If pair ( $node_0, node_1$ ) is a crossed pair Then Begin
        Break up pair ( $node_0, node_1$ );
        Reconstruct  $t$ ;
      End
    Else
      Mark bases  $node_0$  and  $node_1$  as single bases;
    End
  Else Begin
     $node_1$  = A random single node from  $t$ ;
    Mark  $node_0$  and  $node_1$  as a base pair;
    Reconstruct  $t$ ;
  End;
End;
Return ( $t$ );
End.

```

Algorithm 1. Mutation Function

Function Crossover (t_0, t_1 : Tree)

```

Begin
  CrossLength = Random (Length ( $t_0$ ));
   $node_0$  = A random node from  $t_0$ ;
   $node_1$  = A random node from  $t_1$ ;
  SubTree $_0$  = MoveDown ( $node_0$ , CrossLength);
  SubTree $_1$  = MoveDown ( $node_1$ , CrossLength);
  If sub-trees have exteriorly pair Then
    Break them up;
    Swap (SubTree $_0$ , SubTree $_1$ );
  Return ( $t_0, t_1$ );
End.

```

Algorithm 2. Crossover Function

This function includes six inner functions; "Random(L)" returns a random value L' ($0 < L' < L$), the "Length (Sequence)" function that counts the number of terminal nodes of Sequence, "PopulationGeneration (Pop, Number, S)" that generates Number trees of length of S using ESL-TAG, "GetRand ()" that returns a random value between [0,1], our fitness function which is called "ComputScore (Pop)" function by using thermodynamic roles computes the Pop's members energy value, and finally the "RouletteWheel (N, P_i, P_j)" via the roulette wheel method copies $(N * PopulationNumber) / 100$ members of P_i into P_j .

Table 1. The Two Dimensional Table Represents Energy of Two Consecutive Pairs

	A:U	C:G	G:C	U:A	G:U	U:G
A:U	-0.9	-2.1	-1.7	-0.9	-0.5	-1.0
C:G	-1.8	-2.9	-2.0	-1.7	-1.2	-1.9
G:C	-2.3	-3.4	-2.9	-2.1	-1.4	-2.1
U:A	-1.1	-2.3	-1.8	-0.9	-0.8	-1.1
G:U	-1.1	-2.1	-1.9	-1.0	-0.4	-1.5
U:G	-0.8	-1.4	-1.2	-0.5	-0.2	-0.4

4. Results

To assess our algorithm its accuracy has been evaluated and compared with three common thermodynamic based algorithms: Hotknot (21), Pknots (9), and the STAR (1). To compare PreRkTAG with these algorithms, we measured the accuracy of these algorithms on a set of RNAs from the RNaseP data base. Since genetic algorithms have random behavior, every execution of our algorithm on a given sequence may result in different structure. Hence, the PreRkTAG algorithm was applied three times on every sequence and the best structure of each run was stored. Then one of these three with the minimum energy was reported as the predicted structure. Because of the proposed evolution algorithms, in many cases, the mean of the three energy values were close to the picked minimum energy. The RNaseP data base includes natural structures of this set, therefore in this study the accuracy is measured by considering two different parameters; sensitivity and specificity of folded sequences. Sensitivity is defined by the number of base pairs that are predicted correctly per number of base pairs in natural structure. The specificity is equal to the number of true predicted base pairs per number of predicted base pairs. Predictors were employed on 24 sequences from RNaseP data base, which are tabulated in Table 2. The comparison of sensitivity between our algorithm and the other algorithms are given in Table 3. Also Table 4 depicts results of comparing specificity of these algorithms. Unfortunately because of the length of the test sequences, which are fairly long, the Pknots algorithm could not obtain results within a reasonable timeframe on our hardware. Therefore, in Table 3 and Table 4 the Pknots results are not shown. We also provide running time complexity of the examined algorithms in the Table 5. As mentioned, for every RNA sequence, the PreRkTAG has been applied three times and we indicated outcome of the one with minimum energy. For calculating the running time of the algorithm, we followed the same idea and considered the running time of the one with minimum energy. It is well known that choosing suitable parameter values for genetic algorithms is very difficult (24). There is an inverse relationship between the number of generations and the size of population, so that large population size converges fast and

Function PreRkTAG (*S*: RNA Sequences)

```

Begin
  PNumber = Random (Length(S));
  PopulationGeneration (P0, PNumber, S);
  i = 0; stop-condition = False ;
  While (stop-condition = False) Do Begin
    Emembers = {∅} ;
    ComputScore(Pi) ;
    Copy ten percent of high scores Pi into Pi+1;
    For each t in Pi Do Begin
      If GetRand() > Mutation Probability Then Begin
        t' = Mutation (t) ;
        Emembers = Emembers ∪ t';
      End;
    End;
    For each t0 in Pi Do Begin
      If GetRand() > Crossover Probability Then Begin
        t1 = A random member of Pi ;
        (t'0, t'1) = Crossover(t0, t1);
        Emembers = Emembers ∪ (t'0, t'1);
      End;
    End ;
    ComputScore(Emembers) ;
    RouletteWheel (90, Pi ∪ Emembers, Pi+1);
    If Three previous iterations had equal best score Then
      stop-condition = True ;
    If i exceeds flip flops limitation Then
      stop-condition = True ;
    i++;
  End;
End.

```

Algorithm 3. PreRkTAG Algorithm

requires the consideration of a few number of generations. On the other hand, large population size requires more evaluations per generation so as to avoid the decline of the rate of convergence. Generally, the optimal size of population depends on the domain complexity of the feasible solutions. In addition, the probabilities of performing evolutionary functions play a critical role to avoid trapping in local minimums. Small values for the crossover and mutation rates yield to copying most of the individuals to the next population directly which is beneficial in terms of decreasing the running time complexity as it is unnecessary to compute fitness values for most of the individuals. However, the small rates may

Table 2. List of Tested RNA Sequences

RNA's Name	Length	Seq ID
<i>A. Ambivalens</i>	262	AF192349
<i>A. Fulgidus</i>	229	AE000782
<i>A. Truei</i>	320	AF004373
<i>B. Bufo</i>	397	AF044737
<i>C. Paradoxa</i>	351	X89853
<i>D. Janneae</i>	406	AF004372
<i>E. 1B</i>	252	AF192351
<i>E. RH</i>	301	AF192353
<i>G. Gorilla</i>	320	L08685
<i>H. Morrhuae</i>	475	U42981
<i>L. Muta</i>	331	AF004376
<i>M. Formicicum</i>	301	AF121774
<i>M. Musculus</i>	288	L08802
<i>M. Sedula</i>	304	AF121773
<i>M. Thermoautotrophicum</i>	293	U42986
<i>P. Purpurea</i>	383	U38804
<i>S. Carlsbergensis</i>	359	L12746
<i>S. Pastorianus</i>	365	AF186225
<i>S. Uvarum</i>	369	L12749
<i>S. Versatilis</i>	286	NA
<i>U. N2 fixer TP1</i>	294	NA
<i>U. N2 fixer WC1</i>	294	NA
<i>V. Cholerae</i>	399	AE004310
<i>Y. Pestis</i>	377	NA

Table 3. Comparison of Sensitivity Various Algorithms

RNA's Name	PreRkTAG	HotKnots	STAR
<i>A. Ambivalens</i>	0.86	0.43	0.33
<i>A. Fulgidus</i>	0.83	0.70	0.59
<i>A. Truei</i>	0.95	0.62	0.54
<i>B. Bufo</i>	0.31	0.51	0.65
<i>C. Paradoxa</i>	0.92	0.41	0.51
<i>D. Jeanneae</i>	0.67	0.19	0.44
<i>E. 1B</i>	0.88	0.37	0.17
<i>E. RH</i>	0.76	0.13	0.20
<i>G. Gorilla</i>	0.51	0.07	0.38
<i>H. Morrhuae</i>	0.98	0.59	0.06
<i>L. Muta</i>	0.71	0.54	0.31
<i>M. Formicicum</i>	0.54	0.62	0.34
<i>M. Musculus</i>	0.87	0.68	0.53
<i>M. Sedula</i>	0.91	0.95	0.79
<i>M. Thermoautotrophicum</i>	0.94	0.22	0.43
<i>P. Purpurea</i>	0.96	0.68	0.37

<i>S. Carlsbergensis</i>	0.81	0.89	0.75
<i>S. Pastorianus</i>	0.85	0.46	0.88
<i>S. Uvarum</i>	0.87	0.89	0.81
<i>S. Versatilis</i>	0.71	0.41	0.31
<i>U. N2 fixer TP1</i>	0.91	0.56	0.77
<i>U. N2 fixer WC1</i>	0.96	0.88	0.50
<i>V. Cholerae</i>	0.71	0.91	0.87
<i>Y. Pestis</i>	0.87	0.66	0.84

Table 4. The Obtained Specificity Value by the Algorithms

RNA's Name	PreRkTAG	HotKnots	STAR
<i>A. Ambivalens</i>	0.71	0.34	0.31
<i>A. Fulgidus</i>	0.87	0.79	0.51
<i>A. Truei</i>	0.87	0.65	0.67
<i>B. Bufo</i>	0.11	0.17	0.40
<i>C. Paradoxa</i>	0.89	0.29	0.44
<i>D. Jeanneae</i>	0.45	0.04	0.49
<i>E. IB</i>	0.65	0.23	0.21
<i>E. RH</i>	0.71	0.11	0.07
<i>G. Gorilla</i>	0.45	0.08	0.41
<i>H. Morrhuae</i>	0.93	0.65	0.03
<i>L. Muta</i>	0.68	0.52	0.44
<i>M. Formicicum</i>	0.43	0.78	0.21
<i>M. Musculus</i>	0.84	0.76	0.41
<i>M. Sedula</i>	0.81	0.73	0.64
<i>M. Thermoautotrophicum</i>	0.85	0.16	0.39
<i>P. Purpurea</i>	0.89	0.61	0.31
<i>S. Carlsbergensis</i>	0.70	0.72	0.65
<i>S. Pastorianus</i>	0.71	0.38	0.45
<i>S. Uvarum</i>	0.61	0.73	0.65
<i>S. Versatilis</i>	0.55	0.29	0.17
<i>U. N2 fixer TP1</i>	0.79	0.56	0.23
<i>U. N2 fixer WC1</i>	0.91	0.91	0.38
<i>V. Cholerae</i>	0.63	0.80	0.82
<i>Y. Pestis</i>	0.87	0.66	0.84

cause to converge to local minimums. Considering the behavior of genetic algorithms according to these rates shows that high probabilities for performing the mutation function allow for a greater diversity of solutions at very little expense. However, these high probabilities change the concept of genetic algorithms and convert them into random search procedures. The probability of performing the crossover function also has its own difficulties. The high value of the crossover rate helps the evolution progress and causes to it to jump out of local minimums while increasing the chance of converging to the global minimum. On the other hand, it

increases the running time complexity dramatically. A theoretical investigation of the optimal size of population, crossover and mutation rates has been carried out by Goldberge (24) who has found that the best values for the population size is a value base on the length of the solutions. Goldberge showed that a practical crossover rate is in the range of [0.6, 0.95] and for the mutation rate should be in [0.001, 0.2]. In order to evaluate PreRkTAG algorithm, we used Goldberge's investigation for determining our GA's control parameter values. We tested different values for the size of population, crossover and mutation rates on a training set of RNAs. We empirically

Table 5. Comparison of Running Time of the Examined Algorithms in Minutes

RNA's Name	PreRkTAG	HotKnots	STAR
<i>A. Ambivalens</i>	4.575	5.445	6.562
<i>A. Fulgidus</i>	5.123	4.790	5.505
<i>A. Truei</i>	5.622	5.720	6.784
<i>B. Bufo</i>	7.130	7.816	8.011
<i>C. Paradoxa</i>	6.503	5.963	7.976
<i>D. Jeanneae</i>	7.544	6.933	8.715
<i>E. IB</i>	4.222	4.636	6.637
<i>E. RH</i>	5.297	6.296	5.100
<i>G. Gorilla</i>	5.412	5.738	5.486
<i>H. Morrhuae</i>	9.035	8.435	10.033
<i>L. Muta</i>	5.525	6.773	5.631
<i>M. Formicicum</i>	5.473	5.312	6.336
<i>M. Musculus</i>	5.006	5.979	6.580
<i>M. Sedula</i>	6.633	6.499	5.101
<i>M. Thermoautotrophicum</i>	5.148	5.439	6.857
<i>P. Purpurea</i>	7.598	7.796	8.279
<i>S. Carlsbergensis</i>	7.335	6.399	6.955
<i>S. Pastorianus</i>	7.743	7.265	6.569
<i>S. Uvarum</i>	6.519	6.836	7.334
<i>S. Versatilis</i>	5.251	6.482	5.359
<i>U. N2 fixer TP1</i>	5.177	5.773	5.489
<i>U. N2 fixer WC1</i>	5.901	6.789	6.111
<i>V. Cholerae</i>	7.561	7.853	6.763
<i>Y. Pestis</i>	7.104	6.790	7.811

determined the best value for “population size” as a function of length of RNA sequences. In addition, we found that the best “mutation rate” is 0.2 and “crossover rate” is 0.8.

5. Discussion

In this manuscript, we proposed a new genetic algorithm for RNA secondary structure prediction based on a subclass of Tree Adjoining Grammars. This algorithm gives accurate results in a reasonable amount of time and has specificity and sensitivity comparable with existing tools such as Hotknot, Pknots, and the STAR algorithms. Because of the capability of regular ESL-TAG to cover a large class of RNAs' knotted structures, the time complexity of the parsing algorithm damages its integrity, where in our study, we modified implementation of the parsing algorithm to reduce running complexity.

Since using thermodynamics rules has been demonstrated as the most successful approach for prediction of single RNA sequence, we used these rules in fitness function. We designed novel ideas for the evolutionary functions which yielded to accurate predictions close to the natural structure. The eligibility of tree adjoining grammar for

capturing RNAs with knotted structure and also its efficient data structure led us to employ this representation of RNA in the genetic algorithm. It is strongly believed this data structure could be improved and extended toward covering bigger class of knotted structures and a heuristic algorithm like genetic algorithm can use it to predict complicated knotted structures. The algorithm, PreRkTAG, was implemented and tested on a set of 24 different RNA sequences and compared with three well-known algorithms.

The results of the algorithm show that sensitivity of seventeen out of twenty-four samples and specificity of eighteen out of the twenty-four samples are higher than three well-known algorithms. These results show the practicality of the algorithm. This algorithm handles certain class of pseudoknots captured by ESL-TAG and in future work, we intend to extend the algorithm for larger class of pseudoknots.

Acknowledgements

We would like to express our sincere thanks to Professor C.W. Pleij for his kindly helps and guidance on dis-

cussions that led to an improved understanding of the problem characteristics. The authors are also very grateful to Professor D.H. Turner for his guidance on solving our biophysical problems. Finally we thank Professors F.H.D. van Batenburg, M. Zuker, R. Markham, and E. Rivas for their helps for comparing our method vs. theirs. The authors would like to thank the anonymous referees for their helpful suggestions.

Authors' Contribution

Initial idea of the research was from Torabi Dashti, Ahrabian, Nowzari-Dalini, and Zare-Mirakabad. All authors participated in algorithm design, and the structure and organization of the manuscript. Torabi Dashti and Aghaepour implemented the algorithms and tested on different data sets. All authors contributed to read and approved the final manuscript.

Financial Disclosure

None declared.

Funding/Support

The authors would like to acknowledge the financial support of University of Tehran for this research under grant number 6103010/1/02.

References

- van Batenburg FH, Gulyaev AP, Pleij CW. An APL-programmed genetic algorithm for the prediction of RNA secondary structure. *J Theor Biol.* 1995;**174**(3):269-80.
- Hu YJ. GPRM: A genetic programming approach to finding common RNA secondary structure elements. *Nucleic Acids Res.* 2003;**31**(13):3446-9.
- Lyngso RB, Pedersen CN. RNA pseudoknot prediction in energy-based models. *J Comput Biol.* 2000;**7**(3-4):409-27.
- Markham NR, Zuker M. DINAMelt web server for nucleic acid melting prediction. *Nucleic Acids Res.* 2005;**33**(Web Server issue):W577-81.
- Mathews DH, Turner DH. Prediction of RNA secondary structure by free energy minimization. *Curr Opin Struct Biol.* 2006;**16**(3):270-8.
- Rivas E, Eddy SR. A dynamic programming algorithm for RNA structure prediction including pseudoknots. *J Mol Biol.* 1999;**285**(5):2053-68.
- Ruan J, Stormo GD, Zhang W. An iterated loop matching approach to the prediction of RNA secondary structures with pseudoknots. *Bioinformatics.* 2004;**20**(1):58-66.
- Zuker M. Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic Acids Res.* 2003;**31**(13):3406-15.
- Eddy SR, Durbin R. RNA sequence analysis using covariance models. *Nucleic Acids Res.* 1994;**22**(11):2079-88.
- Zare-Mirakabad F, Sadeghi M, Ahrabian H, Nowzari-Dalini A. RNA-Comp: a new method for RNA secondary structure alignment. *MATCH Commun. Math Comput Chem.* 2009;**61**(3):789-816.
- Gorodkin J, Stricklin SL, Stormo GD. Discovering common stem-loop motifs in unaligned RNA sequences. *Nucleic Acids Res.* 2001;**29**(10):2135-44.
- Gutell RR, Power A, Hertz GZ, Putz EJ, Stormo GD. Identifying constraints on the higher-order structure of RNA: continued development and application of comparative sequence analysis methods. *Nucleic Acids Res.* 1992;**20**(21):5785-95.
- Tahi F, Gouy M, Regnier M. Automatic RNA secondary structure prediction with a comparative approach. *Comput Chem.* 2002;**26**(5):521-30.
- Akutsu T. Dynamic programming algorithms for RNA secondary structure prediction with pseudoknots. *Discrete App Math.* 2000;**104**(1-3):45-62.
- Reeder J, Giegerich R. Design, implementation and evaluation of a practical pseudoknot folding algorithm based on thermodynamics. *BMC Bioinformatics.* 2004;**5**:104.
- Uemura Y, Hasegawa A, Kobayashi S, Yokomori T. Tree adjoining grammars for RNA structure prediction. *Theoret Comp Sci.* 1999;**210**(2):277-303.
- Zuker M, Jaeger JA, Turner DH. A comparison of optimal and suboptimal RNA secondary structures predicted by free energy minimization with structures determined by phylogenetic comparison. *Nucleic Acids Res.* 1991;**19**(10):2707-14.
- Walter AE, Turner DH, Kim J, Lyttle MH, Muller P, Mathews DH, et al. Coaxial stacking of helices enhances binding of oligoribonucleotides and improves predictions of RNA folding. *Proc Natl Acad Sci U S A.* 1994;**91**(20):9218-22.
- Cai L, Malmberg RL, Wu Y. Stochastic modeling of RNA pseudoknotted structures: a grammatical approach. *Bioinformatics.* 2003;**19**(Suppl 1):i66-73.
- Shapiro BA, Navetta J. A massively parallel genetic algorithm for RNA secondary structure prediction. *J Supercomp.* 1994;**8**(3):195-207.
- Ren J, Rastegari B, Condon A, Hoos HH. HotKnots: heuristic prediction of RNA secondary structures including pseudoknots. *RNA.* 2005;**11**(10):1494-504.
- Kato Y, Seki H, Kasami T. On the generative power of grammars for RNA secondary structure. *IEICE Trans Inf Syst.* 2005;**88**(1):53-64.
- Joshi AK, Levy LS, Takahashi M. Tree adjunct grammars. *J Comp Sys Sci.* 1975;**10**(1):136-63.
- Goldberg DE. *Genetic algorithms in search, optimization, and machine learning.* Addison-Wesley. 1989
- Freier SM, Kierzek R, Jaeger JA, Sugimoto N, Caruthers MH, Neilson T, et al. Improved free-energy parameters for predictions of RNA duplex stability. *Proc Natl Acad Sci U S A.* 1986;**83**(24):9373-7.