

RNAshapes 2.1.5 manual

Peter Steffen¹, Björn Voß², Marc Rehmsmeier³, Jens Reeder¹ and Robert Giegerich¹

¹Faculty of Technology, Bielefeld University, 33594 Bielefeld, Germany,

²Institute of Biology II, Experimental Bioinformatics, Freiburg University,
Schänzlestr. 1, 79104 Freiburg, Germany,

³Center for Biotechnology (CeBiTec), Bielefeld University, 33594 Bielefeld, Germany

April 15, 2008

Contents

1	Introduction	1
2	The RNAshapes interface by example	2
2.1	Shape representative analysis	2
2.2	Shape probabilities	4
2.3	Consensus shapes analysis	6
2.4	Additional options	7
3	Options	8
3.1	Sequence analysis modes	9
3.2	Additional modes (use with any of the above)	10
3.3	Analysis control	10
3.4	Input/Output	11
3.5	Additional interactive mode commands	13

1 Introduction

This manual describes RNAshapes, a software package that integrates three RNA analysis tools based on the abstract shapes approach [6]: the analysis of shape representatives [1], the calculation of shape probabilities [7], and the consensus shapes approach [4]. This new package is completely reimplemented in C, and outruns the original implementations significantly in runtime and memory requirements. Additionally, we added a number of useful features like suboptimal folding with correct dangling energies, structure graph output, shape matching, and a sliding window approach.

In the following, we will shortly review the notion of abstract shapes, and explain where its power comes from. We will then provide an overview of the problems that can be approached in the new way. Section 2 gives an exemplified introduction to the RNAshapes interface. Section 3 describes all program options in detail.

The abstract shapes approach. An RNA shape is an abstract representation of an RNA secondary structure. It is inspired by the dot-bracket representation known from the Vienna RNA package [3]. Consider the following sequence and two secondary structures from its folding space in dot-bracket representation:

```
AUCGGCGCACAGGACAUCCUAGGUACAAGGCCGCCCGUU
..(((.(...(((....)))..(((....))))))..
..(((....(((....)))..(((....)))..)))..
```

The shapes approach offers five abstraction levels – or *shape types* – ordered in their degree of abstraction. Common to all levels is that they abstract from loop and stack lengths, where unpaired regions are represented by an underscore and stacking regions by a pair of squared brackets. This is the least abstract shape type 1, so the two example secondary structures become:

```
-[[-[-]-[-]]]-
-[[-]-[-]]-
```

The succeeding shape types gradually increase abstraction, ending in type 5, where no unpaired regions are included and nested helices are combined. In this type, our example structures are both represented as:

```
[[][]]
```

For a detailed description of all shapes types see Section 3. These abstractions form the basis of all applications of RNA abstract shape analysis. In the following we give a short overview of the main applications, all integrated in the RNASHAPES package.

Shape representative analysis. Current RNA folding algorithms either calculate a single, minimum free energy prediction, or a huge number of suboptimal structures, most of which are quite similar and therefore redundant. With shapes, we abstract from the concrete secondary structures and only consider classes of structures that fall into different shapes. The *shape representative* (in short: *shrep*) of a shape is the structure with the minimum free energy inside a shape class.

Shape probabilities. In [7], we extended the shapes approach to the computation of shape probabilities. The probability of a shape is the sum of the probabilities of all structures that fall into this shape. Several analyses indicate that this approach is quite effective. For example, an analysis of a conformational switch shows the existence of two shapes with probabilities approximately $\frac{2}{3}$ vs. $\frac{1}{3}$, whereas the analysis of a micro RNA precursor reveals the hairpin shape with a probability near to 1.0 [7].

The new implementation contains three approaches for probability analysis, suitable for different input sizes: Complete probability analysis, sampling shapes probability analysis, and fast high probability shape analysis (see Section 2).

Consensus shapes. The well-known Sankoff algorithm [5] for simultaneous RNA sequence alignment and folding is currently considered an ideal, but computationally over-expensive method. Available tools implement this algorithm under various pragmatic restrictions. In [4], we proposed to redefine the consensus structure prediction problem in a way that does not imply a multiple sequence alignment step. For a family of RNA sequences, our method RNACAST explicitly and independently enumerates the near-optimal abstract shape space, and predicts as the consensus an abstract shape common to all sequences. For each sequence, it delivers the thermodynamically best structure that has this common shape. Since the shape space is much smaller than the structure space, and identification of common shapes can be done in linear time (in the number of shapes considered), the method is essentially linear in the number of sequences. Our evaluation showed that the new method compares favorably with available alternatives [4]. It is particularly useful on sequences with low conservation, where methods based on sequence alignment cannot be employed. We have now integrated RNACAST into the RNASHAPES package.

2 The RNASHAPES interface by example

In the following, we give an exemplified introduction to the RNASHAPES interface. All sequence files used here can be found in the directory `examples`.

2.1 Shape representative analysis

As an example sequence, we use the alanine tRNA of *Natronobacterium pharaonis* (gb: AB003409.1/96-167). It is contained in sequence file `pharaonis.seq` in the `examples` directory. With option `-f`, we

3

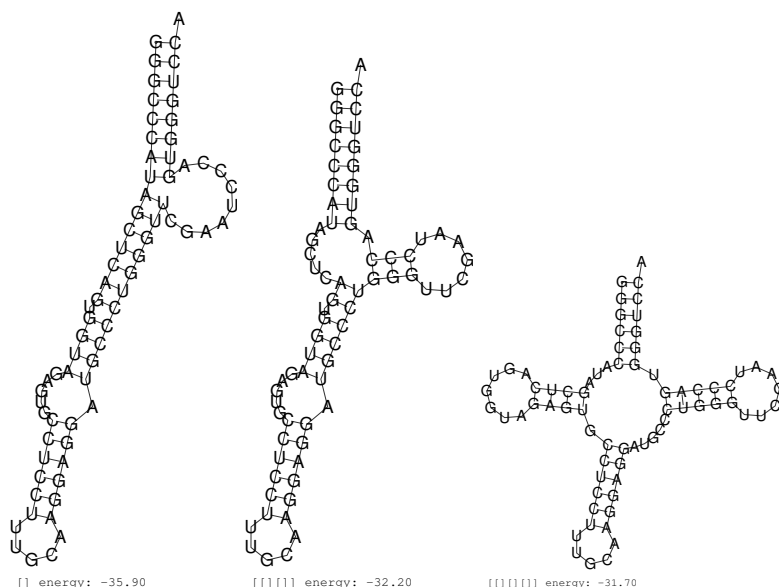


Figure 1: Postscript structure graphs generated with option **-g**

The sequences in the input file are processed one after another, and the results are printed together with the corresponding sequence descriptor.

We can also feed RNASHAPES by standard input, which is useful for directly processing sequence output generated by other programs. For example, the call

```
cat random.fasta | RNASHAPES
```

is equivalent to the **-f** call above. Finally, we can also enter sequences interactively by calling RNASHAPES without any input sequence (but possibly with parameters):

```
RNASHAPES -e 5
Interactive mode. Try 'RNASHAPES -h' for more information.

Input sequence (upper or lower case); :q to quit, -h for help.
.....1.....2.....3.....4.....5.....6.....7.....8
```

The interactive mode supports the following features:

- Direct input of sequences (or alternatively by cut-and-paste).
- Input history with keyboard keys UP and DOWN (library `editline` required).
- Complete interface to all program options. Instead of a sequence, simply type one or several RNASHAPES commands. Examples:
 - **-h** shows the command overview.
 - **-H <option>** shows a detailed help for the given option.
 - **-c 20** sets the energy range to 20%.
 - **-C -f ires.fasta** switches to consensus shapes mode and starts the analysis with file `ires.fasta`.
 - **:s** shows the current settings.

2.2 Shape probabilities

RNASHAPES offers a number of options to calculate structure and shape probabilities. The probability of a shape is the sum of the probabilities of all structures that fall into this shape.

```

RNashapes -q -f pharaonis.seq
GGGCCCCAUAAGCUCAGUGGUAGAGUGCCUCCUUUGCAAGGAGGAUGCCCUUGGUUCGAAUCCCAAGUGGGUCCA
-35.90 0.9897545 [ ]
-32.20 0.0089891 [ [ ] ]
-31.70 0.0012534 [ [ [ ] ] ]
-27.70 0.0000029 [ [ [ [ ] ] ] ]

```

As a second example, we use the pheS-pheT-Attenuator of *E. coli* (embl:V00291.1/3682-3754). It is known to switch from a translationally inactive to a translationally active conformation under specific conditions. These two conformations correspond to two valleys in the structure landscape that are separated by a saddle point (energy barrier). In terms of shape analysis, this means that two shapes with reasonable probability should be present. The corresponding experiment delivers the following results (example file `attenuator.fasta`):

```

RNAsnpes -q -f attenuator.fasta
> pheS-pheT-Attenuator of E. coli (embl:V00291.1/3682-3754)
      ATCCAGGAGGCTAGCGCGTGAGAAGAGAAACGGAACAGCGCCTGAAAGCCTCCCACTGGAGGCTTTTTTTTG
-21.20  ...(((((((.....)))))).)).(((((((.....)))))).).... 0.5381946  [] []
-20.93  .(((((((((((.....)))))).)).)).(((((((((((.....)))))).)).).... 0.3243870  []
-19.33  ...(((((((((((.....)))))).)).)).(((((((((((.....)))))).)))).... 0.0975740  [ [] []
-19.73  ..((.....))..(((((((.....)))))).)).(((((((((((.....)))))).)))).... 0.0388666  [ [] []
-17.30  ...(((((((.....)))))).(((((((((((.....)))))).)))).... 0.0008492  [ [] [] []
-15.13  ..((.....)).(((((((((((.....)))))).)).)).(((((((((((.....)))))).)))).... 0.0001033  [ [] [] []
-14.20  ...(((((((.....)))))).(((((((((((.....)))))).)))).... 0.0000243  [ [] [] []
-13.10  ..((.....)).(((((((((((.....)))))).)).)).(((((((((((.....)))))).)))).... 0.0000010  [ [] [] []

```

As stated above, RNASHapes offers several options to control probability calculation, where option **-q** is the computational most expensive one. This is due to the fact, that this analysis performs an analysis of the complete folding space. Although the computational effort depends only on the size of the shape space (which is much smaller than the complete folding space), it requires a substantial amount of computer main memory. In our experience, option **-q** can be used with sequences up to 250 bases on a computer with 2GB main memory. Additionally, this restriction depends not only on the sequence length, but also on the structural properties of the sequence. To get an impression of the expected running time for a certain calculation, we can use the option **-B** that shows a progress bar:

A slightly more efficient variant of option **-q** above is option **-p**. This performs the same analysis, but without calculation of the corresponding shreps. This option should work with sequences up to a length of around 300 bases (for the presentation, we only use the "short" attenuator sequence here):

```

RNAsnp -p -f attenuator.fasta
> pheS-pheT-Attenuator of E. coli (embl:V00291.1/3682-3754)
ATCCAGGAGGCTAGCGCGTGAGAAGAGAAAAACGGAAAAACAGCGCCTGAAAGCCTCCACAGTGGAGGCTTTTTTG
0.5381946  []
0.3243870  []
0.0975740  [[]]
0.0388666  [][]
0.0008492  [[]]
0.0001033  [[]]
0.0000243  [[]]
0.0000010  [][]

```

```

RNashapes -P 5 -f attenuator.fasta
> pheS-pheT-Attenuator of E. coli (embl:V00291.1/3682-3754)
      ATCCAGGAGGCTAGCGCGTGAGAAAGAGAAAAACGAGCGCCTGAAAAGCCTCCCACTGGAGGCTTTTTTTTG
-21.20  ...(((((((.....((((((((.....)))))))))).))))).((((((((.....))))))))). .... 0.5381930  [] []
-20.93  .(((((((((((.....((((((((.....)))))))))).))))).))))).))))).))))). .... 0.3243860  []
-19.33  ..(((((((((((.....((((((((.....)))))))))).))))).))))).))))).))))). .... 0.0975737  [[] []]
-19.73  ..((((.....)).....((((((((.....)))))))))).))))).((((((((.....))))))))). .... 0.0388665  [] [] []
(Only 4 of 5 probabilities calculated. To get more results, increase energy range (-e or -c).)

```

The last, and most powerful approach for long sequences, is the sampling shapes probability analysis. The sampling shapes approach works in the same manner as Ding and Lawrence’s Sfold program. In each step of the recursive backtracing procedure, base pairs and the structural element they belong to are sampled according to their probability, which is obtained from the partition function. For each sample, we calculate its corresponding shape. The shape probability then results from its frequency in the sample space. The sampling shapes approach is activated with option **-i**, together with the desired number of samples:

```

RNashapes -i 1000 -f attenuator.fasta
... samples ...
Results for 1000 iterations:

-21.20 ...(((((((((.....)))).....)))).(((((((.....)))))))..... 0.5050000  []
-20.93 .(((((((((((.(((((((.....)))).....)))).(((((((.....)))))))..... 0.3560000  []
-19.33 ...(((((((.....(((((((.....)))).....)))).(((((((.....)))))))..... 0.0990000  [[]]
-19.73 ..((.....))..(((((((.....)))).....))..(((((((((((.....)))))))..... 0.0390000  [][]
-13.60 .....((.....)).....(((((((.....)))).....))..(((((((((((.....)))))))..... 0.0010000  [][] []

```

Finally, we can also calculate the probabilities of individual structures by adding the option **-r** to the program call. This option can be used with any analysis mode of RNAsHapes. For example, with our first analysis from above, we receive the following result:

```

RNAsnpes -f pharaonis.seq -e 5 -r
                GGGCCCAUAGCUCAGUGGUAGAGUGCCUCCUUUGCAAGGAGGAUGCCUGGGUUCGAAUCCAGUGGGGUCCA
-35.90  (0.4012148)  ((((((((((((((((((.((((((.....((((((.....)))))).)))))))))).)))))))).  []
-32.20  (0.0009912)  ((((((((((((((((((.((((((.....((((((.....)))))).)))))))))).)))))))).  [[[]]
-31.70  (0.0004404)  ((((((((((((((((((.((((((.....)))))).)))))))))).)))))))).  [[[]]]

```

2.3 Consensus shapes analysis

6

```

RNashapes -C -f ires.fasta
1) Shape: [[] []] Score: -223.50 Ratio of MFE: 0.99
> EMCBCG
CCUUUGCAGGCAGCGGAAAUCCCCACCUGGUAACAGGUGCCUCUGCGGCCAAAAGCCACGUGUAUAAGAUACACCUGCAAAGG
-34.10 ((((((((((.....((((.....)))))))).))))((.....))..((((.....))))))))) [[] []] R = 2
> MNGPOLY
CCUUUGCAGGCAGCGGAAUCCCCACCUGGUGACAGGUGCCUCUGCGGCCGAAAGCCACGUGUGUAAGACACACCUGCAAAGG
-39.10 ((((((((((.....((((.....)))))))).))))((.....))..((((.....))))))))) [[] []] R = 2
> FDI251473
GCACGCAAGCCGCGGGAACUCCCCUUGGUAACAAGGACCCGCGGGGCCAAAAGCCACGUUCUCUGAACCUUGCGUGU
-34.10 ((((((((((.....((((.....)))))))).))))((.....))..((((.....))))))))) [[] []] R = 2
> FMDVALF
GCAUGAUGGCGUGGGAAACUCCCCUUGGUAACAAGGACCCACGGGGCCAAAAGCCACGUCCUCACGGACCCAUGGC
-34.70 ((((((((((.....((((.....)))))))).))))((.....))..((((.....))))))))) [[] []] R = 2
> FAN133359
GCAUGACGGCCGUGGGAACUCCUCCUUGGUAACAAGGACCCACGGGGCCAAAAGCCACGCCACACGGGCCCGUCAUGU
-41.90 ((((((((((.....((((.....)))))))).))))((.....))..((((.....))))))))) [[] []] R = 1
> PIFMDV2
GCAUGUUGGCCGUGGGAACACCUCUUGGUAACAAGGACCCACGGGGCCGAAAGCCAUUGCUAACGGACCCAACAUUGU
-39.60 ((((((((((.....((((.....)))))))).))))((.....))..((((.....))))))))) [[] []] R = 1

```

The consensus shapes are printed in order of their score. Here, the score is the sum of the shrep energies. The "ratio of MFE"-value specifies the ratio of this score to the sum of the minimum free energies. A ratio near 1.0 means a good conservation, a lower ratio means less conservation. The R-value is the rank of the shape in the shape space of the corresponding sequence.

As in all other analysis modes, the shape type can be controlled with option **-t**. For consensus shape analysis, we generally prefer to work with the less abstract level 3:

```

RNashapes -C -f ires.fasta -t 3

```

To get more results, we can also increase the energy range with options **-e** and **-c**, for example:

```

RNashapes -C -f ires.fasta -t 3 -e 10

```

We propose to use the output of the consensus shapes analysis as input for RNAforester [2], a multiple RNA structure alignment program. Use output type **-o f** together with option **-C** to generate suitable input for RNAforester. For example:

```

RNashapes -C -f ires.fasta -o f | RNAforester -m

```

Note that with output type **-o f** only the result for the first consensus is printed (otherwise RNAforester would not work properly). Use the shape match option **-m** to get alternative results. RNAforester is now part of the Vienna RNA package and can be downloaded at <http://www.tbi.univie.ac.at/~ivo/RNA/>.

2.4 Additional options

Sliding window mode. Apart from consensus shapes, each analysis mode can be used with a sliding window mechanism. Here, the complete input sequence is processed by individual calculations of subsequences of the specified window size.

For example, the file V00291 contains the E.coli thrS, infC, rplT, pheS, pheT and himA genes (embl:V00291.1) and has a length of 7784 bases. To start the probability analysis with a window size of 73, we call (this will take some time, maybe stop the calculation with **ctrl-c**):

```

RNashapes -q -f V00291 -w 73
>V00291 V00291.1 12-SEP-2005
1
gattcagtttatgctgctgtaaatccgctcgagtaaacctttcagacgcacggtgatgttatcagttgttct
-8.80 .....(((((((.....)))))))).(((((((.....)))))))). 0.6994365 []
-8.20 ((((((((((.....((((.....)))))))).))))((.....))..((((.....))))))))) 0.2448787 [[] []]
-7.50 ((((((((((.....((((.....)))))))).))))((.....))..((((.....))))))))) 0.0425042 []
-5.70 ((((((((((.....((((.....)))))))).))))((.....))..((((.....))))))))) 0.0101577 [[] []]
-4.90 ((((((((((.....((((.....)))))))).))))((.....))..((((.....))))))))) 0.0014580 [[] []]

```

```

-4.40 .....(((((((.....)))))).....(((((((.....)).(.....))......))) 0.0007761 [] [] []
-4.30 .....((...(((((((.....))))))...((...)).)...(((((((.....))...... 0.0005654 [] [] []
-3.30 ((((((((((.....)))))).....(((((.....)).(.....))....))) 0.0001179 [] [] []
-3.60 ((((((((((.....))))))...((...)).)...(((((((.....))......))) 0.0000690 [] [] []
-2.40 ((((((((((.....)))))).....(((((.....)).(.....))....))) 0.0000288 [] [] []
-1.30 ((((((((((.....)))))).....(((((((.....)).(.....))......))) 0.0000046 [] [] []
-1.00 ((.(((((((.....))))(((((.....)))).....)).(((((.....)).(.....)).... 0.0000021 [] [] []
-1.30 ((.(((((((.....))))(((((.....)))).....)).(((((.....)).(.....)).... 0.0000010 [] [] []

```

This is the result for the first “window”, the bases 1-73. After this, the window is moved by one base, and the calculation continues with bases 2-74, and so on. Note that these succeeding calculations are faster than the first one, since only a single column of the dynamic programming matrices has to be calculated in each window step. The window step size can be changed with option **-W**.

The result for each window is the same as if we would calculate the corresponding subsequence individually. For example, the window 3682-3754 gives the same result as our pheS-pheT-Attenuator analysis from above, as this is exactly the same sequence.

It would be nice to extend the window approach in a way that the results for the individual windows are merged to a result for the complete sequence. Currently, this is not implemented in RNASHapes.

Suboptimal folding. RNASHapes offers a complete suboptimal folding mode. Its implementation is based on a nonambiguous RNA grammar and handles dangling energies correctly. It is activated with option **-s**:

```

RNASHapes -f pharaonis.seq -s -e 5
GGGCCCAUAGCUCAGUGGUAGAGUGCCUCCUUGCAAGGAGGAUGCCCGGGUUCGAAUCCAGUGGGUCCA
-31.10 ((((((((((.....(((((((.....)))))).....))))))((.....)).)))))))). [] []
-31.10 ((((((((((.....(((((((.....)))))).....))))))((.....)).)))))))). [] []
-32.20 ((((((((((.....(((((((.....)))))).....))))))((.....)).)))))))). [] []
.... (69 results more) ...

```

Shape matching. To see only those structures, that fall into a certain shape, we can use the “shape match” option **-m**. For example, to see all clover-leaf structures in the energy range 5 kcal/mol above the mfe, we call:

```

RNASHapes -f pharaonis.seq -s -e 5 -m '[[[] []]]'
GGGCCCAUAGCUCAGUGGUAGAGUGCCUCCUUGCAAGGAGGAUGCCCGGGUUCGAAUCCAGUGGGUCCA
-31.60 ((((((((((.....)))))).....(((((((.....)))))).....)))))).....(((((((.....)))))))). [] [] []
-31.60 ((((((((((.....)))))).....(((((((.....)))))).....)))))).....(((((((.....)))))))). [] [] []
-31.70 ((((((((((.....)))))).....(((((((.....)))))).....)))))).....(((((((.....)))))))). [] [] []
-31.70 ((((((((((.....)))))).....(((((((.....)))))).....)))))).....(((((((.....)))))))). [] [] []

```

The “shape match” option is also available in all other analysis modes of RNASHapes.

Output control. RNASHapes offers a number of options to control the program output. The first option, **-S**, splits structures into smaller parts. This is especially useful when working with long sequences where the program output can be quite confusing for manual inspection. For example:

```

RNASHapes -f pharaonis.seq -S 50
-35.90 1 50
GGGCCCAUAGCUCAGUGGUAGAGUGCCUCCUUGCAAGGAGGAUGCCUG
(((((.....(((((((.....)))))).....))))))
51 72
GGUUCGAAUCCAGUGGGUCCA
))).....)))))))). []

```

The option **-O** can be used to “fine-tune” the format of the printed results, for example when we are only interested in parts of the result, or when results of RNASHapes should be used as input for other programs. See Section 3 for a detailed description.

3 Options

The following section gives a complete description of the RNASHapes command line interface.

-h Display command option overview

-H option Display detailed information on **option**

This displays the corresponding section of the RNASHAPES manual for the given command line option.

-v Show version

This shows the version number of RNASHAPES.

3.1 Sequence analysis modes

-a Shape folding (standard mode)

RNA folding based on abstract shapes. This is the standard mode of operation when no other options are given. It calculates the shapes and the corresponding shreps based on free energy minimization. The energy range can be set by **-e** and **-c**. When not specified, the energy range is set to 10% of the minimum free energy.

-s Complete suboptimal folding

Complete suboptimal folding of RNA. This mode uses a non-ambiguous grammar that also handles dangling bases of multiloop components in a non-ambiguous way. The energy range can be set by **-e** and **-c**. When not specified, the energy range is set to 10% of the minimum free energy.

-p Shape probabilities

Shape probability mode. This option calculates the shape probabilities based on partition function. The probability of a shape is the sum of the probabilities of all structures that fall into this shape. On a computer with 2GB main memory, sequences up to a length of 300 bases can be processed with this mode.

-q Shape probabilities (including shreps)

Shape probability mode. Calculates the shape probabilities based on partition function. This is the same as **-p**, and additionally, the corresponding shreps with their minimum free energies are calculated. Note that this mode is slightly slower than **-p** and can be used with sequences up to a length of 250 bases.

-P value Shape probabilities for mfe-best shapes

Shape probability mode. This mode first calculates the best **value** shapes based on free energy minimization. In a second step, it calculates the probability for each of these best shapes. This mode has lower memory requirements than modes **-p** and **-q** and can be used for longer sequences (up to 500 bases). The energy range must be specified with **-e** or **-c** in order to get the desired number of results.

-i value Sampling with **value** iterations

Probabilistic sampling based on partition function. This mode combines stochastic sampling with a-posteriori shape abstraction. A sample from the structure space holds M structures together with their shapes, on which classification is performed. The probability of a shape can then be approximated by its frequency in the sample.

Sequences up to a length of around 1500 can be handled with this mode. In our experience, 1000 iterations are sufficient to achieve reasonable results for shapes with high probability.

-C Consensus shapes (RNAcast)

For a family of RNA sequences, this method independently enumerates the near-optimal abstract shape space, and predicts as the consensus an abstract shape common to all sequences. For each sequence, it delivers the thermodynamically best structure which has this common shape. Since the shape space is much smaller than the structure space, and identification of common shapes can be done in linear time (in the number of shapes considered), the method is essentially linear in the number of sequences. Input for RNAcast must be provided in multiple fasta format.

We propose to use the output of the consensus shapes analysis as input for RNAforester [2], a multiple RNA structure alignment program. Use output type **-o f** together with option **-C** to generate suitable input for RNAforester. For example:

```
RNAshapes -f test.fasta -C -o f | RNAforester -m
```

Note that with output type **-o f** only the result for the first consensus is printed (otherwise RNAforester would not work properly). Use the shape match option **-m** to get alternative results. RNAforester is now part of the Vienna RNA package and can be downloaded at <http://www.tbi.univie.ac.at/~ivo/RNA/>.

3.2 Additional modes (use with any of the above)

-r Calculate structure probabilities

This calculates the probability of every computed structure. It can be combined with any sequence analysis mode. Note that this option increases processing time of modes **-a**, **-s** and **-C**.

-w value Specify window size

Beginning with position 1 of the input sequence, the analysis is repeatedly processed on subsequences of the specified size. After each calculation, the results are printed out and the window is moved by the window position increment (**-W**), until the end of the input sequence is reached.

-W value Specify window position increment (use with **-w**) (default: 1)

This specifies the increment for the window analysis mode (**-w**).

-m shape Match shape (use with **-a**, **-s**, **-p**, **-q**, or **-C**)

Specify a shape for the corresponding mode of operation. For example, with options **-p -m '[]'** the probability of shape `[]` is computed.

3.3 Analysis control

-e value Set energy range (kcal/mol)

This sets the energy range for shape folding (**-a**), complete suboptimal folding (**-s**), probability analysis with **-P**, and consensus shapes analysis (**-C**). **value** is the difference to the minimum free energy for the sequence.

-c value Set energy range (%) (default: 10)

This sets the energy range as percentage value of the minimum free energy. For example, when **-c 10** is specified, and the minimum free energy is -10.0 kcal/mol, the energy range is set to -9.0 to -10.0 kcal/mol.

-t value Specify shape type (1-5) (default: 5)

The shape type is the level of abstraction or dissimilarity which defines a different shape. In general, helical regions are depicted by a pair of opening and closing square brackets and unpaired regions are represented as a single underscore. The differences of the shape types are due to whether a structural element (bulge loop, internal loop, multiloop, hairpin loop, stacking region and external loop) contributes to the shape representation: Five types are implemented. Their differences are shown in the following example:

```
AUCGGCGCACAGGACAUCCUAGGUACAAGGCCGCCCGUU
..(((.(...((....)).(((.....))))))..
```

Type 5: Most abstract - helix nesting pattern and no unpaired regions

[[] []]

Type 4: helix nesting pattern and unpaired regions in external loop and multiloop

[[] []]

Type 3: nesting pattern for all loop types but no unpaired regions

[[[] []]]

Type 2: nesting pattern for all loop types and unpaired regions in external loop and multiloop

[_[[[] []]]

Type 1: Most accurate - all loops and all unpaired

-[[_[]_[]]]_-

-F value Set probability cutoff filter (use with **-p**, **-q** or **-P**)

This option sets a barrier for filtering out results with very low probabilities during calculation. The default value here is 0.000001, which gives a significant speedup compared to a disabled filter. Note that this filter can have a slight influence on the overall results. To disable this filter, use option **-F 0**.

-T value Set probability output filter (use with **-p**, **-q** or **-P**)

This option sets a filter for omitting low probability results during output. Unlike **-F**, this option does not have any influence on probabilities beyond this value.

-M value Set maximum loop length (default: 30) (use **-M n** for unrestricted)

This option sets the maximum lengths of the considered internal and bulge loops. The default value here is 30. Note that this restriction can have a very slight influence on the calculated structure and shape probabilities. For unrestricted loop lengths, use option **-M n**. This will increase calculation times and memory requirements.

-y value Set minimal shape length

This option sets the minimal shape length. Subshapes smaller than **value** are omitted from the analysis.

-l Allow lonely base pairs

In default mode, RNASHapes only considers helices of length 2 or longer. With option **-l**, lonely base pairs are also included.

-u Ignore unstable structures (use with **-a**, **-s** or **-C**)

This option filters out closed structures with positive free energy.

3.4 Input/Output

-o value Specify output type (1-4,f) (default: 2)

Specifies the output type. Output type 1 mimics RNAfold and RNAsubopt. Type 2 is the default RNASHapes output. Type 3 is similar to type 2, but without parentheses and with only a single space between results. This output type can be used for exporting results as a comma separated text-file to other applications like Microsoft Excel. Type 4 is a colored variant of type 2. Additional output types can be defined with option **-O**.

In consensus shapes analysis (**-C**), output type **f** can be used to generate suitable input for RNAforester (a multiple RNA structure alignment program; see **-C** for details).

-O string Specify output format string

The option **-O** can be used to "fine-tune" the format of the printed results, for example when we are only interested in parts of the result, or when results of RNASHAPES should be used as input for other programs. The syntax is as follows:

```
TYPE{FORMAT}...TYPE{FORMAT}
```

where TYPE specifies the result element:

```
D: structure in dot-bracket notation
S: shape string
E: energy
P: shape probability
R: structure probability (option -r)
C: shape rank (option -C)
V: verbatim output, independent of result element
```

FORMAT is the C-format string that shall be used to print the corresponding result element. Typical C-format strings are `%.2f` for a floating point number with two decimal places and `%s` for a string. For example, to print only the structure followed by its energy, we can use `-O 'D{%s\t}E{%.2f}V{\n}'`. The symbol `'\n'` performs a line break, the symbol `'\t'` a tabulator. An ANSI escape sequence can be used with symbol `'\e'` (see Example 4 below).

The standard output types (option **-o**) are defined as follows:

```
1) 'D{%s }E{%.2f} }R{%.7f} }P{%.7f }S{%s}C{ R = %d}V{\n}'
2) 'E{%-8.2f}R{%.7f} }D{%s }P{%.7f }S{%s}C{ R = %d}V{\n}'
3) 'E{%.2f }R{%.7f }D{%s }P{%.7f }S{%s}C{ %d}V{\n}'
4) 'E{%-8.2f}R{%.7f} }D{\e[1;31m%s\e[0m }
   P{\e[1;30m%.7f\e[0m }S{%s}C{ R = %d}V{\n}'
```

-S value Specify output width for structures

This splits the structure strings into parts of the specified length. This option is useful when displaying results for long sequences that would otherwise not fit onto the screen.

-# value Print only the first value results

This option specifies the total number of results to be printed. When this number is reached, the program terminates. Note that this option does not reduce calculation time or memory requirements (except for modes **-s** and **-i**).

-g value Generate structure graphs for first value structures

This generates postscript structure graphs for the first value structures computed for a sequence. If multiple sequences are given, value graphs are generated for each sequence.

The filenames of the structure graphs consist of several parts:

1. When the input sequence is given in fasta format, the first 12 characters of the sequence description are taken. White-spaces and special characters are removed. When no description is available, "rna" is chosen as standard name.
2. The sequence position in window mode (option **-w**).
3. The running number of the result.

For example, the first result of a sequence called "xyz" at position 7 in window mode will be saved in file xyz_7.1.ps.

-L Highlight uppercase characters in structure graphs

Used with option **-g**, this generates postscript structure graphs where all uppercase characters in the input sequence are highlighted. This option is useful for marking interesting regions of the input sequence.

-N Do not include additional information in graph output file

In standard operation, the postscript structure graph generation (option **-g**) generates files with shape, energy, and shape probability (if available) included at the bottom. Use this option to suppress this.

-A Omit samples in output

Omit output of samples in sampling mode (**-i**)

-f file Read input from file

Let RNAsHapes load its input data from **file**. **file** can contain a plain single sequence, or multiple sequences in fasta format. When given multiple sequences, each sequence is processed separately in the order of input.

Valid characters in an input sequence are "ACGU" and "acgu". "T" and "t" will be converted to "U". Other letters are mapped to "N" and will not be paired. All other characters are ignored.

-B Show progress bar (use with **-p**, **-q** or **-P**)

Setting this option activates a progress bar. This is useful when experimenting with options **-p** and **-q**, to get an impression of the expected running time.

-z Enable colors (in interactive mode: disable colors)

This option enables colored output. In interactive mode, this is the default setting, so use **-z** to disable colors here.

-Z Enable colors for dotbracket and shape strings

This option colors dotbracket and shape strings in the result output, such that corresponding structural elements have the same color in both representations.

-D string Convert dotbracket-string to shape (choose type with **-t**)

Convert a dotbracket-string into a shape. Choose the shape type with option **-t**. The default shape type is 5. For example:

```
RNAsHapes -D '(((((((.....))))))...(((.....)))' -t 4
_[]_[]_
```

-U Start graphical user interface

This option starts the graphical user interface included in the RNAsHapes distribution. It requires Java 1.4.2 or later (download from <http://java.sun.com/>). Note that the RNAsHapes distribution for Microsoft Windows includes a slightly different user interface. It does not require Java and additionally, it offers an interactive visualization of the calculated RNA structures.

3.5 Additional interactive mode commands

:s Show current configuration

This command shows the current settings in an interactive session.

:d Reset configuration

This command sets all settings to their default values.

:e string Execute system command

Command **:e** executes a system command. For example, we can use the command **:e gv rna.1.ps** to open a structure graph file created with option **-g** (on a unix machine with gv installed).

:q Quit

This command quits an interactive RNASHAPES session.

References

- [1] R. Giegerich, B. Voß, and M. Rehmsmeier. Abstract Shapes of RNA. *Nucleic Acids Res.*, 32(16):4843–4851, 2004.
- [2] M. Höchsmann, B. Voß, and R. Giegerich. Pure multiple RNA secondary structure alignments: A progressive profile approach. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 1(1):53–62, 2004.
- [3] I.L. Hofacker, W. Fontana, P.F. Stadler, L.S. Bonhoeffer, M. Tacker, and P. Schuster. Fast folding and comparison of RNA secondary structures. *Monatshefte f. Chemie*, 125:167–188, 1994.
- [4] J. Reeder and R. Giegerich. Consensus shapes: an alternative to the Sankoff algorithm for RNA consensus structure prediction. *Bioinformatics*, 21(17):3516–3523, 2005.
- [5] D. Sankoff. Simultaneous solution of the RNA folding, alignment and protosequence problems. *SIAM J. Appl Math*, 45(5):810–825, 1985.
- [6] P. Steffen, B. Voß, M. Rehmsmeier, J. Reeder, and R. Giegerich. RNASHAPES: an integrated RNA analysis package based on abstract shapes. *Bioinformatics*, 2005. Epub ahead of print.
- [7] B. Voß, R. Giegerich, and M. Rehmsmeier. Complete probabilistic analysis of RNA shapes. 2005. Manuscript under review.