# RNA Secondary Structure Prediction

This unit details the steps for predicting the secondary structure of an RNA sequence using free energy minimization (Mathews et al., 1999). Two protocols are given—one for the computer program RNAstructure (see Basic Protocol) and another for the *mfold* server (see Alternate Protocol). The *mfold* server is a World Wide Web adaptation of the *mfold* package for Unix computers. This package is available for use on Unix platforms and is described elsewhere (Zuker et al., 1999; *http://www.rpi.edu/~zukerm/*). RNAstructure is for personal computers (PCs) with Windows 95, Windows 98, or Windows NT. The *mfold* server is accessed via the Web using a standard browser (e.g., Microsoft Internet Explorer or Netscape Navigator). Both programs yield equivalent predictions. The minimum free energy structure and a set of suboptimal structures with similar free energies are predicted. In addition, a triangular plot called the energy dot plot that shows the superposition of all possible base pairs within a prescribed free energy increment is predicted. The choice of protocol is a matter of convenience for the user. RNAstructure requires a PC, whereas the *mfold* server requires accessing the software across the Internet.

## PREDICTING A SECONDARY STRUCTURE WITH RNAstructure

*BASIC PROTOCOL*

The last few years have seen large advances in the power of the personal computer. Processors have become faster and the standard amounts of memory (i.e., RAM) larger. Furthermore, the 32-bit Microsoft Windows operating systems no longer have the segmented memory limitation of older operating systems. These advances make the PC capable of more significant calculations, such as RNA secondary structure prediction. Table 11.2.1 shows the length of time required for three different personal computers to predict RNA sequences of five lengths.

This protocol lists the steps for structure prediction on a PC using the computer program RNAstructure. It assumes some basic familiarity with Windows.

### Materials

Personal computer with Windows 95, Windows 98, or Windows NT (Microsoft)
RNAstructure (free of charge at the Turner Lab Web site,
*http://rna.chem.rochester.edu*)

*NOTE:* It is recommended that a computer with a 100 MHz or faster processor be used. Also, memory use increases with the length of the sequence being analyzed; therefore

**Table 11.2.1** Time in Minutes for Secondary Structure Prediction by RNAstructure[a]

| Computer | RNA sequence length (nt) | | | | |
|---|---|---|---|---|---|
| | 77 | 268 | 433 | 631 | 1542 |
| PIII | 0.02 | 0.12 | 0.29 | 0.71 | 6.40 |
| PII | 0.02 | 0.33 | 0.75 | 1.65 | 12.89 |
| P90 | 0.05 | 1.29 | 3.01 | 6.72 | 55.00 |

[a]This table displays the folding times for five RNA sequences from length 77 to 1542 nt on three different personal computers. PIII is a 450 MHz Pentium III computer with 192 MB RAM; PII is a 233 MHz Pentium II computer with 64 MB RAM; and P90 is a 90 MHz Pentium computer with 16 MB RAM. The sequences for lengths 77, 268, 433, 631, and 1542 nt are RR1664 tRNA (Sprinzl et al., 1998), *Bacillus stearothermophilus* SRP RNA (Larsen et al., 1998), IVS LSU Group I intron from *Tetrahymena thermophila* (Damberger and Gutell, 1994), *Saccharomyces cerevisiae* A5 Group II intron (Michel et al., 1989), and small subunit rRNA from *E. coli* (Gutell, 1994), respectively.

Contributed by David H. Mathews, Douglas H. Turner and Michael Zuker

32 MB of RAM is recommended for predicting structures for sequences of ~500 nt, although less memory will work. RNAstructure will also work on an Apple Macintosh computer emulating Windows. RNAstructure is known to work on a Macintosh G3 computer using Virtual PC 2.1, available from Connectix, on the web at *http://www. connectix.com*.

### Download and install RNAstructure

This section briefly covers the downloading and installation of the RNAstructure sofware for users not familiar with downloading and installing software.

Use a browser such as Microsoft Internet Explorer or Netscape Navigator to access the Turner Lab Web site and follow the links to download RNAstructure. The RNAstructure homepage is at *http://rna.chem.rochester.edu/RNAstructure.html*. Click on "Register to Download NOW." Registration is required so that all users can be notified of periodic upgrades to the program; the registration information is used for no other purpose. After submitting the registration information, choose a download site and save RNAstructure.zip to a local computer.

To extract the program files for RNAstructure, a decompression program such as PKZip or WinZip is necessary. Links to the companies that distribute these shareware programs are given on the Turner Lab Web site. Unzip the files to a temporary folder on your computer and double-click on *install.exe*. This will start the installation program. After installing the program to the chosen directory, the zip file and installation files can be safely deleted from your computer.

### Overview of RNAstructure

Start RNAstructure either from its folder or from the Windows Start button. RNAstructure is composed of modules that accept a sequence, predict the sequence's structure, and then output the result. All the modules are selected in the File menu. Note that throughout this unit, Windows convention for indicating menu options is used, for example, File|New Sequence refers to the "New Sequence" menu option that appears on the "File" menu.

Online help is available anytime during the execution of RNAstructure using the Windows help system. Simply choose Help|Help. The online help will start at a table of contents. RNAstructure help also contains an index of keywords.

Three modules represent the core of RNA secondary structure prediction. The sequence editor, File|New Sequence, or File|Open Sequence, allows the user to input the nucleotides in the RNA sequence. The RNA Fold Single Strand module calculates the lowest free energy structure and a set of suboptimal structures. The Draw module displays the predicted structures. A fourth module, Dot Plot, will calculate the energy dot plot for a sequence (Zuker, 1989). Other modules extend the functionality of the prediction algorithm. These are described in the RNAstructure online help.

### Input the sequence

RNAstructure requires a specific format, the *.seq* file, for sequences. The sequence editor module, illustrated in Figure 11.2.1, is used to enter a sequence and save it in the proper format. A sequence can be input in one of three ways: (1) reading directly from a GenBank file, (2) copying and pasting from a separate document, or (3) typing the sequence manually. It is important that most nucleotides be in capital letters because lowercase nucleotides are forced to be single stranded by the secondary structure prediction algorithm.
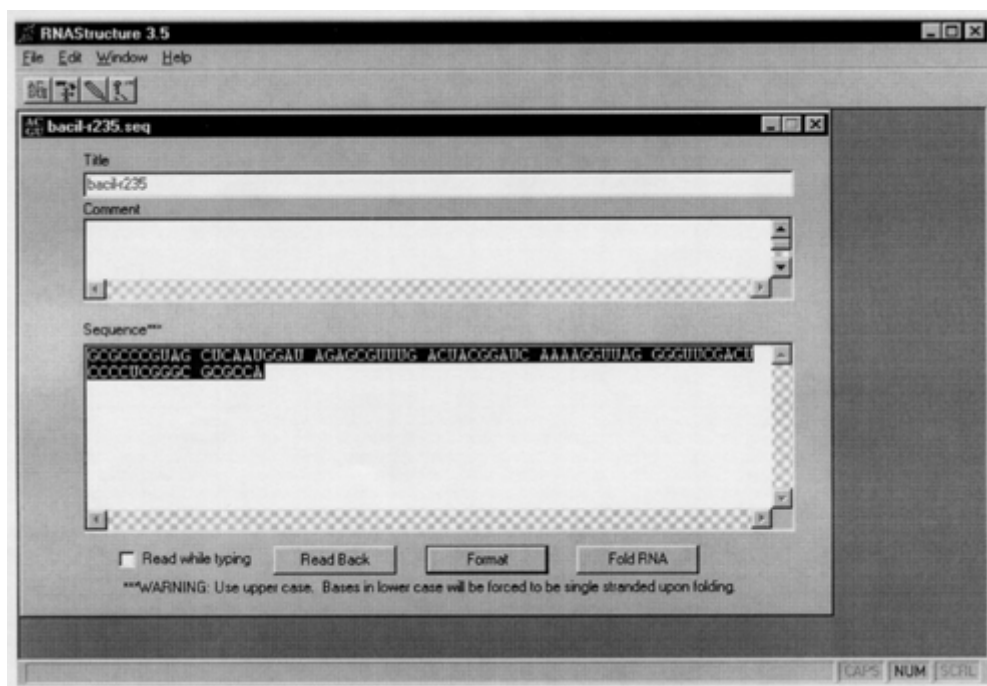
**RNA Secondary
Structure
Prediction**

**11.2.2**

Supplement 2

Current Protocols in Nucleic Acid Chemistry

**Figure 11.2.1**   A screen shot of the sequence editor module of RNAstructure.

If the sequence is saved as a GenBank file, open the sequence editor with File|Open Sequence. This will open the standard Windows Open File dialog. Choose the GenBank file and the sequence editor will open and display the sequence. Save the sequence as a *.seq* file using File|Save.

To enter a sequence either manually or by copy-and-paste, open the sequence editor with File|New Sequence. This will open the sequence editor with all the fields blank. In the first field, "Title," enter a name that will appear with the predicted structure for the sequence. Comments may be placed in the second field. In the third field, "Sequence," enter the RNA sequence from 5′ to 3′. Sequences can be pasted from other programs, such as word processing programs, using Edit|Paste, or typed directly into RNAstructure. When finished, save the sequence as a *.seq* file using File|Save.

### *Features of the sequence editor module*
To check a sequence, RNAstructure will audibly read back the nucleotides in a sequence. If the check box in front of "Read while typing" in the lower left-hand corner is checked, the program will read nucleotides as they are entered. Once a sequence is entered, clicking on the button labeled "Read Back" will produce an audible reading of the sequence.

In an RNA sequence, spaces are ignored and T is treated as U by the prediction algorithm. X can be used for unknown nucleotides or four consecutive Xs can be used to omit a section of the sequence. X does not base pair or stack on other base pairs.

Clicking the button labeled "Format" will format the sequence in lines with 6 columns of 10 nucleotides each for easier viewing of the sequence. Clicking the button "Fold RNA" closes the Sequence Editor module and opens the RNA Fold Single Strand module.

### *Predict an RNA secondary structure*
Open the RNA Fold Single Strand module either by clicking the "Fold RNA" button on the sequence editor window or by choosing File|RNA Fold Single Strand. A window

labeled "Fold RNA Sequence" opens as illustrated in Figure 11.2.2. This window obtains information from the user that is needed for predicting a structure.

Specify a sequence by clicking the "Sequence File" button. (Note that if the window was reached directly from the Sequence Editor module, the sequence last edited will already be specified.) The standard Windows Open File dialog box appears. Choose the *.seq* file that contains the sequence of interest. The Fold RNA Sequence window will now contain the name of the sequence next to the "Sequence File" button. Also, the remaining fields will have default parameters.

The next parameter to specify is the CT file. This file is created by RNAstructure to save the base pairing information. A default name appears to the right of the "CT File" button. It can be changed by clicking on the "CT File" button and specifying a different file in a standard Windows Save File Dialog box.

A check box appears to the left of the label, "Generate Save File." The default, with the box checked, is to create a save file. This file stores free energy data that is calculated by the prediction algorithm. It can be used to quickly calculate a different set of suboptimal structures using the Refold module. A Save file is also necessary to calculate a dot plot.

Finally, three parameters, "Max % energy difference," "Max number of structures," and "Window size," control the output of suboptimal secondary structures. These parameters can be left at their default values. Increasing "Max % energy difference" and "Max number of structures" will increase the number of suboptimal structures that are output. Increasing "Window size" will increase the difference in the suboptimal structures and can therefore reduce the number of suboptimal structures generated.

At this point, optional folding constraints can be entered using the Force menu. Four types of constraints are currently allowed. A base can be forced double-stranded or single-stranded with Force|Double or Force|Single, respectively, according to experimental constraints that can be determined by enzymatic cleavage (see *UNIT 6.1*). A base pair can
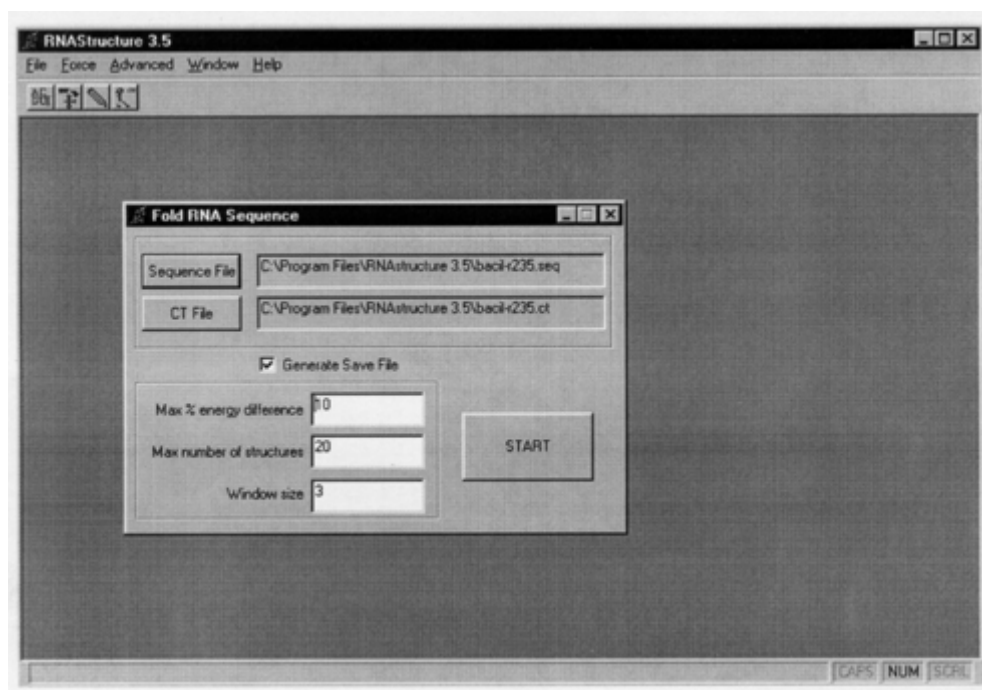


**Figure 11.2.2** A screen shot of the RNA fold single-strand module of RNAstructure.

be specified with Force|Pair. Also, Us can be forced into G-U pairs according to experimental evidence from flavin mononucleotide (FMN) photocleavage (Burgstaller and Famulok, 1997; Burgstaller et al., 1997). The chosen constraints are displayed by Force|Current or erased with Force|Reset. The online help manual describes the use of constraints.

Start the calculation by clicking the button labeled "Start." A progress bar appears on the screen to give the approximate percent of the calculation completed.

### *Display the results*
When the calculation is completed, a dialog box with two buttons is displayed. Choose "Draw Structures" to display the predicted structures. If "Exit" is chosen, the structures may be displayed later by choosing the Draw module from File|Draw and selecting the CT file to be displayed.

The Draw module is now open. Figure 11.2.3 shows a drawing of a predicted tRNA structure. One secondary structure is displayed on the screen at a time. The structure can be enlarged or reduced by specifying the percent zoom on the Zoom dialog box, which is opened by selecting View|Zoom. Alternatively, pressing the control key and the right arrow will enlarge the structure and pressing the control key and the left arrow will reduce the structure. If the window is too small to display the entire structure, scroll bars appear.

Suboptimal structures are selected for display with the Structure dialog box that is opened by selecting View|Structure Number. Alternatively, structures can be chosen by pressing the control key and either the up or down arrow to change to structures with higher or lower free energy, respectively. The upper left hand corner of the Draw module window always specifies the structure number and free energy of the currently displayed structure.

Structures can be printed or copied to the clipboard to be pasted into other programs. File|Print prints the currently displayed structure. Two options are available, "Original



**Figure 11.2.3** A predicted RNA secondary structure as drawn by RNAstructure.

Scale" or "Scale to One Page," in a window that opens. "Original Scale" will print the nucleotides at the same size no matter what the length of the sequence. For long sequences, the structure will span multiple pages. "Scale to One Page" will shrink or enlarge the structure to fit on one page. Edit|Copy copies the currently displayed structure to the clipboard as a bitmap to be pasted into other programs.

### *Dot plots*
The dot plot contains all possible base pairs that can occur within a prescribed free energy increment from the minimum free energy. A base pair in the dot plot is colored according to the minimum free energy of a structure that contains it. It provides an indication of how well the secondary structure is predicted (Zuker and Jacobson, 1995). An energy dot plot of a predicted structure is calculated with the Dot Plot module, opened with File|Dot Plot. Choose a save file for a previously predicted secondary structure in the standard Windows Open File dialog box.

## PREDICTING A SECONDARY STRUCTURE WITH THE *mfold* SERVER

The *mfold* server takes advantage of the Web to make RNA secondary structure prediction available to a large audience. It requires no special setup and is much more user friendly than the unadorned Unix version of *mfold*. This protocol outlines the steps involved in predicting a secondary structure on the *mfold* server using a standard browser.

### *Materials*

Web browser, e.g., Netscape Navigator or Microsoft Internet Explorer, and access to the Web are required.

*NOTE:* This protocol assumes basic familiarity with the Web.

### *Access the server and enter the required information*
The *mfold* server is accessed at *http://www.rpi.edu/~zukerm/*. This page gathers all the necessary information required to predict an RNA secondary structure. There are many fields for entering information; each is labeled, and many of these labels link to explanations of the information required in that field. This protocol describes what information is required, starting at the top of the page and working down.

Be sure to follow the link that reads "Notice," at the top of the page. This provides information to *mfold* users on accessing results, how long results are saved, and security on the server. Other links at the top of the page lead to an RNA page with helpful links to other RNA-related sites, the Zuker Homepage, or e-mail to Michael Zuker.

A name for the sequence must be provided in the first field. An unnamed sequence will be assigned a name according to the date and time of submission. The next field takes the sequence, which can be pasted from other locations or typed directly. As the caption explains, blanks and nonalphabetic characters are ignored. N can be used to represent a nucleotide that will neither pair nor stack and four Ns can be used to join regions of sequence. Also, T is interpreted as U. Currently, a maximum of 3000 nt is allowed.

The next field is for optional folding constraints. These constraints can be determined experimentally with enzymatic cleavage (see *UNIT 6.1*). Five constraints are possible: (1) force nucleotides to be double stranded, (2) force specific base pairs, (3) force nucleotides to be single stranded, (4) prohibit specific base pairs, and (5) prohibit a nucleotide between positions i and j from pairing to nucleotides between positions k and l. Each of these constraints is described in detail on the page that is linked by "constraint information." Next choose whether the RNA strand is linear or circular; the default is linear.

In the next field, choose the folding temperature. The default and recommended temperature is 37°C because the current thermodynamic parameters are most accurately known at that temperature (Mathews et al., 1999). Structure prediction at any other temperature uses a set of parameters from prior studies (Walter et al., 1994). Therefore, for structure prediction at temperatures close to 37°C, e.g., 30° or 45°C, it is usually better to use the default temperature rather than the exact temperature.

Next, enter the percent suboptimality, upper bound, and window size for suboptimal structure generation. Increasing the percent suboptimality and the upper bound will increase the number of suboptimal structures that are output. Increasing the window size will increase the structural difference in the suboptimal structures and can reduce the number of suboptimal structures.

Choose whether the prediction should be an immediate job (predicted while one waits) or a batch job (an e-mail message is sent when the folding has been completed). Currently, a limit of 500 nt is allowed for immediate jobs. If a batch job is chosen, enter an e-mail address to receive a message when the prediction is complete.

The remaining fields affect the drawing of the predicted structures. These settings are used for generating a compressed file that contains drawings of all the predicted structures and are not used for online viewing of structures. An image resolution of low, medium, or high can be chosen. The structure format can be selected as automatic, bases, or outline. A structure drawn in outline format does not specify the identity of each nucleotide in the structure, whereas bases format shows each nucleotide identity explicitly. Automatic mode will choose either outline or bases depending on the length of the sequence. A base numbering frequency other than default can be chosen. This indicates at what interval the nucleotides should be labeled in the structures that are output. A structure annotation method can be selected (Zuker and Jacobson, 1998). These annotations can be used as an indication of confidence in a predicted base pair (Zuker and Jacobson, 1995).

Finally, start the prediction by following the link labeled "Fold RNA" at the bottom of the page. One will receive notification of the status of the job. The other link at the bottom of the page, labeled "reset form," will return all values to default and erase the sequence.

The RNA secondary structure prediction is fast. A 433-nt Group I intron is folded in <1 min and the whole *E. coli* 16S rRNA, 1542 nt, is folded in <11 min of CPU time (Mathews et al., 1999). The actual time for predicting a structure can vary depending on how busy the server is.

### View the results
When the secondary structure prediction is complete, the output page will be available. For immediate jobs, the *mfold* server moves directly to this page. For batch jobs, the output is available by following the link labeled "View previous foldings" at the top of the *mfold* server homepage. All previous jobs within two days will be listed. One must view the results at the same computer used to submit the job or else remember the exact URL (web address) to access the results.

The top of the output page lists the sequence folded and the suboptimal structures generated. The rest of the output page is composed of links to various methods of displaying the predicted structures. The first method of displaying the information is the energy dot plot, which plots all base pairs that are contained in its structures, within a prescribed free energy increment from the lowest free energy. The approximate lowest free energy possible for a structure that contains a pair is indicated by the color used to plot that pair. The dot plot provides an indication of how well the secondary structure is

determined (Zuker and Jacobson, 1995). Two formats are available. The PostScript format can be viewed only with software designed to display PostScript and can be printed on PostScript-capable printers. The GIF format is native to Web browsers, and following the GIF link will display the dot plot on the screen in an interactive mode. This mode allows the user to resize the plot, change the plot's characteristics, and click on dots to get the exact free energy indicated.

All the predicted structures are available to download in a single compressed file. Choose a compression method and a file format and press "create" to download all structures. The compression is either "tar" or "zip." On a PC, these files can be decompressed with a decompression program such as WinZip, available as shareware at *http://www. winzip.com*. For the Macintosh, decompress the files with a decompression program such as DropStuff with Expander Enhancer from Aladdin Systems, available at *http://www.aladdinsys.com*. On a Unix machine, tar files are decompressed with the tar command. The structures can be formatted as PostScript, GIF, CT, Mac CT, GCG connect, or XRNA ss. A description of these formats is found at the "file formats" link. If in doubt, use the GIF format, because a Web browser can read and display this format. Figure 11.2.4 shows a GIF rendering of a tRNA prediction on the *mfold* server.

Next, structural information in the form of ss-count is available. This indicates the propensity of a nucleotide to be single stranded, and counts the number of times a nucleotide appears single stranded in the set of computed suboptimal structures. This information is available as a text file, at the "(ss-count file)" link or as a plot in GIF or PostScript that is calculated by pressing the "View" button. The online plot can be averaged over several nucleotides or magnified by changing the default values under "View ss-count information."
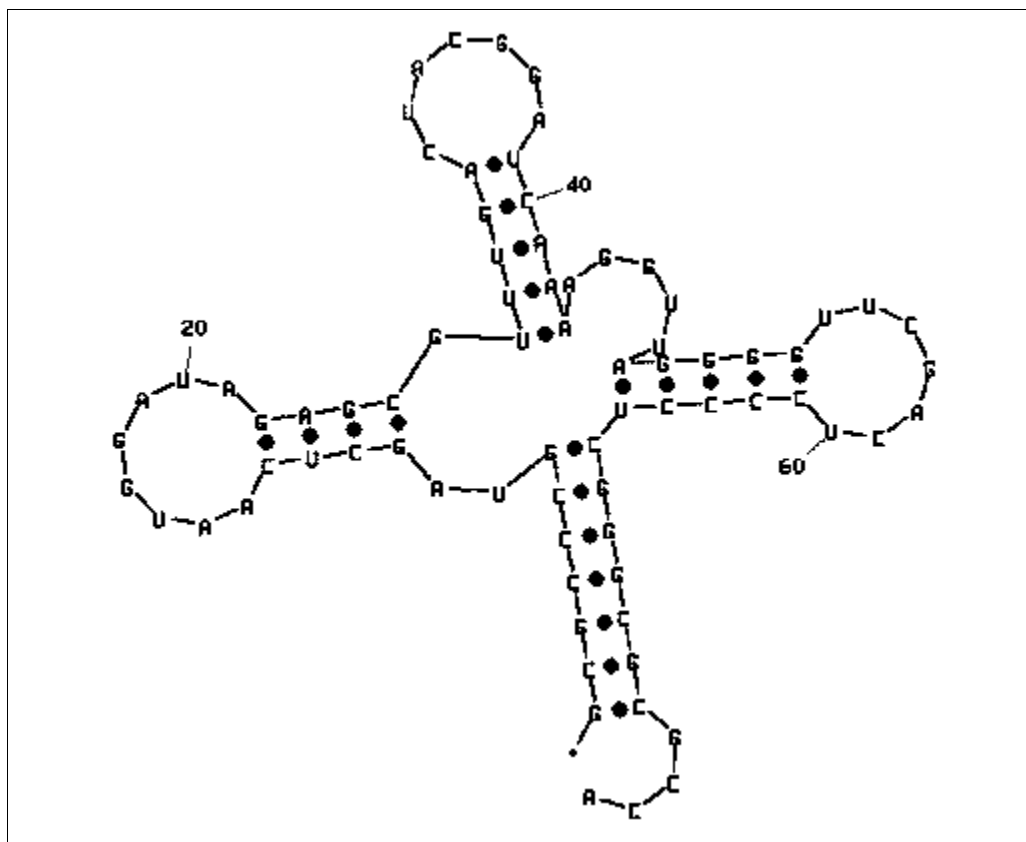


**Figure 11.2.4** A predicted RNA secondary structure as drawn in GIF format by the *mfold* server.

The individual suboptimal structures are available in several formats, under "View Individual Structures." The formats are the same as those available for downloading the whole set of structures. If in doubt, choose the GIF format, because it will be displayed in the Web browser.

Finally, the last method for displaying the predicted structure information is the dot plot folding comparison. This creates a plot of i versus j where each dot indicates a base pair in a suboptimal structure. The plot indicates regions of overlap of the suboptimal structures. It is available as either PostScript or GIF. The GIF format is interactive.

## COMMENTARY

The *mfold* server and RNAstructure predict secondary structures by finding the lowest free energy structure using the same algorithm and parameters (Zuker, 1989; Mathews et al., 1999). These programs calculate the free energy of structures on the basis of nearest-neighbor parameters that are developed from studies of small model systems, examination of the database of known structures, and optimization of the accuracy of the folding algorithm (Mathews et al., 1999; Xia et al., 1998). This method assumes that RNA secondary structure is determined by equilibrium interactions.

### Accuracy

The algorithm is tested by predicting structures for sequences for which secondary structures have been determined by comparative sequence analysis. A recent test was performed with a database of 151,503 nt in 955 structures (Mathews et al., 1999). When RNAs were folded in domains of <700 nt, the lowest free energy structure contained, on average, 73% of known base pairs. For a given sequence, a set of 750 suboptimal structures contains one structure that, on average, has 86% of known base pairs.

### Improving acuracy

The accuracy of secondary structure predictions can be improved with experimental constraints. *UNIT 6.1* discusses the use of enzymes and chemical reagents to probe RNA structures. Enzymatic cleavage can determine bases that are either base paired or single stranded, and flavin mononucleotide (FMN) promotes photocleavage at Us in G·U base pairs (Burgstaller et al., 1997; Burgstaller and Famulok, 1997). Both of these data can be used to constrain the prediction of secondary structures, and a recent study demonstrated that these data improve the accuracy of lowest free energy structures (Mathews et al., 1999). For example, the predicted secondary structure of the td Group I intron from T4 contains 56% of known base pairs. When FMN cleavage data are used to constrain Us during the prediction, the predicted structure contains 83% of known base pairs (Mathews et al., 1999; Burgstaller et al., 1997; Damberger and Gutell, 1994).

Chemical modification accessibility for reagents such as kethoxal and dimethylsulfate cannot be directly applied as a constraint during structure prediction because it not only identifies single-stranded nucleotides, but also nucleotides at the ends of helices. A low free energy structure can be constructed that is consistent with the accessibility data using the program Mix & Match (Mathews et al., 1997).

## Literature Cited

Burgstaller, P. and Famulok, M. 1997. Flavin-dependent photocleavage of RNA at G·U base pairs. *J. Am. Chem. Soc*. 119:1137-1138.

Burgstaller, P., Hermann, T., Huber, C., Westof, E., and Famulok, M. 1997. Isoalloxazine derivatives promote photocleavage of natural RNAs at G·U base pairs embedded within helices. *Nucl. Acids Res*. 25:4018-4027.

Damberger, S.H. and Gutell, R.R. 1994. A comparative database of group I intron structures. *Nucl. Acids Res.* 22:3508-3510.

Gutell, R.R. 1994. Collection of small subunit (16S- and 16S-like) ribosomal RNA structures. *Nucl. Acids Res.* 22:3502-3507.

Larsen, N., Samuelsson, T., and Zwieb, C. 1998. The signal recognition particle database (SRPDB). *Nucl. Acids Res*. 26:177-178.

Mathews, D.H., Banerjee, A.R., Luan, D.D., Eickbush, T.H., and Turner, D.H. 1997. Secondary structure model of the RNA recognized by the reverse transcriptase from the R2 retrotransposable element. *RNA*. 3:1-16.

Mathews, D.H., Sabina, J., Zuker, M., and Turner, D.H. 1999. Expanded sequence dependence of thermodynamic parameters improves prediction of RNA secondary structure. *J. Mol. Biol.* 288:911-940.

Sprinzl, M., Horn, C., Brown, M., Ioudovitch, A., and Steinberg, S. 1998. Compilation of tRNA sequences and sequences of tRNA genes. *Nucl. Acids Res*. 26:148-153.

Walter, A.E., Turner, D.H., Kim, J., Lyttle, M.H., Müller, P., Mathews, D. H., and Zuker, M. 1994. Coaxial stacking of helices enhances binding of oligoribonucleotides and improves predictions of RNA folding. *Proc. Natl. Acad. Sci. U.S.A*. 91: 9218-9222.

Xia, T., SantaLucia, J. Jr., Burkard, M.E., Kierzek, R., Schroeder, S.J., Jiao, X., Cox, C., and Turner, D.H. 1998. Thermodynamic parameters for an expanded nearest-neighbor model for formation of RNA duplexes with Watson-Crick base pairs. *Biochemistry* 37:14719-14735.

Zuker, M. 1989. On finding all suboptimal foldings of an RNA molecule. *Science* 244: 48-52.

Zuker, M. and Jacobson, A.B. 1995. "Well-determined" regions in RNA secondary structure predictions. Applications to small and large subunit rRNA. *Nucl. Acids Res*. 23:2791-2798.

Zuker, M. and Jacobson, A.B. 1998. Using reliability information to annotate RNA secondary structures. *RNA* 4:669-679.

Zuker, M., Mathews, D.H., and Turner, D.H. 1999. Algorithms and thermodynamics for RNA secondary structure prediction: A practical guide. *In* RNA Biochemistry and Biotechnology (J. Barciszewski and B.F.C. Clark, eds.) NATO ASI Series, Kluwer Academic Publishers, Boston.

## Key References

Mathews et al., 1999. See above.

*Derives the thermodynamic parameters used by the secondary structure prediction algorithm and tabulates the accuracy of the algorithm with a large database of structures from sequence comparisons.*

Zuker, 1989. See above.

*Explains the method for predicting suboptimal structures using a dynamic programming algorithm.*

## Internet Resources

http://www.rpi.edu/~zukerm

*Michael Zuker's homepage at the Department of Mathematics at Rensselaer Polytechnic Institute is the home of the mfold online server. It also contains many links to other sites with information on RNA structure.*

http://mfold.burnet.edu.au

*The Macfarlane Burnet Centre for Medical Research maintains a mirror of the mfold server administered by Ewen Bell. Some details of the site's layout differ from the mfold server at Michael Zuker's homepage, but the computations are identical.*

http://rna.chem.rochester.edu

*The Turner Lab homepage is the source for downloading RNAstructure.*

Contributed by David H. Mathews and
    Douglas H. Turner
University of Rochester
Rochester, New York

Michael Zuker
Rensselaer Polytechnic Institute
Troy, New York