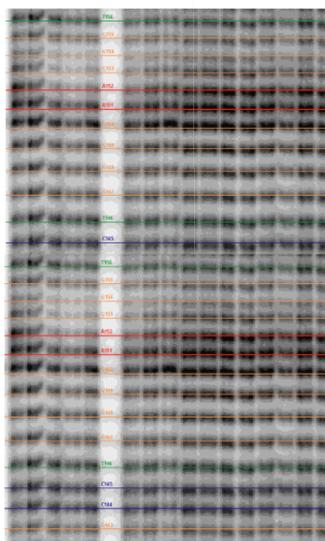


SAFA User's Manual

v1.1

Copyright 2004, Stanford University



SAFA

v.1.1

Semi-Automated Footprinting Analysis

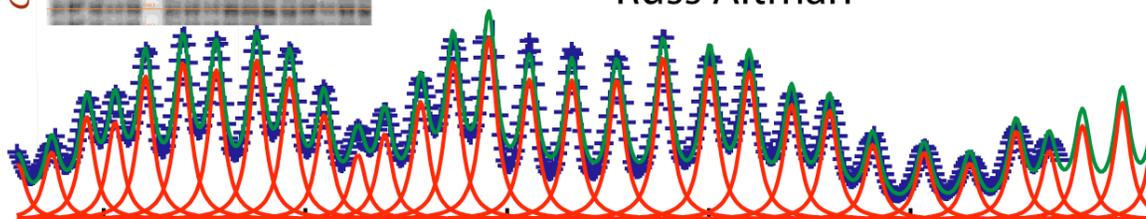
Alain Laederach

Rhiju Das

Sam Pearlman

Dan Herschlag

Russ Altman



Funded by NIH PPG P01-GM-66275. Special thanks to K. Takamoto and M. Brenowitz

Please cite: Das, Laederach, Pearlman, Herschlag, Altman, SAFA: Semi-Automated Footprinting Analysis software for high-throughput quantification of nucleic acid footprinting experiments. RNA 2005, 11: 344-354.

Table of Contents

Table of Contents	2
I. Introduction	3
II. Installation	5
III. Steps in SAFA	6
A. Loading .gel files	6
B. Defining Lanes	7
C. Defining an Anchor Lane	9
D. Aligning the Gel	9
E. Loading a Sequence File	11
F. Choosing cleavage sites and offsets	11
G. Assign Bands	13
H. Quantifying your gel	14
I. Saving Data	16
J. Visualizing and Normalizing Data	17
K. Mapping Data onto Secondary Structure	22
L. Additional Functionality	25
IV. Future Plans	26

I. Introduction

Quantitative analysis of gels from hydroxyl radical footprinting and other structure mapping techniques can provide a great deal of insight into the structural details of RNA molecules. We have developed and implemented a software package (SAFA v0.9b5) that allows rapid quantification of a footprinting gel. By automating many of the steps involved in gel analysis, we estimate that an entire gel with thousands of bands can be quantified in less than 10 minutes, and along with others, we have demonstrated this with numerous gels. In general all the automated features have a manual override, such that even difficult or exceptional gels can be analyzed with the package.

The analysis procedure implemented in SAFA involves several steps:

A) Load Gel

The gel is loaded into the software, and the user is prompted to define the area of the gel to analyze.

B) Define Lanes

The user defines the first lane in the gel manually (i.e. The area of that lane). The computer then guesses the rest of the lane boundaries. The user can always override the computer guess, and erase and add boundaries as desired.

C) Align the Gel

Following lane definition, the gel is aligned. Bands that correspond to the same length product would be aligned vertically in a perfect gel, but temperature imperfections, buffer gradients, and other factors often cause gels to “smile” or “frown.” By drawing lines connecting such bands, anchor lines traversing the gel are assigned. A simple interpolation then transforms the gel image to one where the anchor lines become perfectly horizontal.

D) Assign Bands

Using the aligned gel, only one lane (e.g. a reference ladder) needs to be assigned manually, and the assignments are used as starting guesses for a least squares (Levenberg-Marquardt) fitting of a sum of Lorentzians to the data. This fitting model is further constrained by assuming the width of the band is linearly dependent on the position of the previous and next bands.

E) Quantify

The least squares fit of the sum of Lorentzians to the data can be done all at once, or one lane at a time so that the user can examine the results immediately. Following the completion of the analysis of all lanes, the data may be saved in a text file.

F) Data Visualization

Immediately after the gel has been quantified the user may visualize the individual Lorentzians, the calculated profiles and their fit to the data. The user may choose to visualize the data later, any time after the quantifying step has been completed.

This manual walks you through the different steps involved in analyzing a gel using SAFA. Section II provides detailed installation instructions. Limited support can be provided; report bugs and questions to Alain Laederach (alain@helix.stanford.edu). It is possible to save your progress at any point during a SAFA session by using the *Save Dump* option in the file menu. The dump can be read in again by using *Load Dump* from the same menu. It is recommended that you save after every session so that you can always go back and look carefully at your analysis.

Details of the science behind SAFA, descriptions of the algorithms implemented have now been published. If you use SAFA to generate results that you publish, please cite the SAFA primary reference:

Das, Laederach, Pearlman, Herschlag, Altman, SAFA: Semi-Automated Footprinting Analysis software for high-throughput quantification of nucleic acid footprinting experiments. RNA 2005. *in press*.

II. Installation

SAFA currently requires MATLAB to be installed for it to run on a Macintosh. Obtaining MATLAB is relatively straightforward (<http://www.mathworks.com>). Many universities have special licensing agreements with MathWorks with licensing costs averaging between \$100-\$200 a year. A Windows standalone is now available for download from <http://safa.stanford.edu>. This means that for a windows machine you **do not** need to purchase Matlab. Simply download the Standalone and run the MCRInstaller. Afterwards double clicking on SAFA.exe will launch SAFA. You do not need to read the rest of the installation instructions if you have a PC. We are currently working on a Mac standalone, so check the SAFA webpage regularly for updates.

Installation of MATLAB is relatively easy and detailed instructions can be found on <http://www.mathworks.com/support/>.

SAFA is a series of MATLAB files. It is packaged as a simple zip archive (SAFAv09b5.zip) downloadable from the web site (<http://safa.stanford.edu>). On either Macs or PCs double clicking on the package will unpack the software. It is recommended that SAFA be installed in a separate directory under the MATLAB toolbox directory, such as:

`<MATLAB_Home>/toolbox/SAFAv09b5/`

At this point, SAFA is installed, but MATLAB needs to be made aware of SAFA's location. First start MATLAB. Then, you can permanently add the SAFA installation directory to MATLAB's script path by choosing the *Set Path* menu option under the *File* menu in the MATLAB application. In the *Set Path* window that pops up, click the *Add with Subfolders* button, and navigate to the SAFAv09b5 directory in `<MATLAB_Home>/toolbox/`, select it, and hit *OK*. This will make the SAFA scripts available to MATLAB for use.

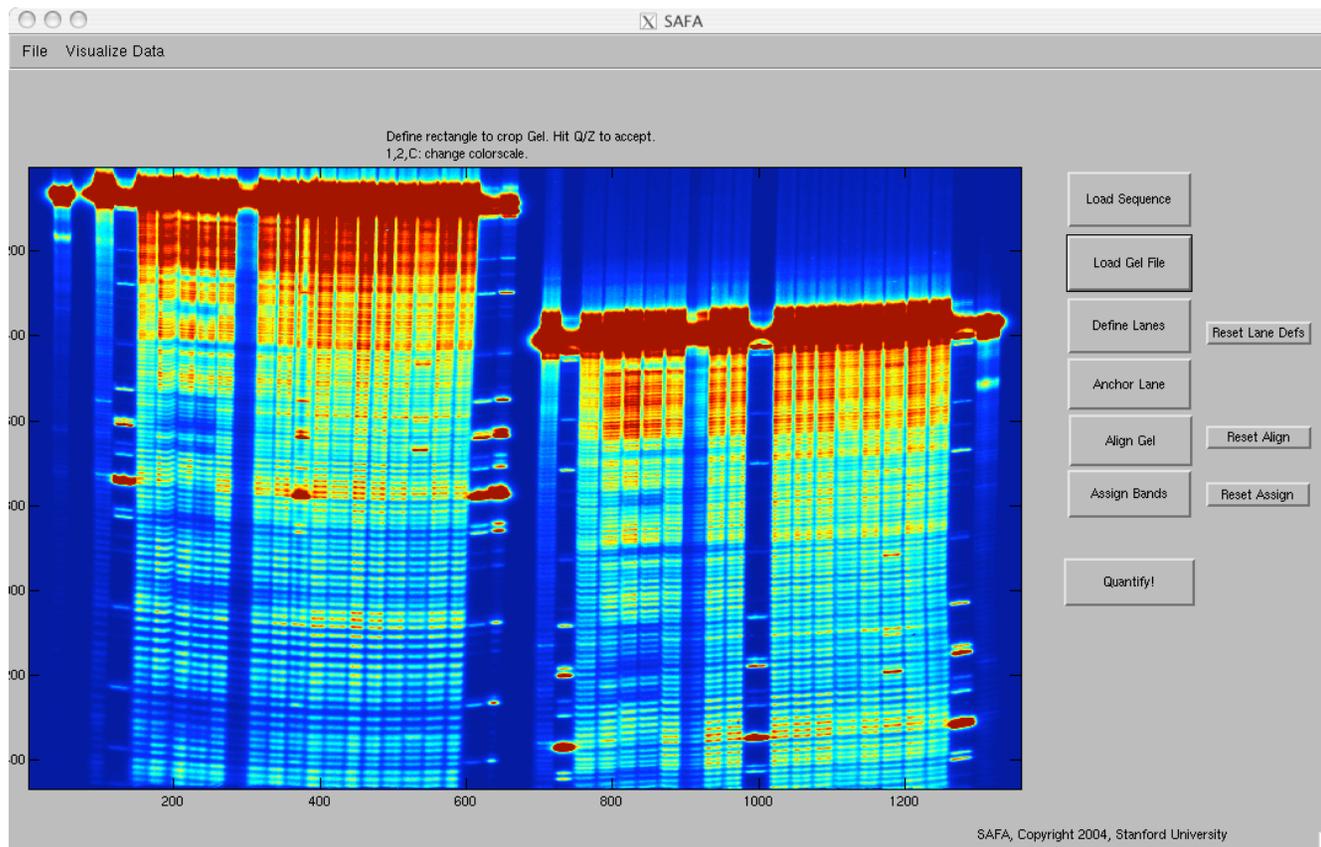
To run SAFA once you have told MATLAB where you unpacked the data, at the MATLAB command prompt type:

`>>SAFA`

III. Steps in SAFA

A. Loading .gel files

The first step in analyzing a gel is loading the Gel file into SAFA. Currently accepted formats are ImageQuant .gel files from Molecular Dynamics. Click the *Load Gel* button and choose the appropriate file. The image of your gel will appear in the main window as shown below:



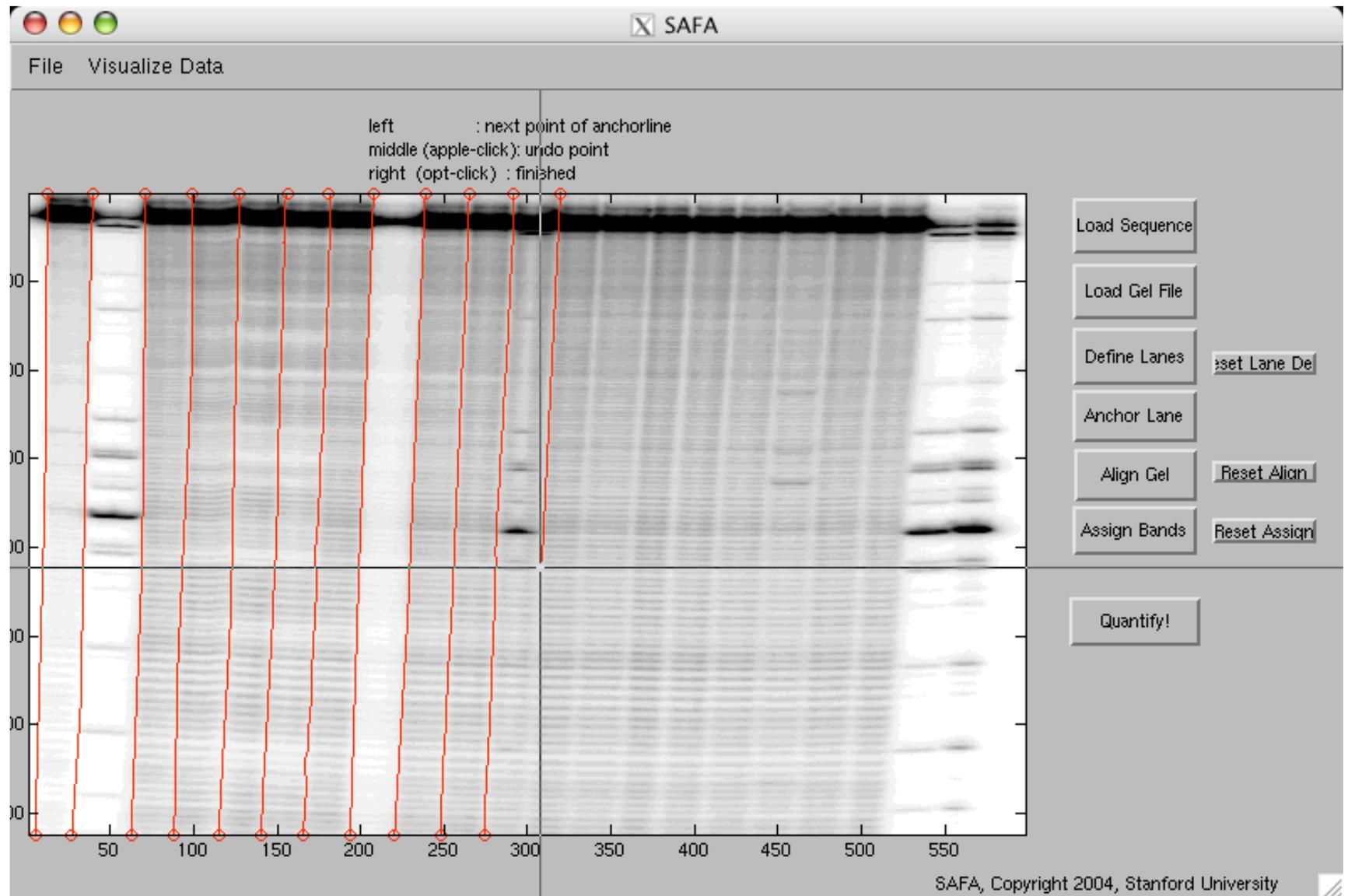
The cursor will become a crosshair, indicating that a function in SAFA is running. A new action (i.e., clicking on another button) should only be started once the cursor is no longer a crosshair, indicating completion of the current function. Future versions will not allow clicking a new function button while another function is running.

Load Gel allows you to crop your gel image. Place the cursor at the uppermost left corner of the gel you want to analyze, click once, and hold the button down. Then drag the mouse to the bottom right corner of the gel, and release the mouse button. You may repeat this as many times as you wish until you are satisfied you have chosen the correct region, and then hitting 'Q' or 'Z' will finalize your choice, and a cropped image of your gel will appear. If you do not wish to crop your gel, simply click at the top left and bottom right corners of the image. Once the gel is loaded properly the gel image is resized to fit inside the window and the cursor crosshairs will disappear. During this process, you may change the view of the gel from grayscale to false color and back again by hitting the 'C' key. Additionally, the contrast may be increased or decreased by hitting the '1' and '2' keys respectively.

It should be noted that SAFA will be quantitatively analyzing the gel image that is read in. To obtain a correct answer with SAFA it is critical that the gel image not be over-exposed. The response in band intensity should be linear for the quantification to be done correctly.

B. Defining Lanes

Lanes are defined by clicking on the *Define Lanes* button. You will be prompted to enter the number of bins per lane. Each lane will be discretized into this number of bins for use in the next steps. The default value of 10 is generally applicable for high quality gels. For lower quality gels with severe 'smiling' or 'frowning' or other band distortions, higher numbers of bins (e.g. 20 or 30) may yield better results at the cost of additional computation time. The program will then ask you to define the left-most lane. Do this by clicking down the left boundary of the left-most lane. The first left-click will define where the red line starts at the top of the gel. Once you have finished drawing a boundary, right clicking will terminate the line. Once you have defined both the left and right boundaries of the first lane, you may hit 'G' to ask the program to guess the next boundary. If the guess is satisfactory, hit 'G' again to define further lane boundaries automatically. Otherwise place the pointer over the last boundary line and hit 'E' to erase the last computer guess and then define the lane boundary manually as with the first lane. The program senses the right side of the gel and when no more boundaries can be added. Hit 'Z' to finish the assignment procedure. SAFA then bins the image. This is a complicated procedure and can take several seconds on slower machines. SAFA finally draws the centerlines of each lane and the lane boundaries turn green. At this point you may proceed to the next step.

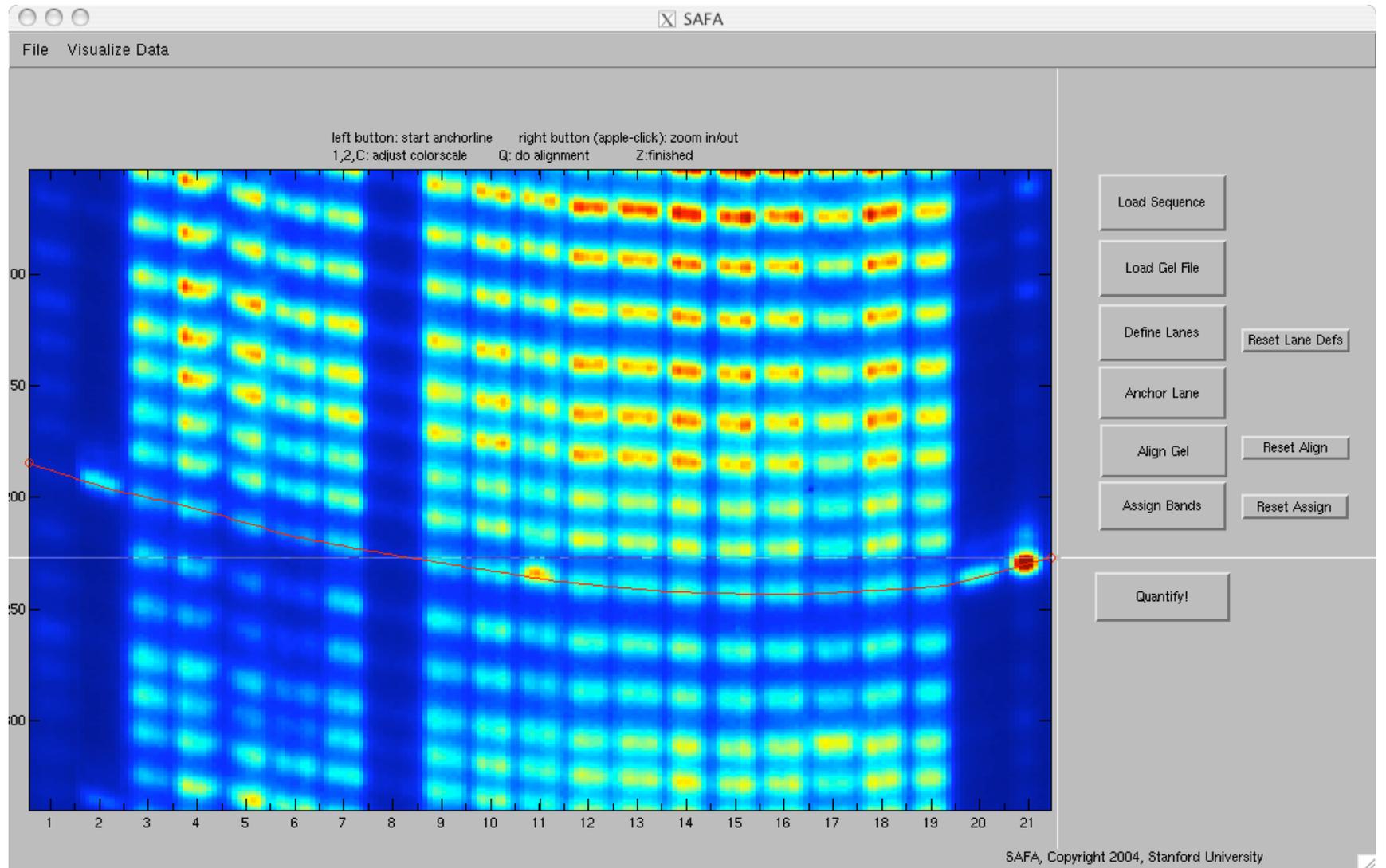


C. Defining an Anchor Lane

Prior to aligning the gel, a lane must be defined that will serve as an anchor lane. This is the lane containing the bands to which corresponding bands in other lanes will be aligned. By clicking on *Anchor Lane*, you can visually select what lane to use. The cursor will become crosshairs and you can click on the lane you wish to use. We recommend using a lane in the middle of the gel with a set of known bands, such as a ribonuclease T1 ladder for RNA gels. Once you have clicked on the lane your choice will be confirmed in the text above the gel image.

D. Aligning the Gel

Aligning the gel by clicking on *Align Gel* activates a routine that will allow you to define anchor bands. This will correct the gel for distortions by defining bands across the gel that correspond to the same product length. This is an essential step for the automated quantification of the gel, as a proper alignment is key to proper band assignment.



The user draws the red line to establish one or more anchor lines for the alignment routine. Interpolation of the profiles is used to stretch or compress the lanes accordingly and effectively straighten the gel. Practically, you should find a band that is easy to see across the entire gel. Hitting 'C' toggles the gel visualization between grayscale and false color. In either mode, the contrast of the gel can be increased or decreased by typing '1' or '2', respectively. Left clicking allows you to draw the anchor line, while middle clicking undoes the point, and right clicking finishes an anchor line, and zooms in and out of the gel when not defining a line. In general two anchor lines are required for gel alignment, one near the top of the gel and another near the bottom. For gels with more distortion or an unusual running pattern, more anchor lines may be needed. 'Q' causes the gel to be aligned, while Z does this and then exits the alignment routine. As many anchor lines as needed may be drawn and 'Q' can be hit as many times as needed as well. If you would like to reset the alignment (i.e. undo all of the changes and remove all the anchor lines), click on *Reset Align*. To erase one particular anchor line during this process, place the cursor near the line in question and hit 'E'.

E. Loading a Sequence File

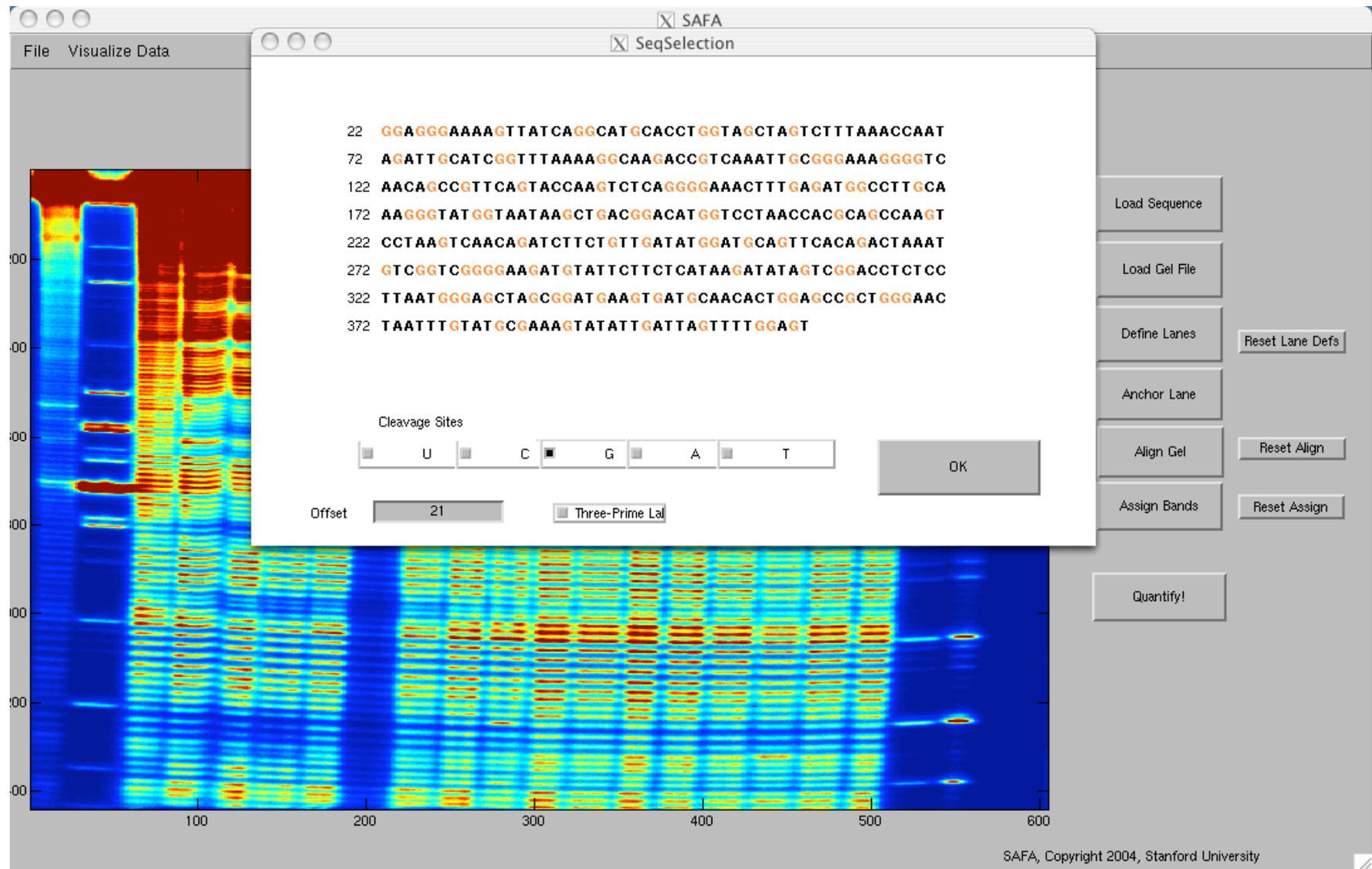
Prior to assigning the bands, the sequence of the molecule needs to be established. A simple routine has been implemented that can read in FASTA format sequence files. The format is very simple:

```
>TtLSU
GGAGGGAAAAGTTATCAGGCATGCACCTGGTAGCTAGTCTTTAAACCAATAGATTGCATCGGT
```

The first line is a comment line and the next lines are the sequence of the molecule.

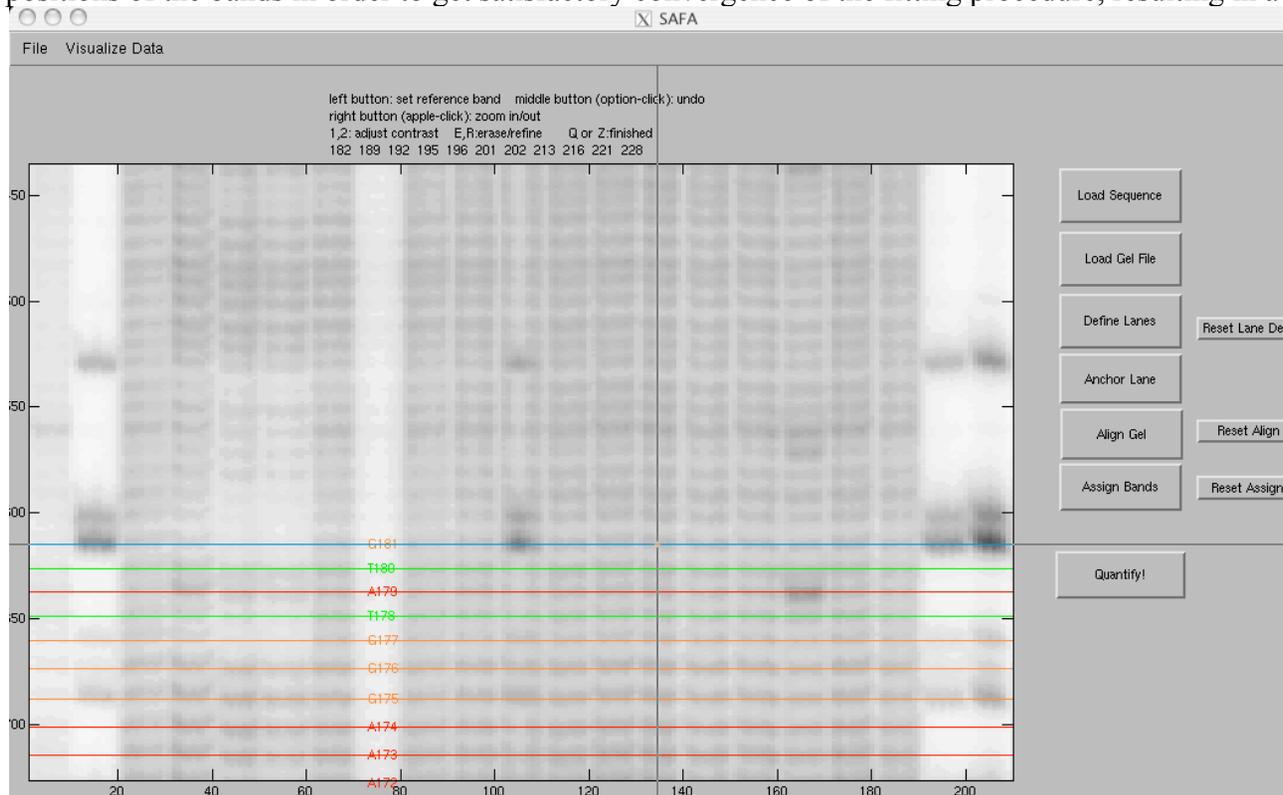
F. Choosing cleavage sites and offsets.

When you read in a .fas file by clicking the *Load Sequence* button, a Sequence Selection window opens. This window will allow you to choose one or more sets of nucleotide cleavage sites, e.g., G sites if you are using a ribonuclease T1 ladder as a reference in an RNA experiment. You should also choose an offset if your favorite numbering scheme does not call the first nucleotide "1." For example, with the standard "L-21" *Tetrahymena* ribozyme sequence above, the first residue is numbered 22, and the offset 21 is entered. Finally, you can select whether the molecule is 5'-labeled (band numbers decrease for bands of faster mobility, going "down" the gel) or 3'-labeled. When finished hit the OK button to finalize your choices.



G. Assign Bands

Once the sequence file is loaded, you will be able to click on *Assign Bands*. Left clicking places the next band on the Ladder as defined by the sequence file and the cleavage sites in the sequence browser window. You will only have to identify the nucleotide positions corresponding to your reference ladder, e.g. G's if you are using a ribonuclease T1 ladder as a reference in an RNA experiment. SAFA will guess and assign the positions of bands between these reference sites by linear interpolation. Each nucleotide is displayed with a particular color corresponding to its residue. Hitting '1' and '2' adjusts the contrast while it is possible to zoom in on the gel by right-clicking. Typing either E or R while placing the cursor erases the line closest to the cursor and allows you to reposition it. This can be useful if your gel did not run homogenously. It is important to make accurate choices for the positions of the bands in order to get satisfactory convergence of the fitting procedure, resulting in accurate data.

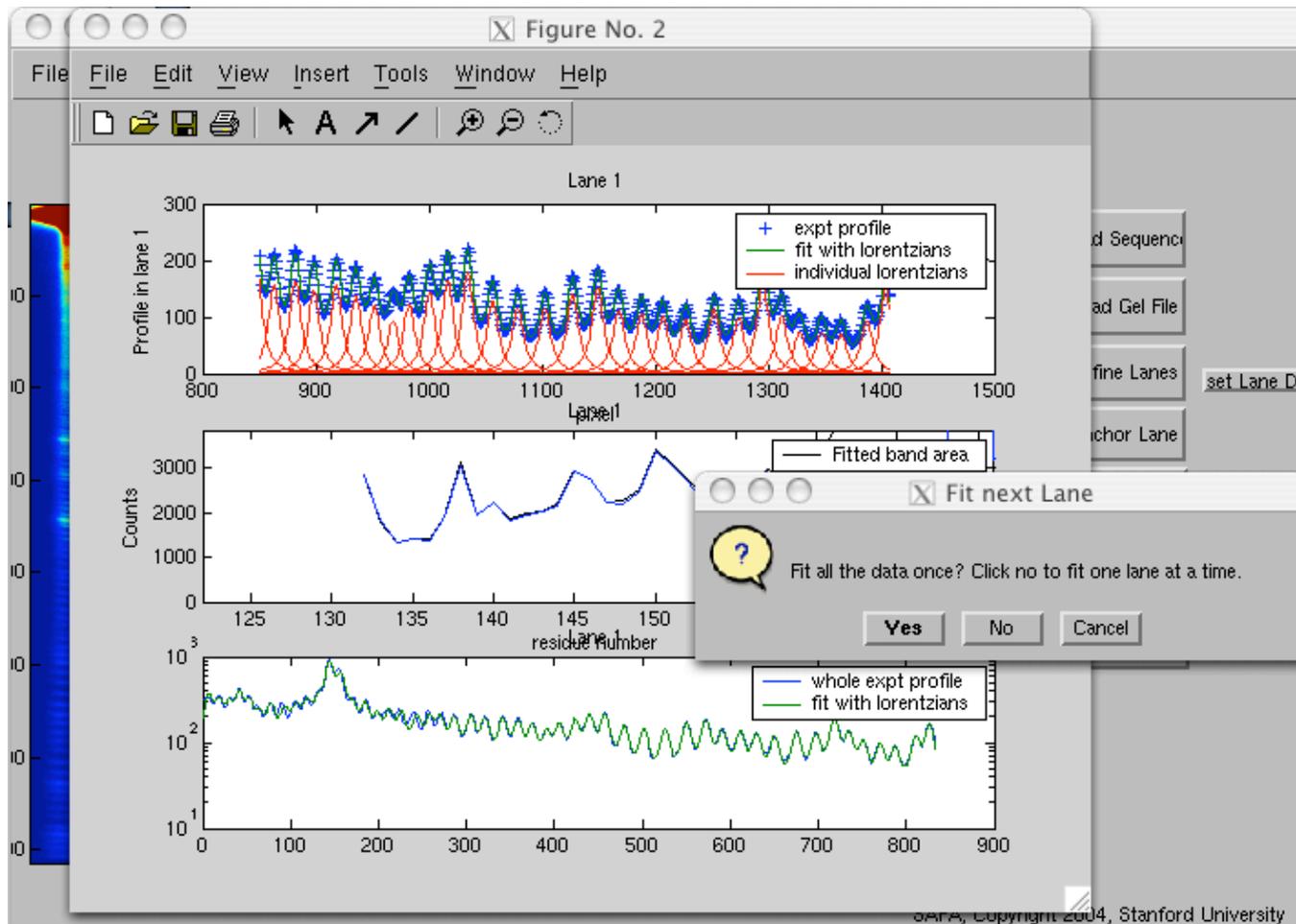


H. Quantifying your gel

Quantifying the gel is accomplished by clicking on *Quantify*. This will begin a least-squares Levenberg-Marquardt optimization procedure that fits a sum of Lorentzians to the data. You will be prompted for the number of bands to fit, and you should estimate how many bands (starting from the bottom of the gel) are clearly resolved from each other – typical RNA gels range from 20 to 50. The positions of these first, “best” bands will be optimized, starting with your position guesses. For any bands higher up the gel, the computer will assume your guesses for the band positions are correct.

SAFA enforces non-negative peak amplitudes and areas in all but the highest bands on the gel. While this method slows the analysis slightly, it provides better results.

Once you have selected a model, SAFA will fit the data and a new window will appear with the fitting plotted:



The top frame shows the Lorentzians in red, the fitted profile in green and the experimental data in blue crosses. The second frame shows the fitted peak areas. In black are the fitted areas for the first, “best” bands whose positions were refined by the quantification procedure, and in blue are the bitted areas for all bands. In general these two lines should agree. The third frame shows the fitted profile in blue and the data in green.

I. Saving Data

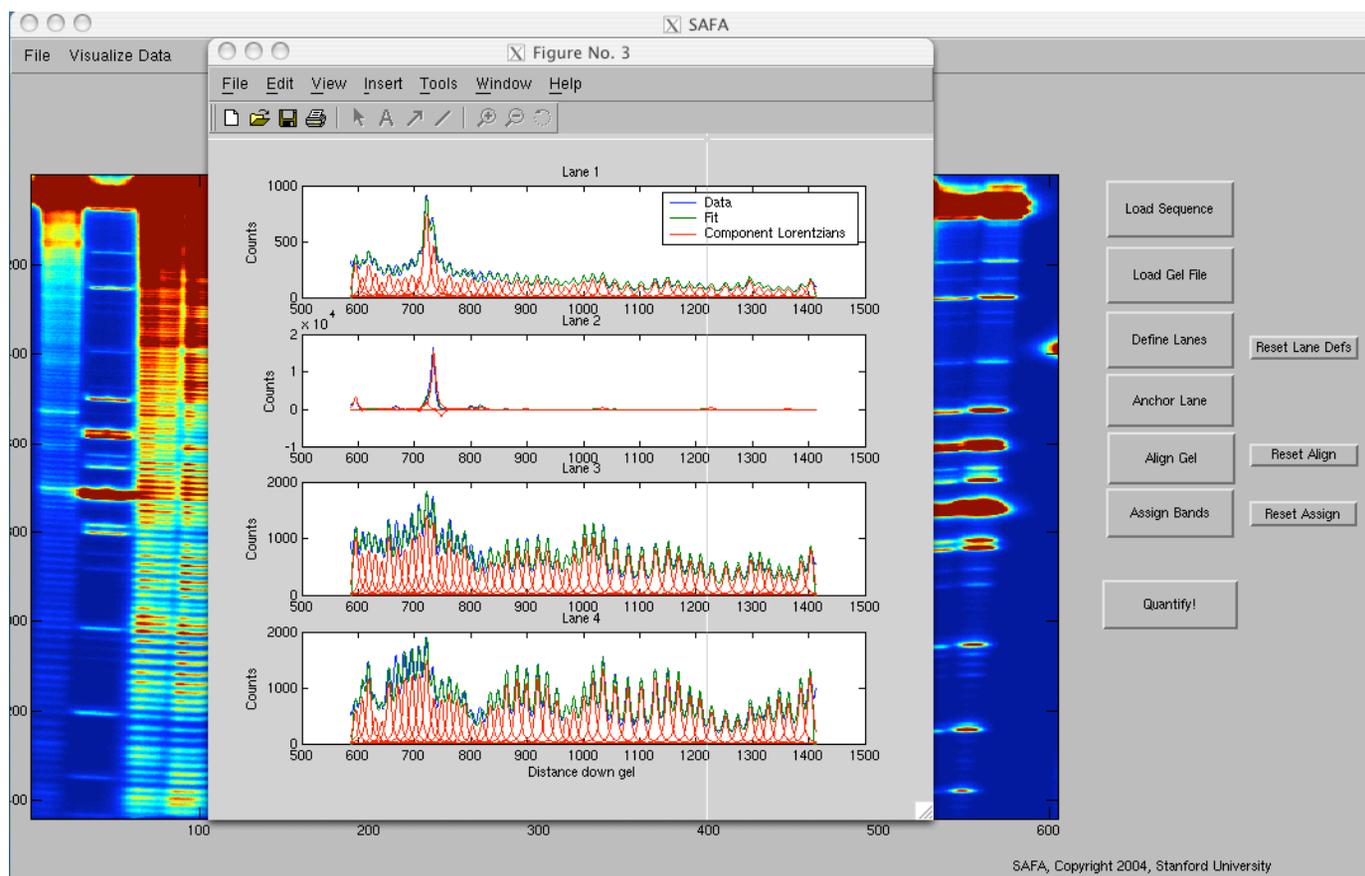
When the fitting procedure is finished, all lanes are fit (provided you did not cancel the fitting before it had completed all the lanes), and you will be prompted as to whether you want save the data in .txt format. The output format is very simple:

134	1.21E+03	-1.10E+04	6.02E+03	1.15E+04	1.25E+04	1.05E+04	1.10E+04	2.32E+03
135	1.29E+03	-1.08E+04	6.78E+03	1.30E+04	1.46E+04	1.24E+04	1.01E+04	2.13E+03
136	1.27E+03	-1.15E+04	7.62E+03	1.05E+04	1.04E+04	9.18E+03	9.51E+03	2.08E+03
137	1.75E+03	-1.08E+04	9.57E+03	7.82E+03	5.15E+03	4.44E+03	8.57E+03	2.00E+03
138	2.79E+03	-1.16E+04	1.04E+04	1.02E+04	8.37E+03	6.72E+03	9.19E+03	2.00E+03
139	1.76E+03	-1.26E+04	7.99E+03	6.60E+03	4.51E+03	3.71E+03	8.26E+03	1.81E+03
140	2.03E+03	-1.41E+04	8.67E+03	7.79E+03	5.30E+03	4.01E+03	9.21E+03	2.06E+03
141	1.54E+03	-3.20E+03	7.66E+03	8.67E+03	7.04E+03	5.86E+03	8.34E+03	1.83E+03
142	1.73E+03	-1.05E+04	9.74E+03	1.27E+04	1.20E+04	1.01E+04	1.24E+04	2.55E+03
143	1.83E+03	-1.27E+04	1.00E+04	1.34E+04	1.31E+04	1.08E+04	1.23E+04	2.40E+03
144	1.93E+03	-1.18E+04	1.19E+04	1.73E+04	1.65E+04	1.40E+04	1.59E+04	3.11E+03
145	2.61E+03	-1.20E+04	1.37E+04	1.82E+04	1.76E+04	1.46E+04	1.66E+04	3.38E+03
146	2.46E+03	-1.31E+04	1.40E+04	1.86E+04	1.80E+04	1.51E+04	1.71E+04	3.48E+03
147	1.98E+03	-1.28E+04	1.25E+04	1.60E+04	1.53E+04	1.31E+04	1.40E+04	2.88E+03
148	1.93E+03	-1.26E+04	1.30E+04	1.58E+04	1.49E+04	1.26E+04	1.37E+04	2.73E+03
149	2.16E+03	-1.18E+04	1.42E+04	1.63E+04	1.37E+04	1.16E+04	1.41E+04	2.88E+03
150	2.80E+03	-3.85E+03	1.74E+04	2.05E+04	1.77E+04	1.52E+04	1.84E+04	3.63E+03
151	2.58E+03	-9.74E+03	1.64E+04	1.79E+04	1.45E+04	1.25E+04	1.51E+04	3.12E+03
152	2.29E+03	-1.10E+04	1.70E+04	1.34E+04	7.93E+03	6.19E+03	8.87E+03	1.96E+03
153	1.91E+03	-9.85E+03	1.16E+04	9.62E+03	5.24E+03	4.18E+03	5.89E+03	1.32E+03
154	1.47E+03	-1.04E+04	7.66E+03	7.28E+03	4.45E+03	3.85E+03	4.06E+03	9.89E+02
155	1.87E+03	-1.08E+04	1.02E+04	1.19E+04	1.04E+04	9.15E+03	8.92E+03	1.84E+03
156	2.16E+03	-1.05E+04	1.29E+04	1.61E+04	1.43E+04	1.21E+04	1.20E+04	2.40E+03

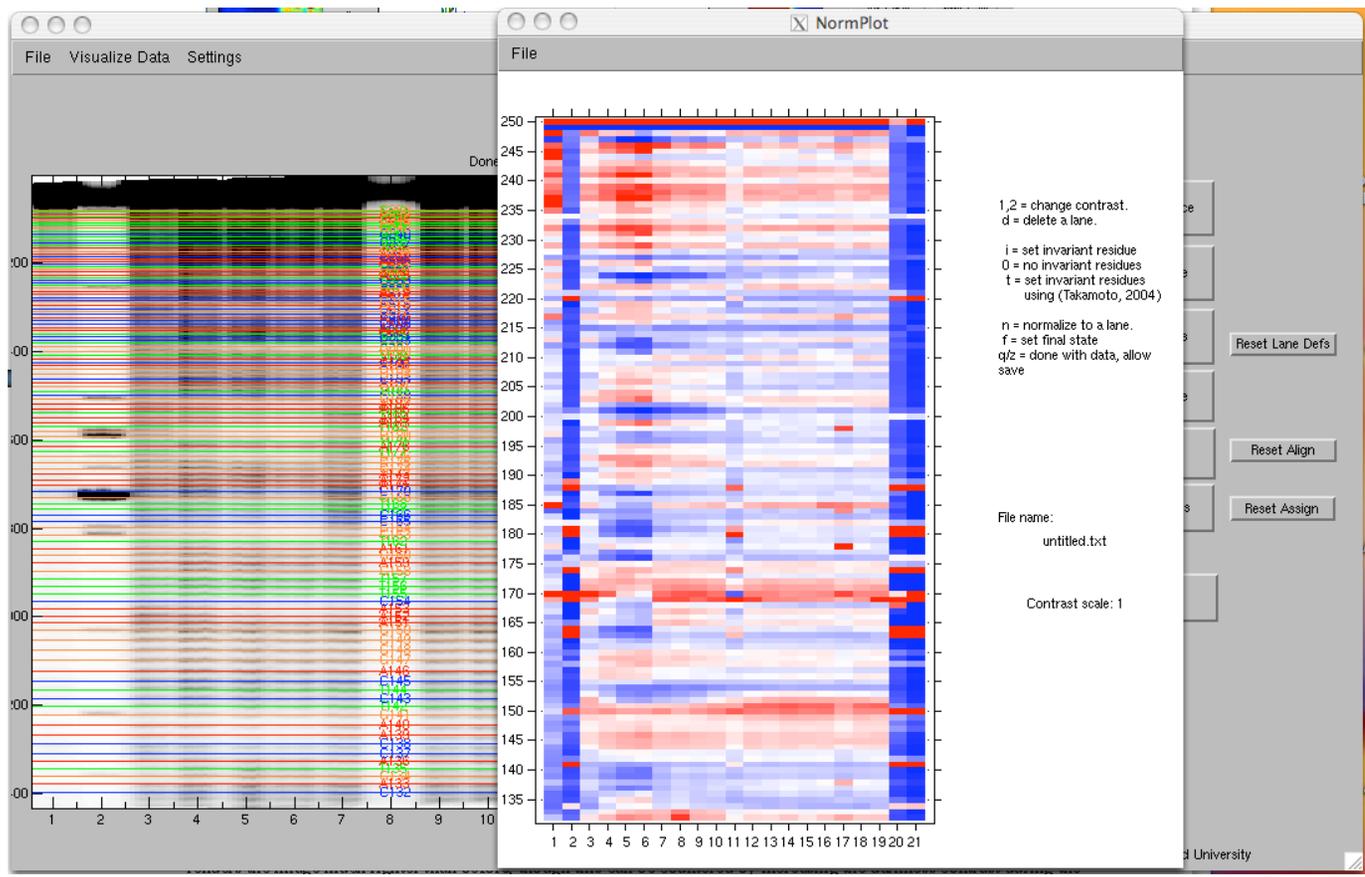
The first column is the position number while the next columns are the peak areas. It is possible to obtain negative peak areas at the extreme bands of the gel. Usually these negative values result when there is no band at that particular position. The data in columns containing large negative numbers, as well as the data in neighboring positions, should be considered with caution.

J. Visualizing and Normalizing Data

After deciding whether or not to save the data, you will be prompted for the number of lanes to plot per page. The fit of the sum of Lorentzians to the experimental data as well as the individual Lorentzian curves will be displayed for each lane. This is a useful format for printing the data for your notebook; you can also save these fits to PDF files using the menu options. After each page, click on the plot window to show the next set of lanes. Once *Quantify* has been called, you may view the data fits by choosing “Plot Fits” under the “Visualize Data” menu item.

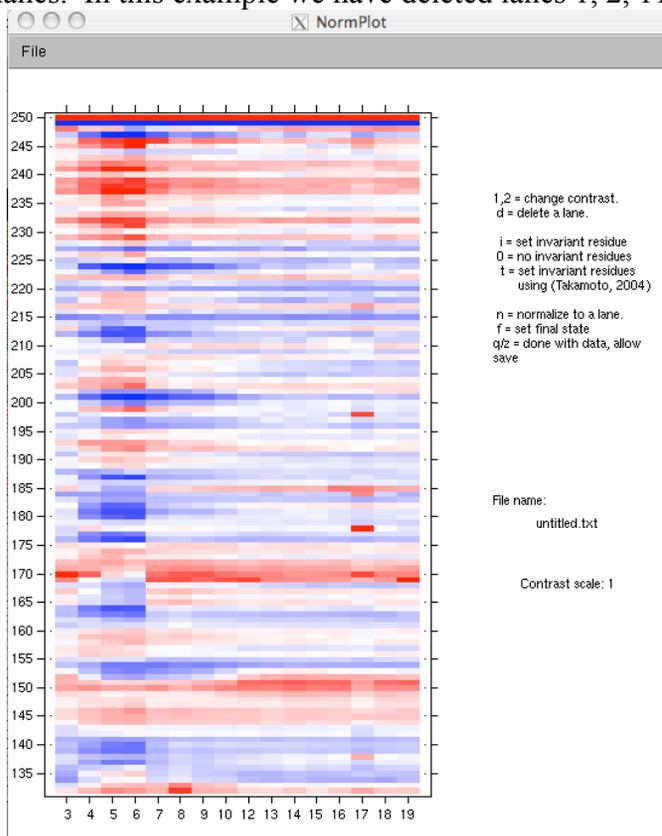


We have implemented a rather simple data normalization utility. This utility is particularly useful for experiments in which the relative intensity of bands must be compared across lanes (e.g. quantitative titrations or time-resolved kinetic experiments.) The basic principle behind normalization is to identify invariant portions (lanes or bands) of the gel, and to use these invariant regions to then normalize the rest of the data [see Takamoto, Chance, and Brenowitz (2004), *Nucleic Acids Res.* 32(15):E119]. To start the normalization utility, click on the ‘Visualize Data’ menu item and select ‘Normalization/Colorplot.’ The following window comes up:



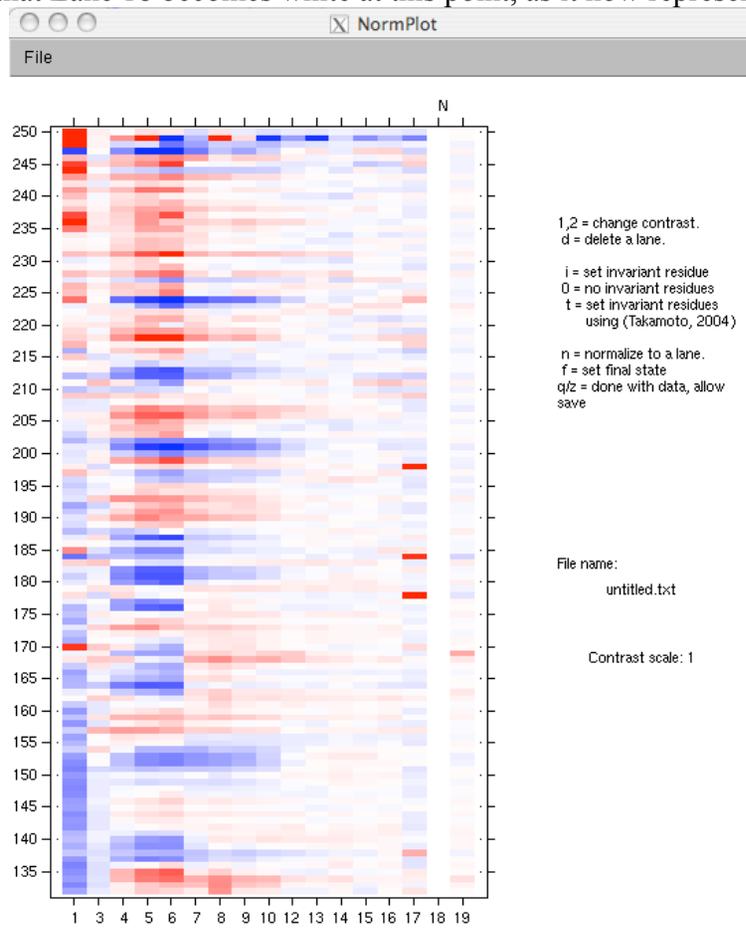
The NormPlot utility allows the user to visualize their entire data set in a manner analogous to microarray data. The data representation follows a standard color scale, with red indicating more counts and blue fewer counts. The data is initially normalized to the average number of counts for the entire gel. When the cursor is a crosshair, the utility is active, and either “q” or “z” needs to be hit in order to access the File menu.

The first step in normalizing the data is to remove non-data lanes (e.g. ladder lanes). This is accomplished by mousing over the lanes and hitting the “d” key to delete lanes. In this example we have deleted lanes 1, 2, 11, 20 and 21:



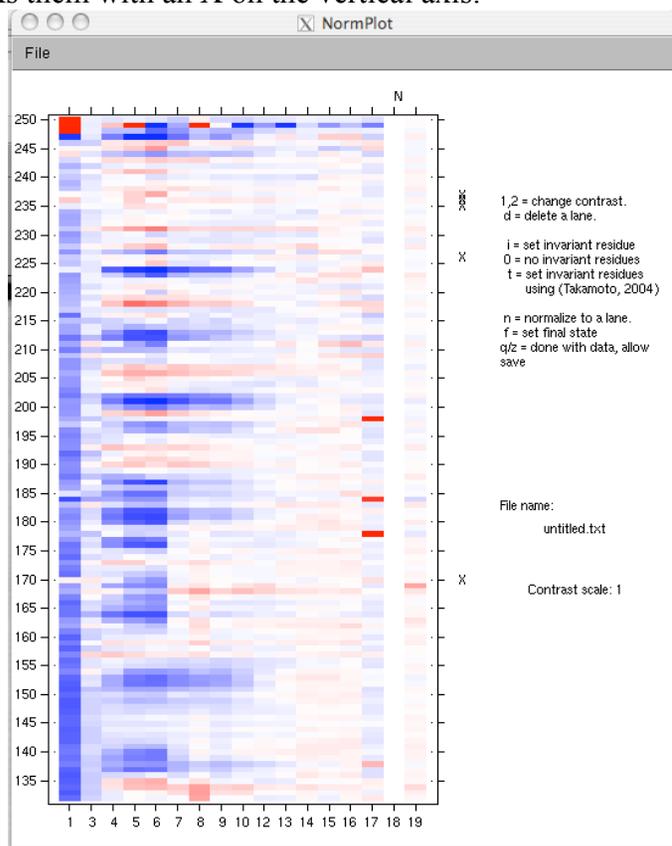
Note that lane numbers are preserved on the horizontal axis. The utility renormalizes the data at every lane deletion.

The data can then be normalized to one or more lanes. For example, in this gel Lane 18 is an “unfolded” lane, meaning that one can assume that the cleavage is relatively uniform in this lane. By mousing over Lane 18 and hitting “n” for normalize to lane, the data is instantly normalized to that lane (note that Lane 18 becomes white at this point, as it now represents a standard state with value 1).



The utility also places an “N” at the top of Lane 18 to indicate the data is normalized to that lane. At this point the data is normalized to a lane and can be exported by hitting “q” or “z” and then clicking on “File” and “Save Data.” The user can also save a copy of the image in TIFF, PDF or Adobe Illustrator format.

In some cases, specifically for titrations, invariant residues must also be identified across bands. NormPlot has an automated utility to perform this function based on the methods described by Takamoto et al. *NAR* 2004. The utility automatically identifies invariant residues in the data and marks them with an X on the vertical axis:



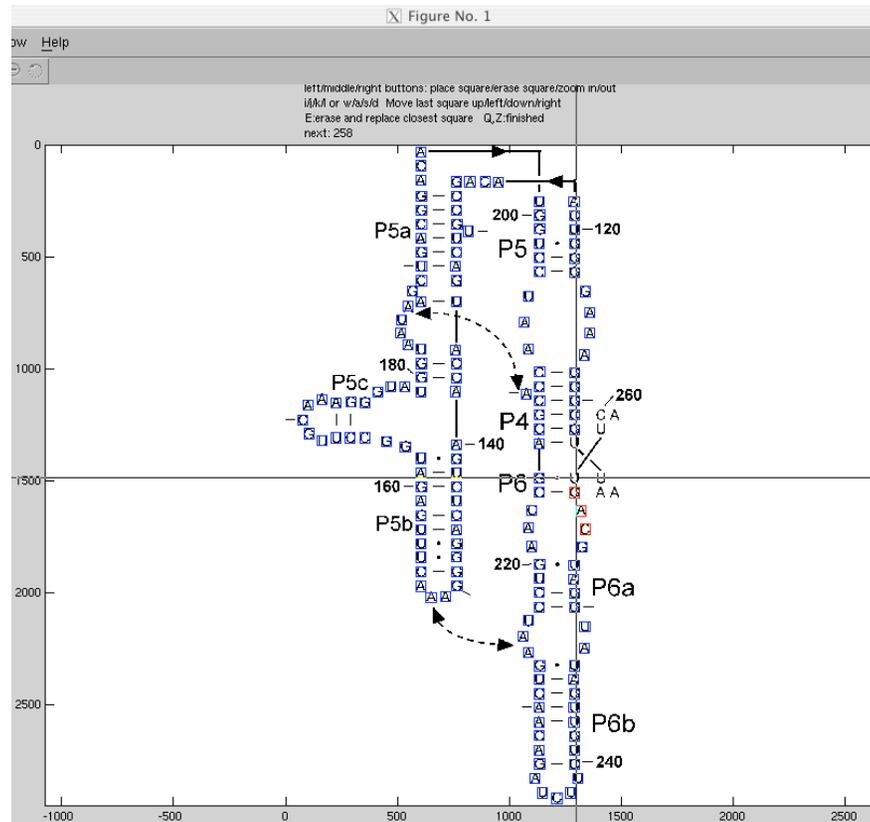
At this point the utility has identified five invariant residues and normalized the data to the relative intensities of the bands at these residues. The user can always add or subtract invariant residues using “i” or “0.” Finally, the user can also set a lane as the final state by mousing over the lane with the crosshairs and hitting “f.” This sets the values of the band intensity in the final lane to a value of one, and scales the rest of the data linearly.

The NormPlot utility will always try to load the latest data stored in SAFA after having hit the quantify button. The user can also read in previously saved data in the .txt format using the File--> Read Data menu option. To reset the application simply select File--> Latest SAFA data. Finally, the data is exported in the same format (tab-delimited text) as the raw SAFA data, however the .norm.txt extension is appended to the filename. It is recommended to abide by this file nomenclature to avoid confusion as to the source of the data.

K. Mapping Data onto Secondary Structure

A prominent use of data from footprinting experiments is to determine the relative solution-accessible surface area of different regions of the molecule being studied. We have developed a way to quickly and clearly map the quantified, normalized data onto the molecule’s secondary structure. This tool can be accessed from the *Secondary Structure Plot* menu item under *Visualize Data* in SAFA. If a gel has been quantified in the current SAFA session, that data will be used initially for secondary structure mapping and a colorplot of the data will prompt the user to choose a lane of data to use for the mapping (though new data can easily be loaded if desired). The basic steps for mapping accessibility onto secondary structure are three-fold:

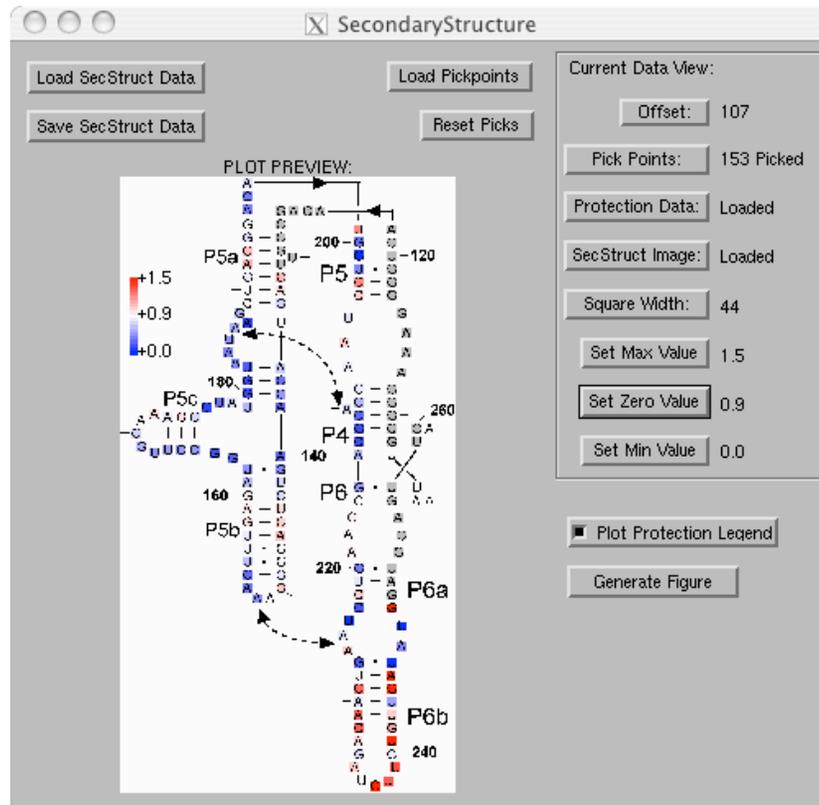
1. Create and then load an image (in JPEG, GIF, or TIFF format) of the molecule’s secondary structure using the *SecStruct Image* button. It is recommended that the image be high-resolution (150 dpi or higher). If desired the image should have some space reserved on one side for a color legend depicting the relative solution accessibility levels.
2. The *Pick Points* button will become active once an image is loaded, and will open a figure in which the user can place squares over each residue in the region of the molecule that has had its solution-accessibility levels quantified. The size of the squares can be set with the *Square Width* button. An offset indicating the position of the first residue that has been assigned a square can be set as well using the *Offset* button. The Pick Points figure is shown below:



3. Load the accessibility data directly from SAFA or from the .txt or .norm.txt files saved from previous SAFA sessions. The format of these files is simple and can be generated from other sources as well: The first column indicates the residue number for that row, and each successive column contains count data for that residue from the footprinting experiment for a lane of the gel.

The application maps accessibility to color in the following manner: A “zero” level of protection can be chosen (initially set to be the average of the minimum and maximum count levels in the data) which will be plotted as white. Each residue that has higher accessibility than this will be colored red with intensity on a linear scale from this zero-level to a maximum value (also user-chosen, defaulting to the maximum value in the data) plotted as solid red. Similarly, residues with accessibility values lower than the zero-value will be colored blue with intensity on a linear scale from the zero-level to a minimum value (also user-chosen, defaulting to the minimum value in the data) plotted as solid blue. The secondary structure image and the points picked for the residue locations on that image can be loaded independently from previously saved data files using the *SecStruct Image* and *Load Pickpoints* buttons, and all fields can be saved and later loaded together from those data files using the *Load/Save SecStruct Data* buttons. The current state is previewed in a figure within the window. The max-, min-, and zero-values may be reassigned by hitting the appropriate *Set Max/Zero/Min Value* buttons.

When the structure image, residue locations, and accessibility data have all been loaded or assigned, the *Generate Figure* button will become active. This generates a full-sized figure showing the accessibility levels plotted on the secondary structure of the molecule. The *Plot Protection Legend* checkbox controls whether the color legend is displayed in the preview plot as well as in the full-sized figure. If a legend is being displayed, its location can be chosen by the user before saving the figure or exporting to the file format of choice. A figure of the main SecondaryStructure window after picking points and loading protection data is shown below:



L. Additional Functionality

- Resetting SAFA: If you wish to clear SAFA of any loaded gels or other data choose *Reset Application* under the File menu.
- Exporting your Gel: After aligning your gel, you may choose *Export Gel* under the File menu. This will save a .tif file for importing into ImageQuant if desired.
- Displaying the square-root of the image: The PC version of MATLAB appears to have a bug in its image display functionality where dark regions of the gel can flip to white, creating a “halo” effect in the gel image. To counter this, the

user has the option of selecting “Render SQRT Image” under the *Settings* drop-down menu. This will fix the problem, but renders the image much lighter than before. This can be countered by increasing the darkness contrast during the loading of a gel, the defining of lanes, or the alignment of the gel. While using the square root of the image may lose some resolution, it can be toggled on and off when needed. This option defaults to “on” under Windows, and “off” for other platforms.

IV. Future Plans

We feel that the ability to immediately see not just the raw data, but extended visualization of the results of quantifying the gel will be an asset to researchers using SAFA. Toward this end we plan to integrate additional modeling and visualization functionality into SAFA, which will allow for mapping the relative exposure of each position on the molecule to solution onto a 3D model of the molecule.