



From SEDFIT to SEDPHAT

SEDPHAT for AUC SV Data Analysis for Equilibrium Systems

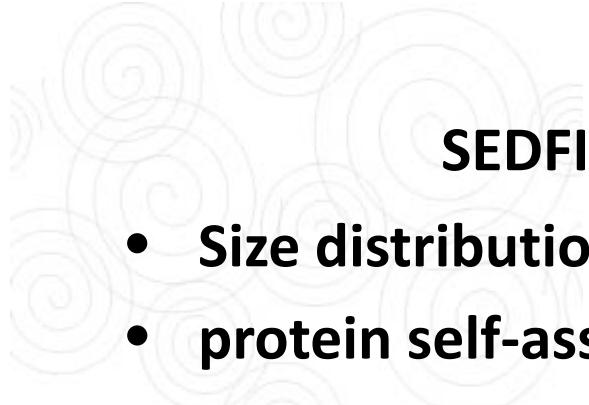
<http://www.analyticalultracentrifugation.com>

Ph.D. Chen Hui-Ju

IBC, AS

chenhj@gate.sinica.edu.tw

02-27855696 ext 5010



SEDFIT

- **Size distribution analysis**
- **protein self-association**
- **size-and-shape distributions**
- **macromolecules**
- **sedimentation velocity ultracentrifugation**
- **Lamm equation modeling and**

SEDPHAT

- multi-site binary and ternary protein interactions
- heterogeneous protein-protein interactions
- multi-protein complexes by multi-signal
- **specific models**

Papers with applications of SEDFIT

- Uncontrolled Zinc- and Copper-Induced Oligomerisation of the Human Complement Regulator Factor H and Its Possible Implications for Function and Disease. *Journal of Molecular Biology*, Volume 384, Issue 5, 31 December 2008, Pages 1341-1352. Ruodan Nan, Jayesh Gor, Imre Lengyel, Stephen J. Perkins
- Specificity and Reactivity in Menaquinone Biosynthesis: The Structure of *Escherichia coli* MenD (2-Succinyl-5-Enolpyruvyl-6-Hydroxy-3-Cyclohexadiene-1-Carboxylate Synthase). *Journal of Molecular Biology*, Volume 384, Issue 5, 31 December 2008, Pages 1353-1368. Alice Dawson, Paul K. Fyfe, William N. Hunter
- Evidence for the Oligomeric State of 'Elastic' Titin in Muscle Sarcomeres. *Journal of Molecular Biology*, Volume 384, Issue 2, 12 December 2008, Pages 299-312. Ahmed Houmeida, Andy Baron, Jeff Keen, G Nasir Khan, Peter J. Knight, Walter F. Stafford III, Kavitha Thirumurugan, Beatrix Thompson, Larissa Tskhovrebova, John Trinick
- The High-Resolution NMR Structure of the Early Folding Intermediate of the *Thermus thermophilus* Ribonuclease H. *Journal of Molecular Biology*, Volume 384, Issue 2, 12 December 2008, Pages 531-539 Zheng Zhou, Hanqiao Feng, Rodolfo Ghirlando, Yawen Bai
- NUTS and BOLTS: Applications of fluorescence-detected sedimentation. *Analytical Biochemistry*, In Press, Uncorrected Proof, Available online 6 December 2008. Rachel R. Kroe, Thomas M. Laue
- *Helicobacter pylori* neutrophil-activating protein promotes myeloperoxidase release from human neutrophils. *Biochemical and Biophysical Research Communications*, Volume 377, Issue 1, 5 December 2008, Pages 52-56. Chung-An Wang, Yen-Chun Liu, Shin-Yi Du, Chio-Wen Lin, Hua-Wen Fu
- The family 52 β-xylosidase from *Geobacillus stearothermophilus* is a dimer: Structural and biophysical characterization of a glycoside hydrolase. *Biochimica et Biophysica Acta (BBA) - Proteins & Proteomics*, Volume 1784, Issue 12, December 2008, Pages 1924-1934

Papers with applications of SEDPHAT

- **Solution Structure of the Complex Formed between Human Complement C3d and Full-length Complement Receptor Type 2.** *Journal of Molecular Biology, Volume 384, Issue 1, 5 December 2008, Pages 137-150.* Keying Li, Azubuike I. Okemefuna, Jayesh Gor, Jonathan P. Hannan, Rengasamy Asokan, V. Michael Holers, Stephen J. Perkins
- **Role of tryptophan-208 residue in cytochrome c oxidation by ascorbate peroxidase from Leishmania major-kinetic studies on Trp208Phe mutant and wild type enzyme.** *Biochimica et Biophysica Acta (BBA) - Proteins & Proteomics, Volume 1784, Issue 5, May 2008, Pages 863-871.* Rajesh K. Yadav, Subhankar Dolai, Swati Pal, Subrata Adak
- The **dimeric assembly** of *Photobacterium leiognathi* and *Salmonella typhimurium* SodC1 Cu,Zn superoxide dismutases is affected differently by **active site demetallation** and **pH**: An analytical ultracentrifuge study. *Archives of Biochemistry and Biophysics, Volume 471, Issue 1, 1 March 2008, Pages 77-84.* B. Catacchio, M. D'Orazio, A. Battistoni, E. Chiancone
- **tRNA-dependent asparagine formation** in prokaryotes: Characterization, isolation and structural and functional analysis of a ribonucleoprotein particle generating Asn-tRNA^{Asn}. *Methods, Volume 44, Issue 2, February 2008, Pages 146-163.* Marc Bailly, Mickaël Blaise, Hervé Roy, Marzanna Deniziak, Bernard Lorber, Catherine Birck, Hubert D. Becker, Daniel Kern
- Characterization and Further Stabilization of Designed Ankyrin Repeat Proteins by Combining **Molecular Dynamics Simulations and Experiments.** *Journal of Molecular Biology, Volume 375, Issue 3, 18 January 2008, Pages 837-854.* Gianluca Interlandi, Svava K. Wetzel, Giovanni Settanni, Andreas Plückthun, Amedeo Caflisch
- On the **quaternary structure of a C-type lectin from Bothrops jararacussu venom – BJ-32 (BjcuL).** *Toxicon, Volume 52, Issue 8, 15 December 2008, Pages 944-953.* F.P. Silva Jr., G.M.C. Alexandre, C.H.I. Ramos, S.G. De-Simone

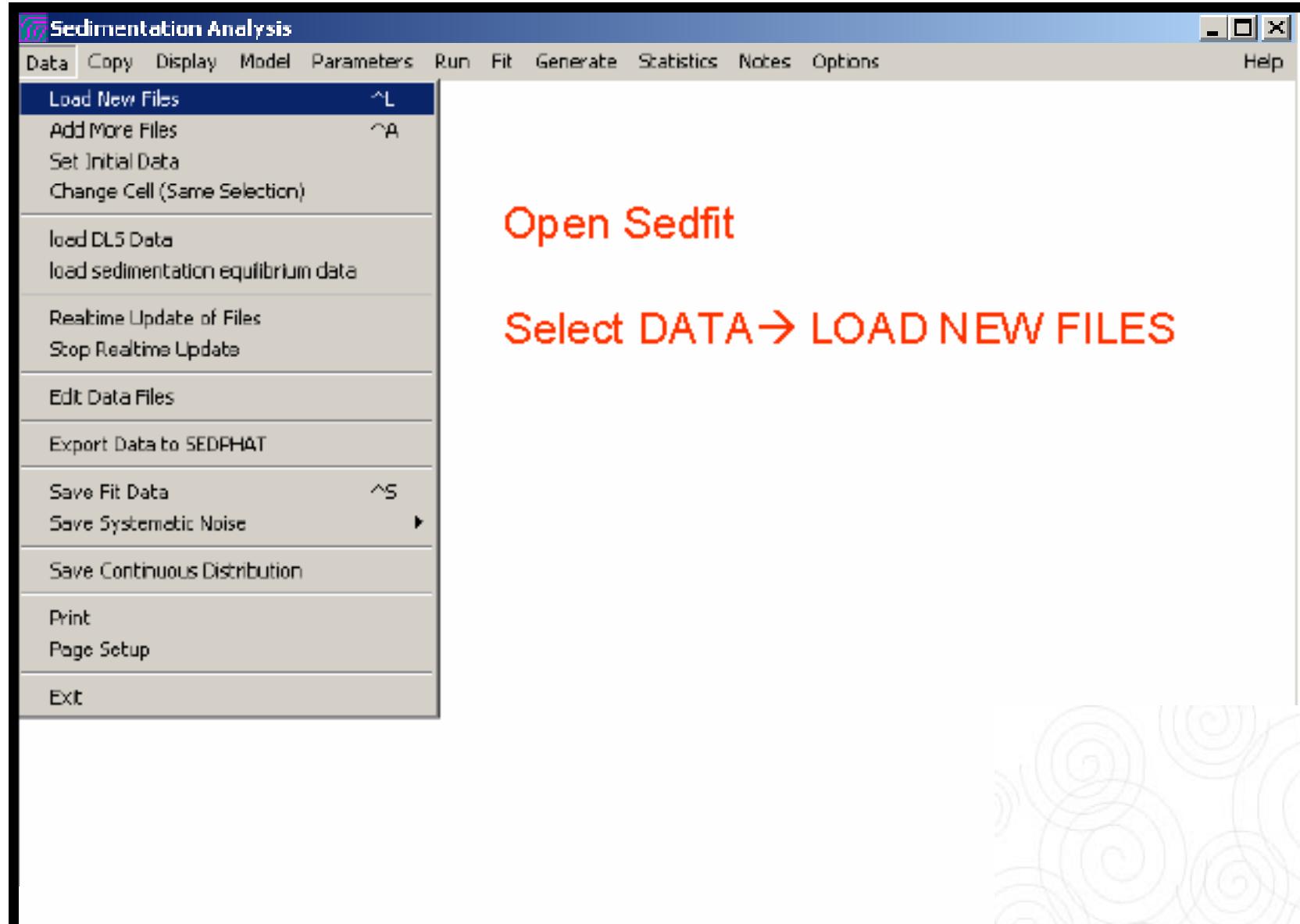
Step by step

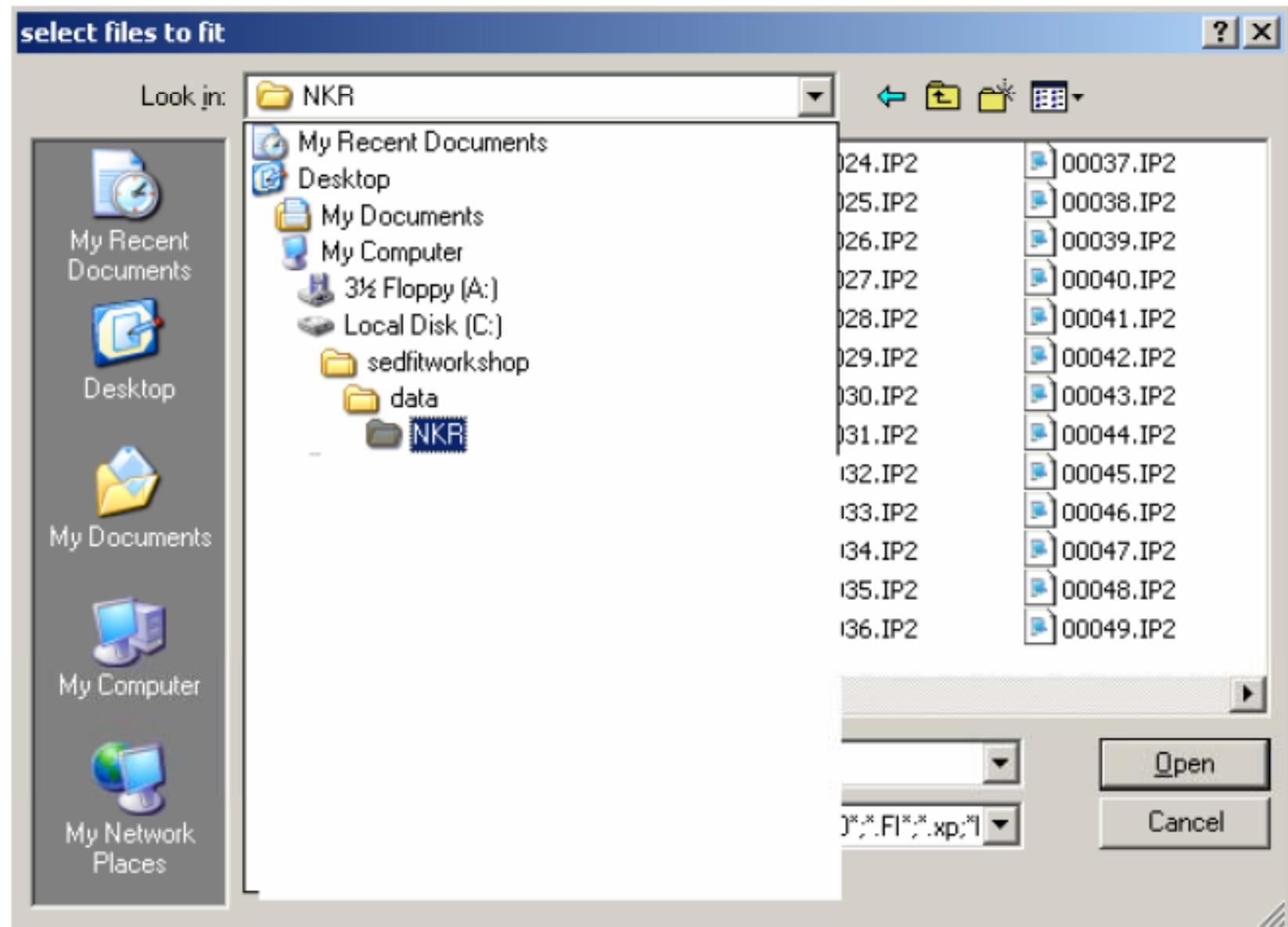
Content

- SedFit C(S) analysis
- SedFit C(M) analysis
- SedFit C(f/fo) analysis
- Transform to SEDPHAT
- SEDPHAT models :
 - (1) Single species & heterogenesis
 - (2) Monomer to dimer

Parameter setup

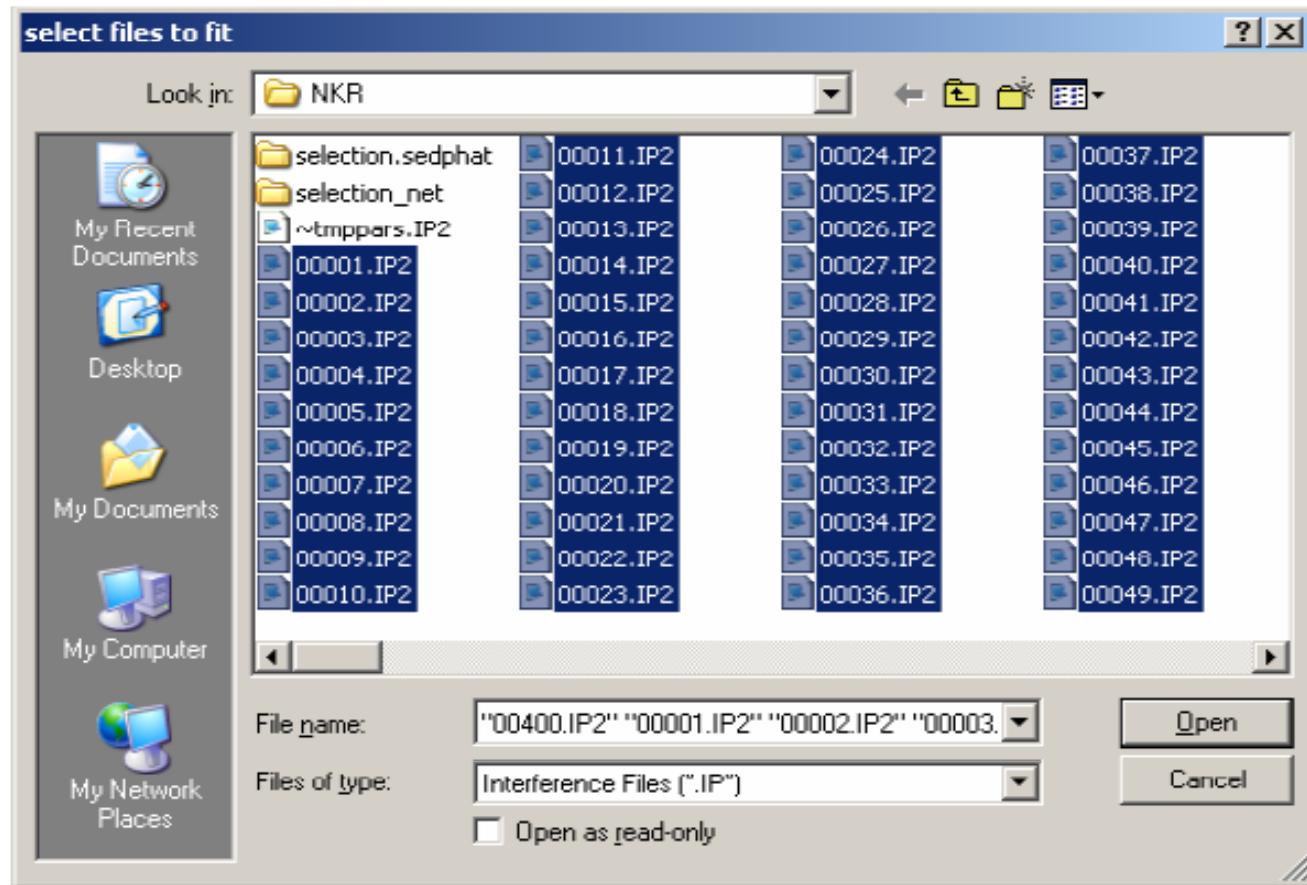
- Buffer density and viscosity
- Protein sample v-bar
- Rmsd
- Z value
- Boundary observation
- Statistics analysis



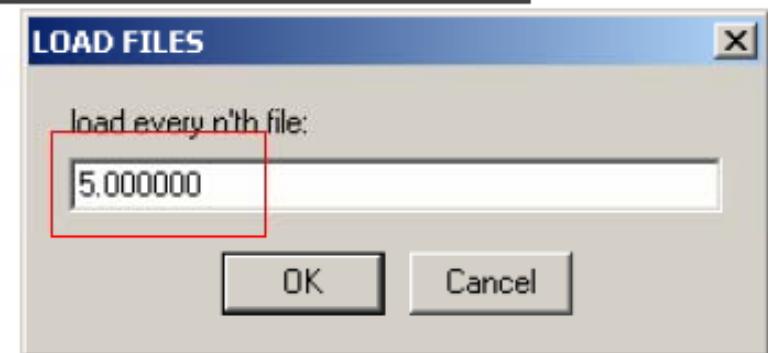


Choose scans within the NKR folder

Select all the interference scans for Cell 2 (highlight scans 1-400)



and load every 5th scan





If the data has been loaded before, Sedfit recognizes this and asks if the user wants to use the same parameters as the previous analysis (Sedfit had automatically saved this information in a temporary parameters file "tmppars.ip2" located with the raw data).

Select "No"

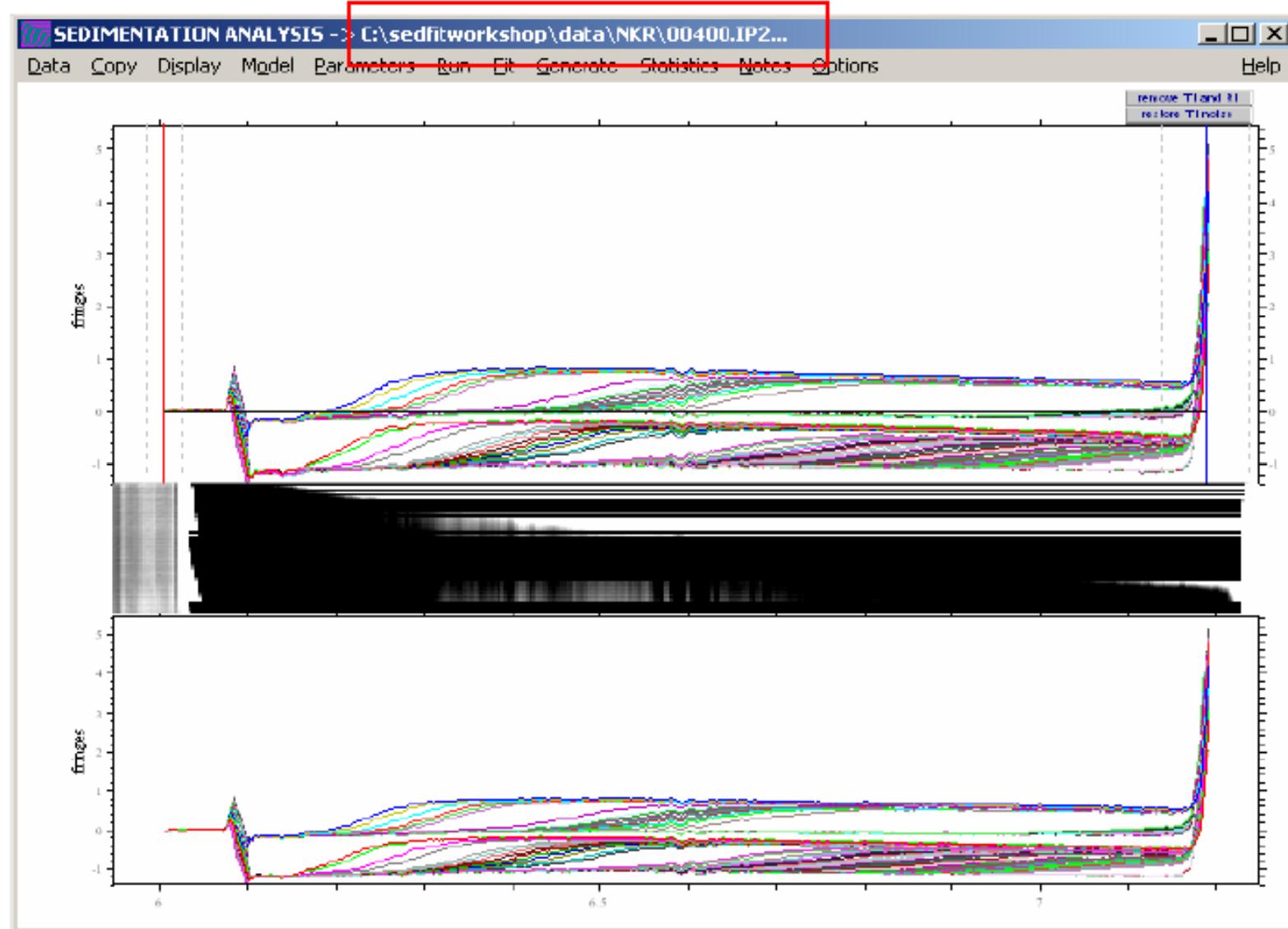


Selecting YES will load all the parameters that we used before (ie. Meniscus, bottom, fitting limits, etc.)

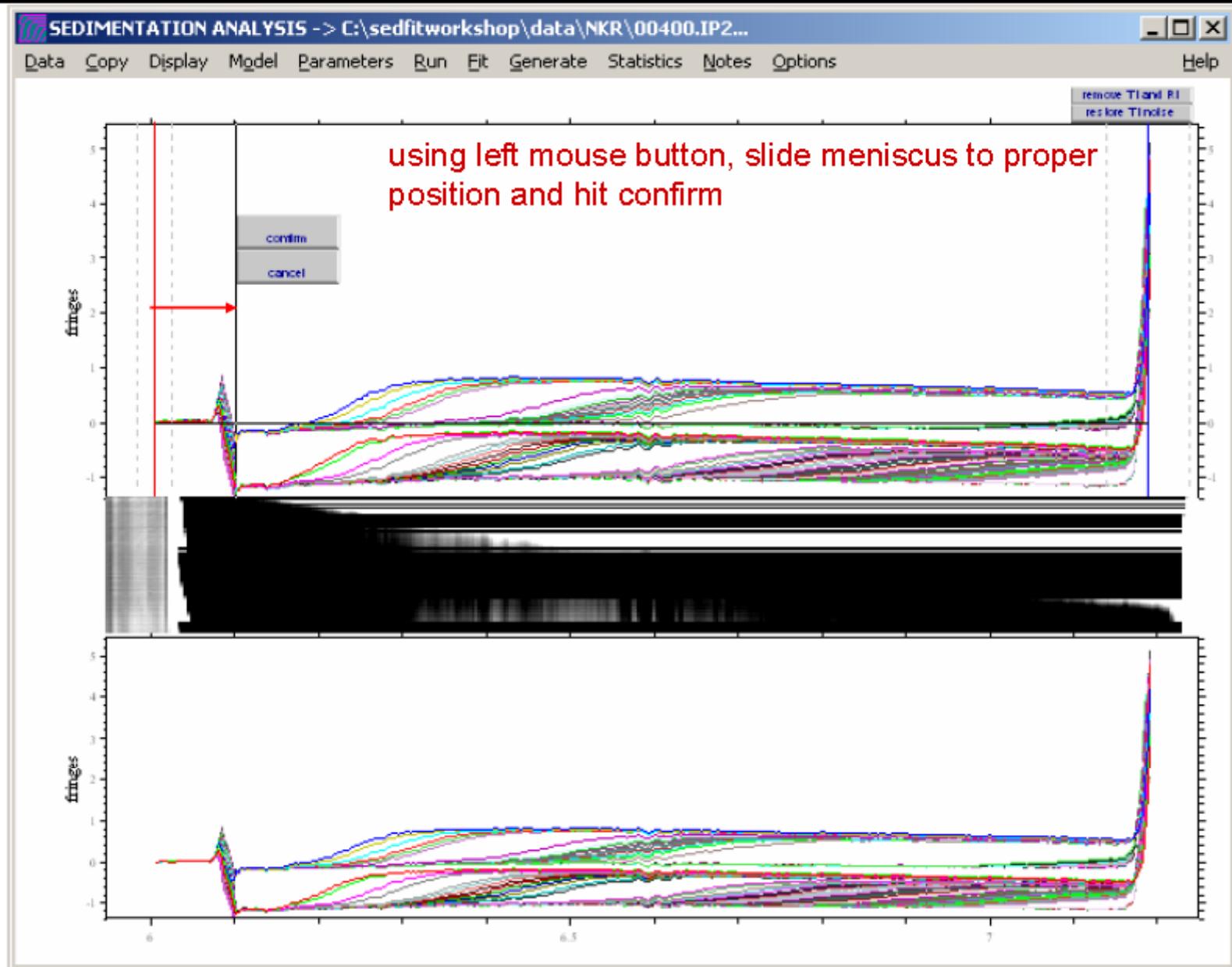
Selecting NO will result in the user needing to set these parameters again.

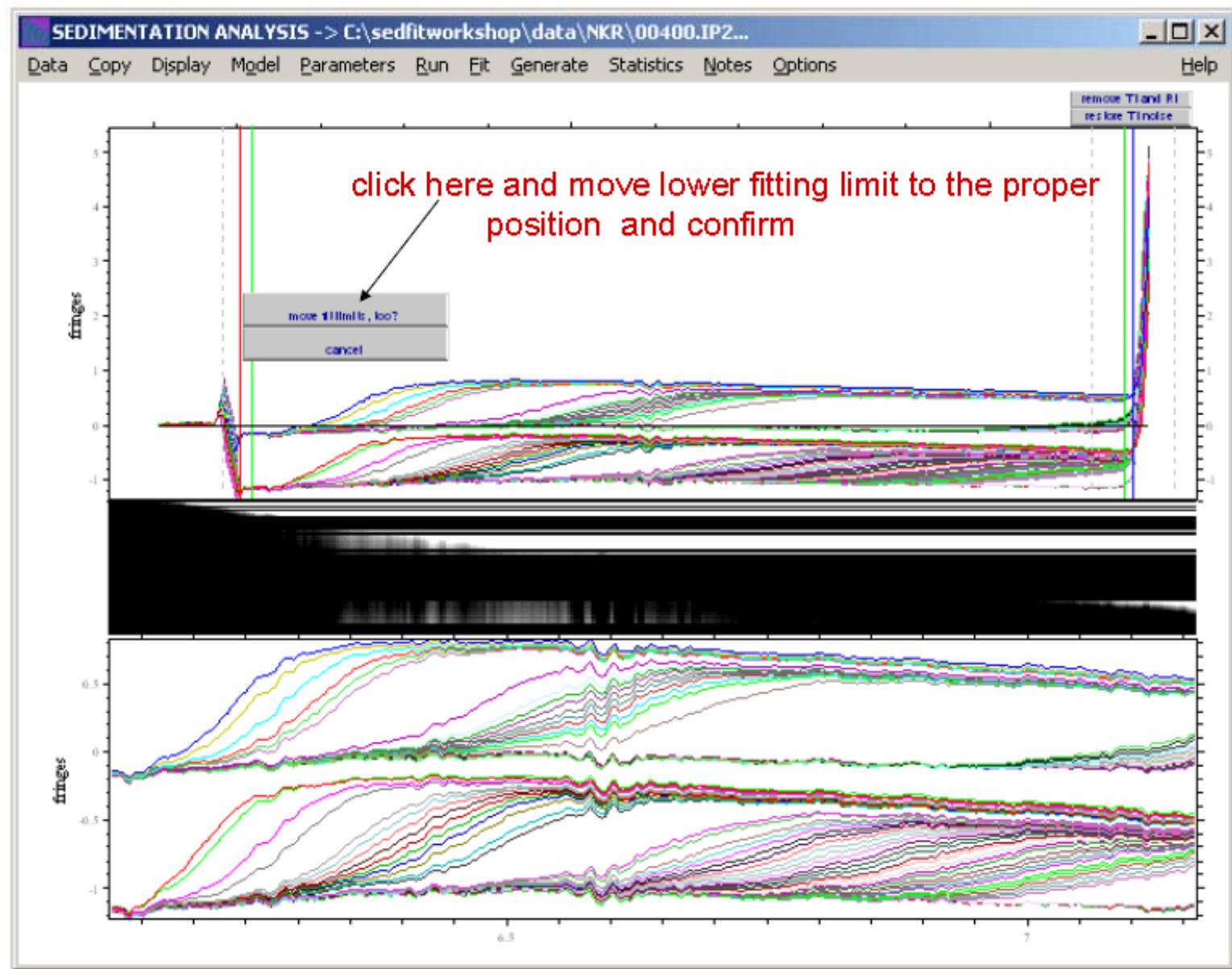
Select YES

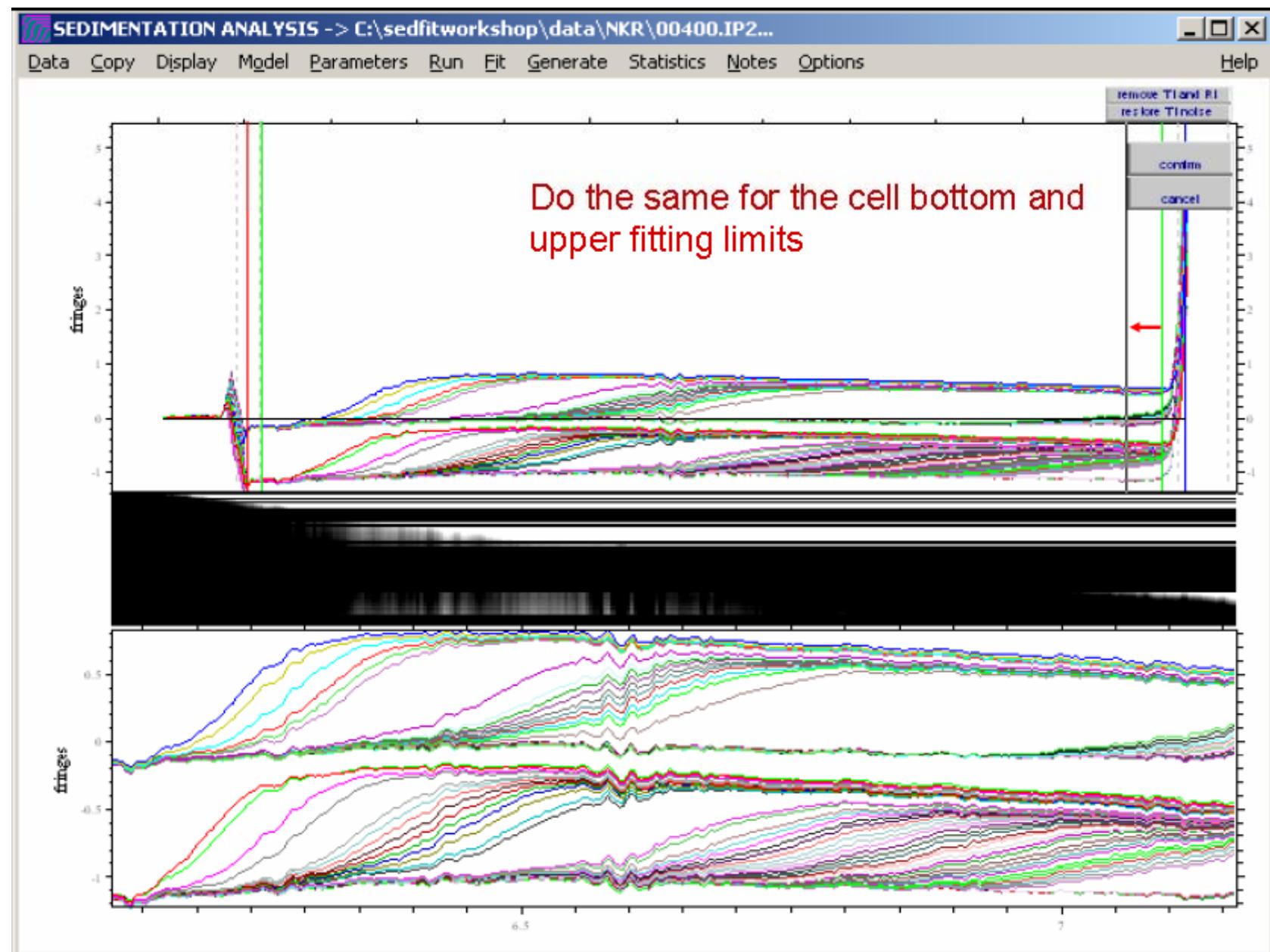
Every 5th scan loaded.

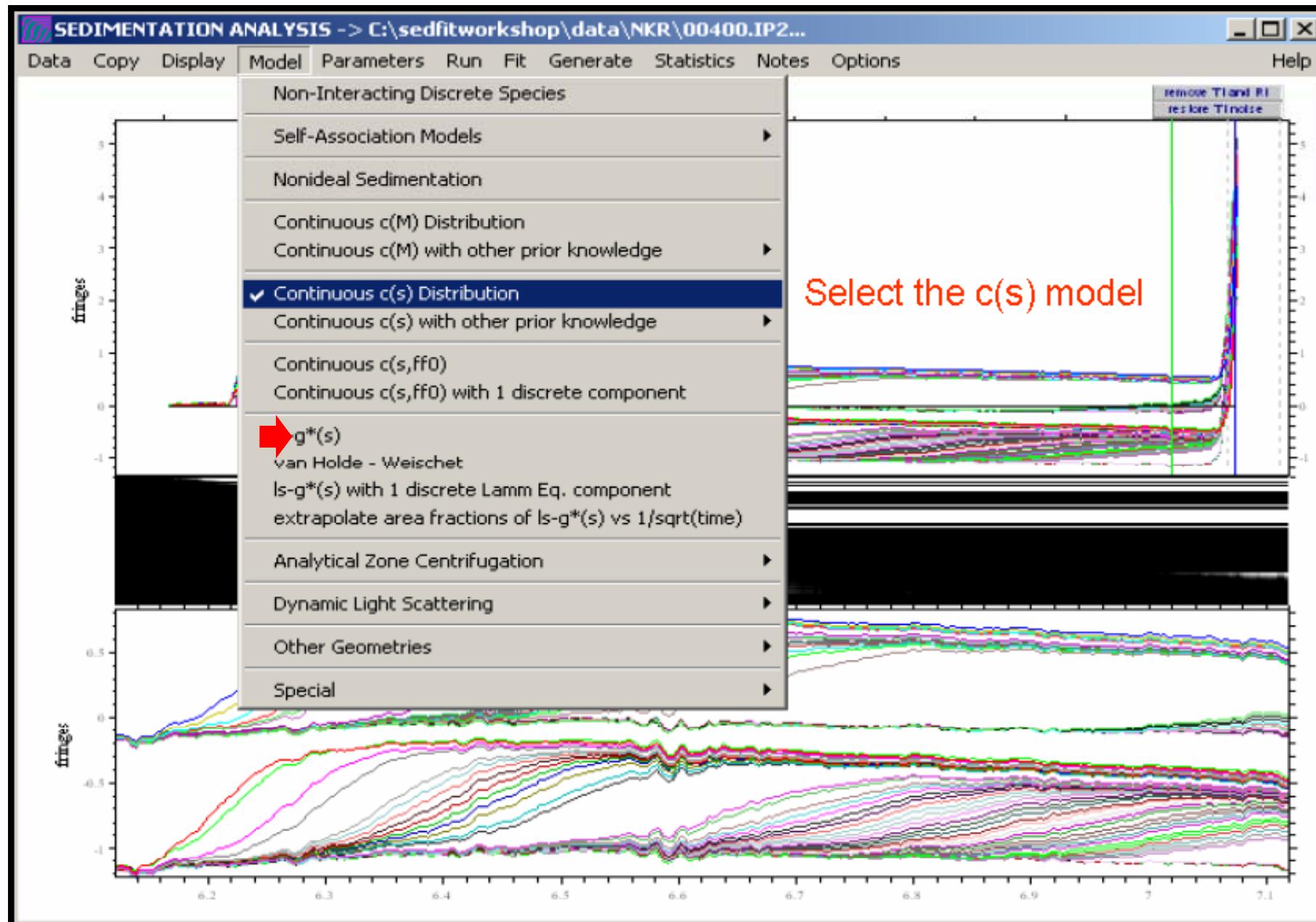


Next, select the meniscus position

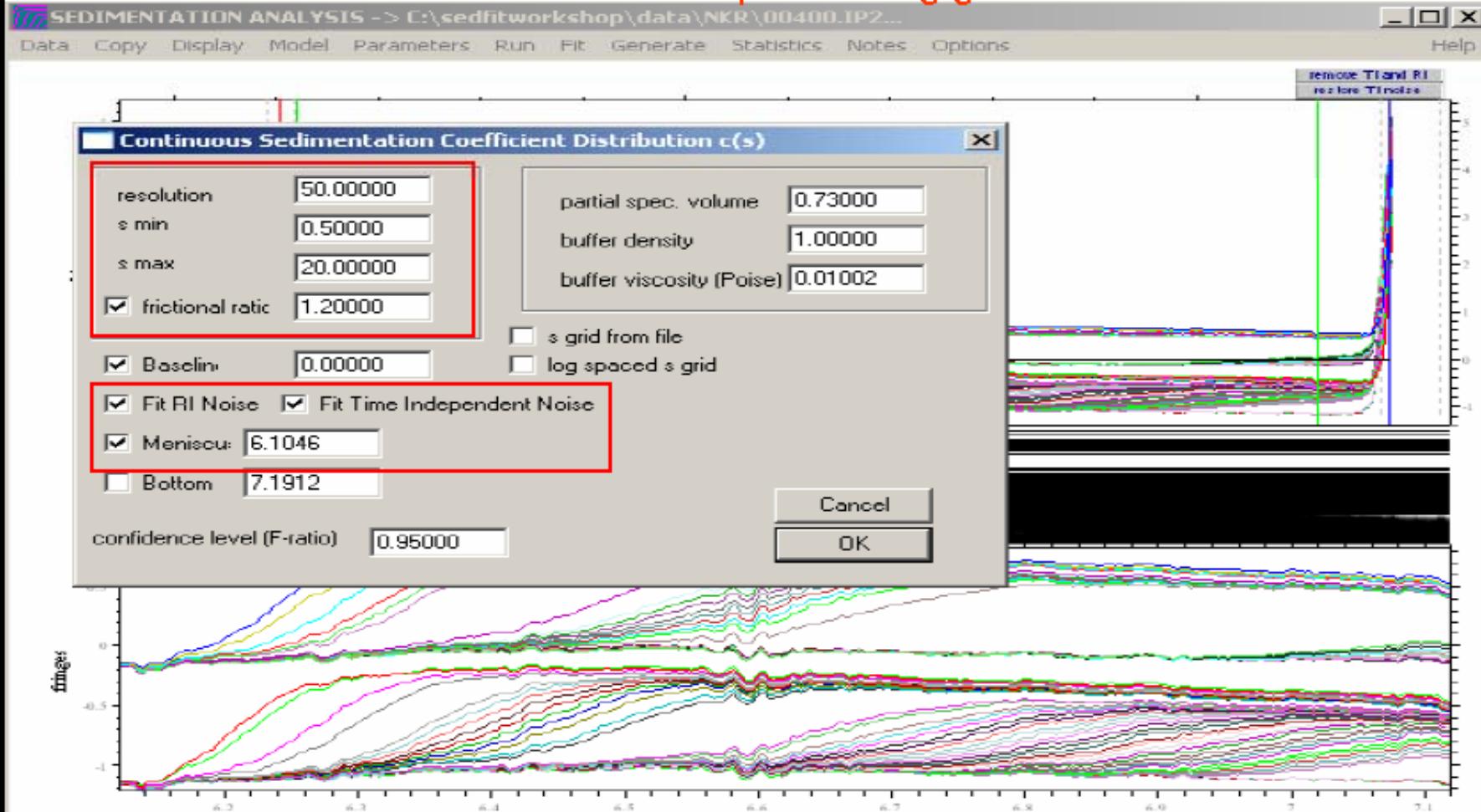








Select PARAMETERS and input starting guesses

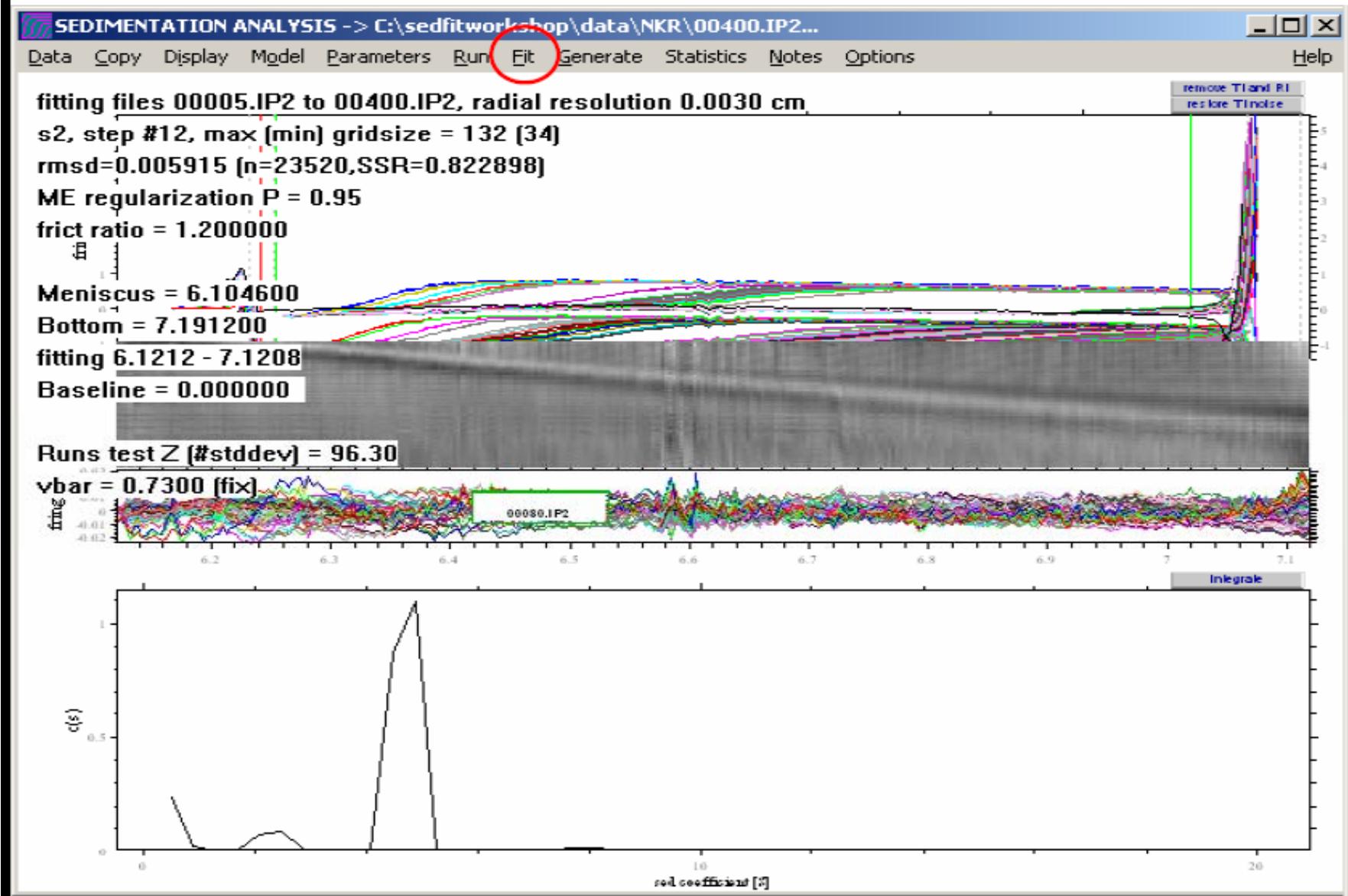


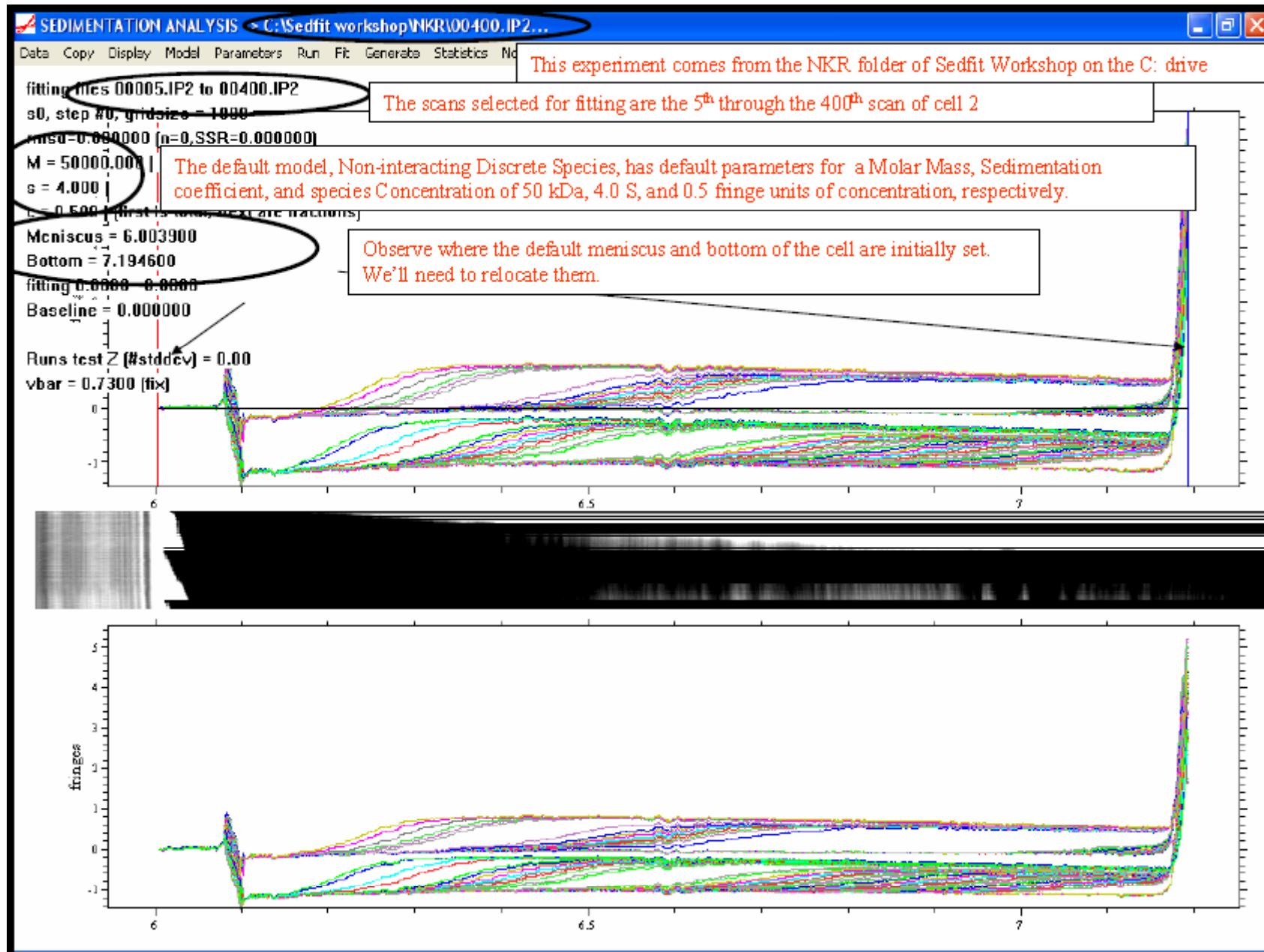
Here we are going to model the experimental data with simulated curves for 50 s-values between 0.5 and 20S.

We will float the frictional ratio Fit for TI and RI noise and Float the meniscus

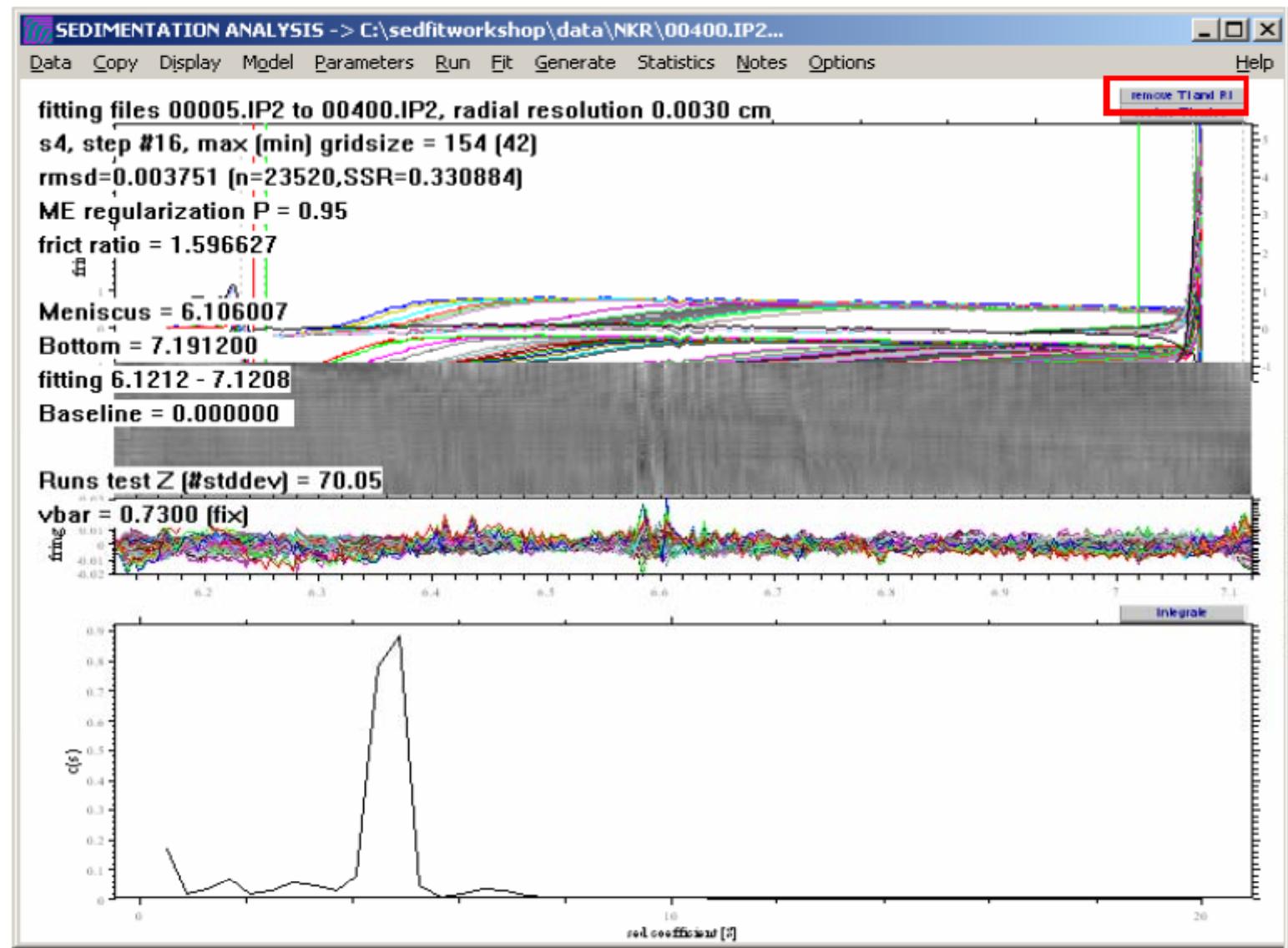
HIT OK

Perform a RUN which optimizes the linear parameters. Then FIT

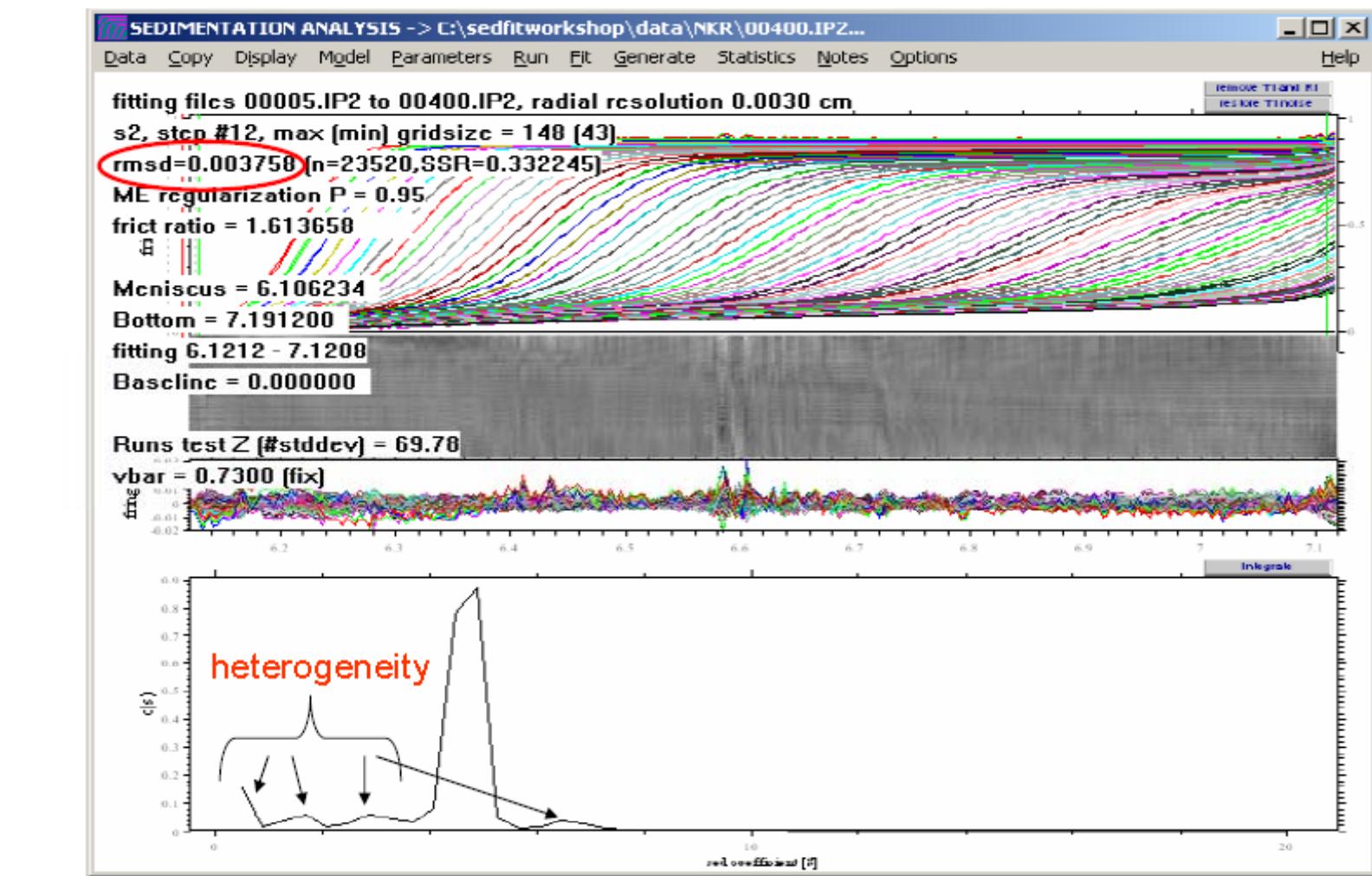


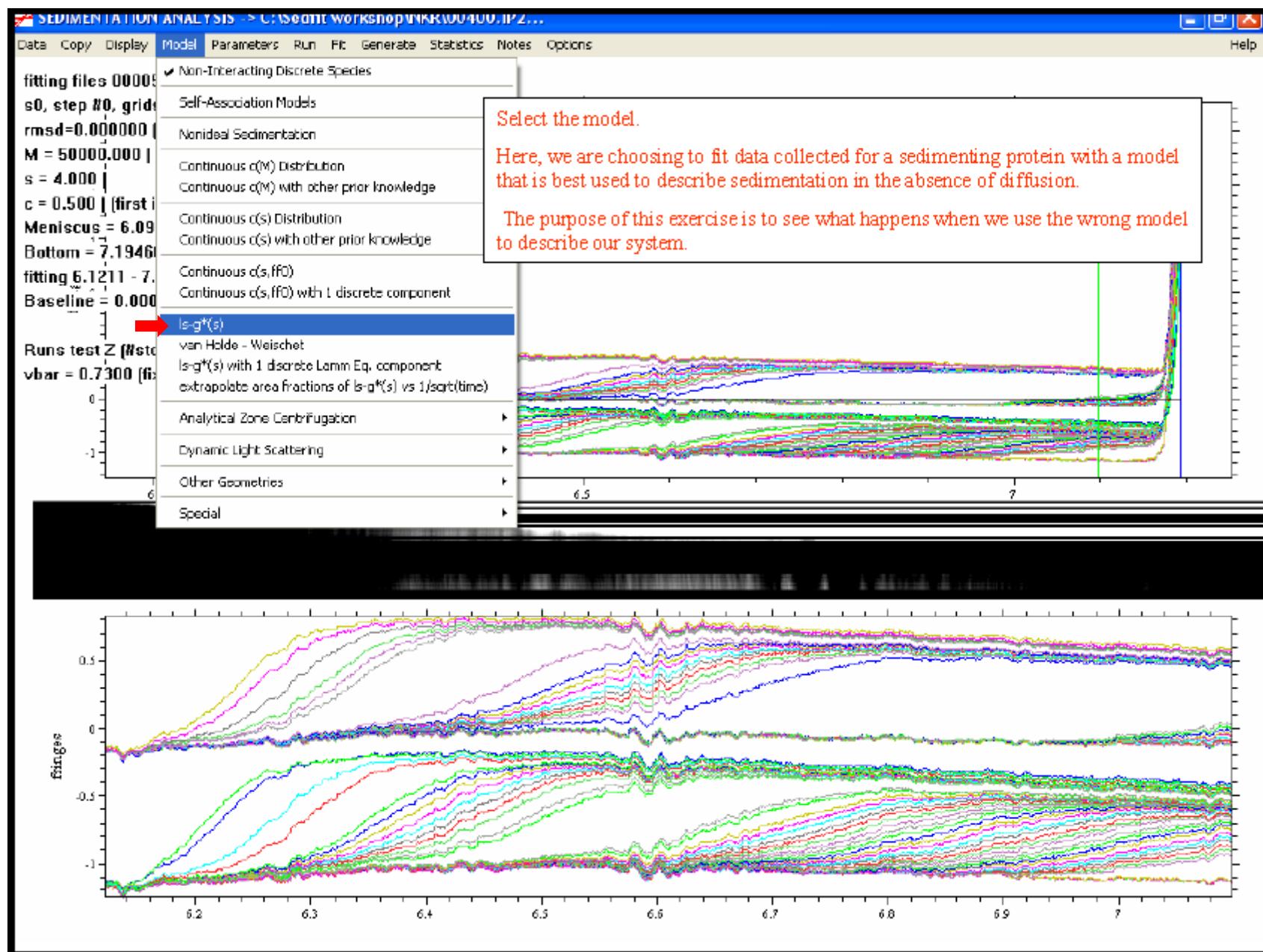


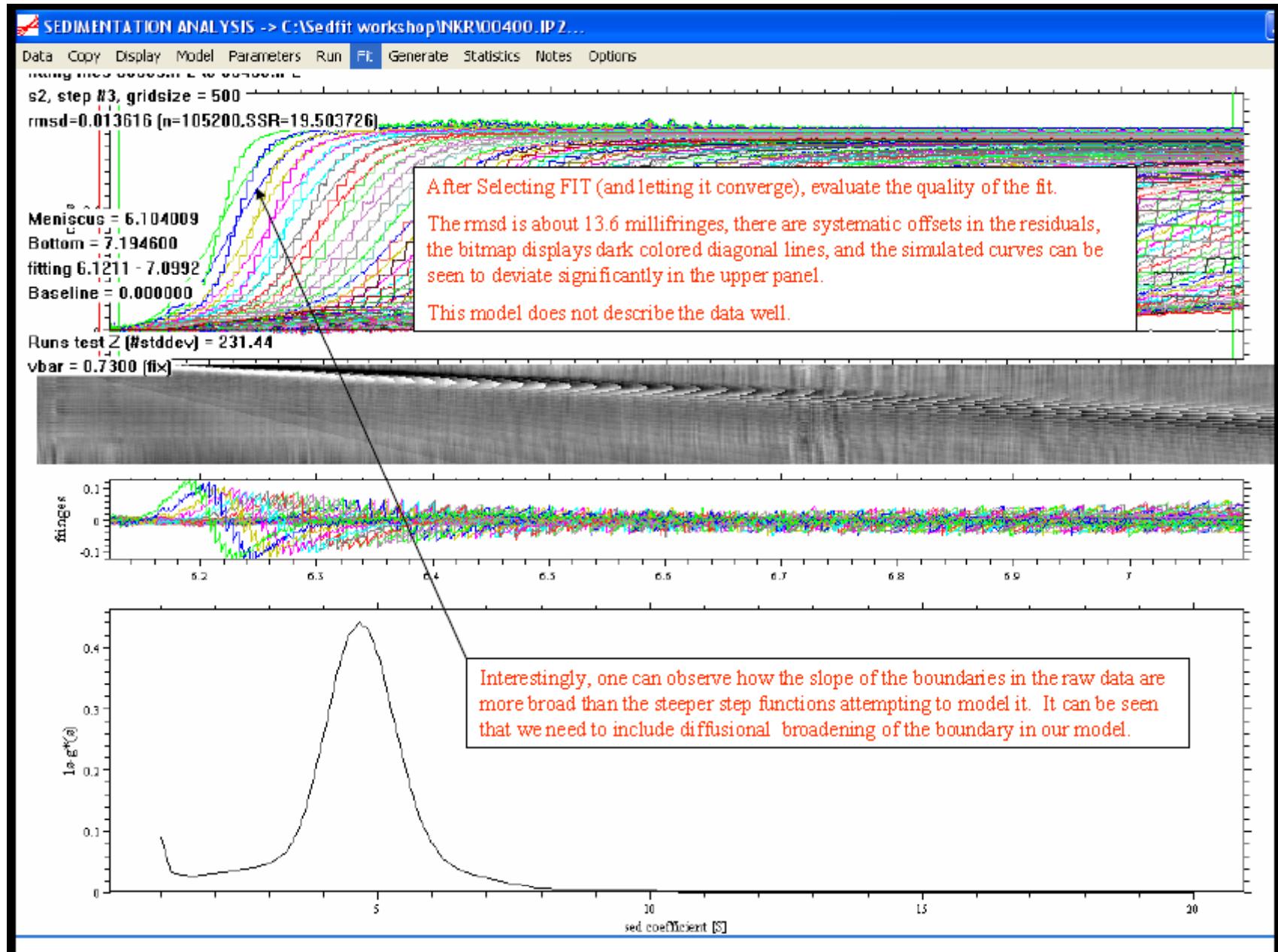
Fitting optimized all parameters. Then subtract the calculated RI and Ti noise.

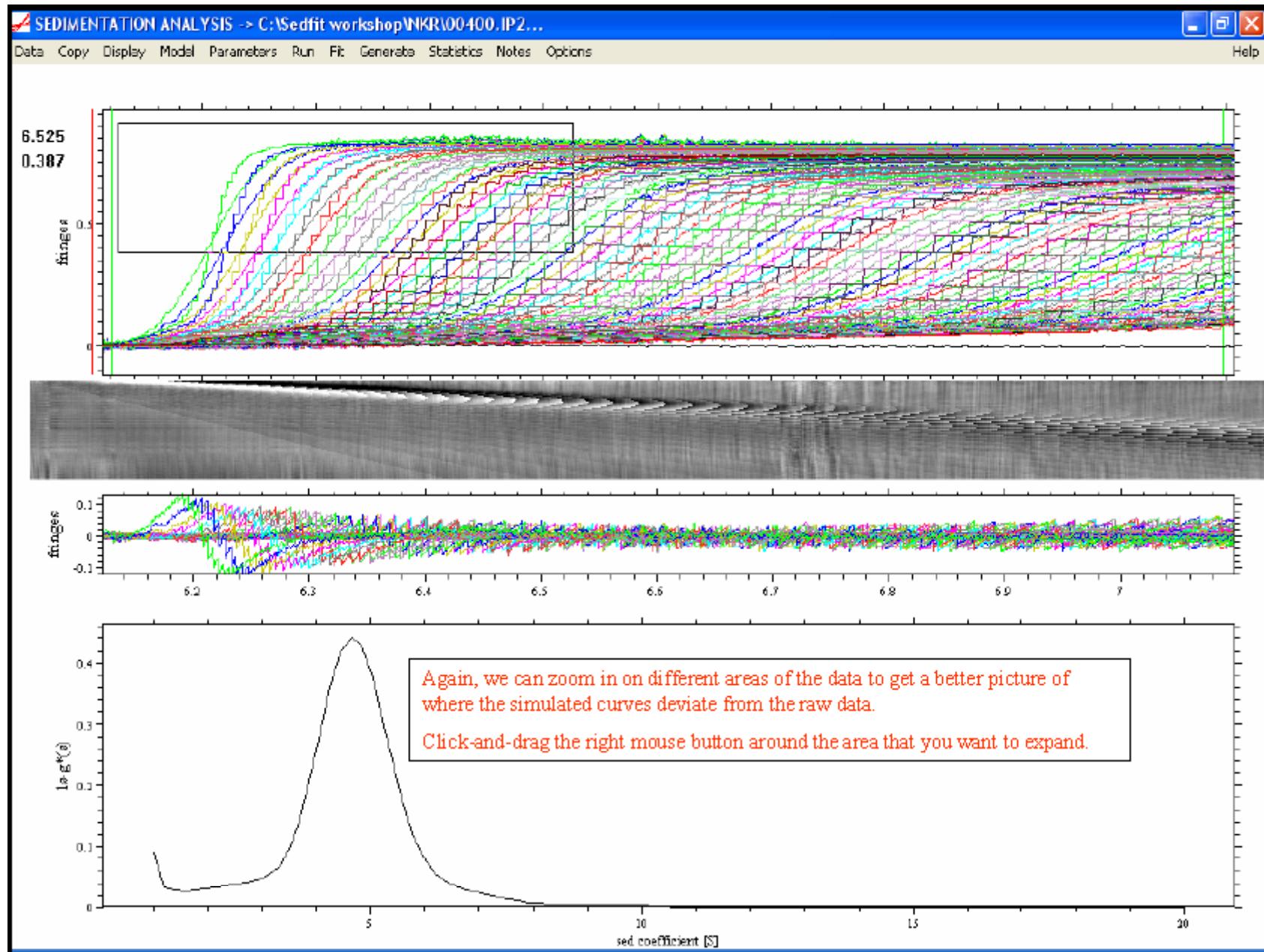


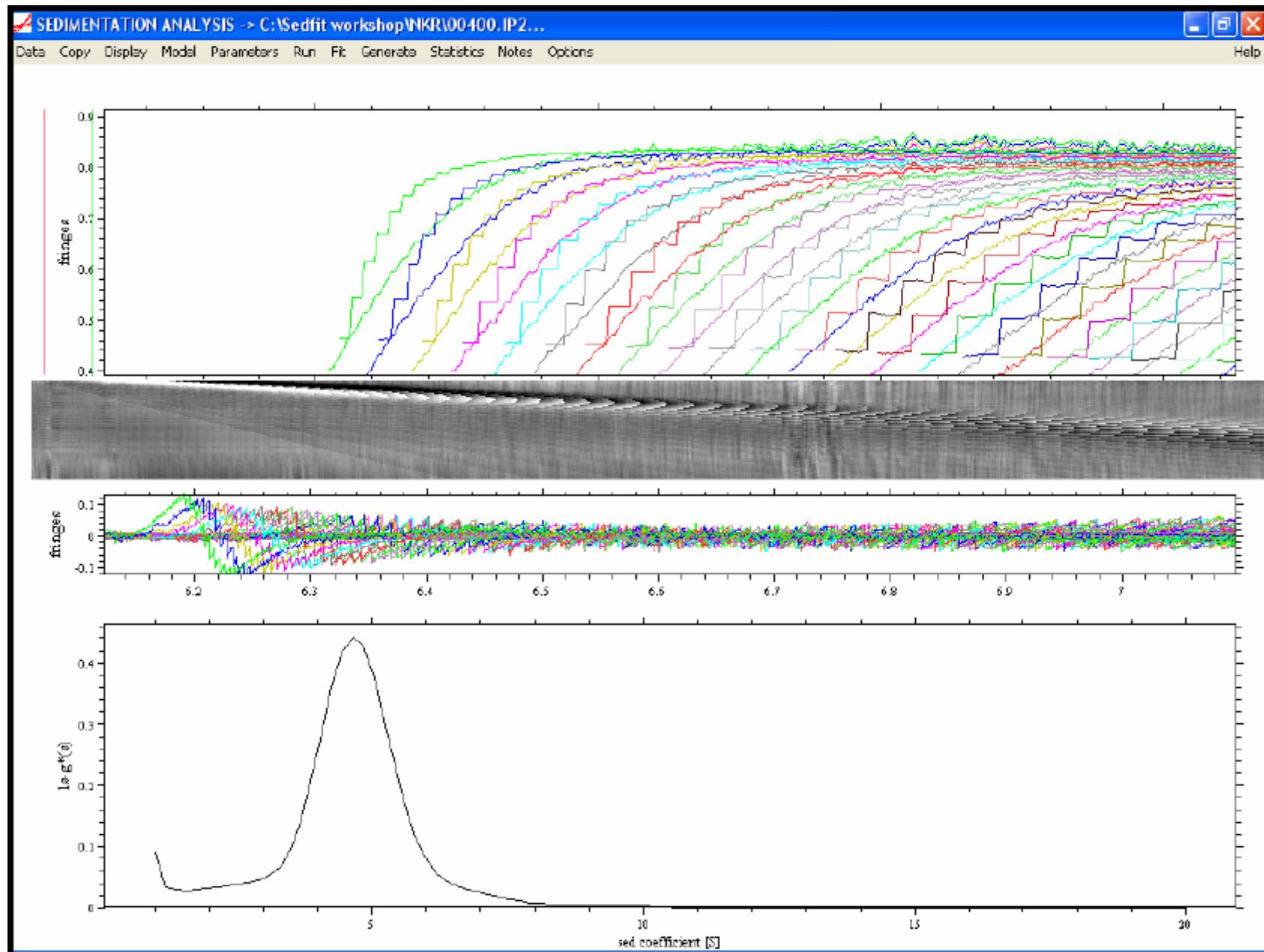
Assess the quality of the fit! Notice root mean square deviation (rmsd) >0.010
Compare to the Is-g*(s) model, and the Non-interacting discrete species model.
Notice the heterogeneity of the sample revealed by the c(s) distribution.

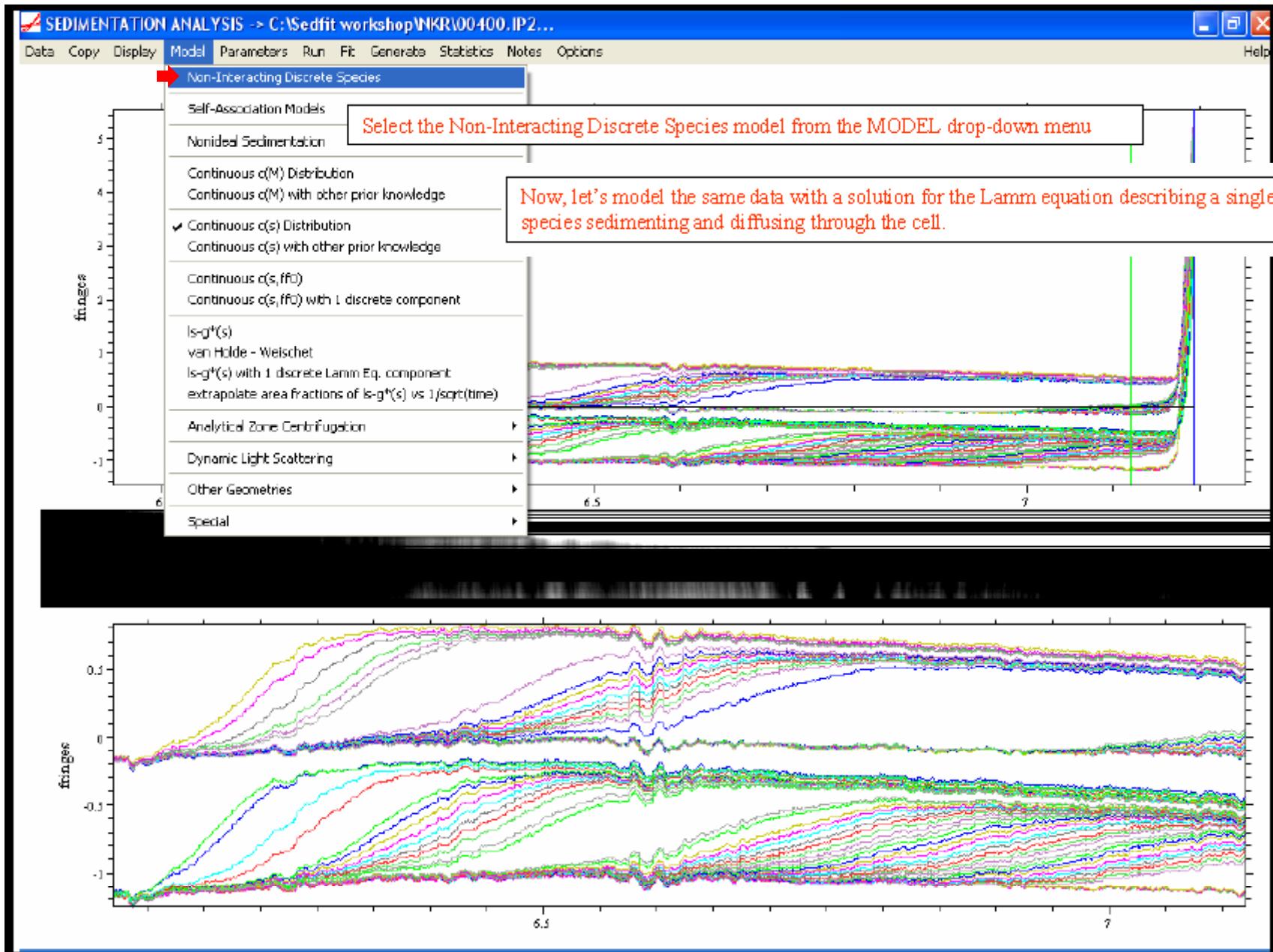


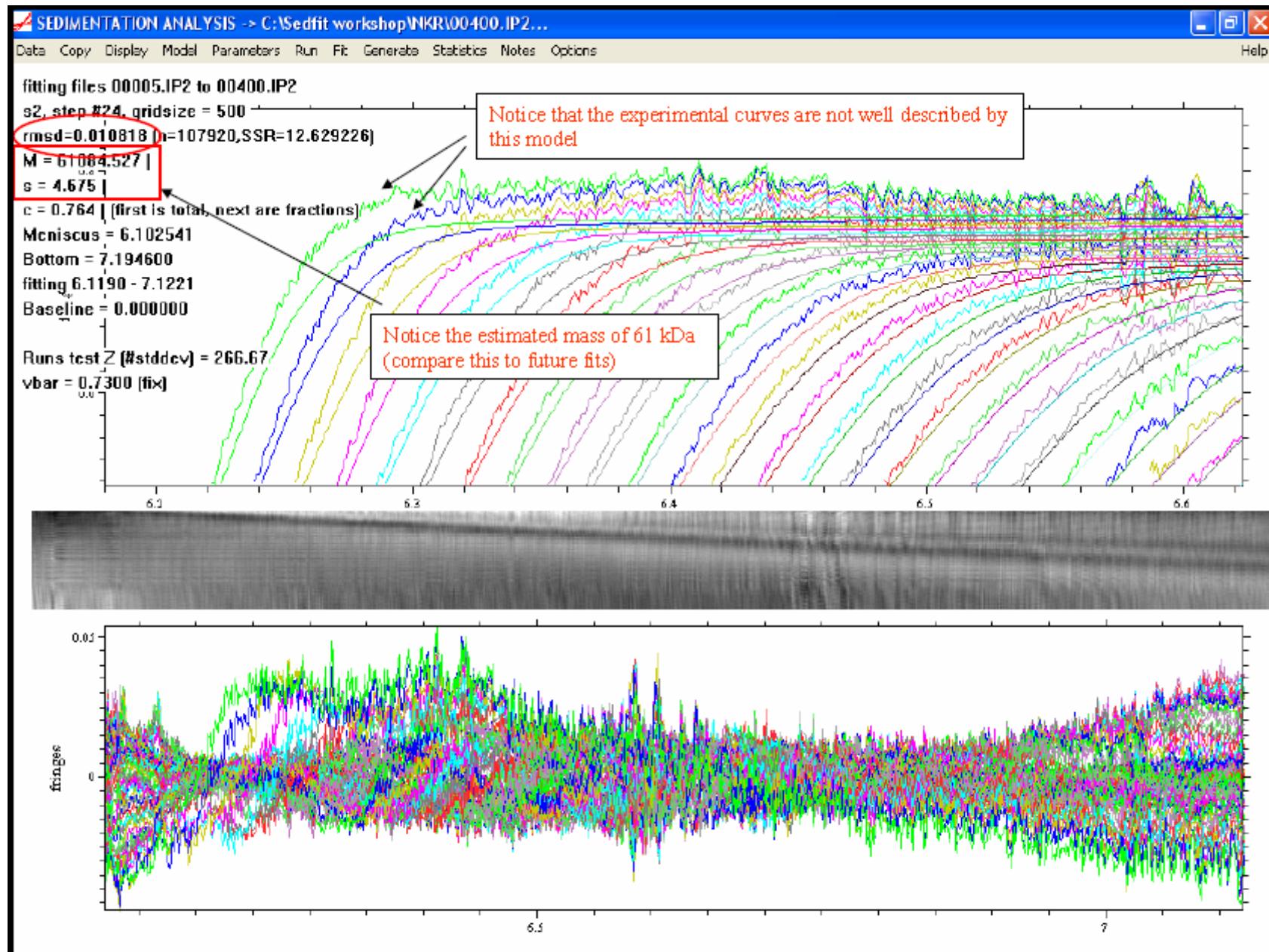




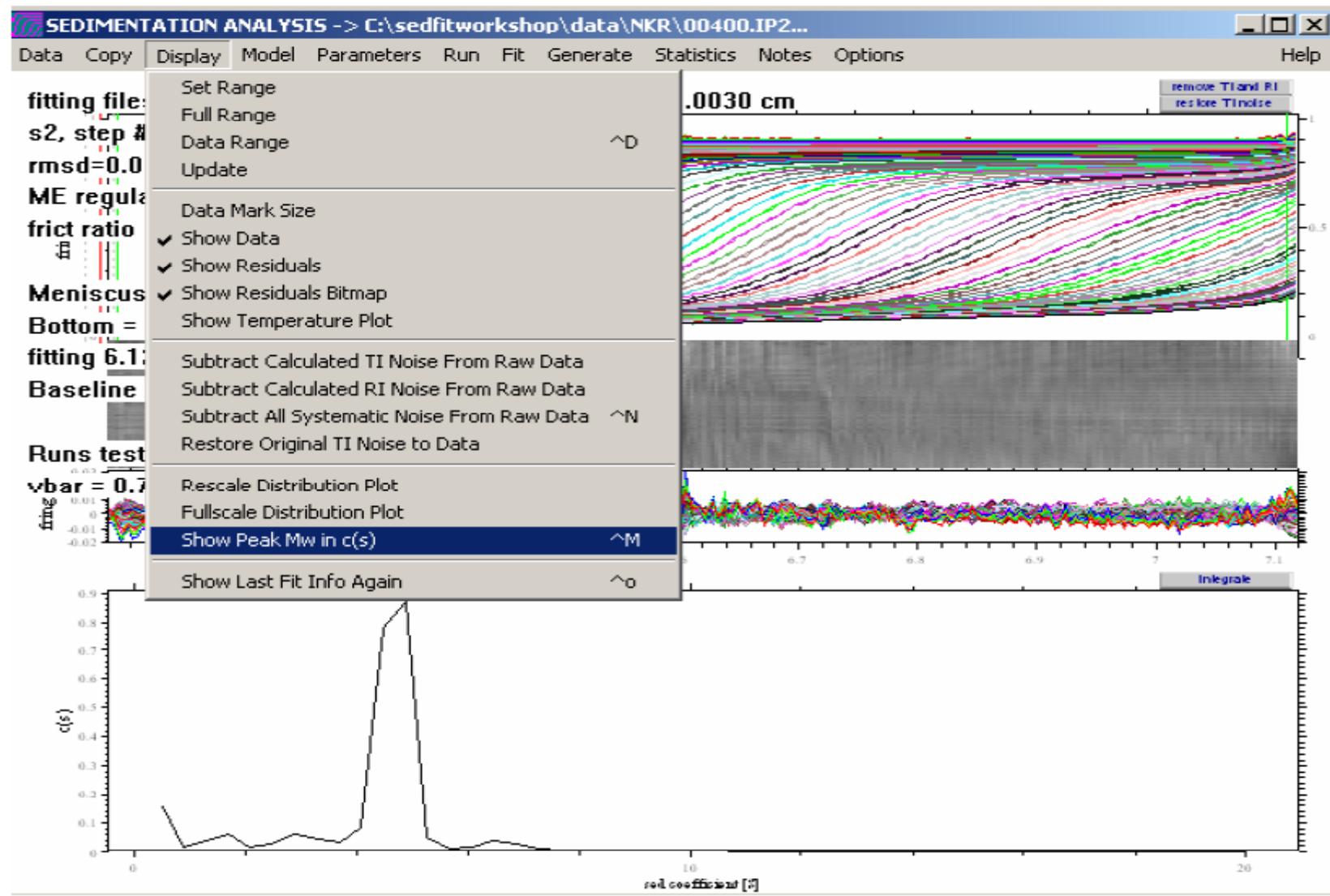




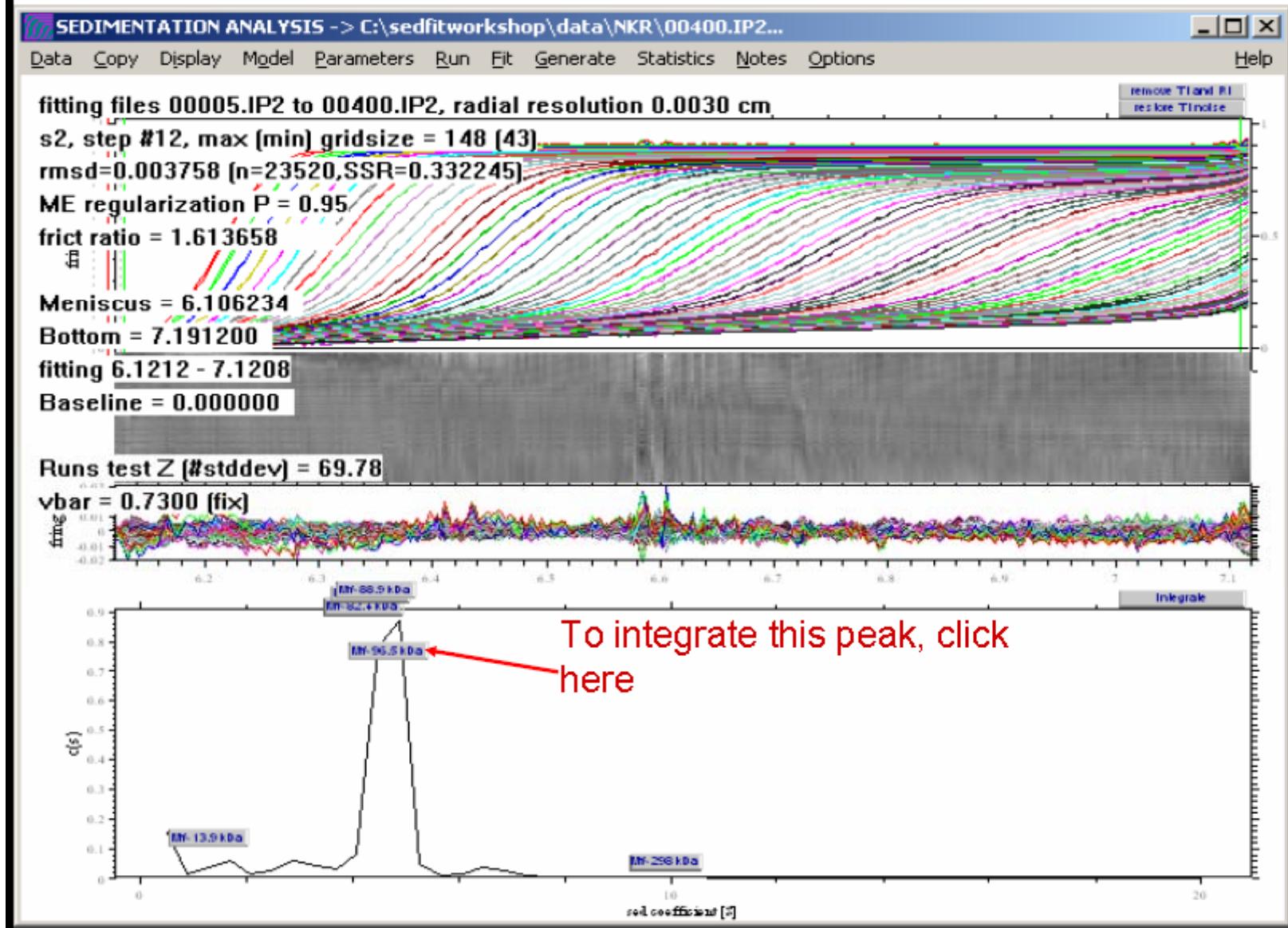




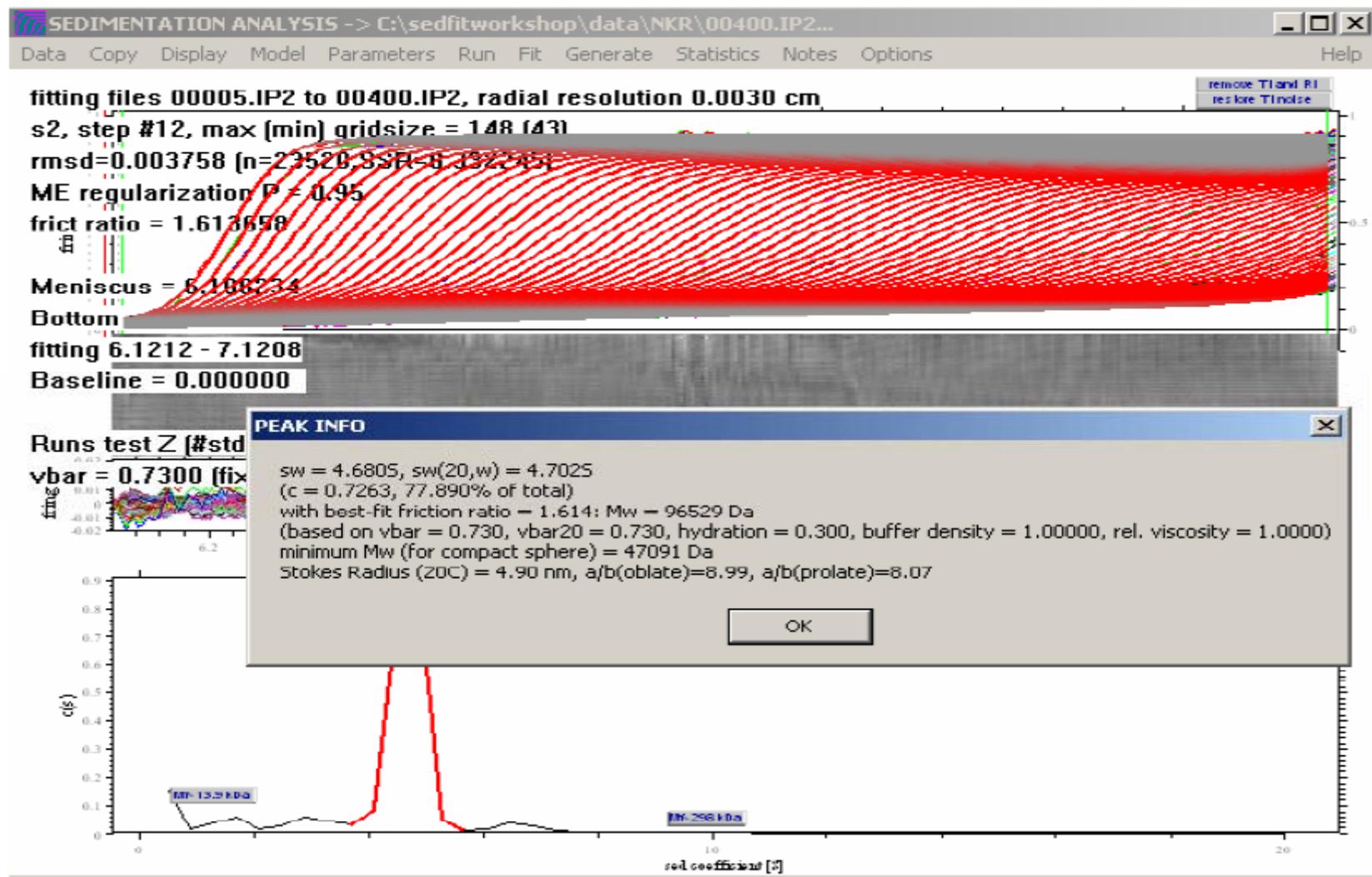
Display the peak MW either from the top menu or by hitting “control M”

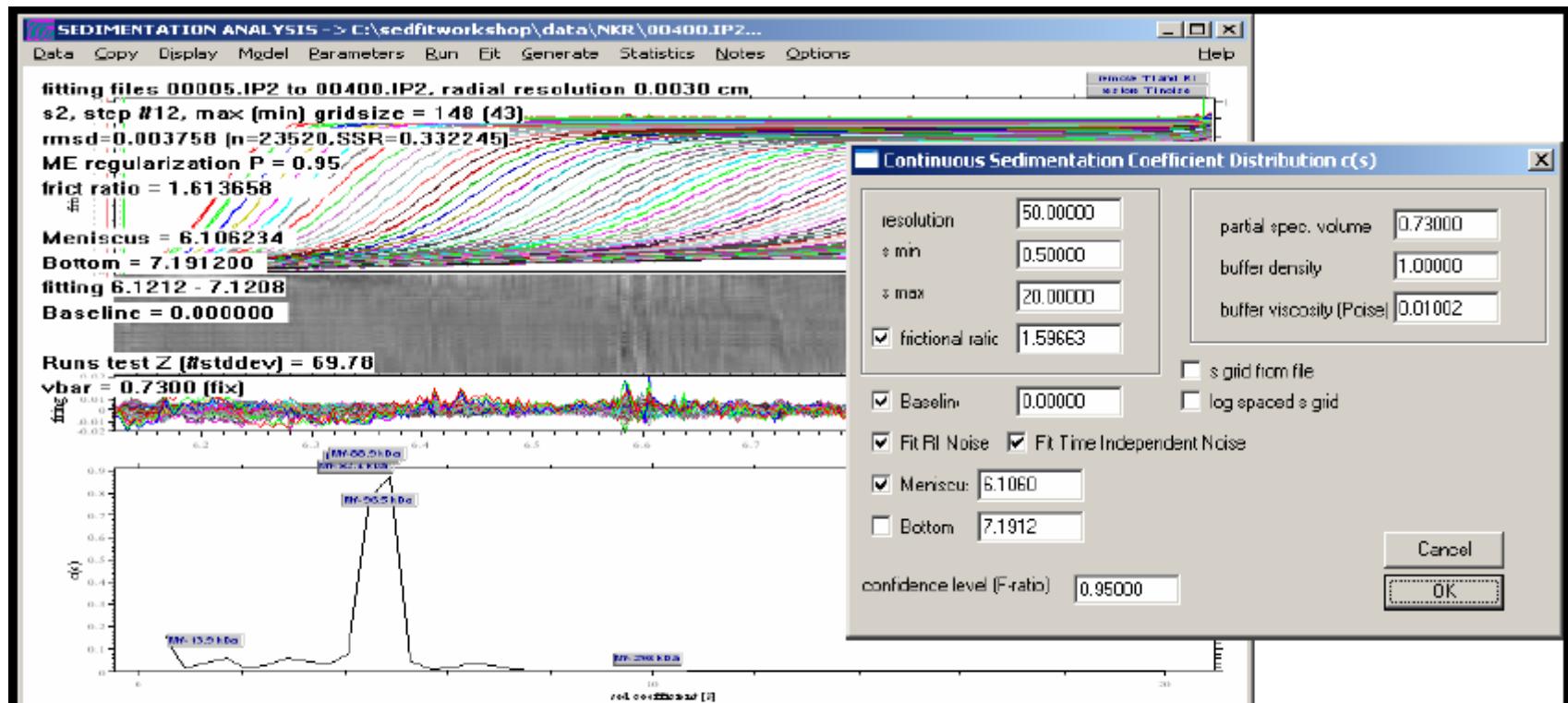


Copy and paste the final cS with MW displayed into Power Point to save



Location of this MW displayed.
Peak information window opens. Copy and paste to ppt.





Try the following:

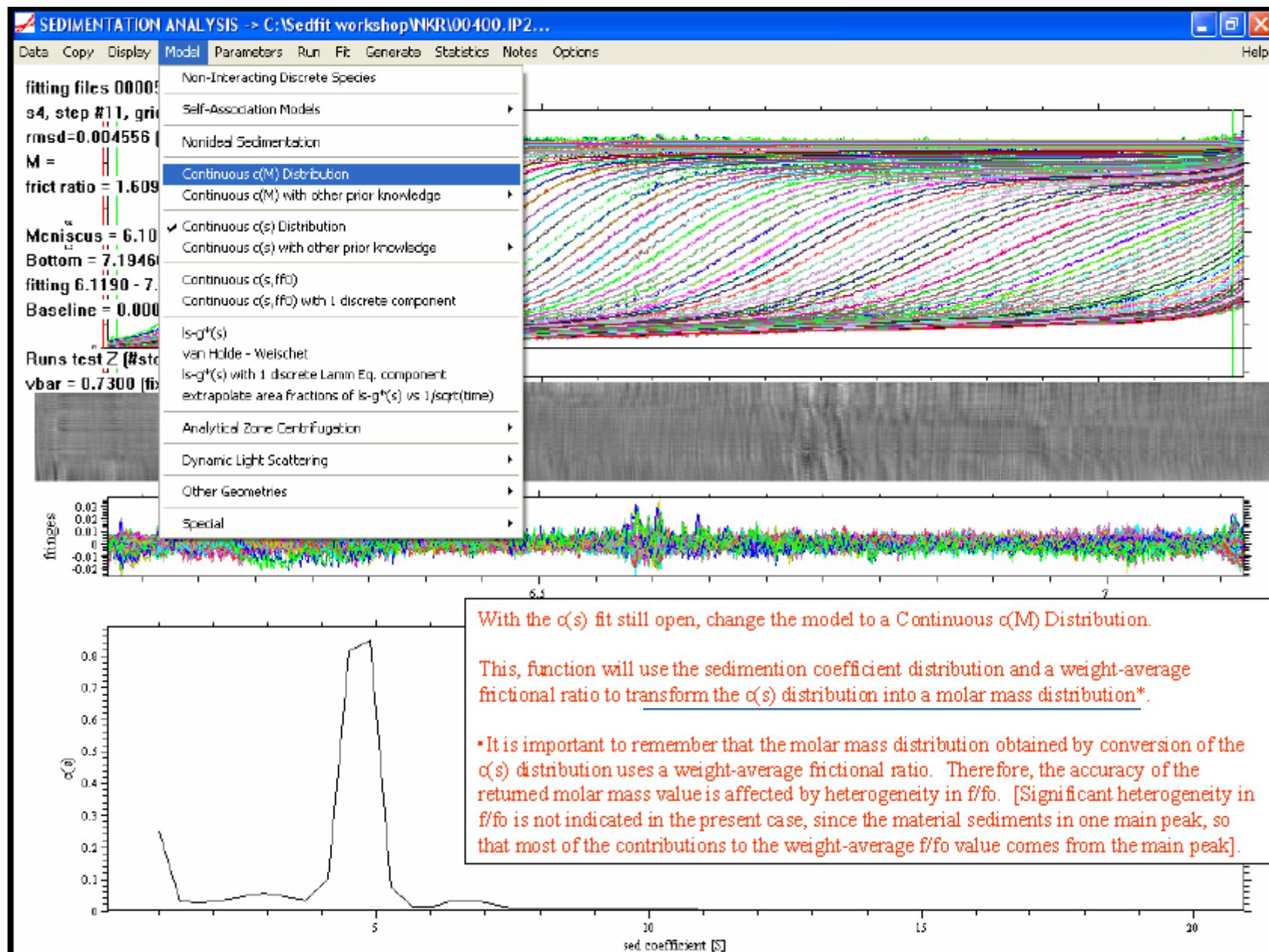
Zoom in the distribution range from 10-20 S (with dragging a rectangle with the right mouse button) – there's not much material out there.

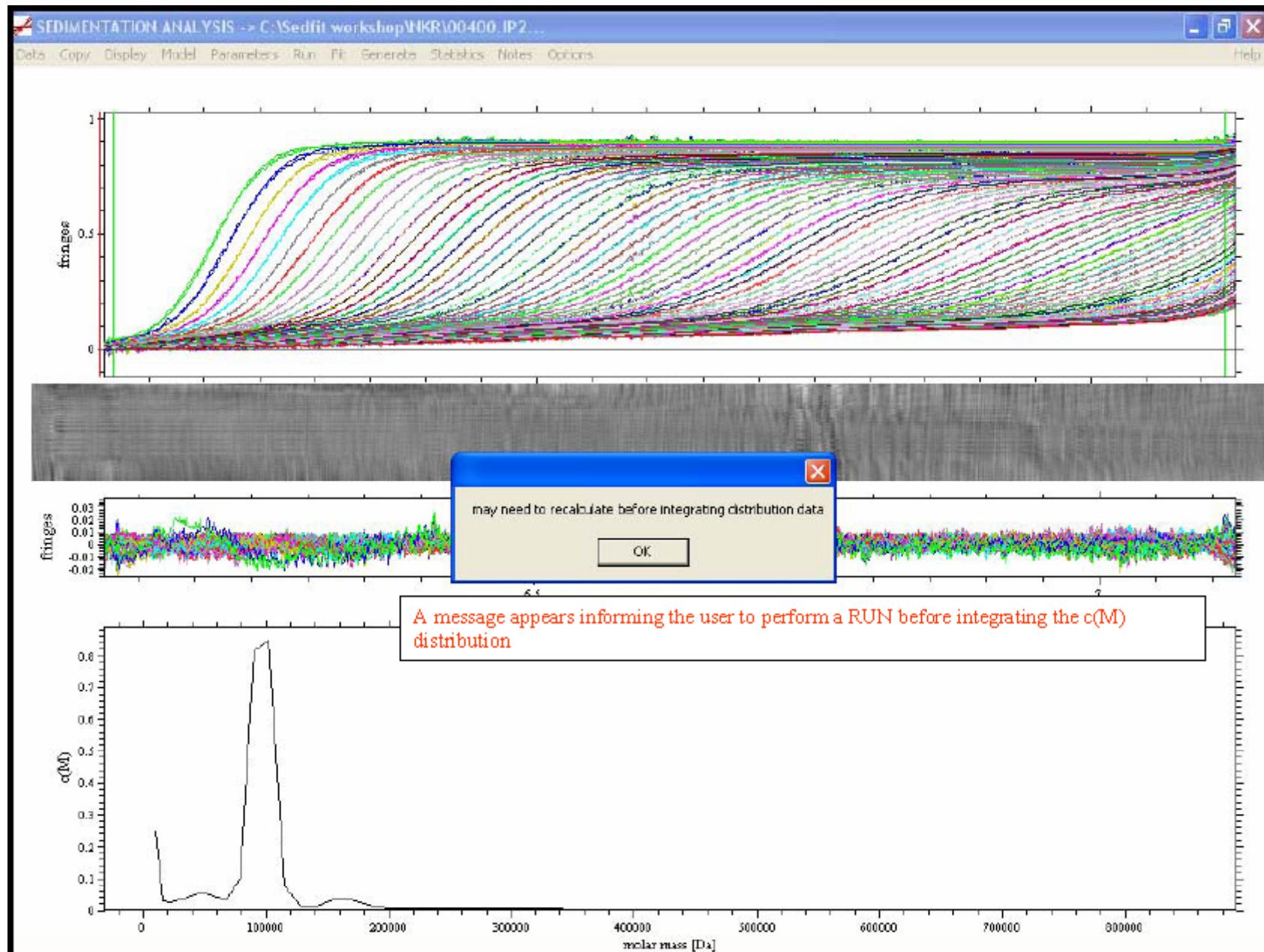
You can further refine the fit by selecting $s_{max} = 15$, and $resolution = 100$

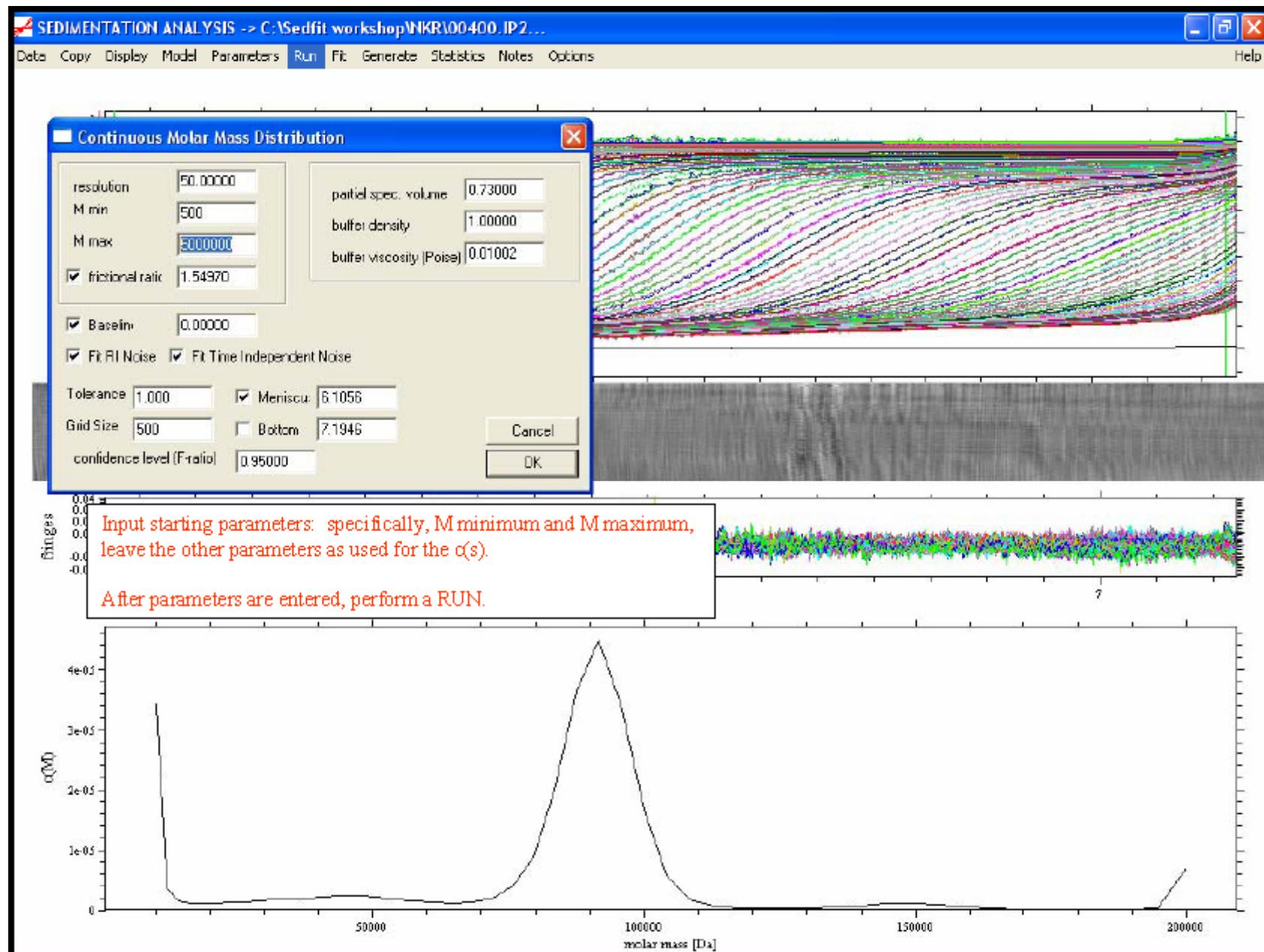
Do a RUN – you'll see a smoother distribution.

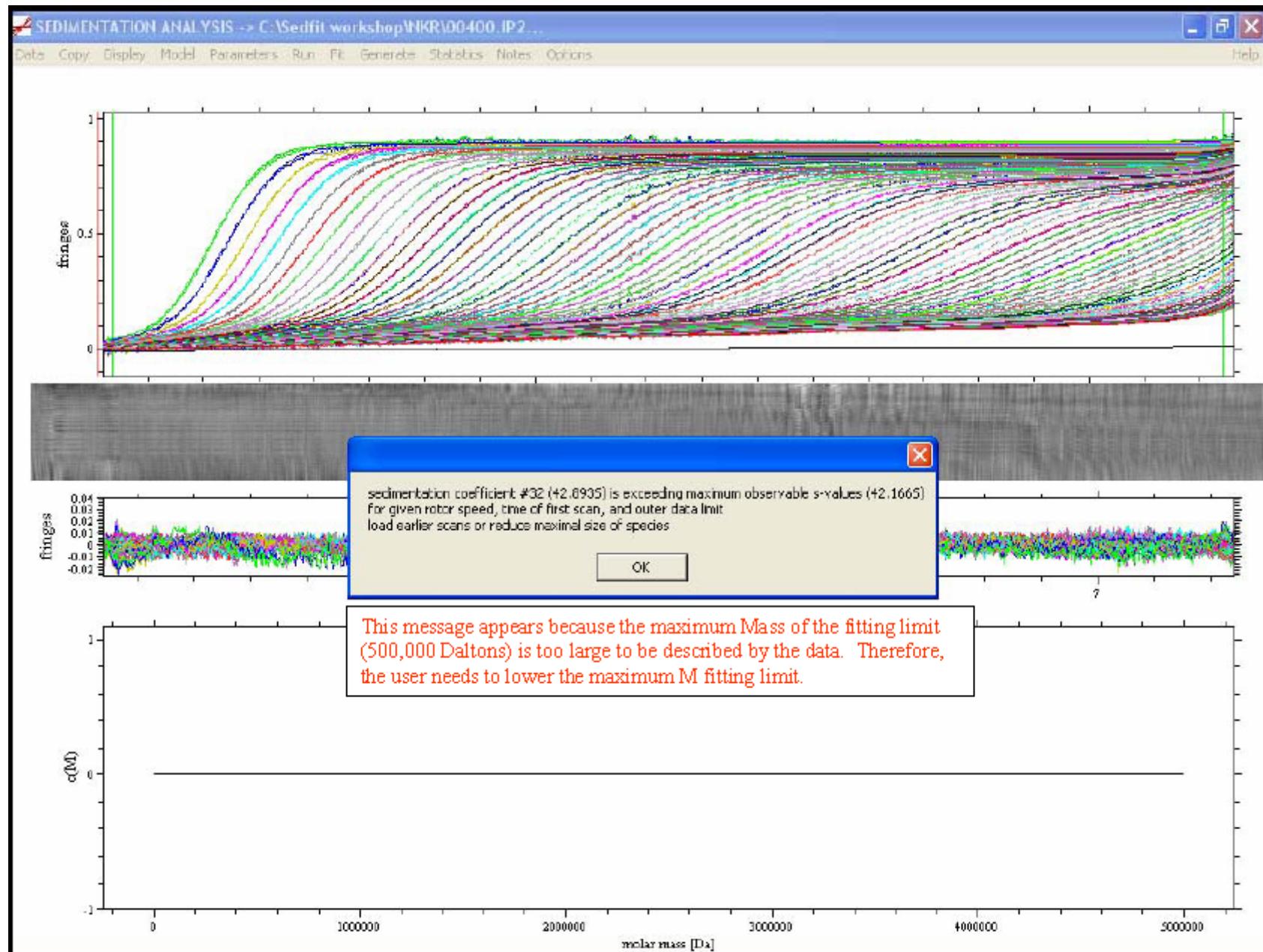
Then,

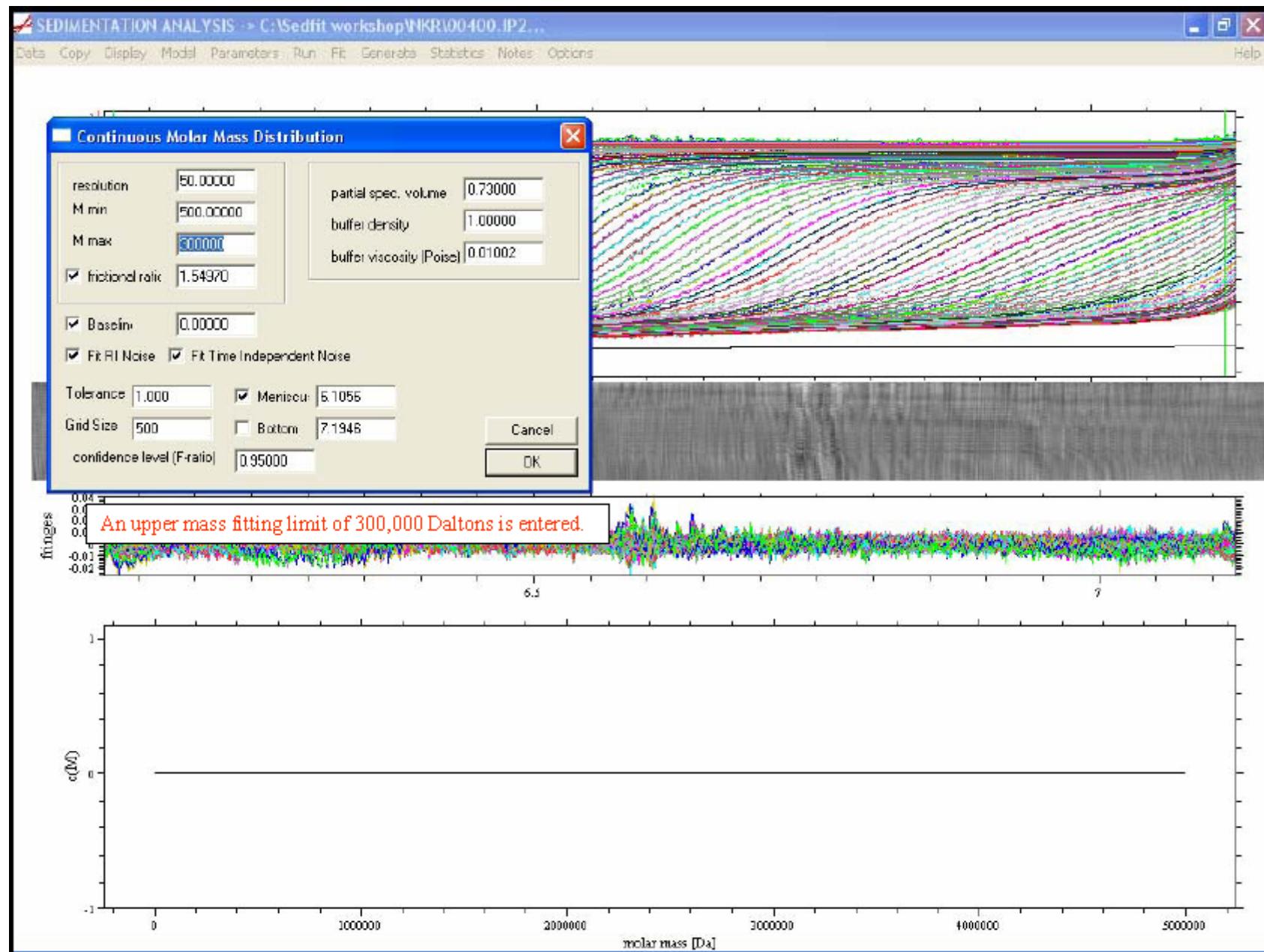
Since the distribution goes up at s_{min} , test what decreasing the lower limit s_{min} does (for example, set $s_{min} = 0.2$). One of the two should be expected: a) a new peak appears and $c(s)$ goes down at the new s_{min} , and $rmsd$ gets slightly better; b) $c(s)$ goes up further at s_{min} . What that means is that s_{min} is low enough to describe all possible sedimentation at the lower end, and is correlated with the baseline (baseline = no sedimentation). This is no problem. In this case, you can simply ignore this partial peak at s_{min} , and you can zoom into the rest of the distribution to see that better. The option b) is what happens here.

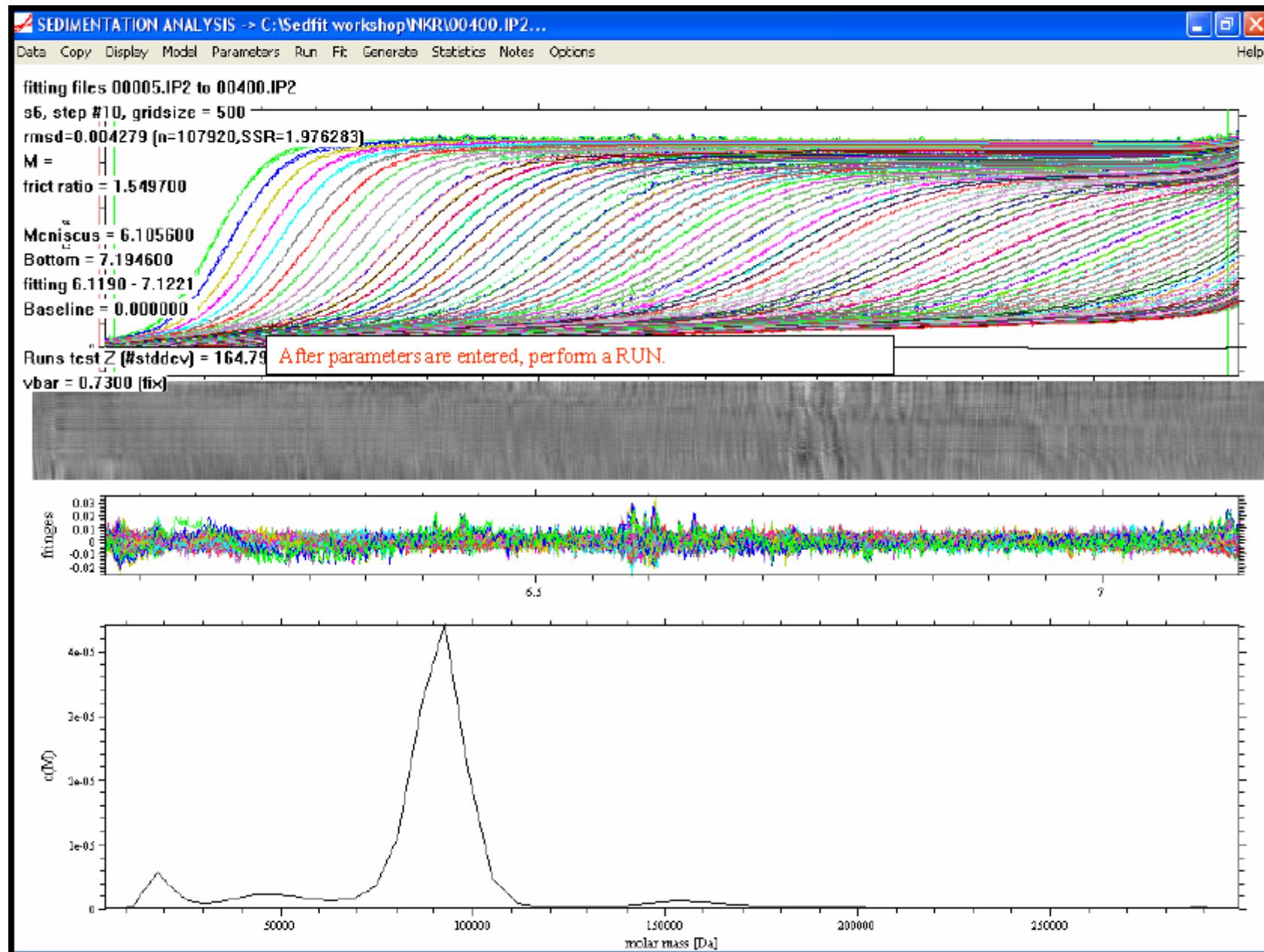


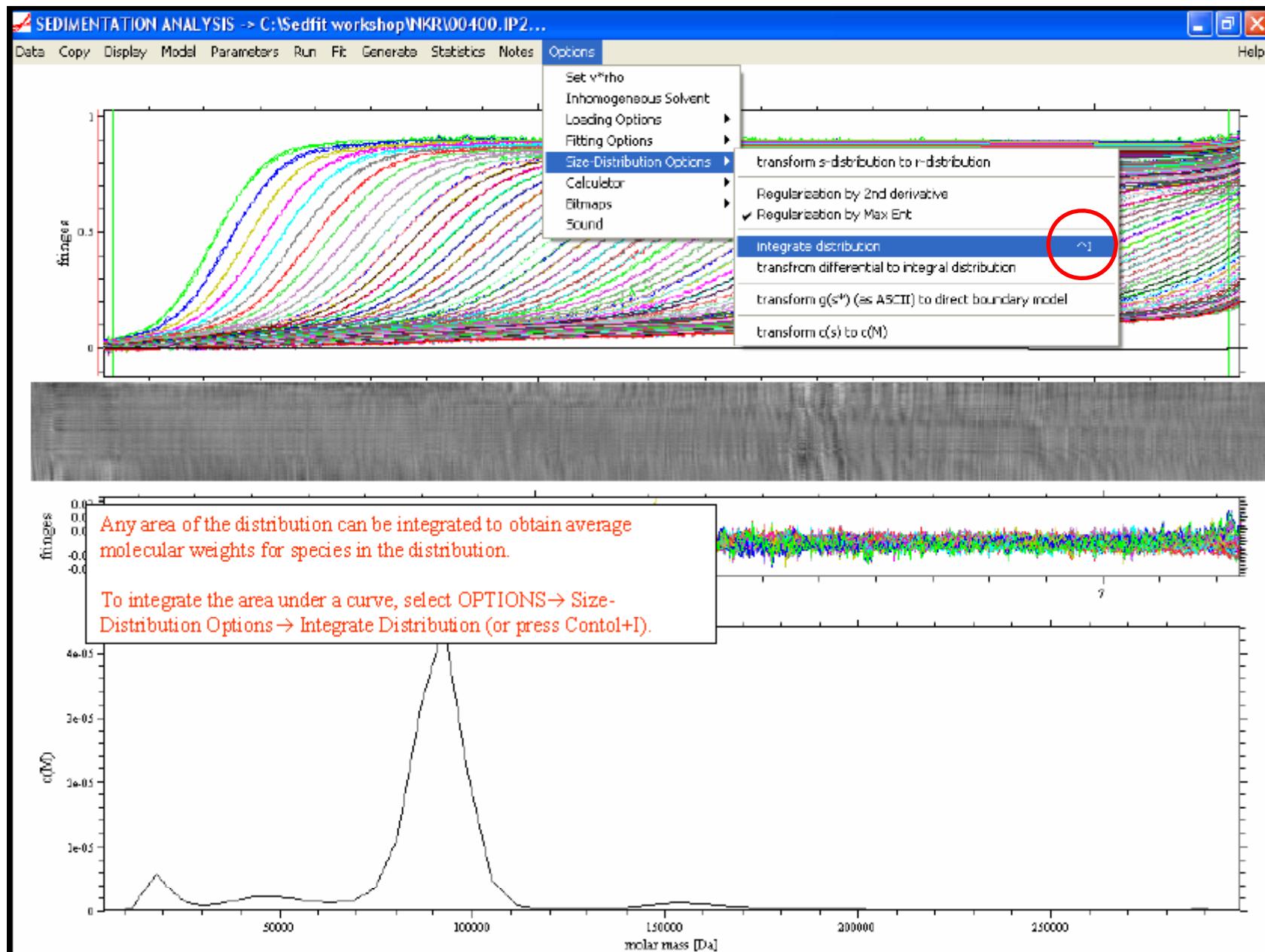


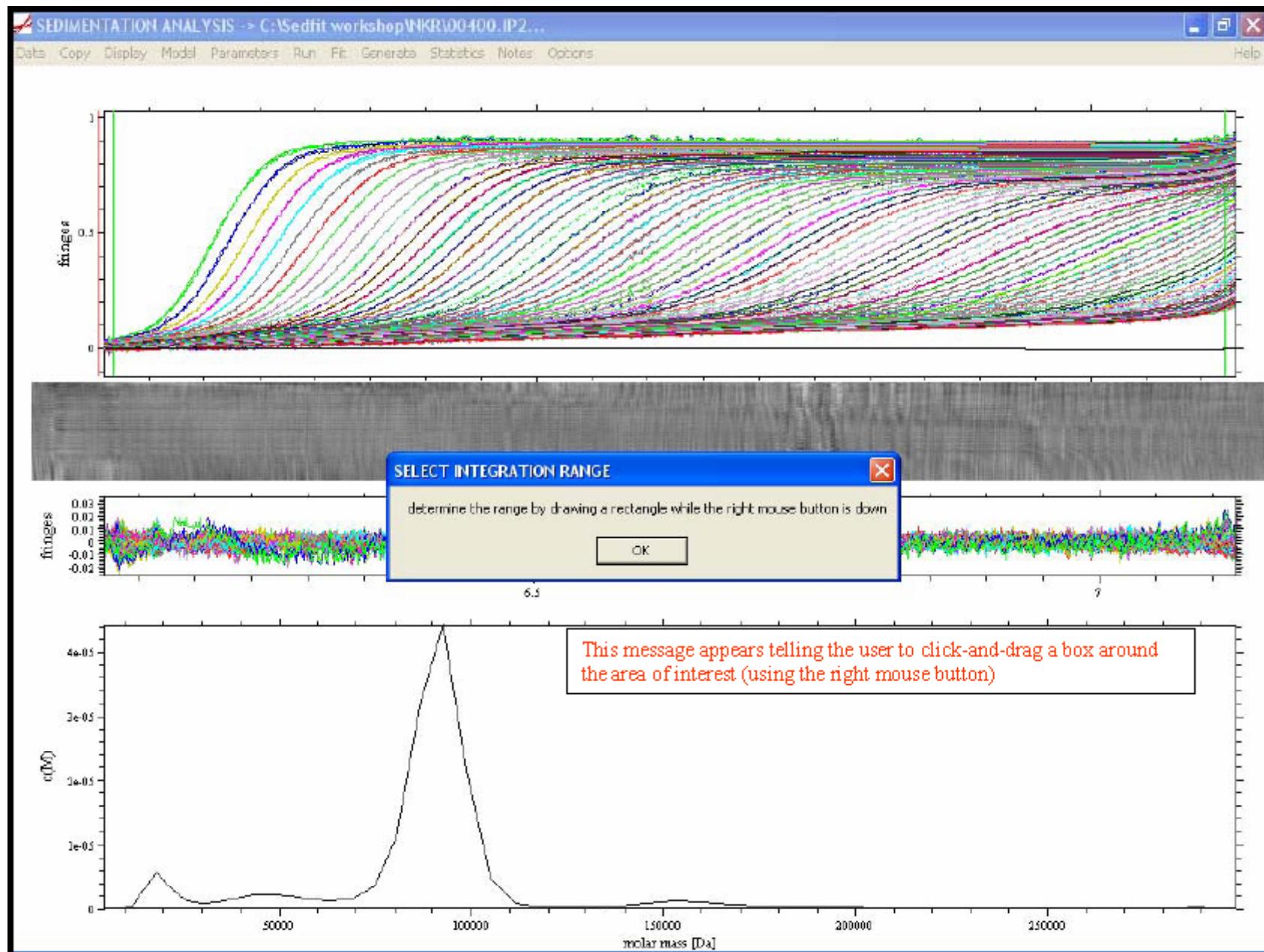


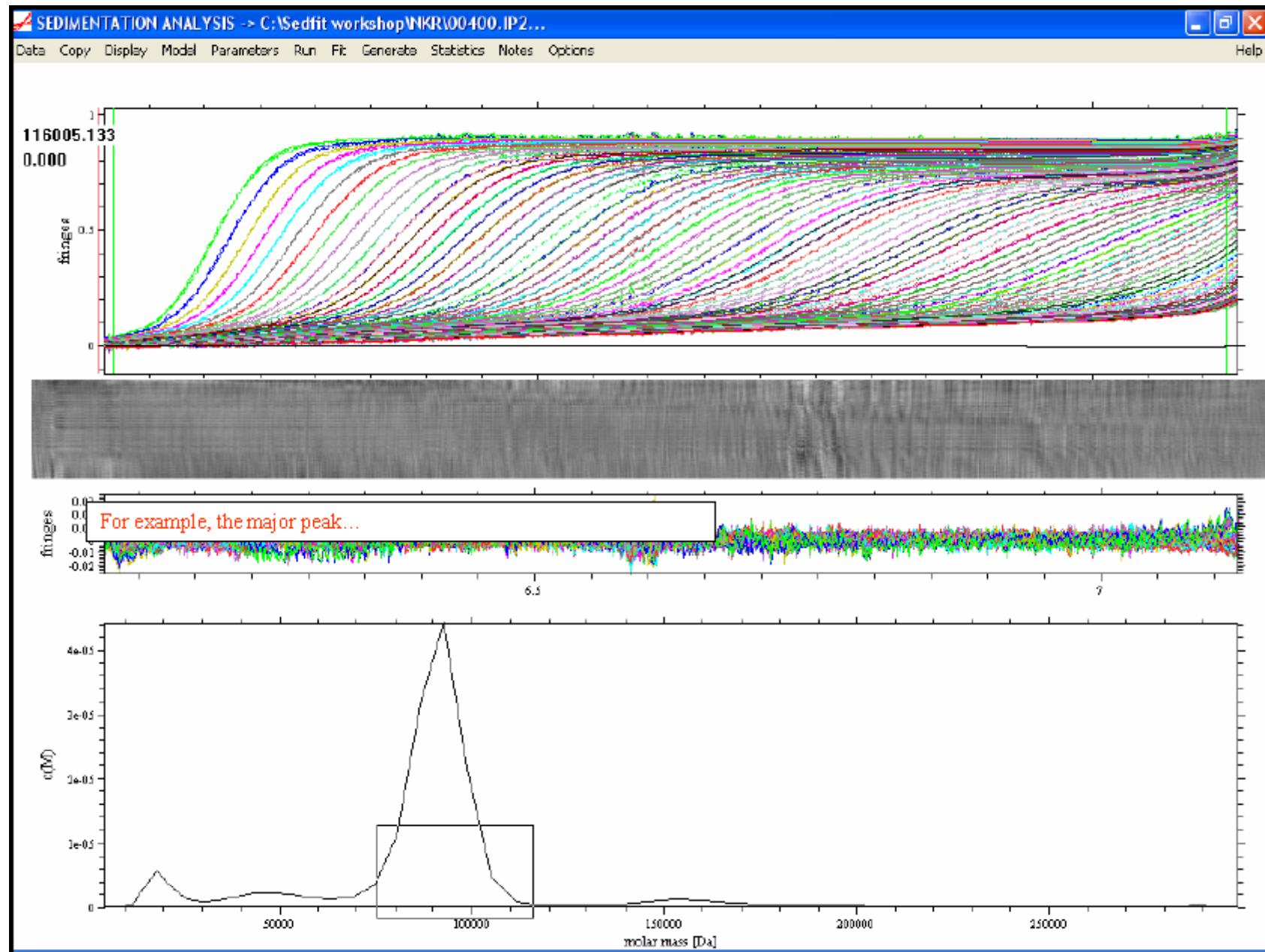


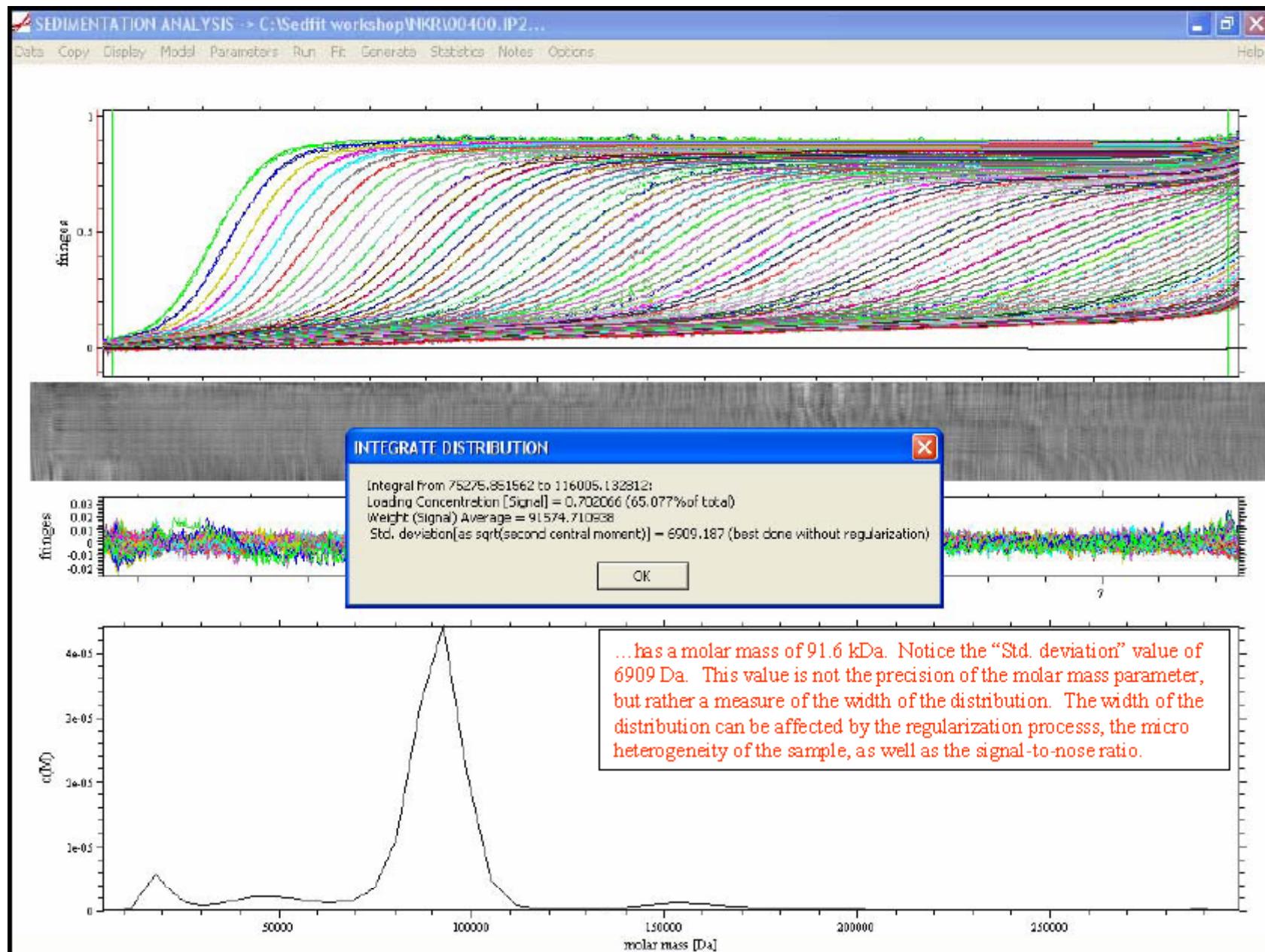


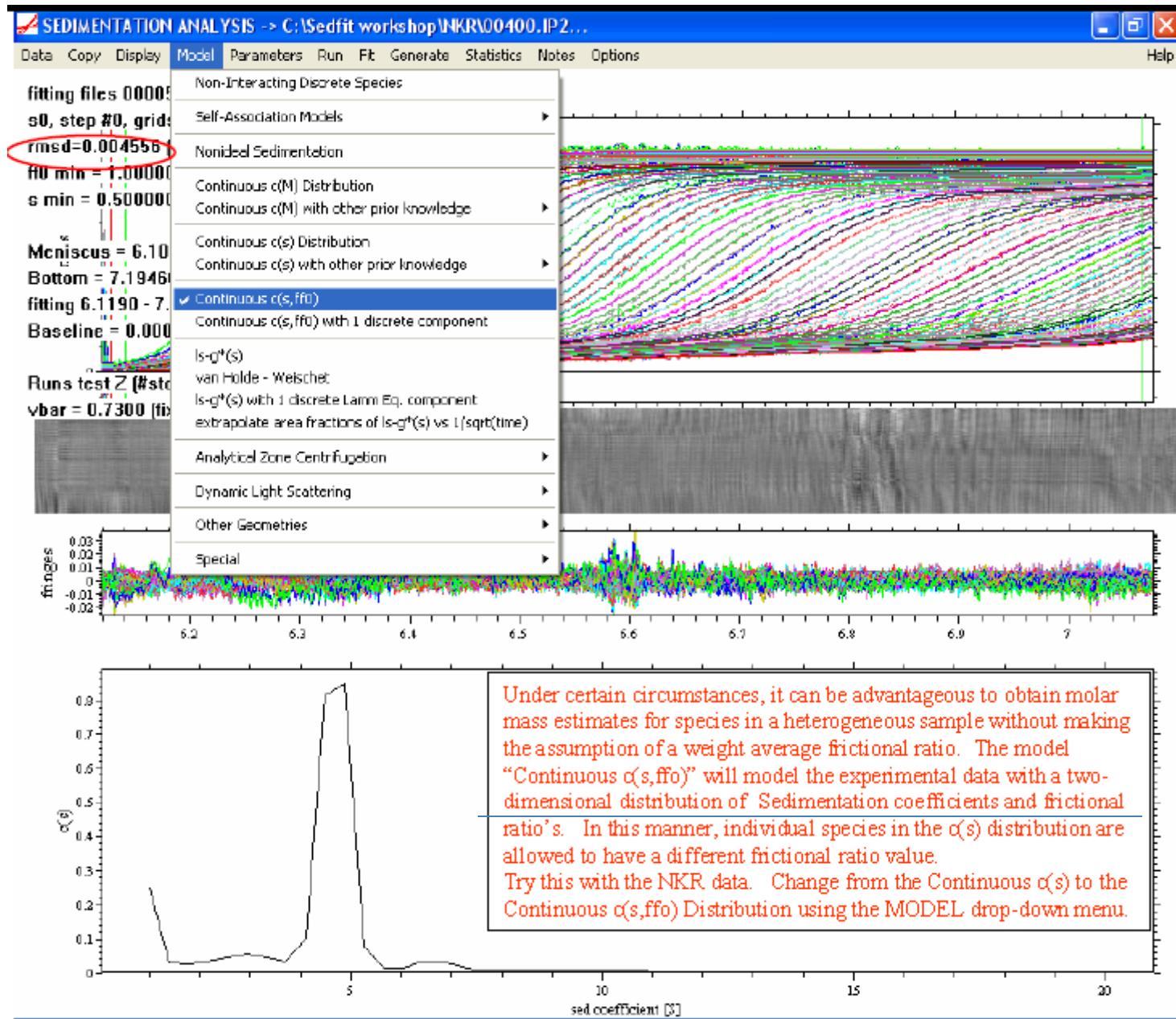


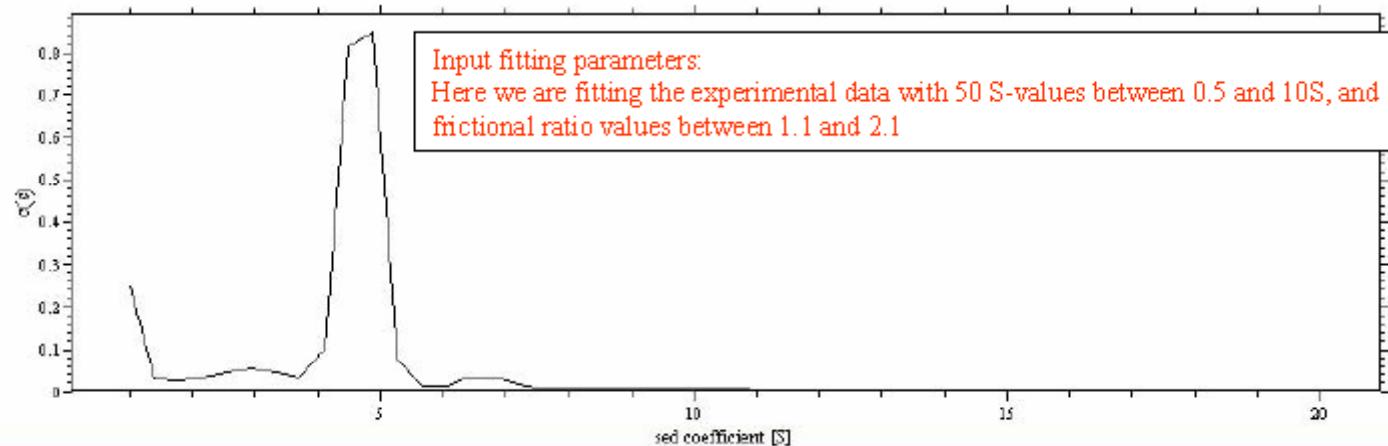
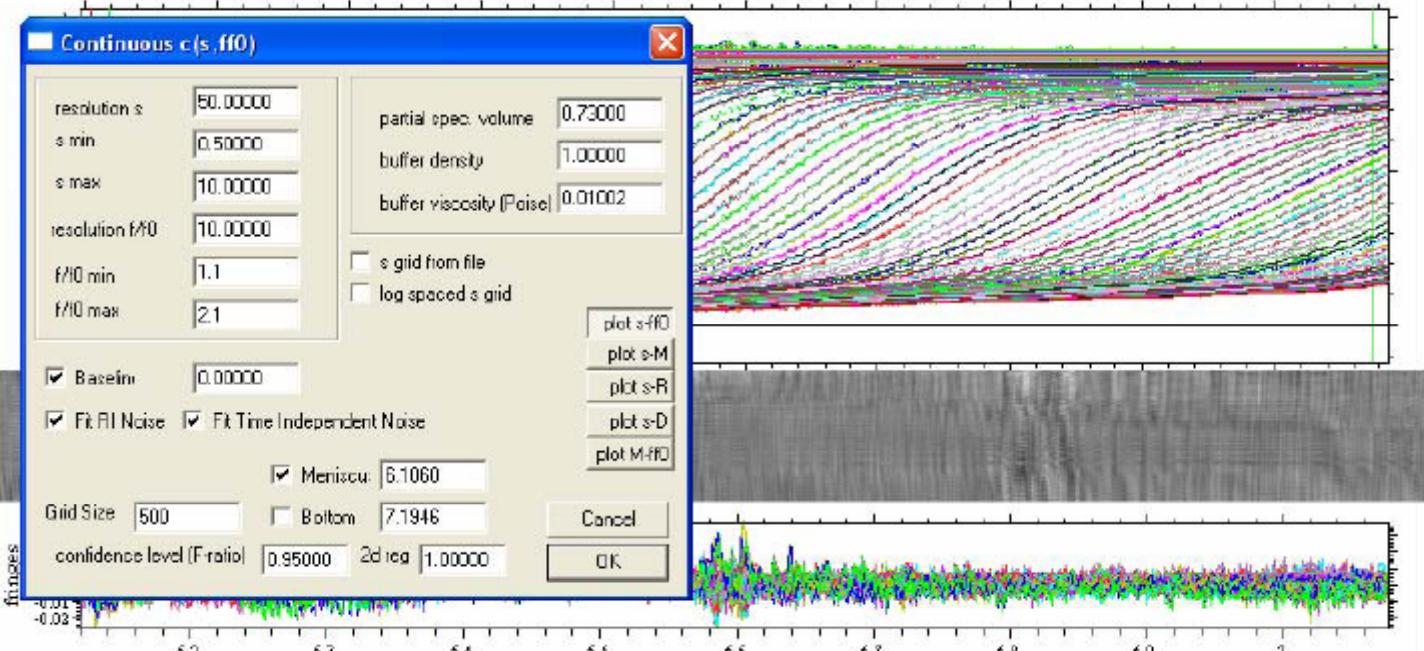


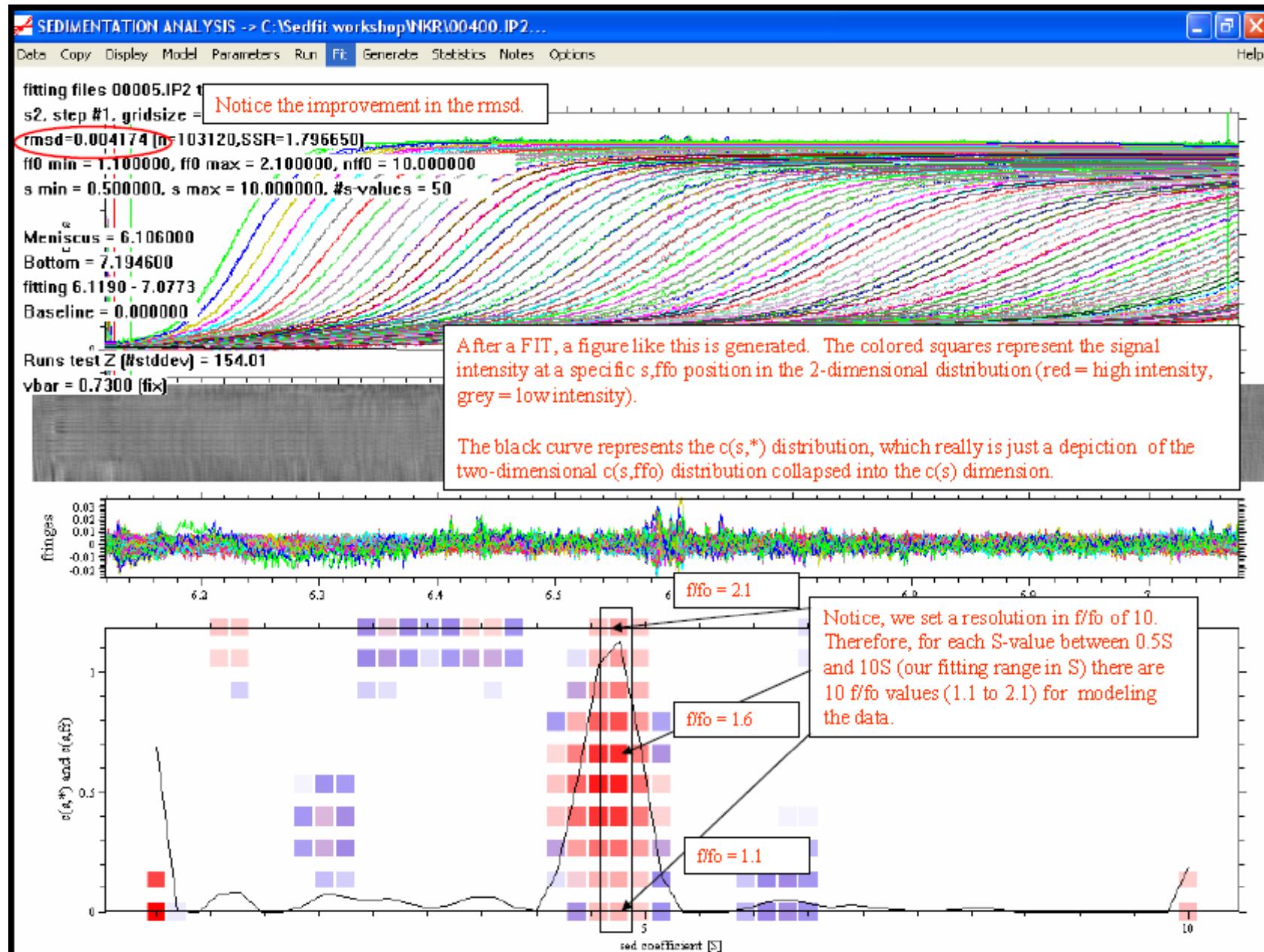


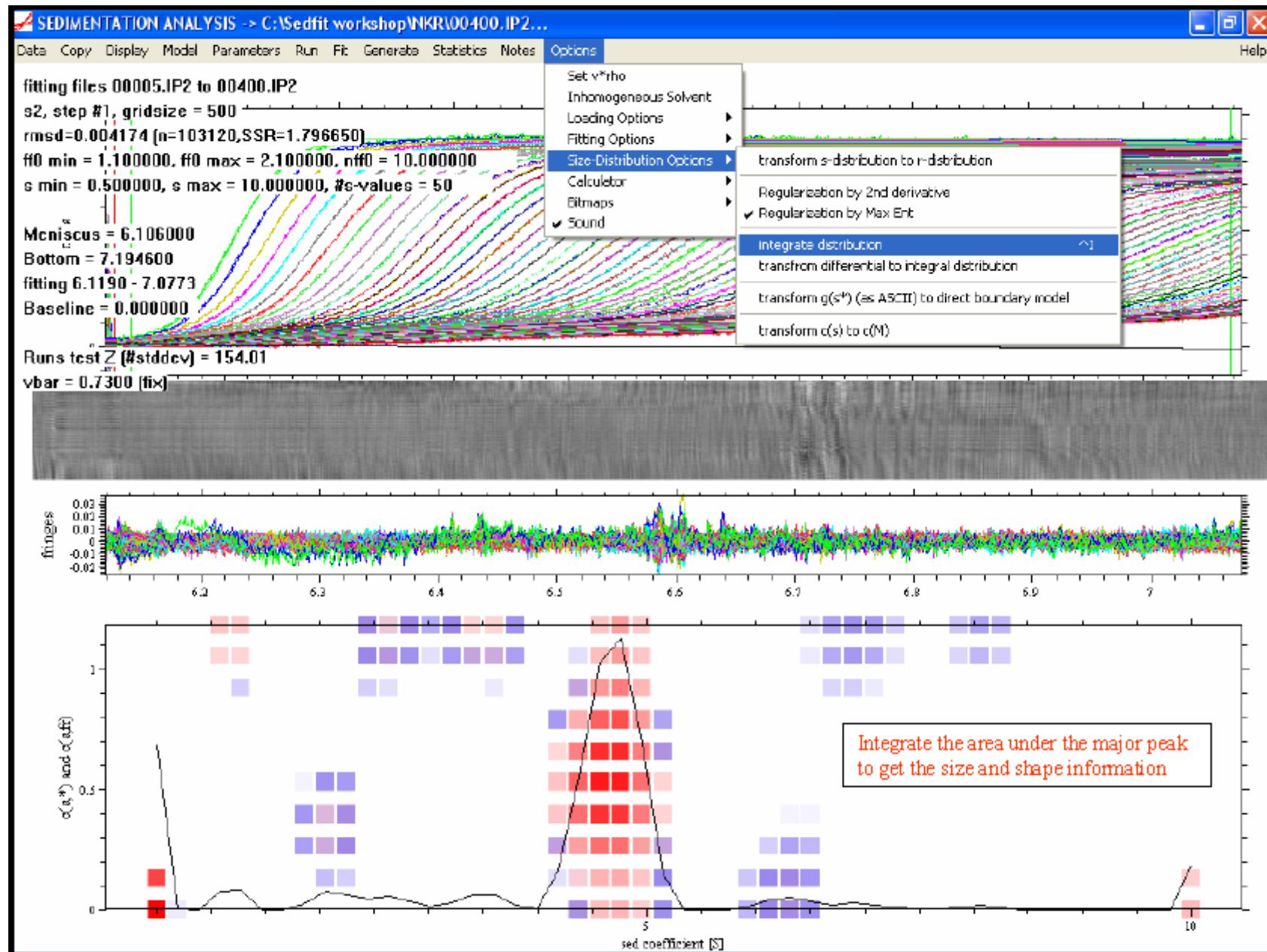


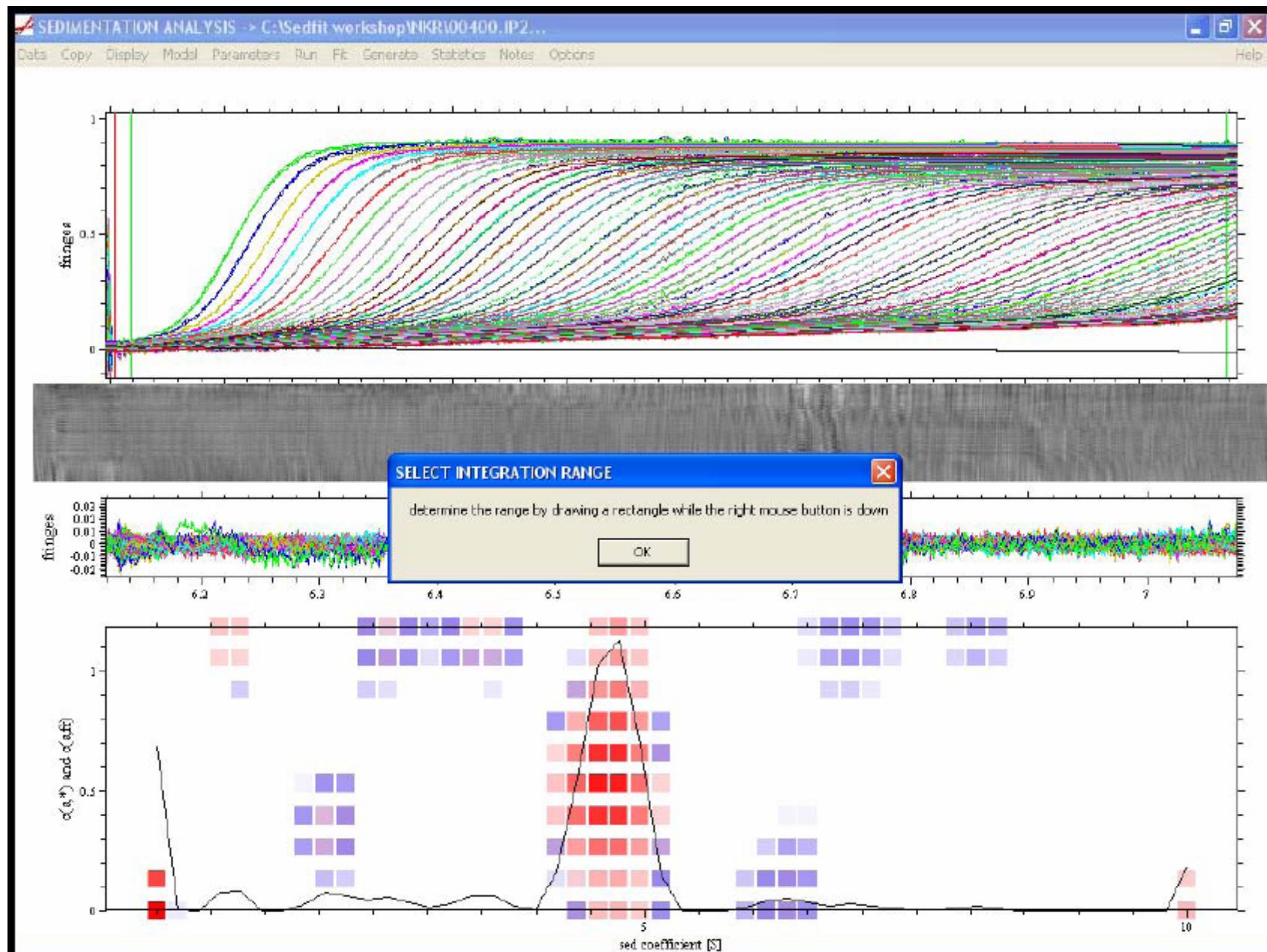


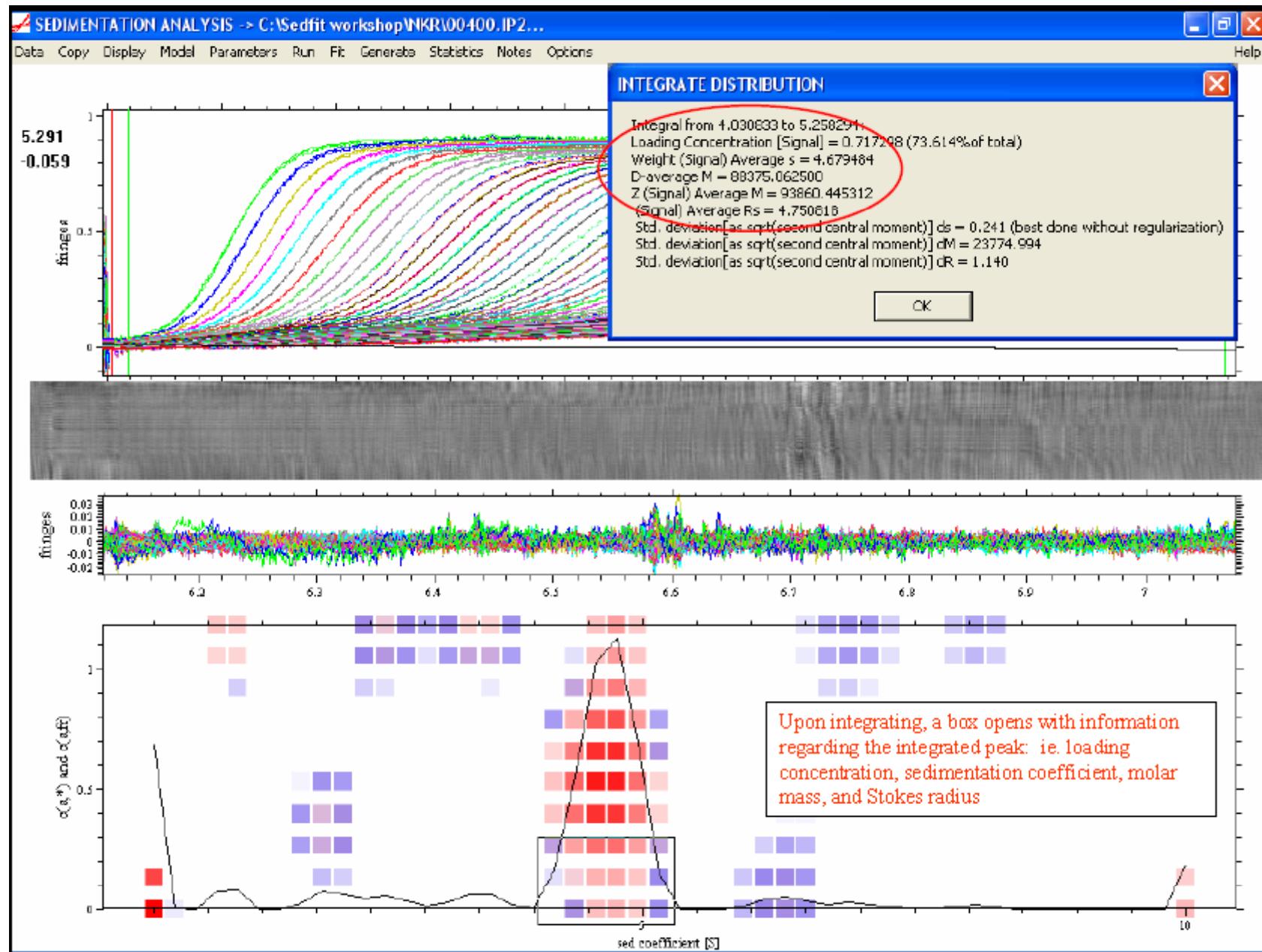


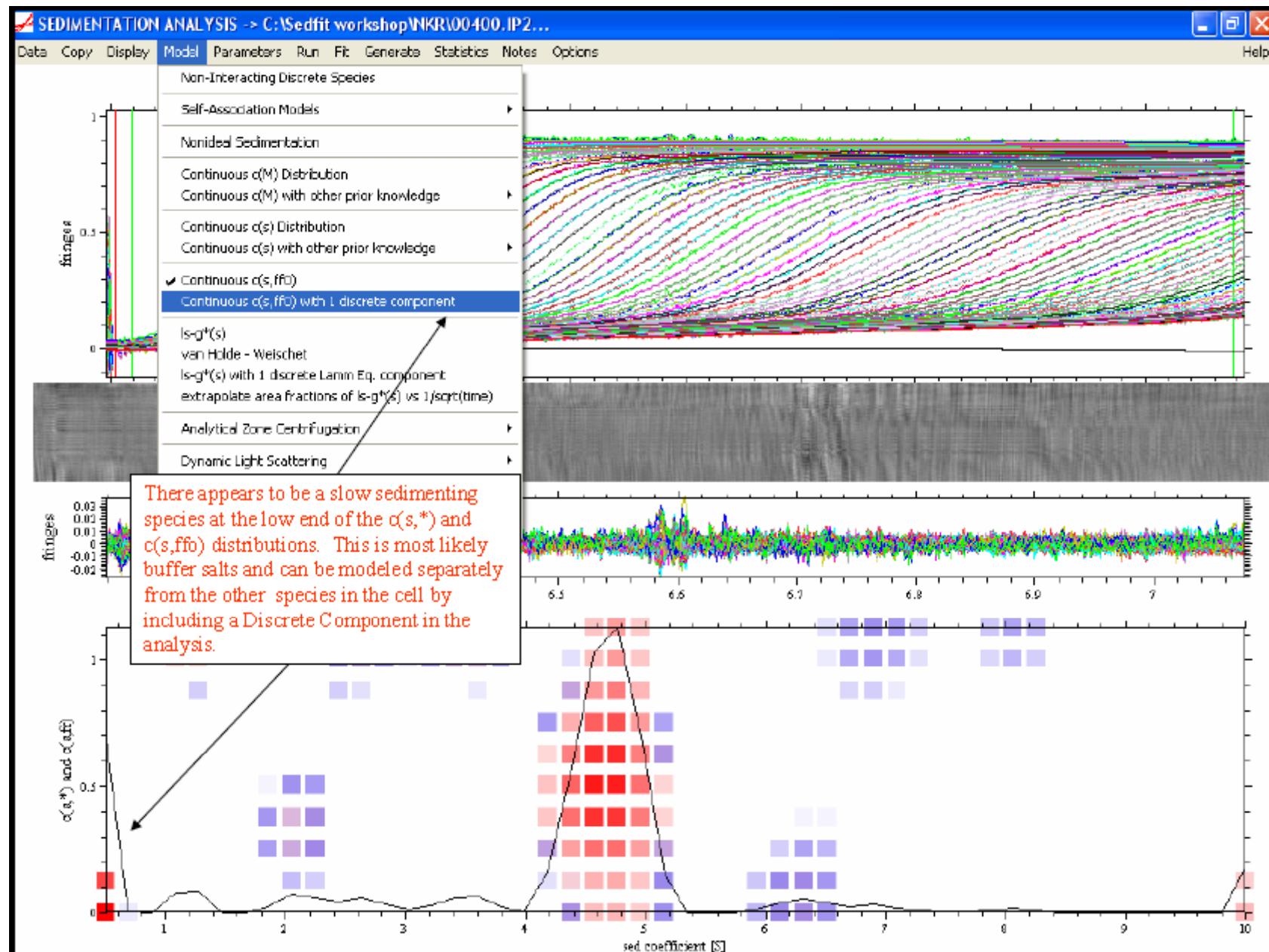


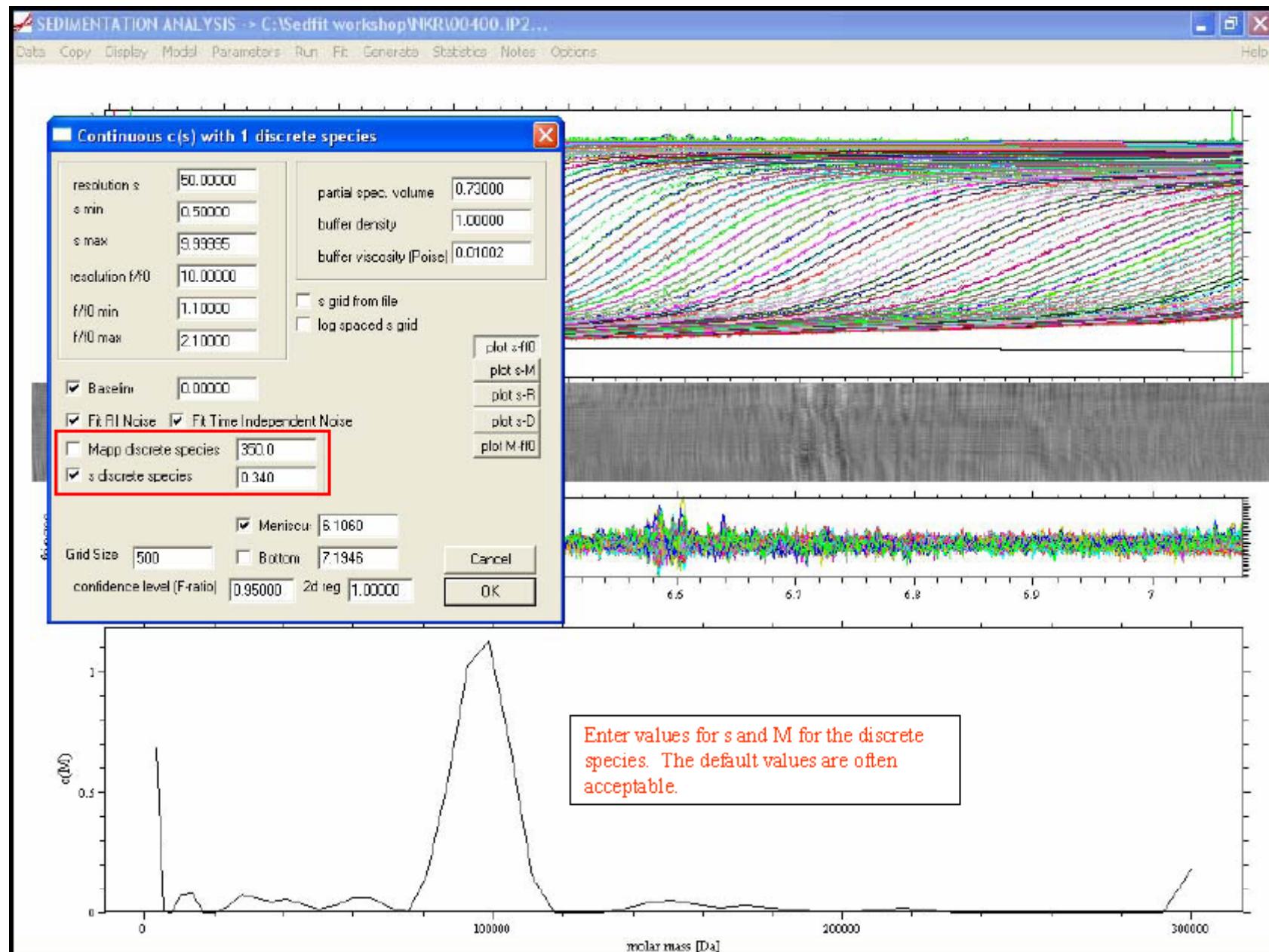


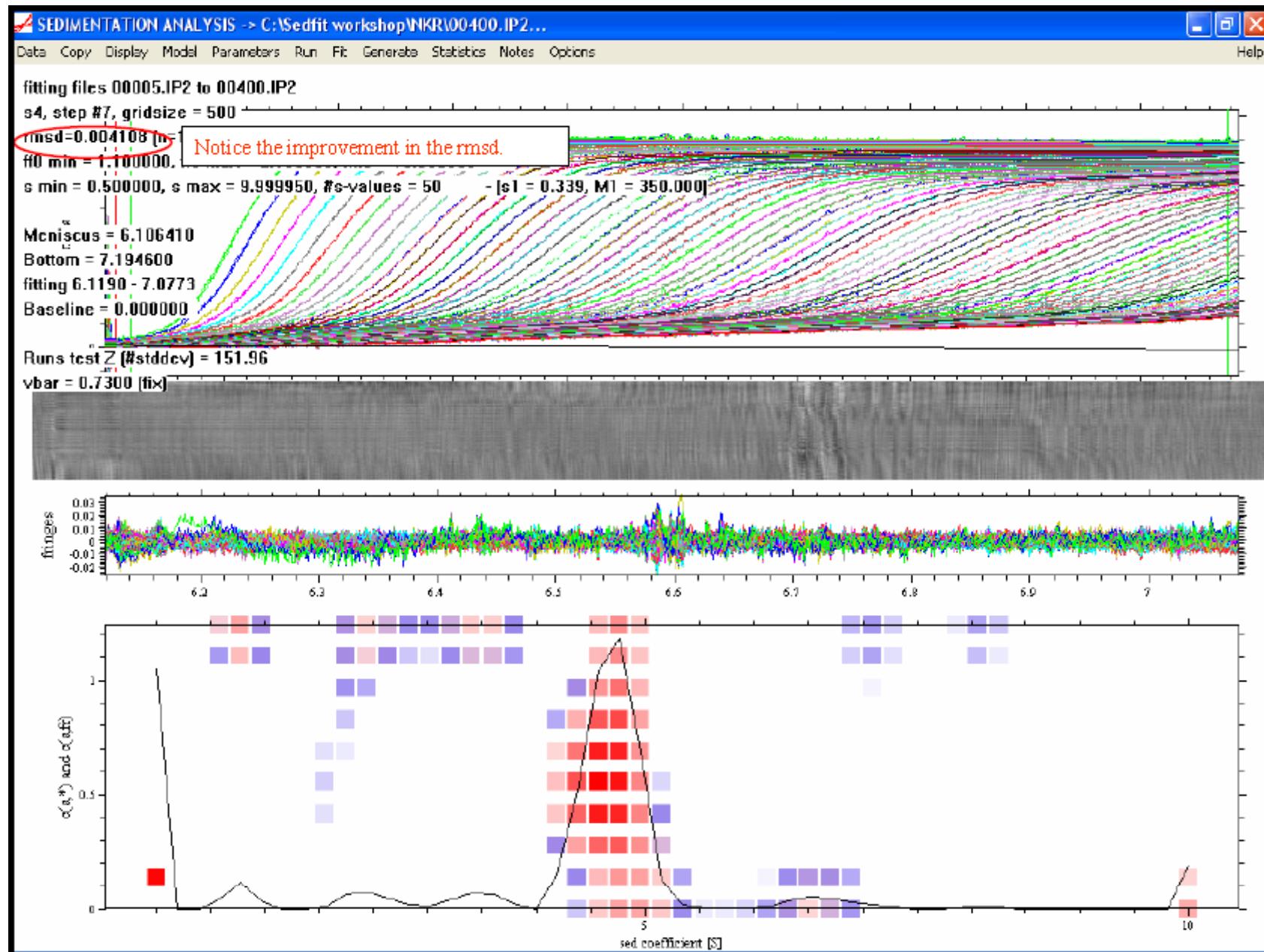


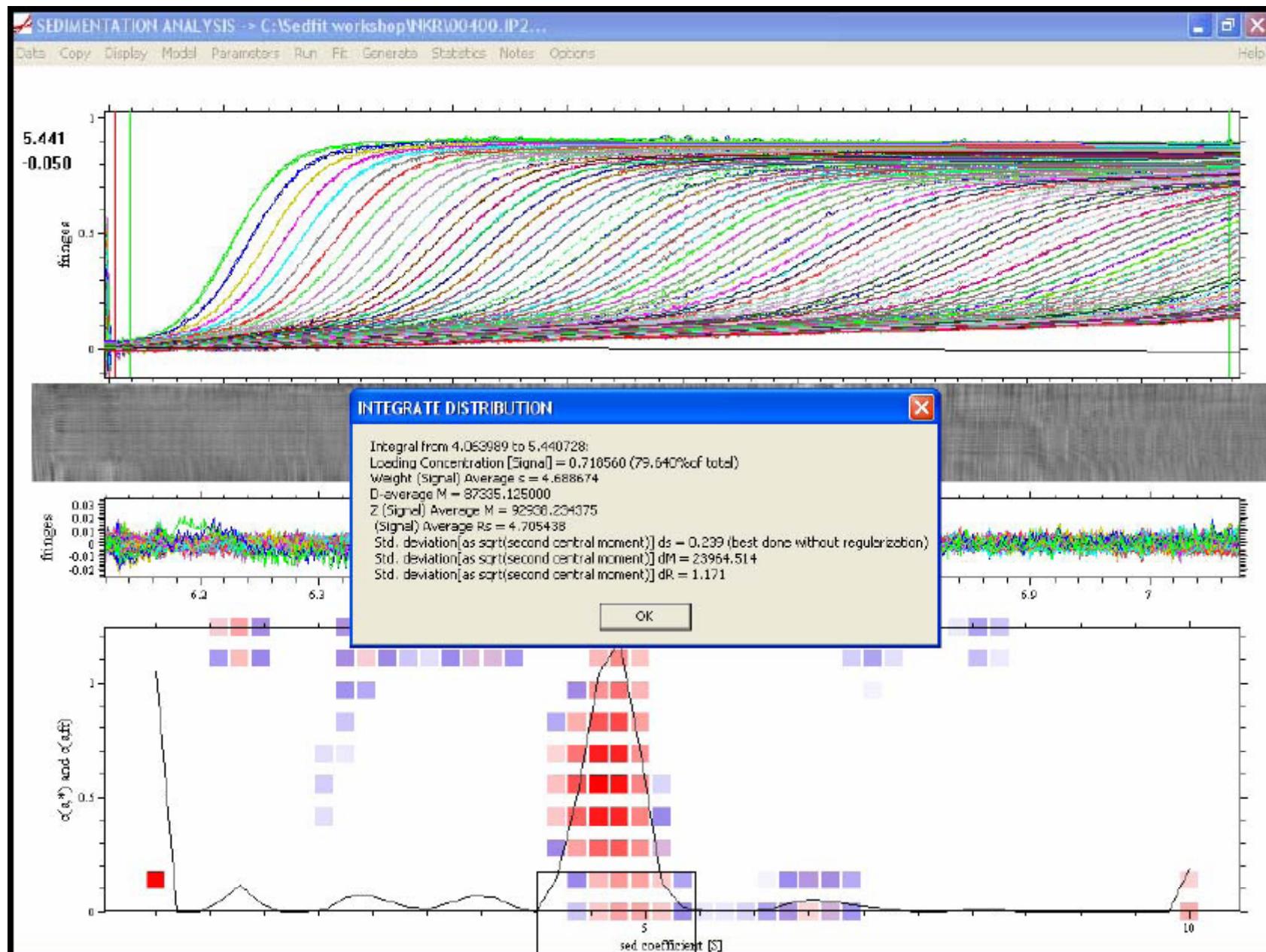


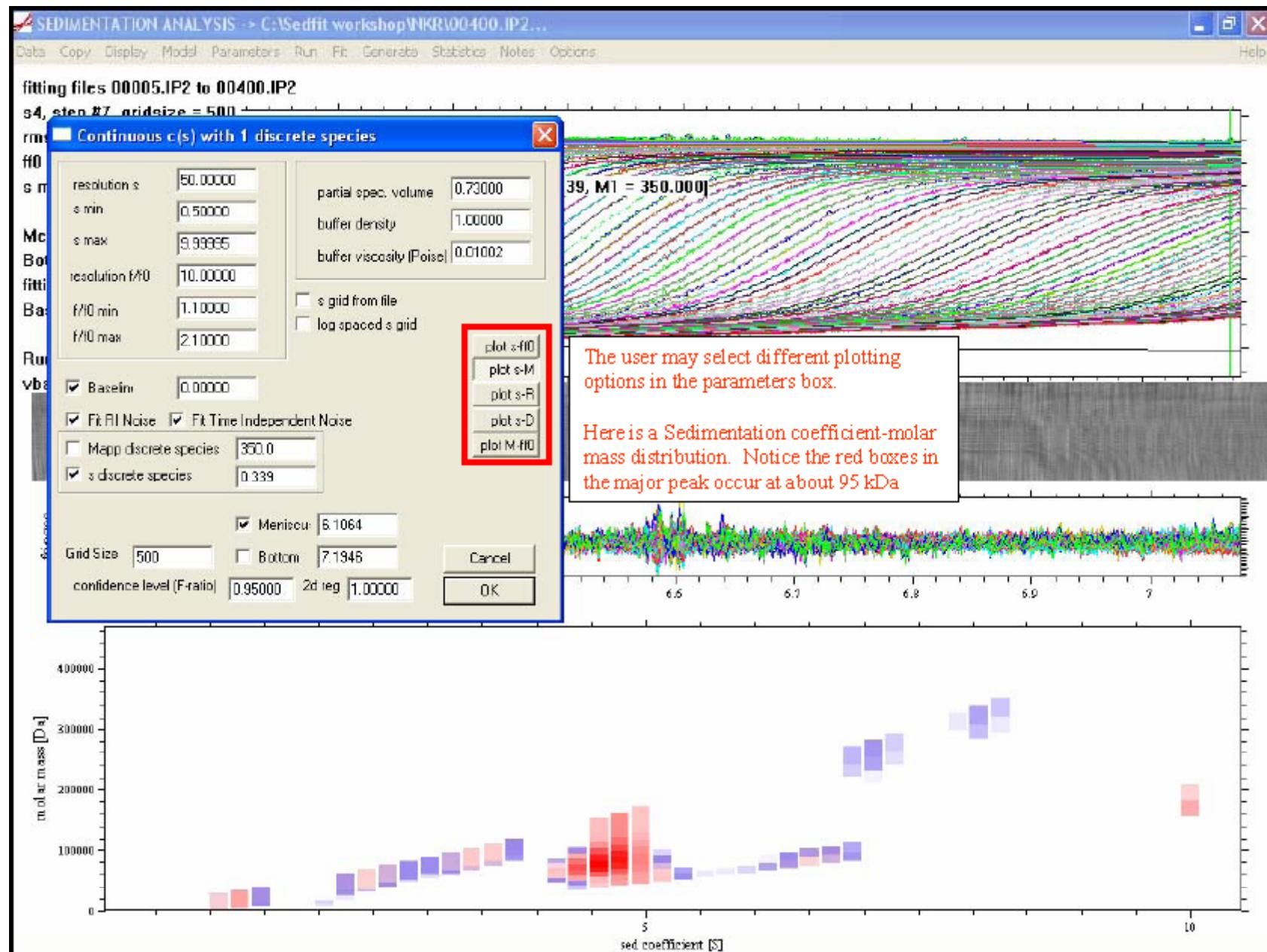


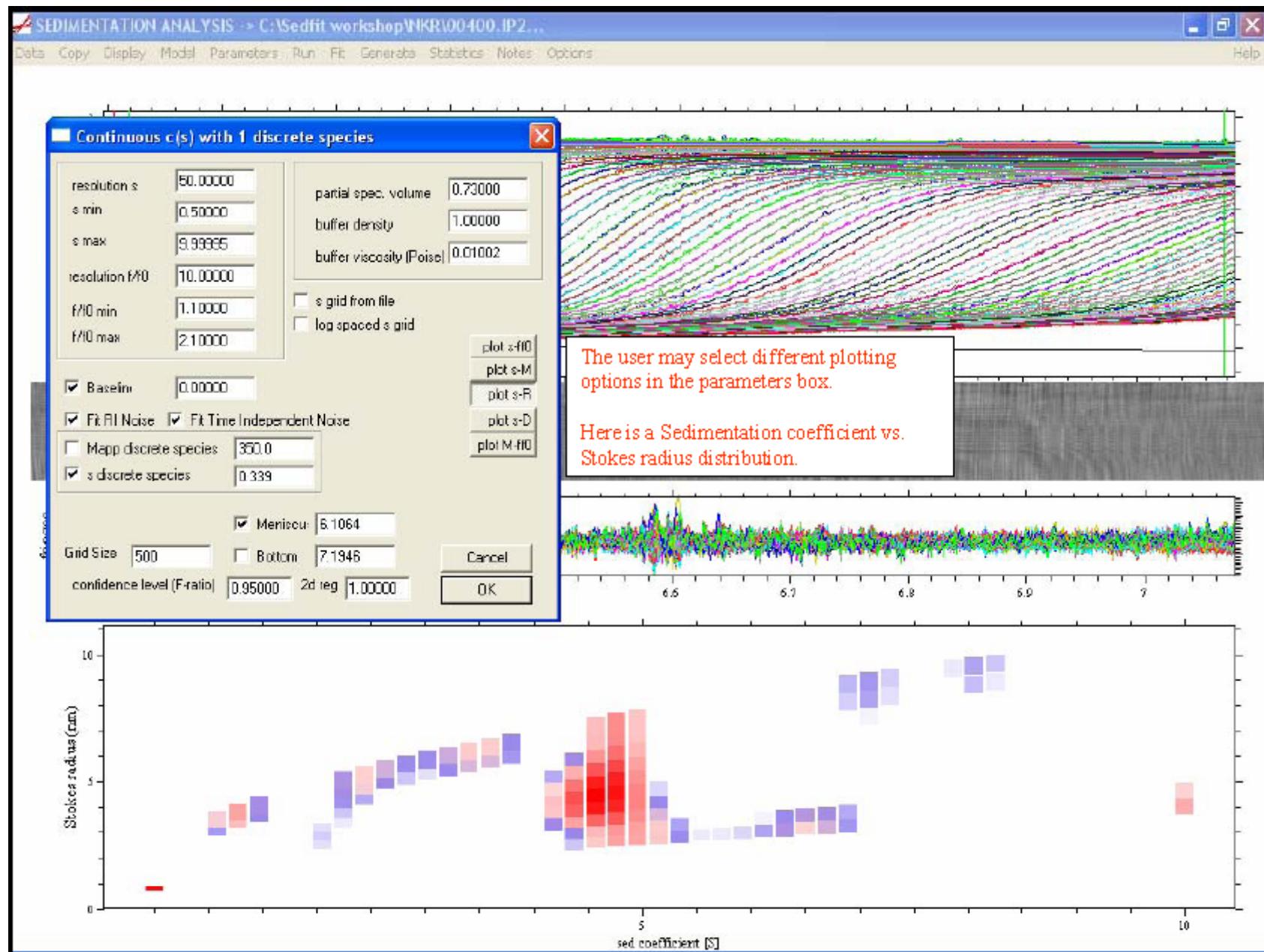


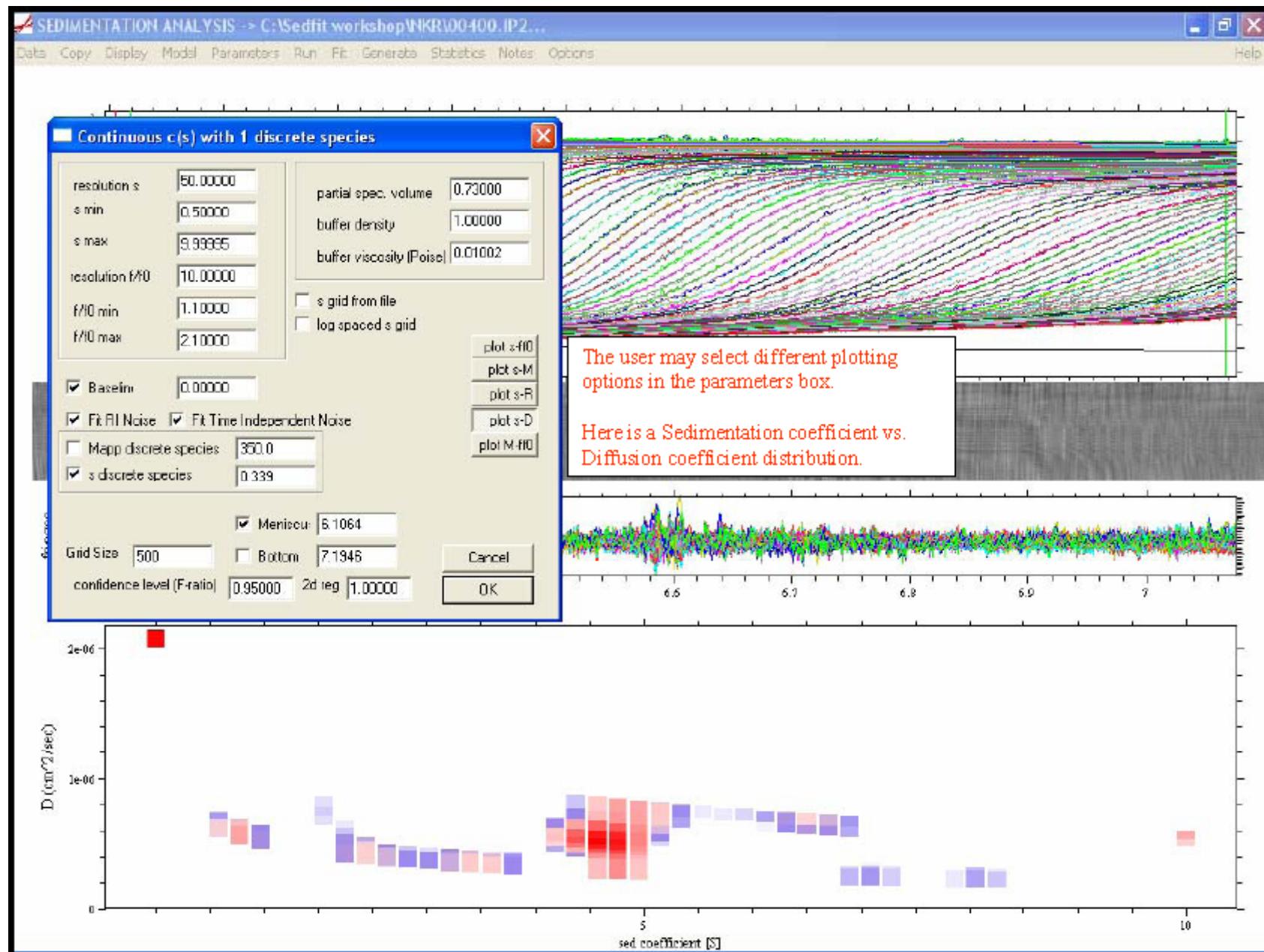


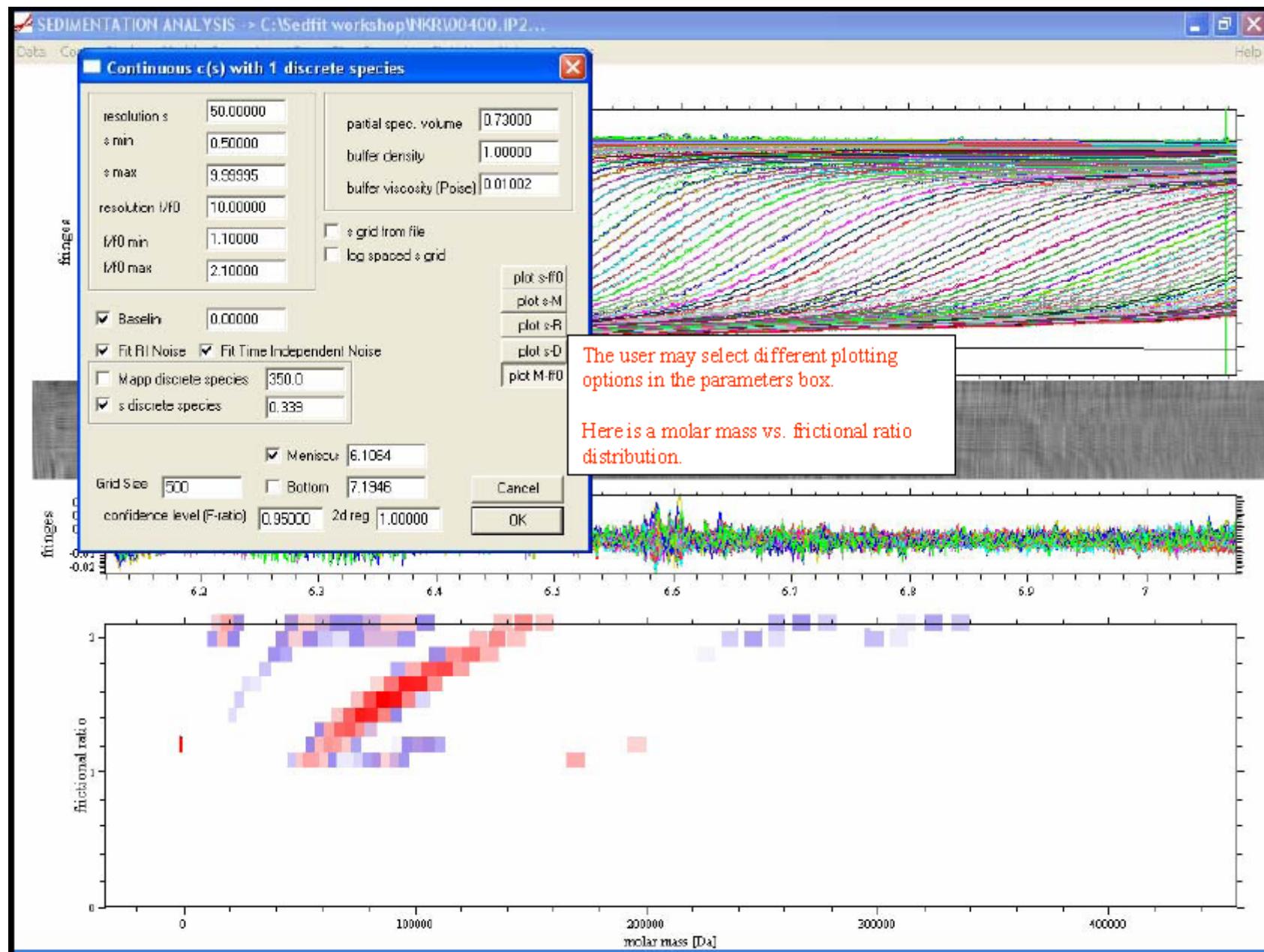


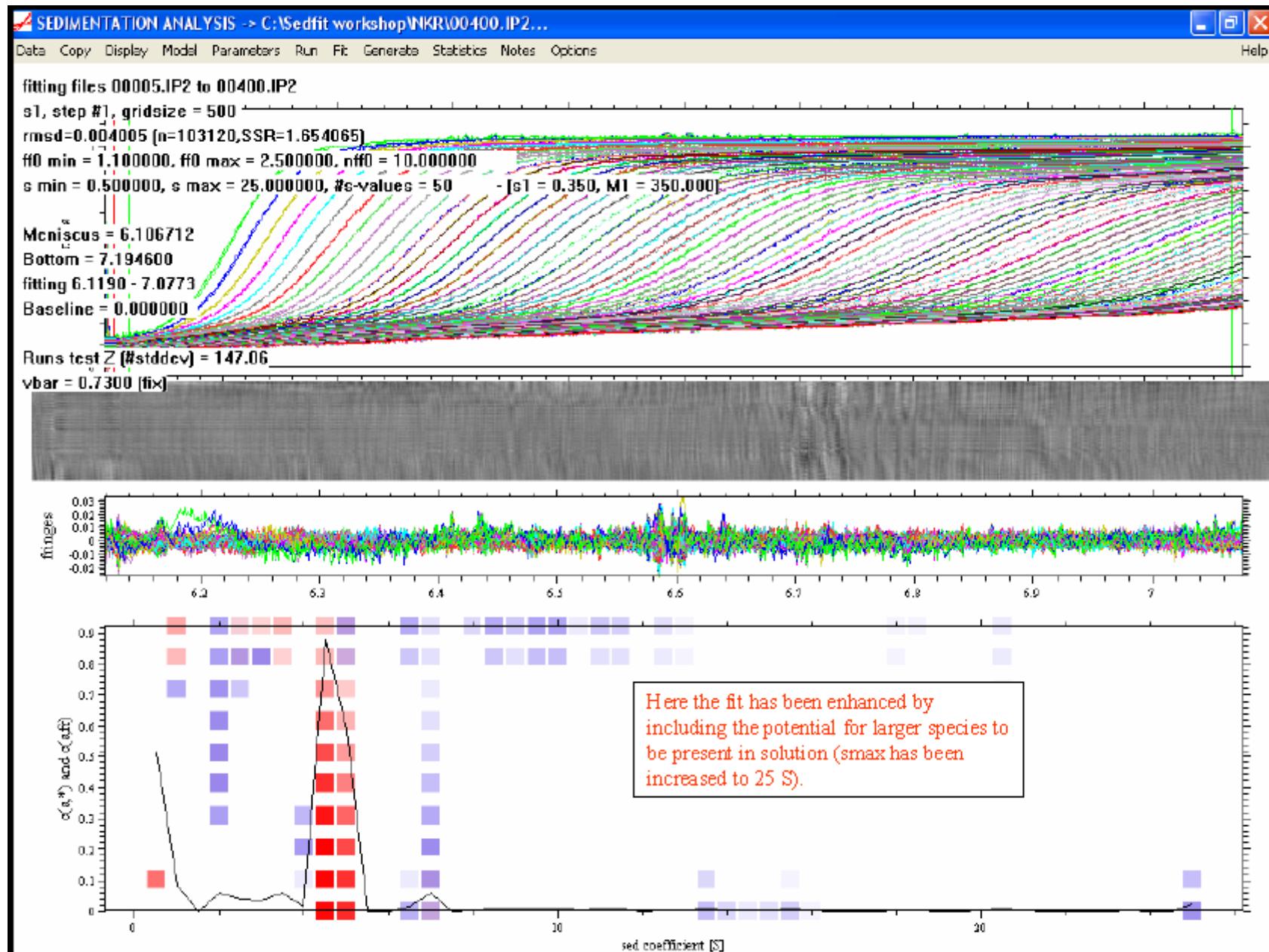


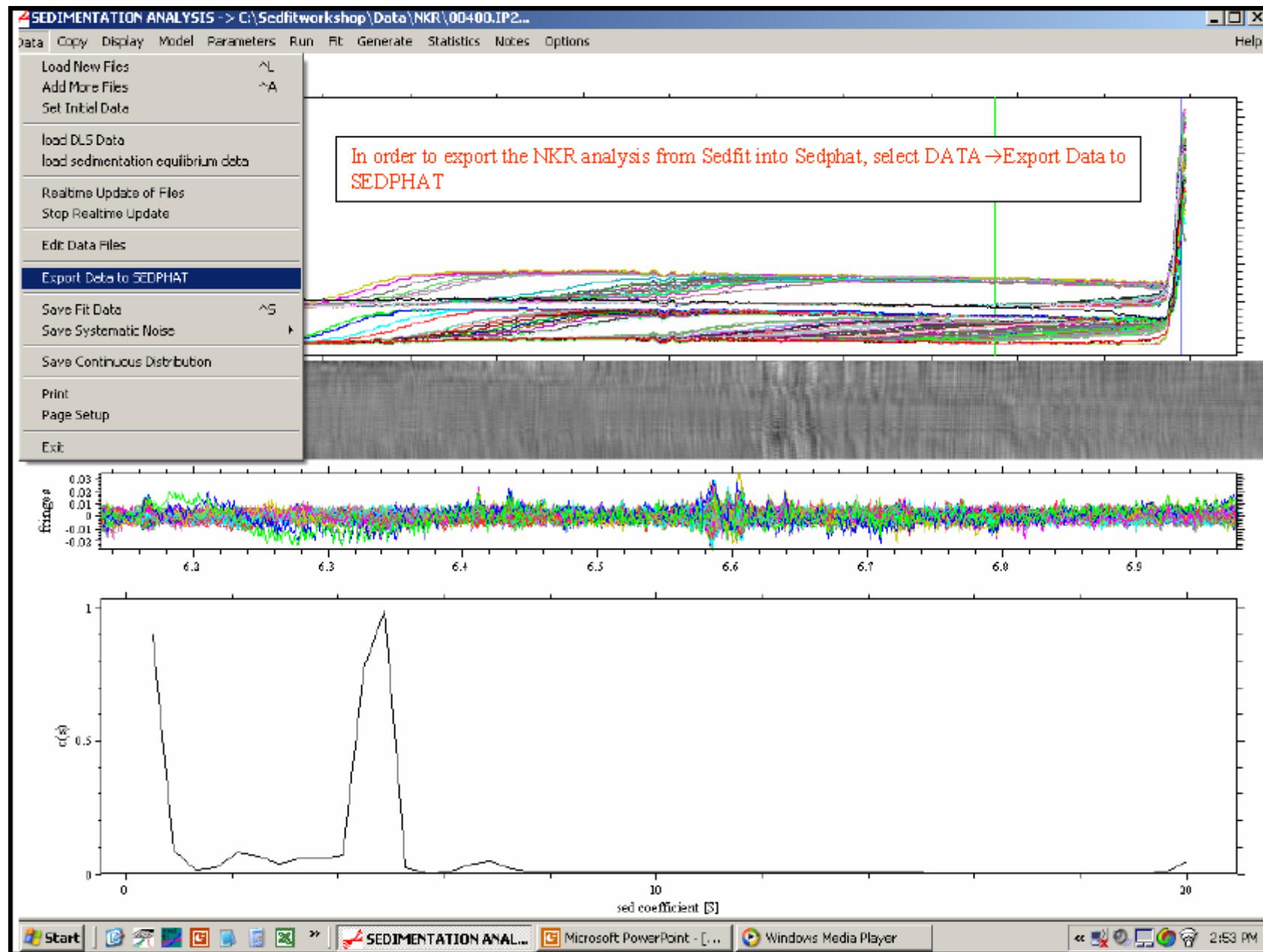


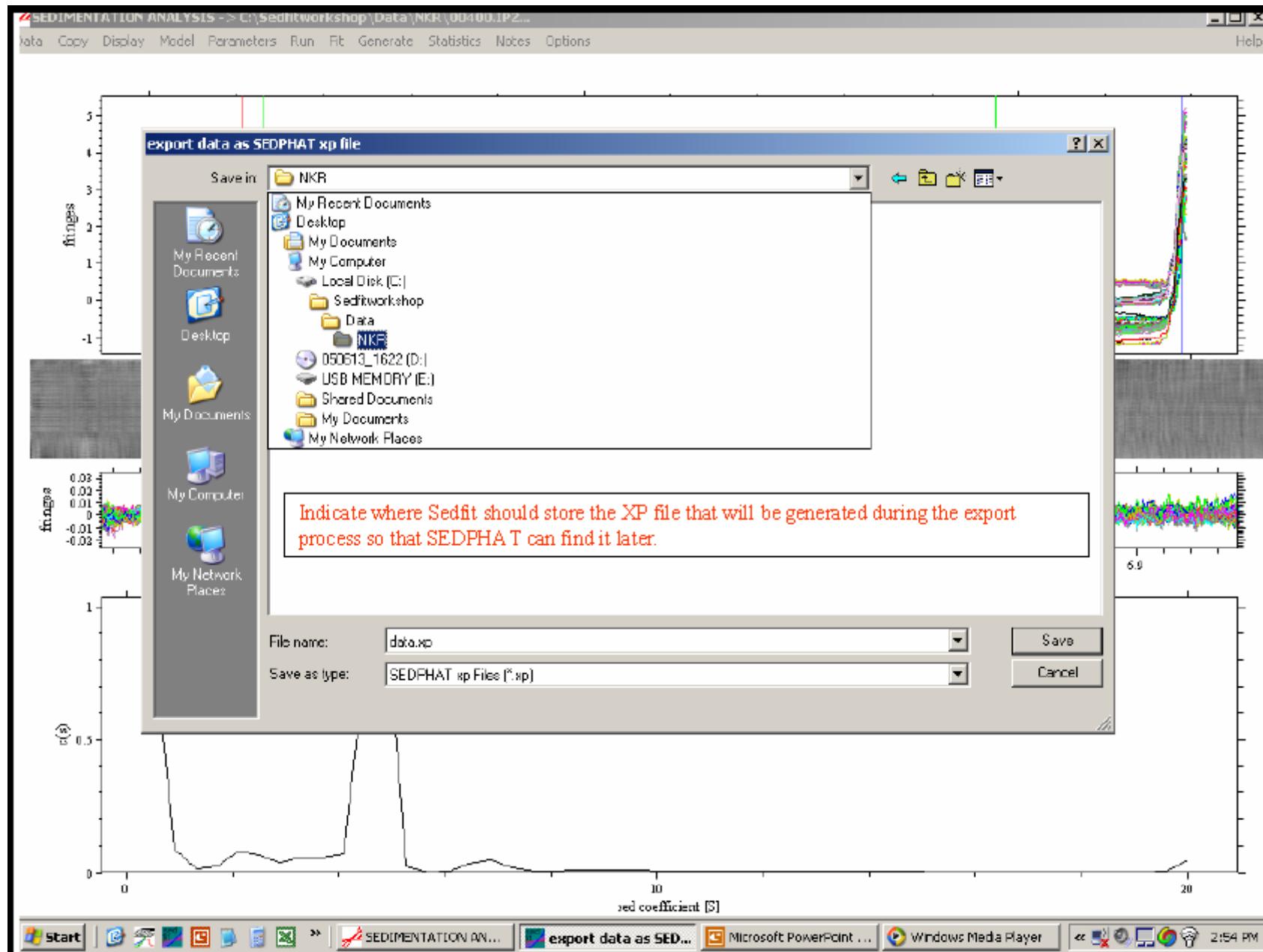


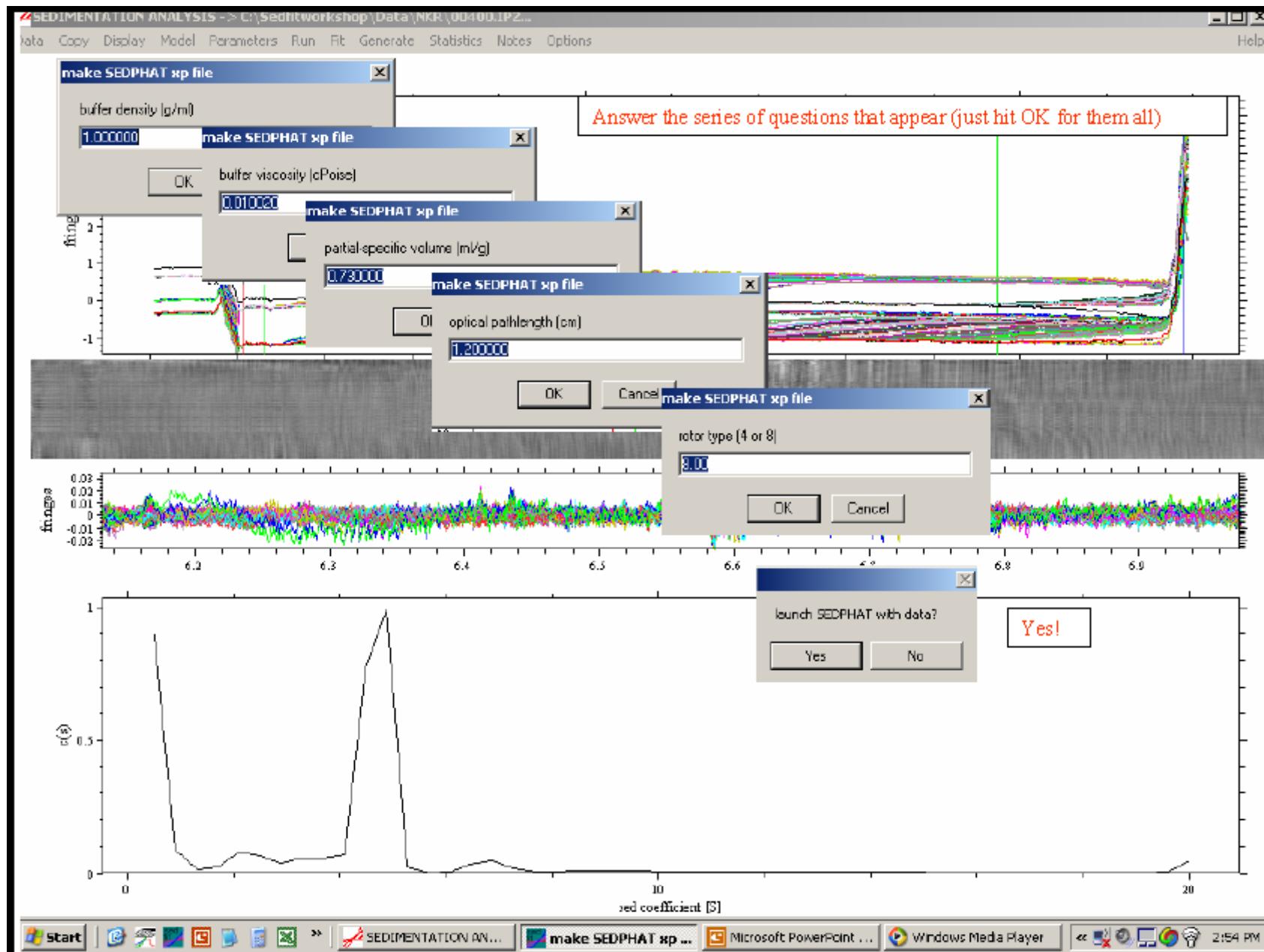


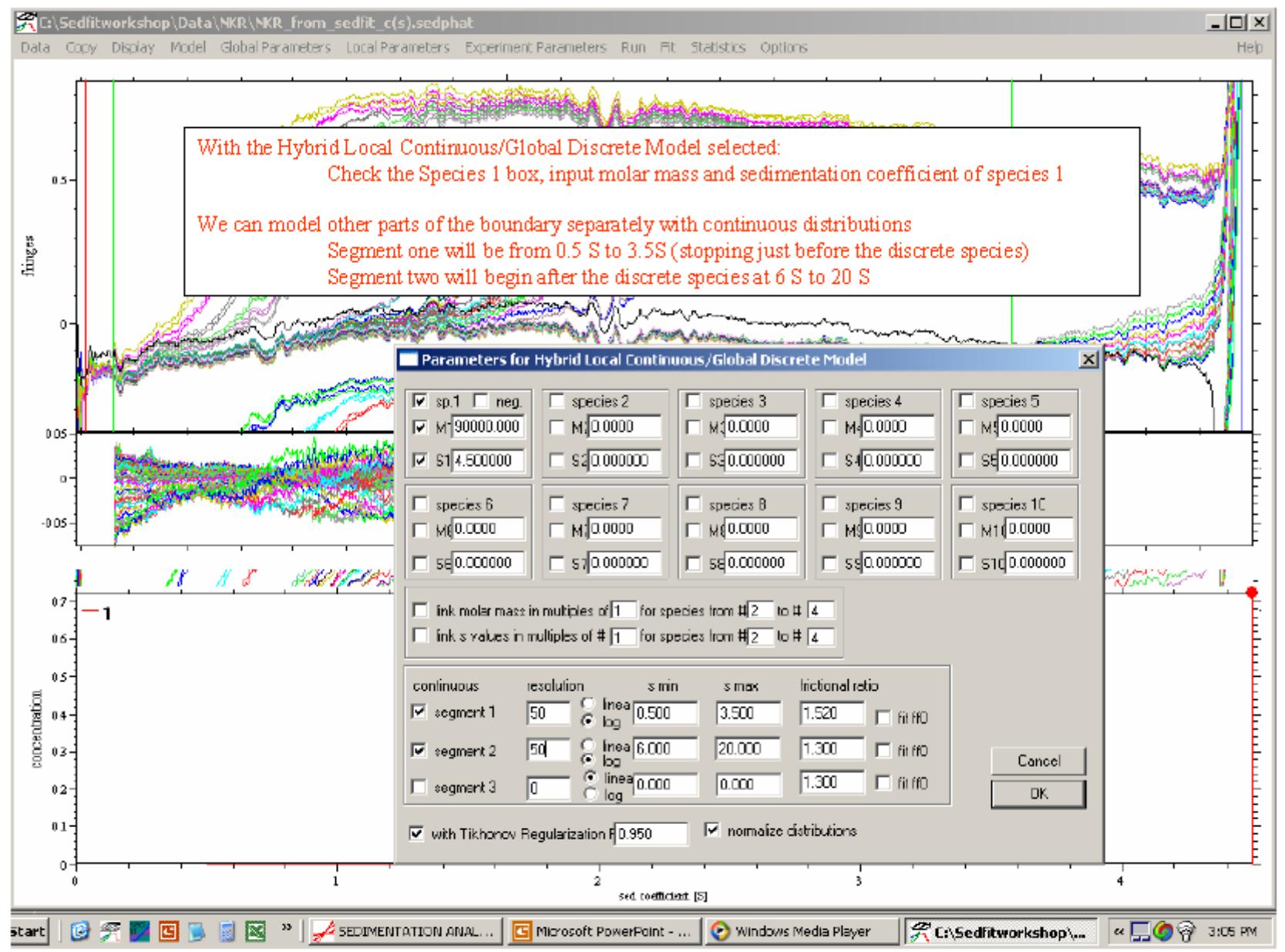


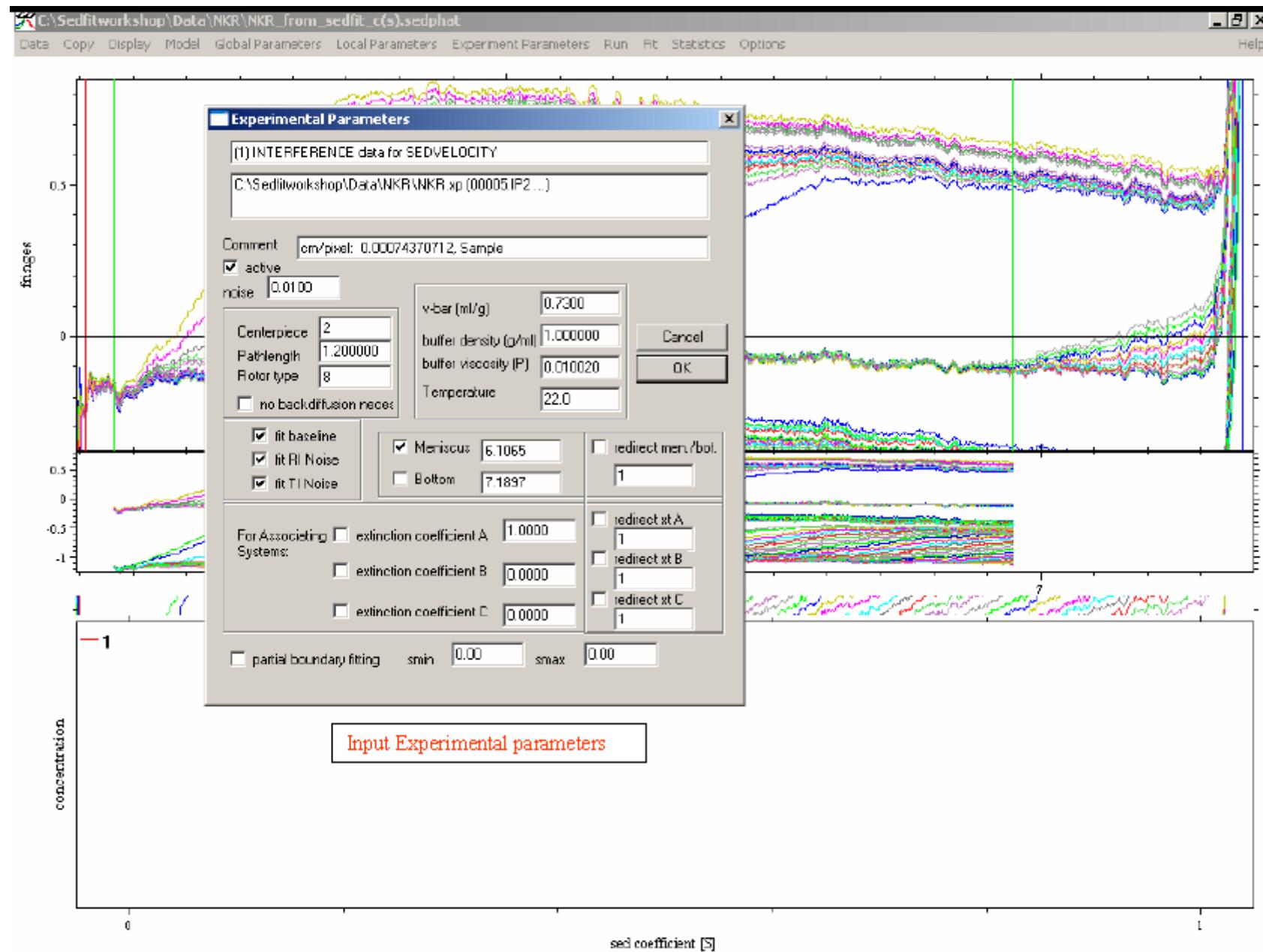


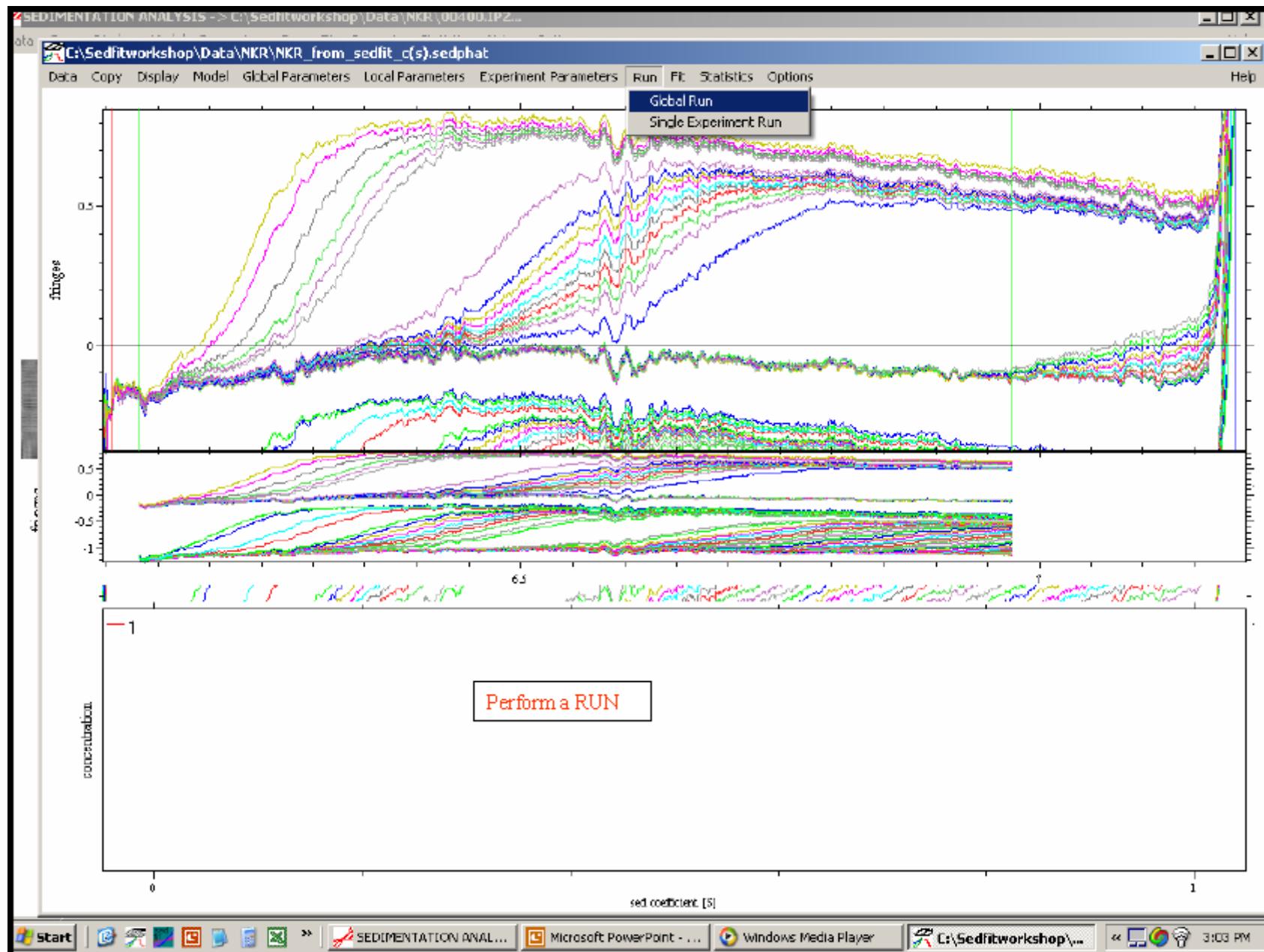


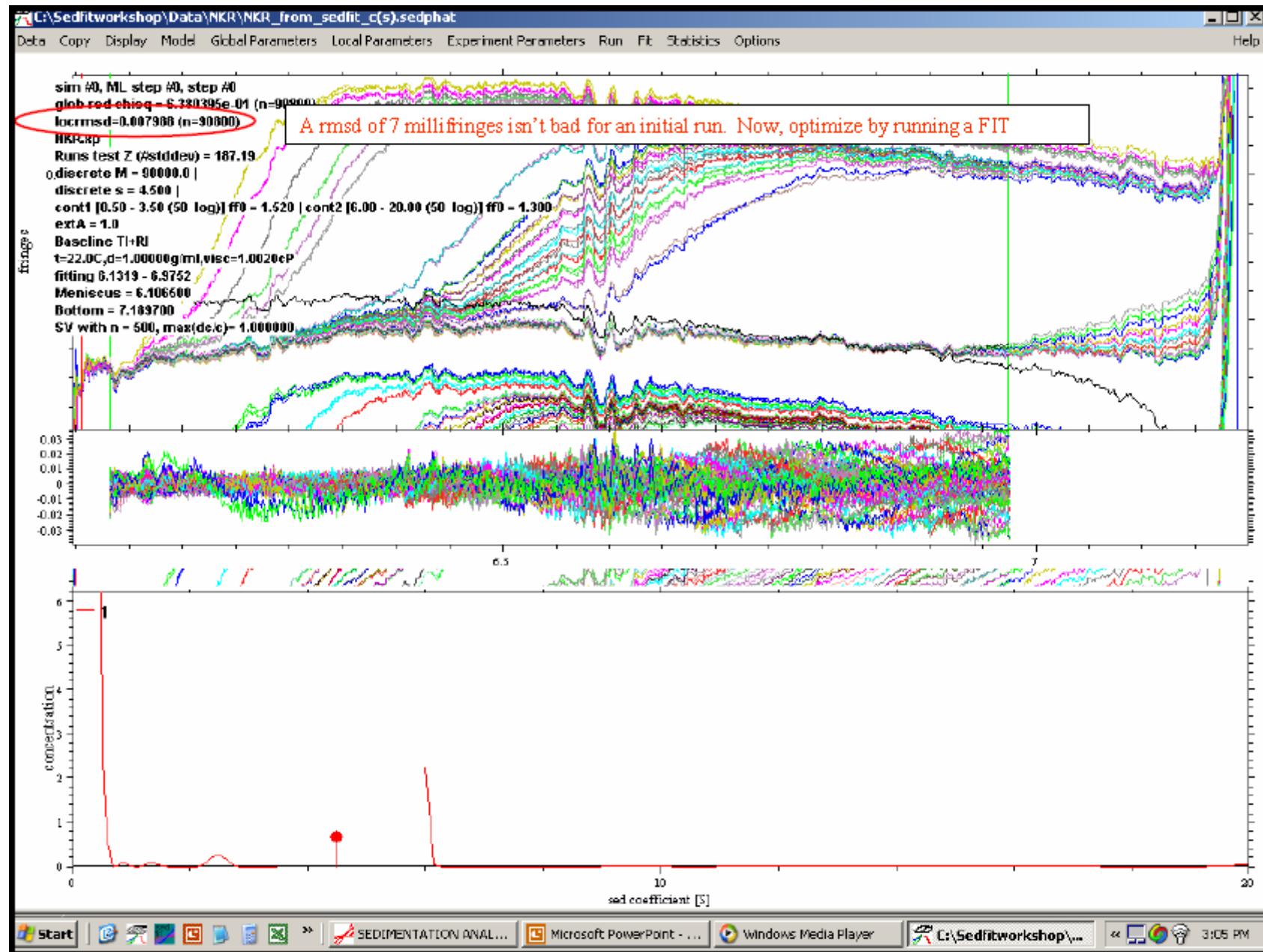


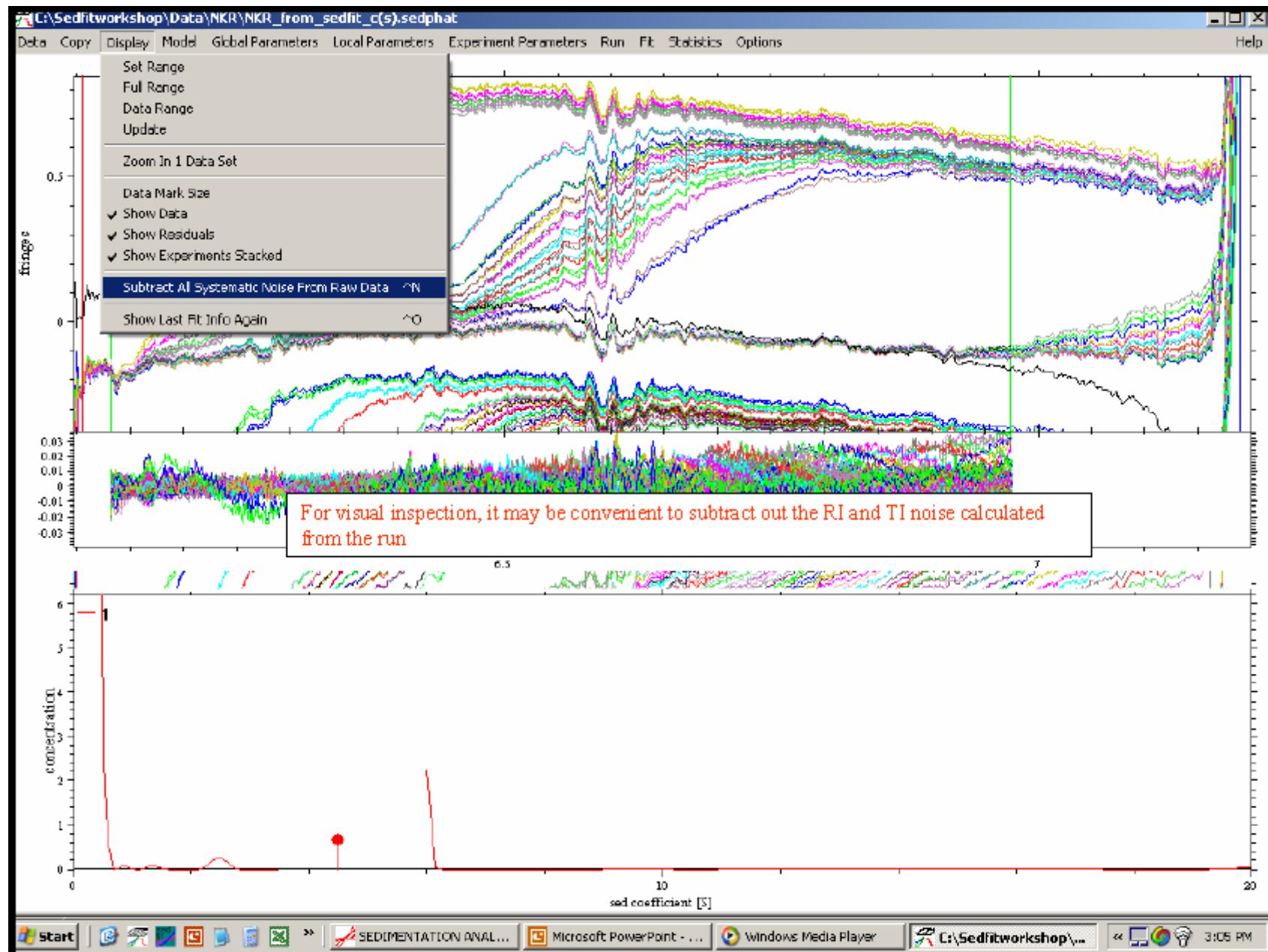


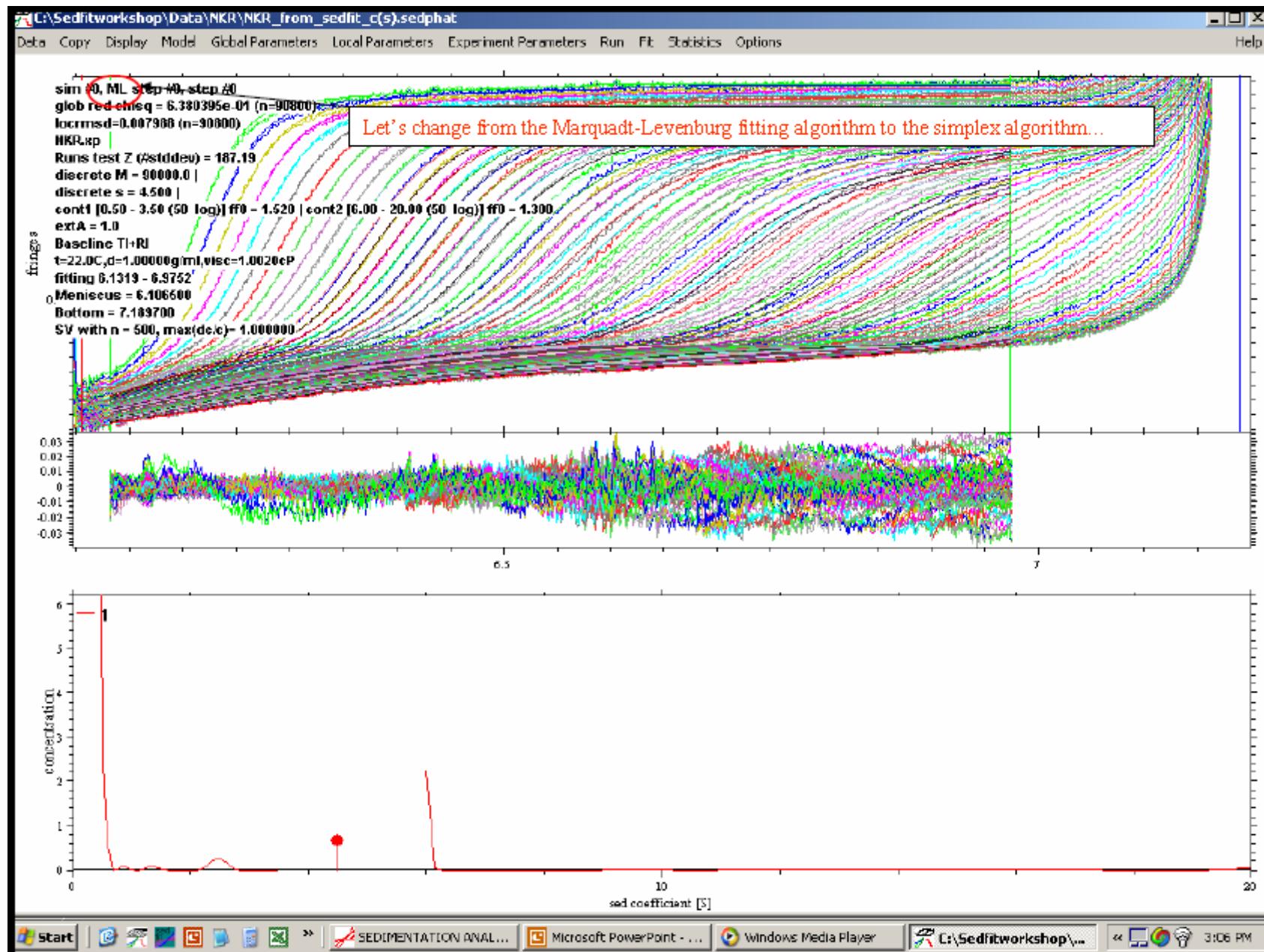


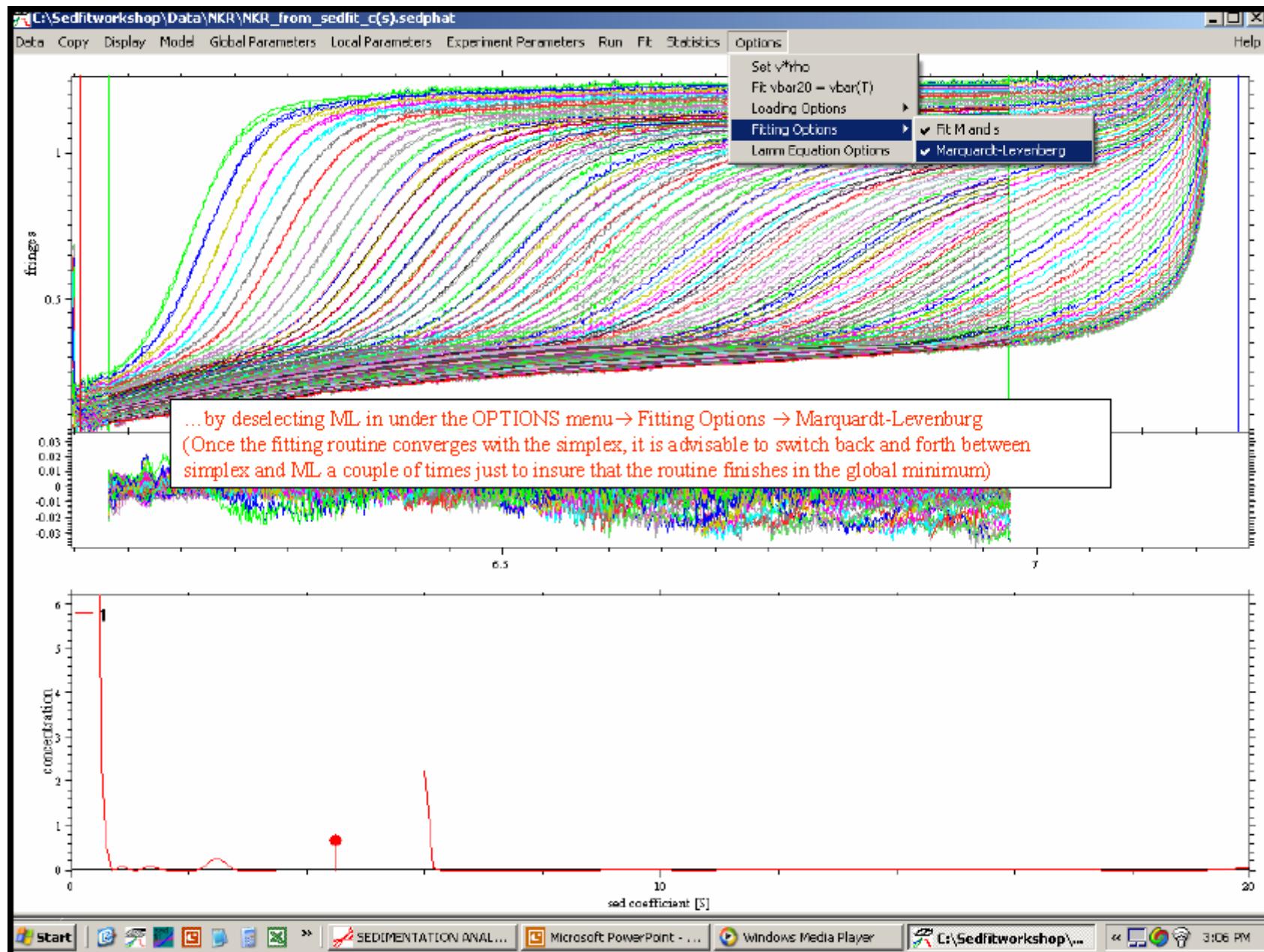


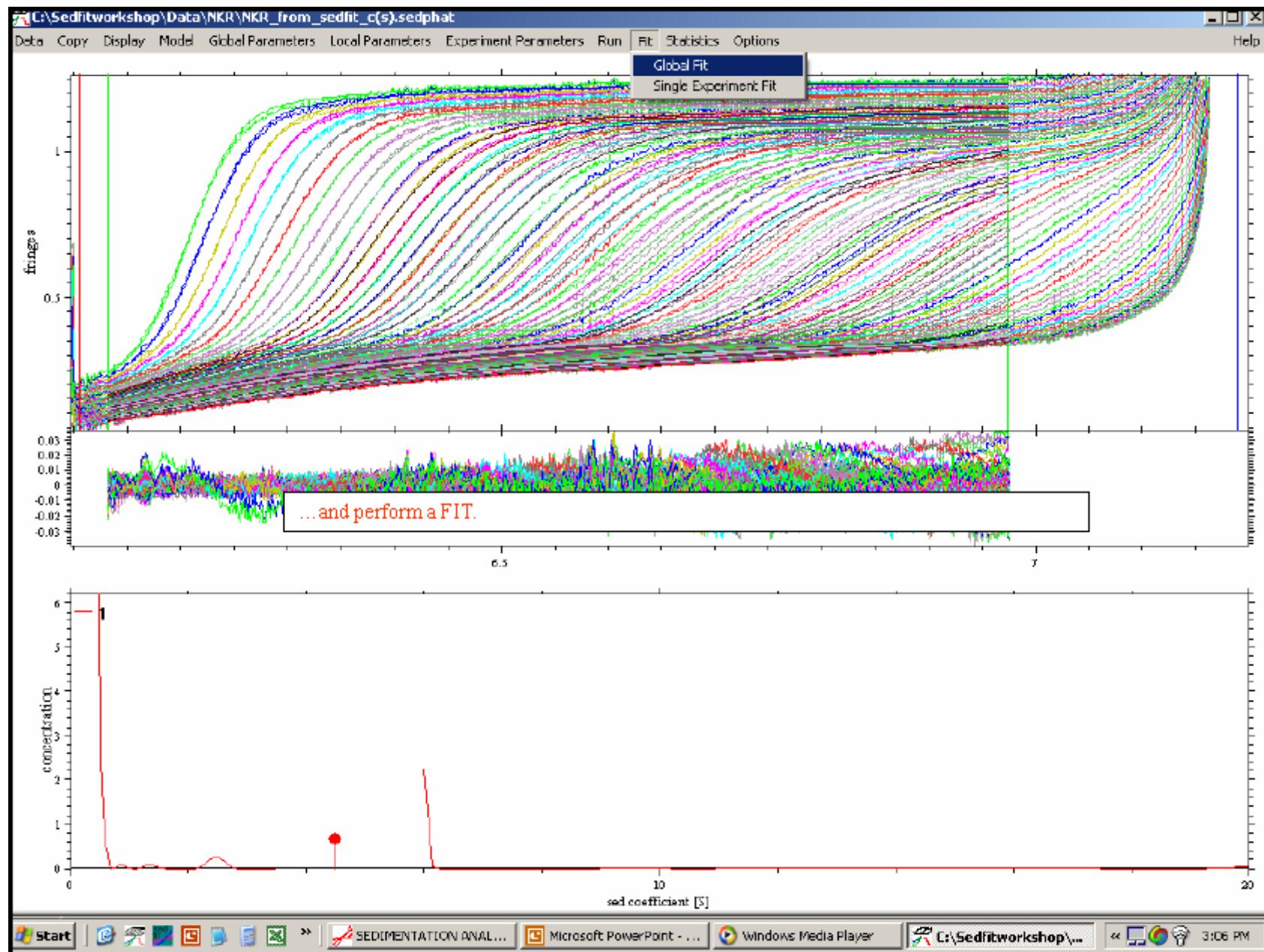


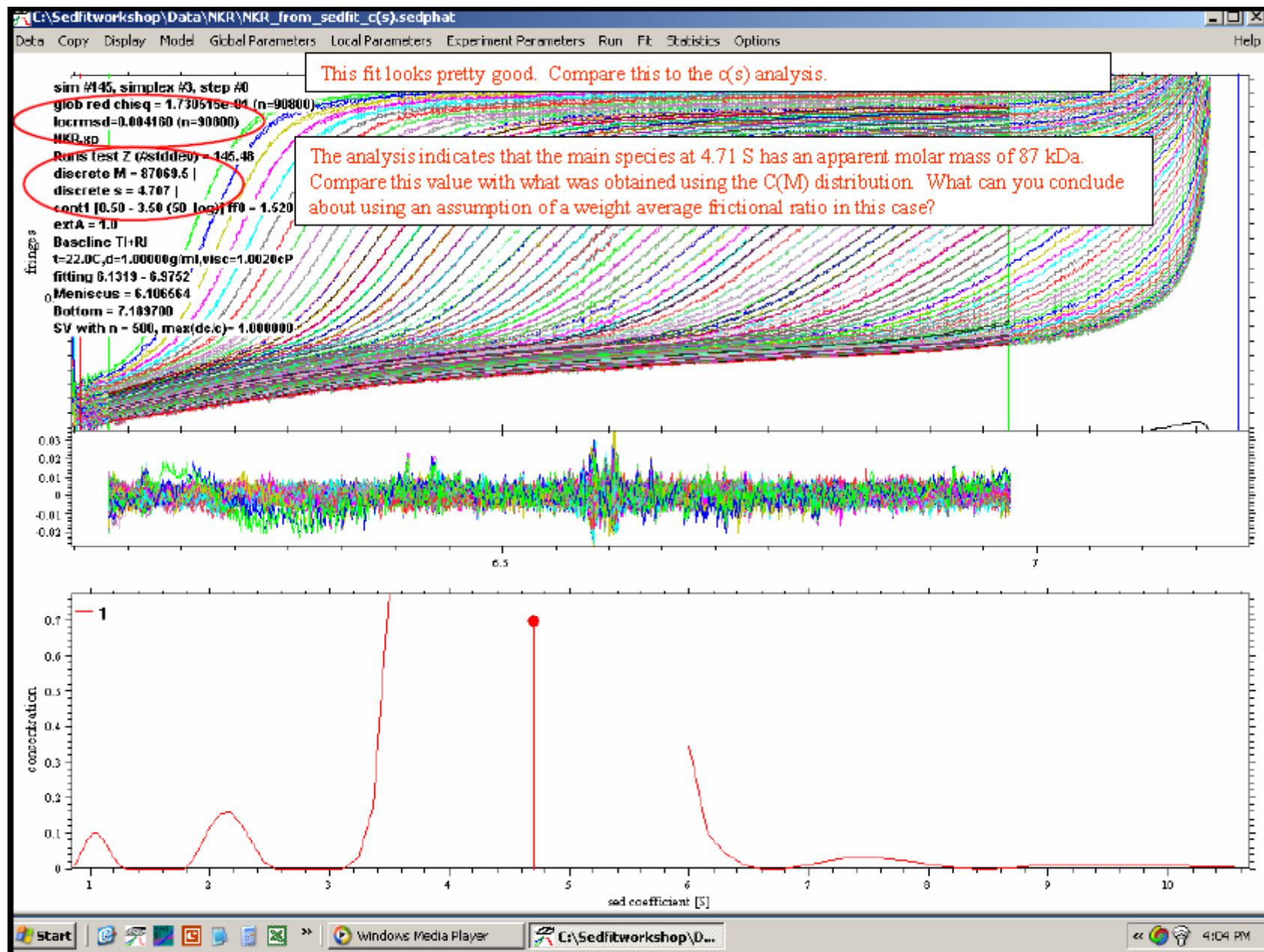


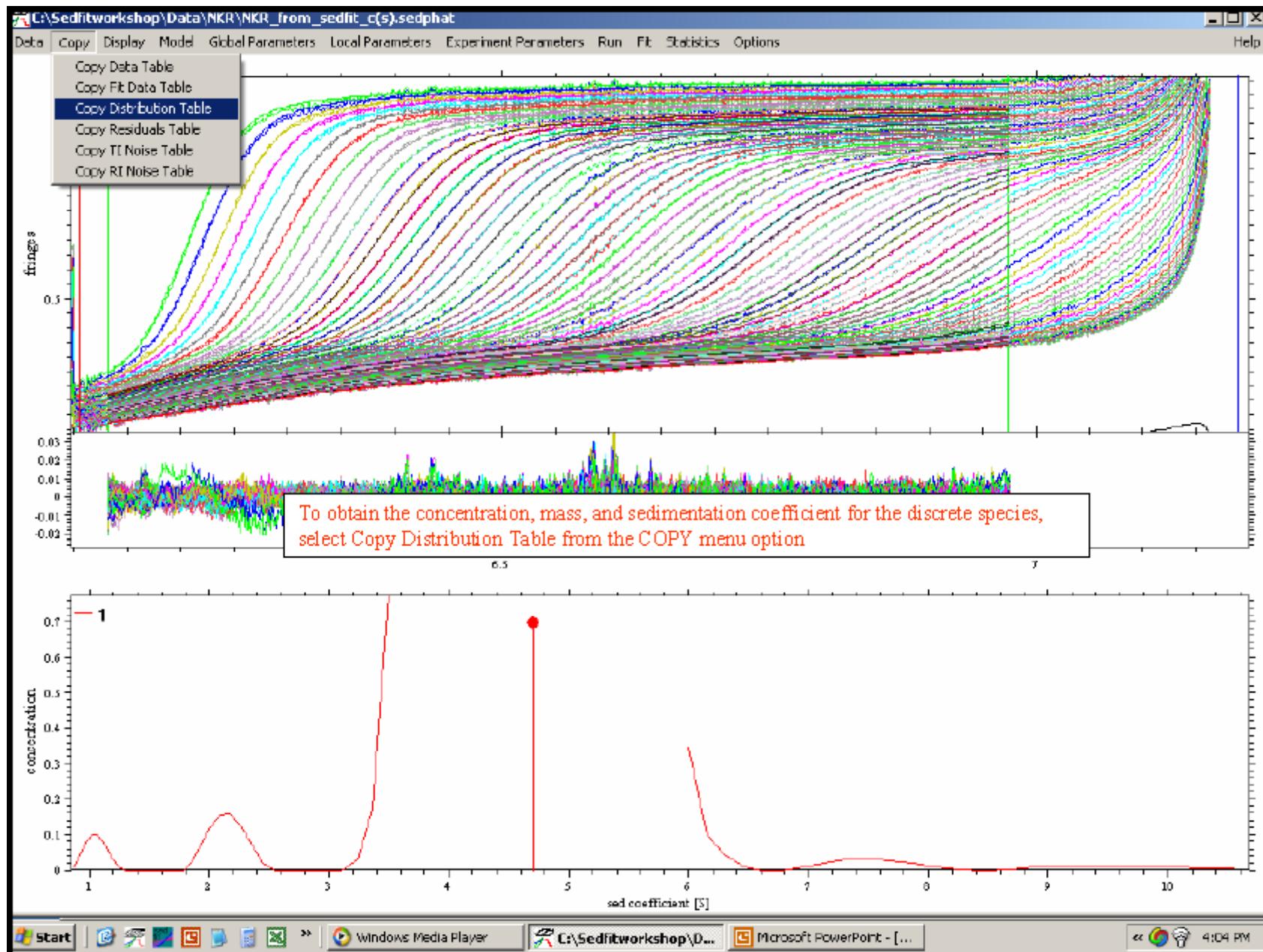


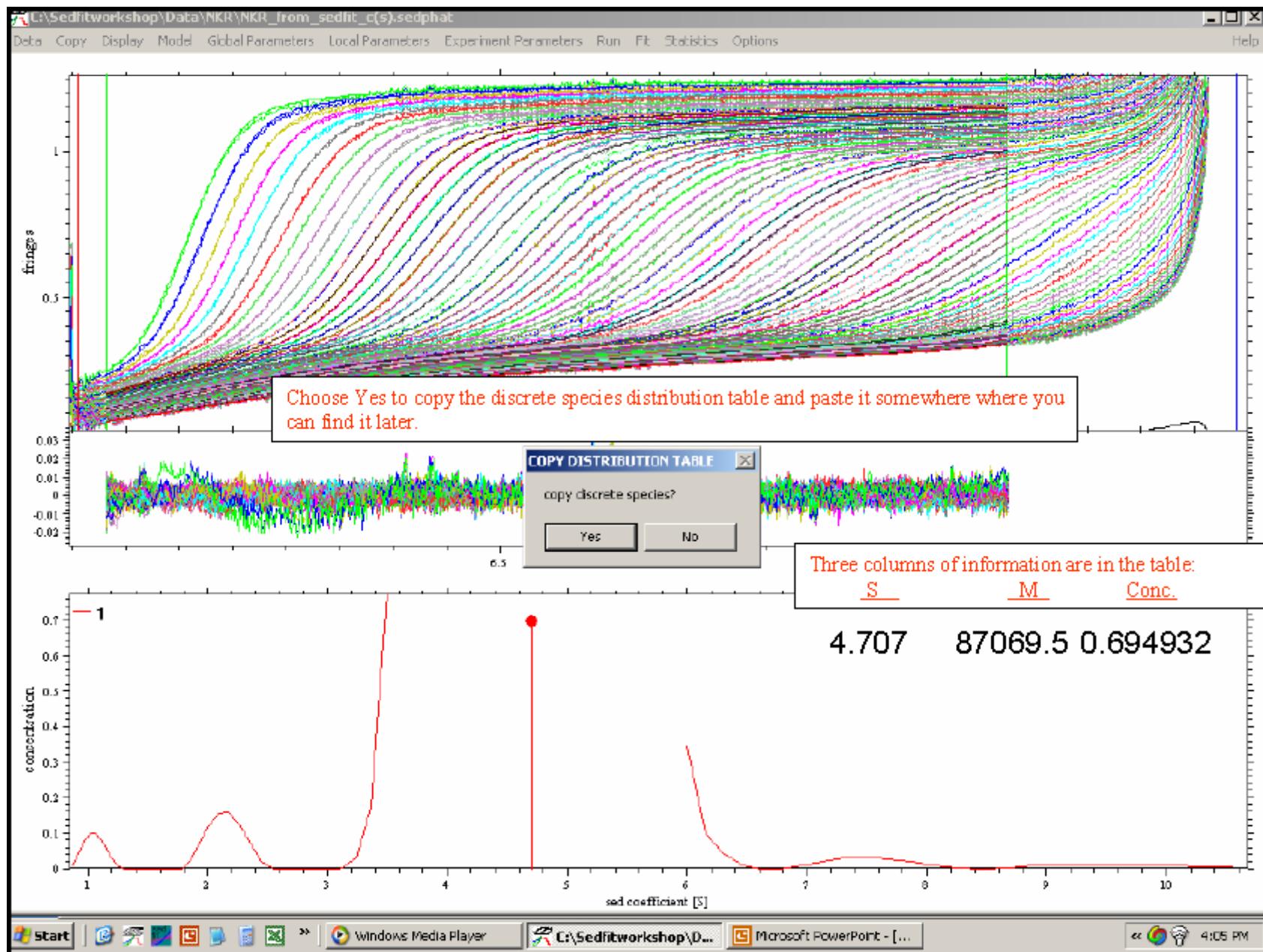








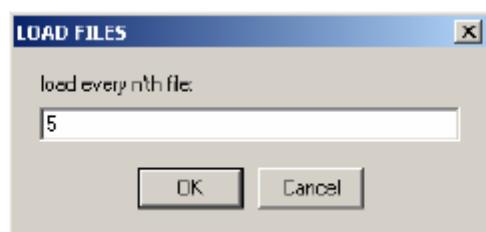
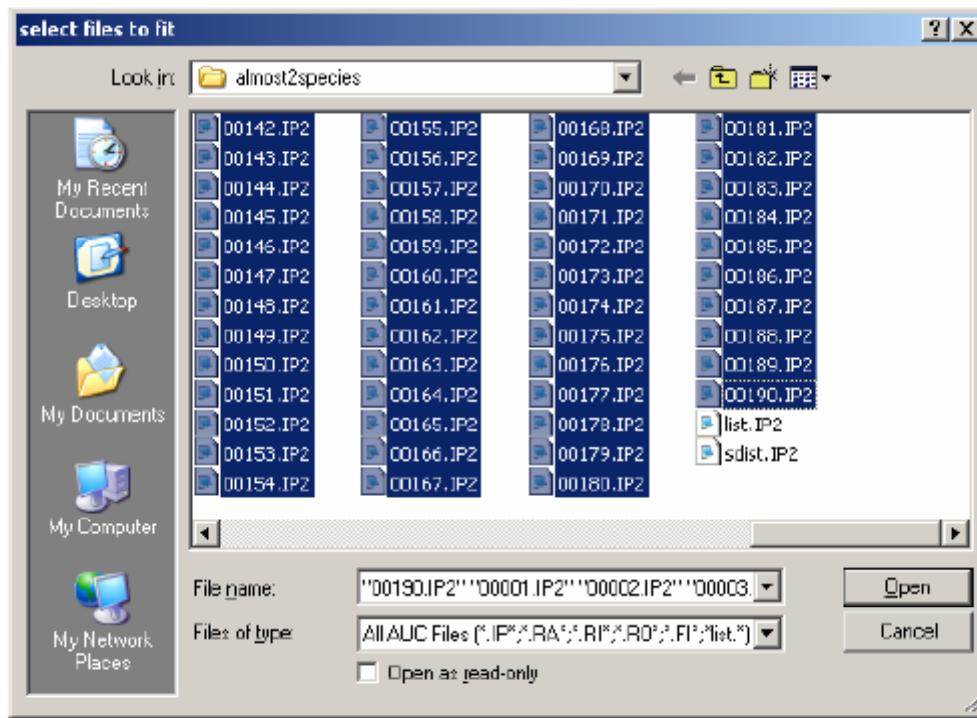




IV-1: Lamm equation modeling with a monomer-dimer system

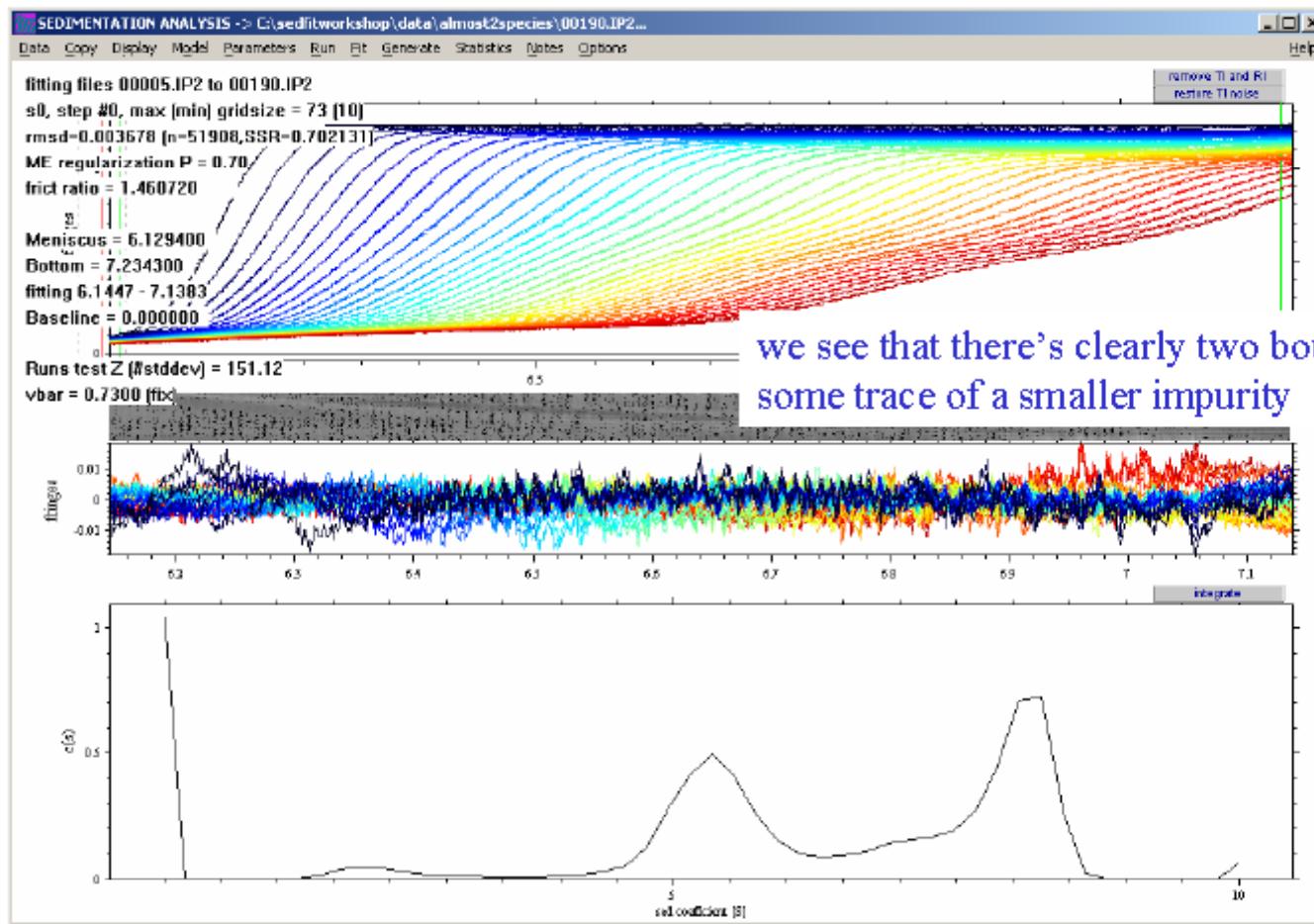
Let's suppose we have a protein ~ 100 kDa that, based on $c(s)$ seems to form dimers. There's also a smaller Mw impurity. The goal of this exercise is to make the transition from $c(s)$ in SEDFIT to set up an analysis of the same data as a interacting monomer-dimer system in SEDPHAT.

open SEDFIT, and load all files in the example “almost2species”



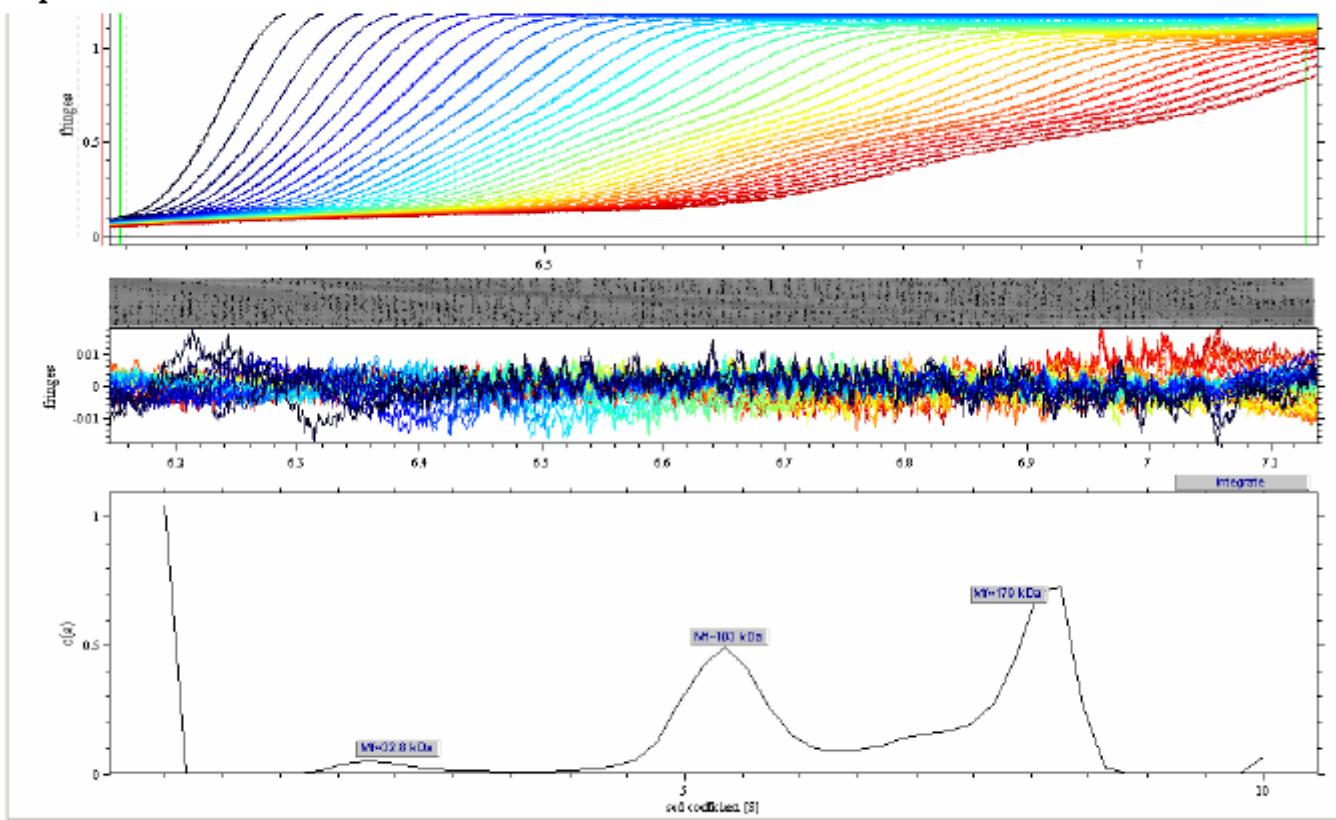


YES, then to a RUN command
then do control-D [data range] and control-N [subtract TI+RI noise]
This will reproduce a previously made standard c(s)
[if you don't get to restore a previously saved fit, see tutorials on c(s)]

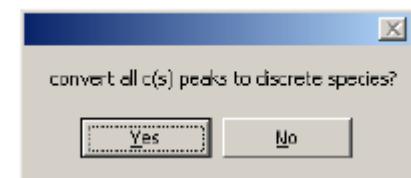
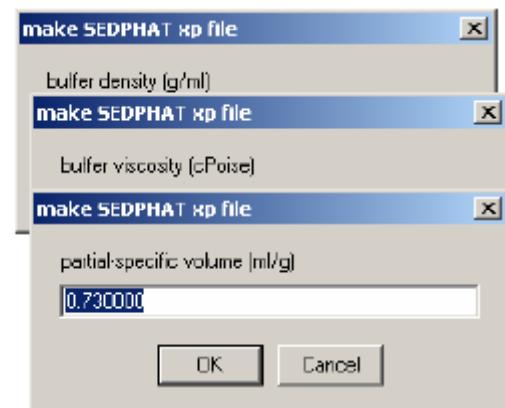
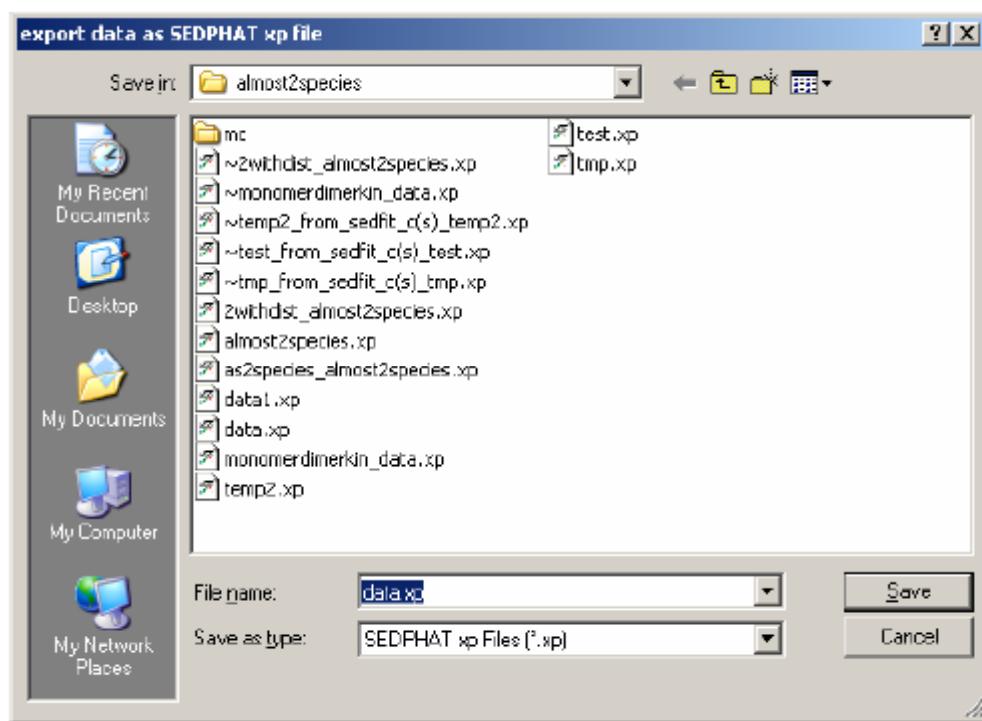


we see that there's clearly two boundaries, and
some trace of a smaller impurity

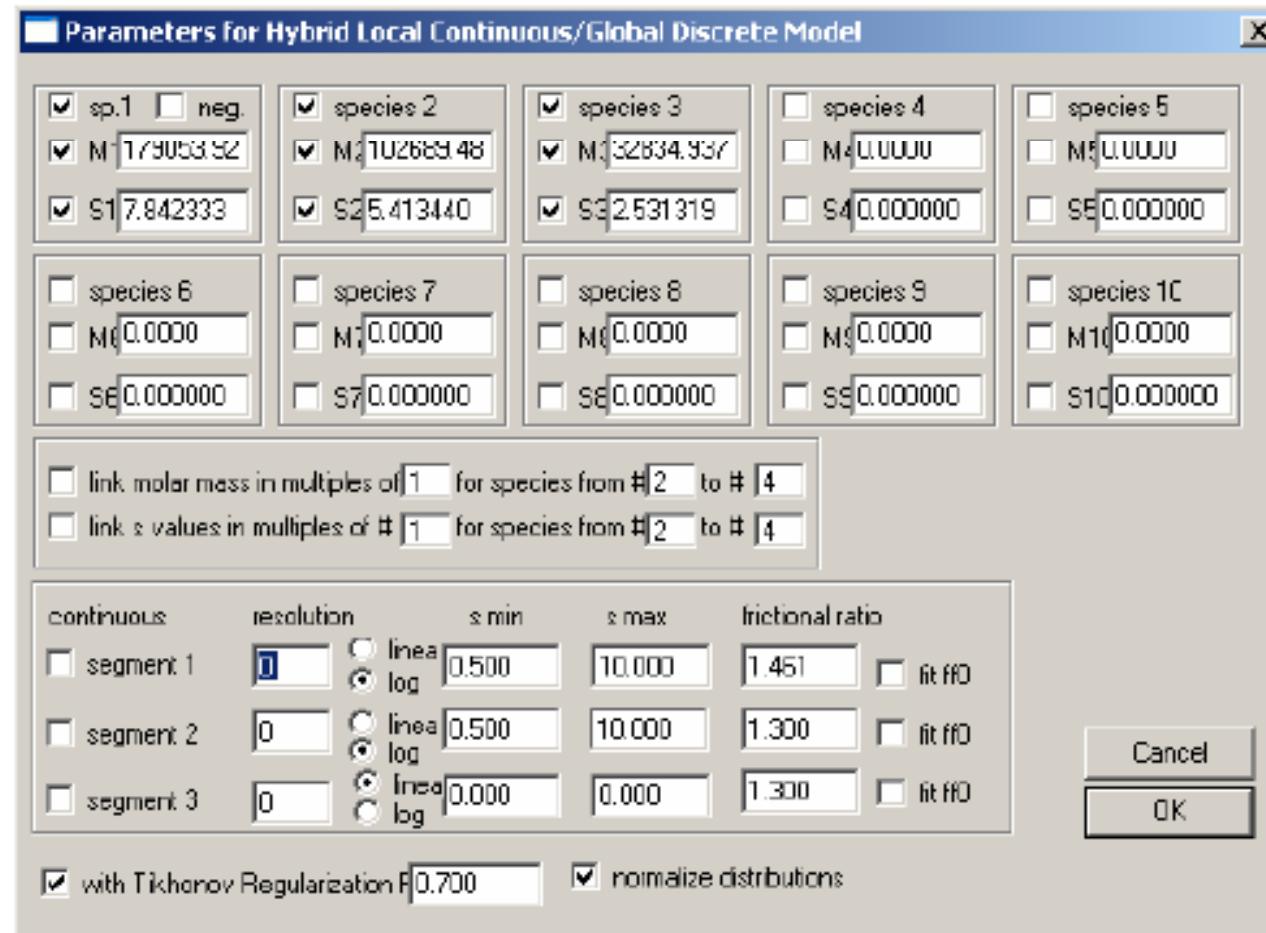
control-M shows the approximate Mw estimates: The 103 kDa should be the monomer, but the 179 kDa for what should be the dimer is too low a number, which means there might be excess boundary spread. Also, there's something in between the two peaks – this region should not be populated if we had really independent monomer and dimer



let's export this. Use the function Data → Export Data to SEDPHAT
Give the xp-file to be generated some filename, and default through all the parameter boxes.

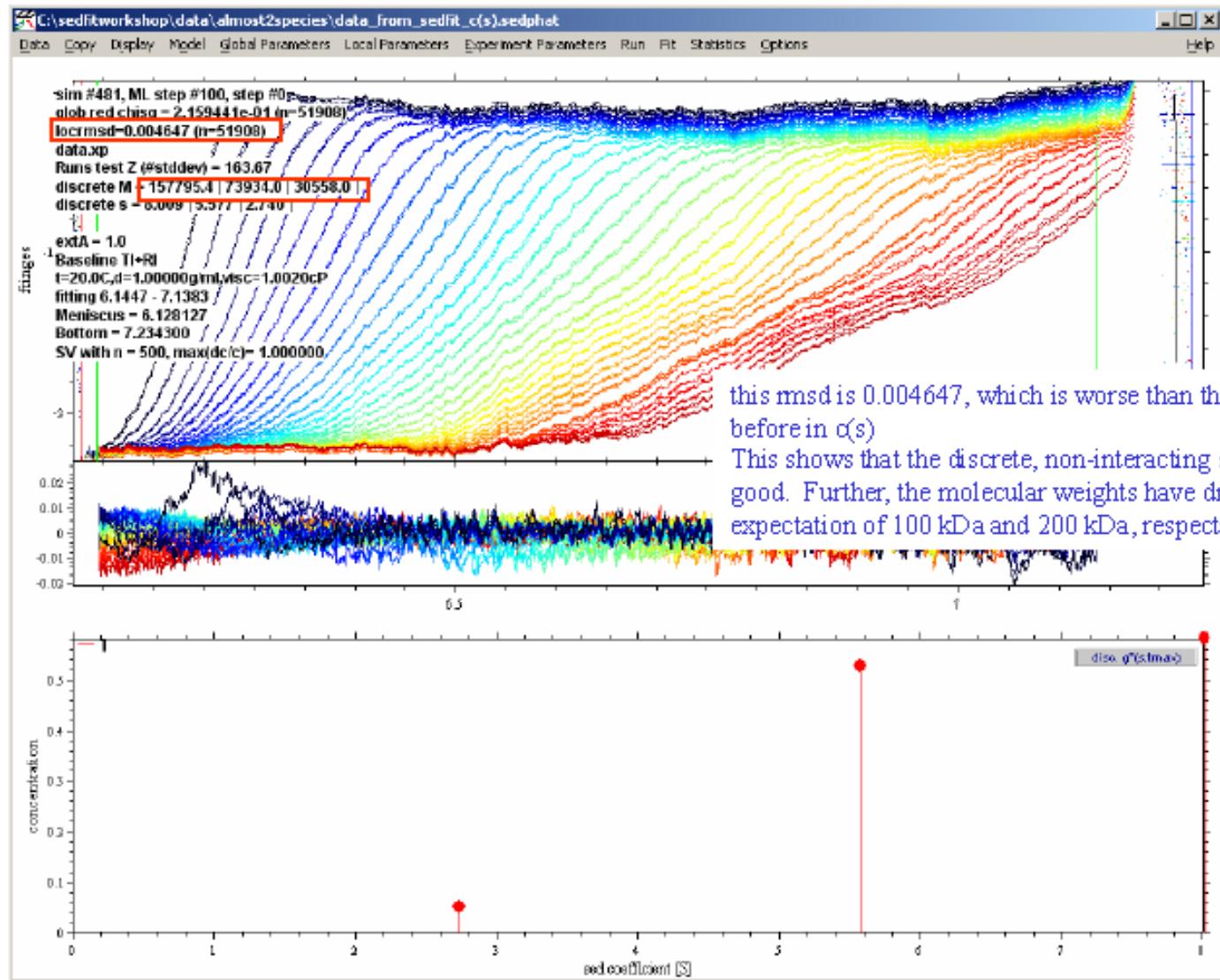


this will use the automatic integration in SEDFTT to generate starting guesses for SEDPHAT model treating all species as discrete species.

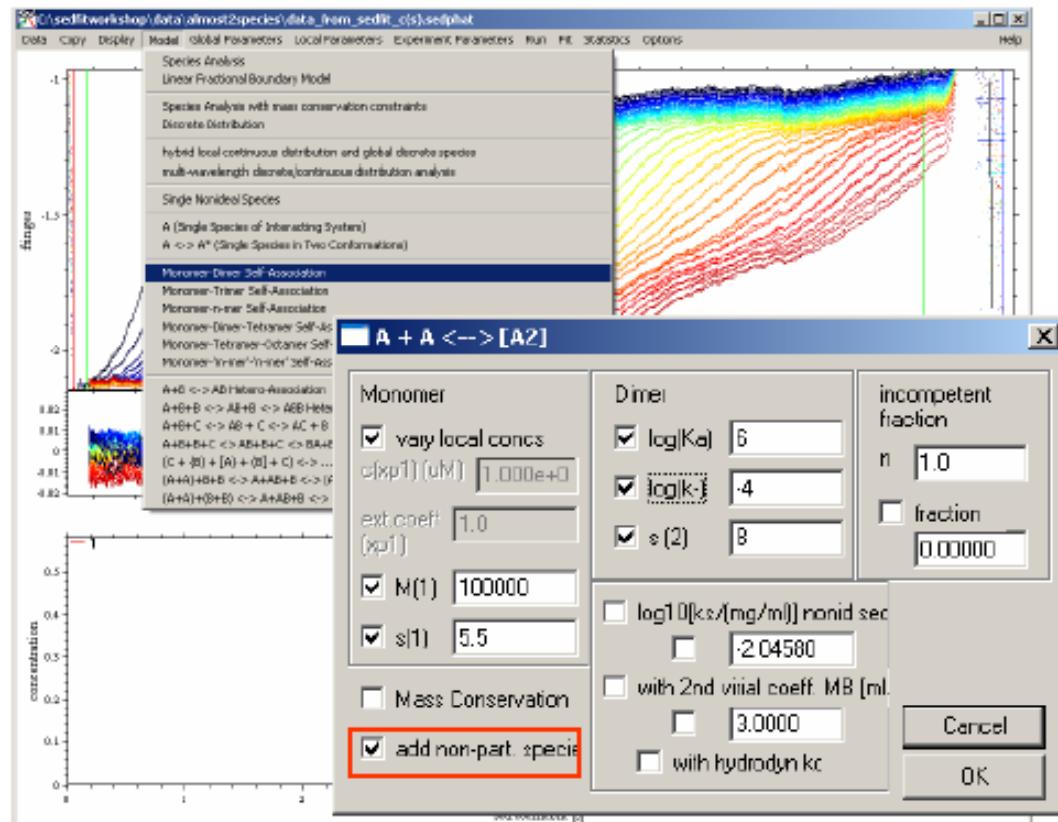


accept, do RUN and then FIT

then fine-adjust the display by using control-D [set data range] and right-double click in the distribution window [full range distribution plot] and control-O [show last fit info]



Let's try the monomer-dimer model in order to account for finite reaction kinetics.



also check 'add non-part. species' box, since we will need to account for the slow sedimenting degradation product when fitting this data.

Then hit OK

we use the following initial estimates:

$M_1 = 100\text{kDa}$ [that should be the Mw from sequence, but we don't know the exact vbar, therefore we'll float it here]

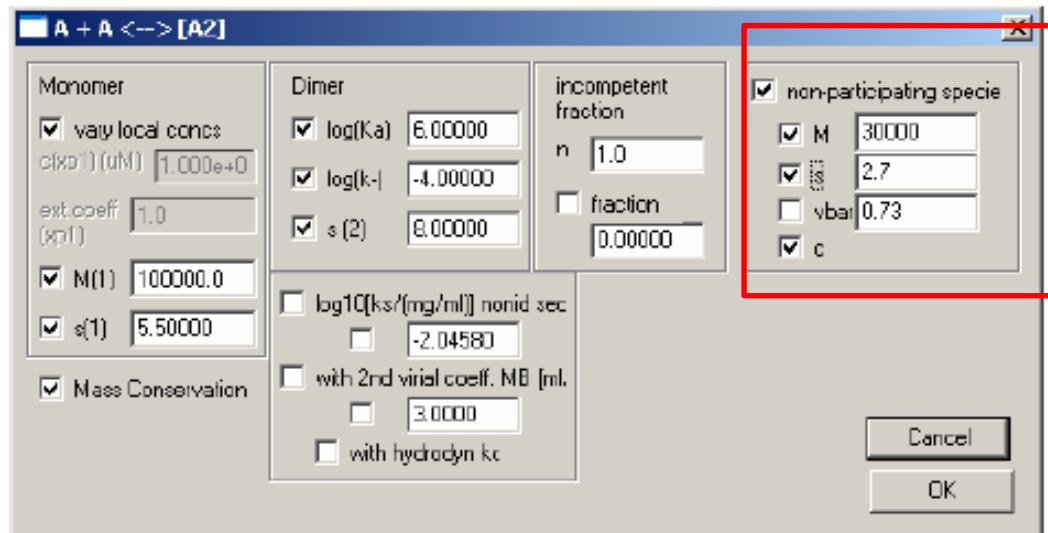
$s_1 = 5.5, s_2 = 8$ [these are the s -values from the discrete model, which should be excellent starting guesses and we'll float]

$\log(Ka) = 6$ [since we loaded $\sim \mu\text{M}$ concentration, and see roughly equal height of monomer and dimer peaks in $c(s)$, we should be close to K_d at these μM concs]

$\log(k-) = -4$ [the reaction must be fairly slow (i.e. -4 or slower) for us to resolve peaks in $c(s)$, but we want to put the starting guess at a value close to the range of -3 to -4 where the SV experiment is sensitive to the value of $\log(k_{off})$. Otherwise the optimization of this parameter would become difficult.]



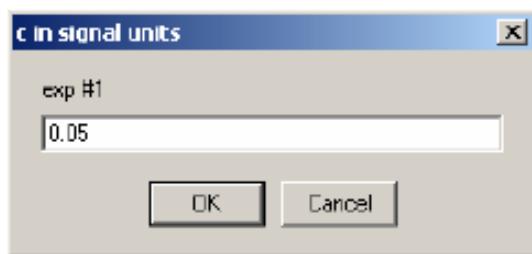
this is no problem, all SV models obey mass balance.



because we had checked the 'add non-part. species' box, now an extended box shows up with a field where we enter the parameters of the non-participating species.

This will be 30kDa and 2.7 S, which are the estimates taken from the discrete species fit

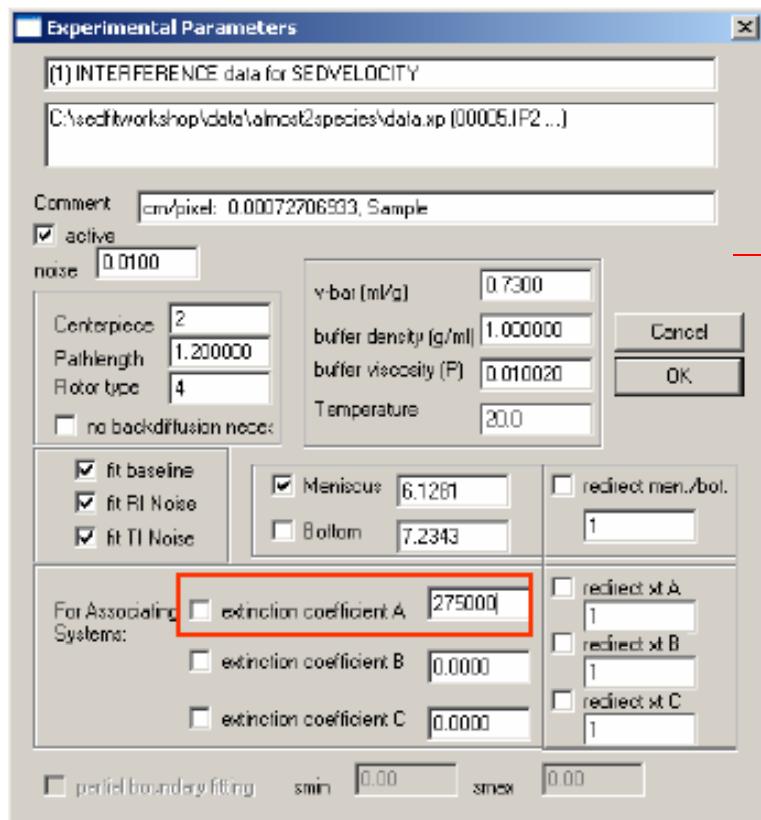
Then hit OK



Now a input box pops up asking us how large the concentration of the non-participating species is. This is in signal units, and from the discrete species fit we can estimate a value of 0.05. (Or this can be obtained also from integrating the previous c(s) around this contaminant peak).

One more set of parameters is required, due to the fact that in the interacting species models we want to deal with molar binding constants. Therefore, we need starting concentrations in molar units, as well as extinction coefficients to relate the molar concentration to signal.

Open the Experimental Parameters:



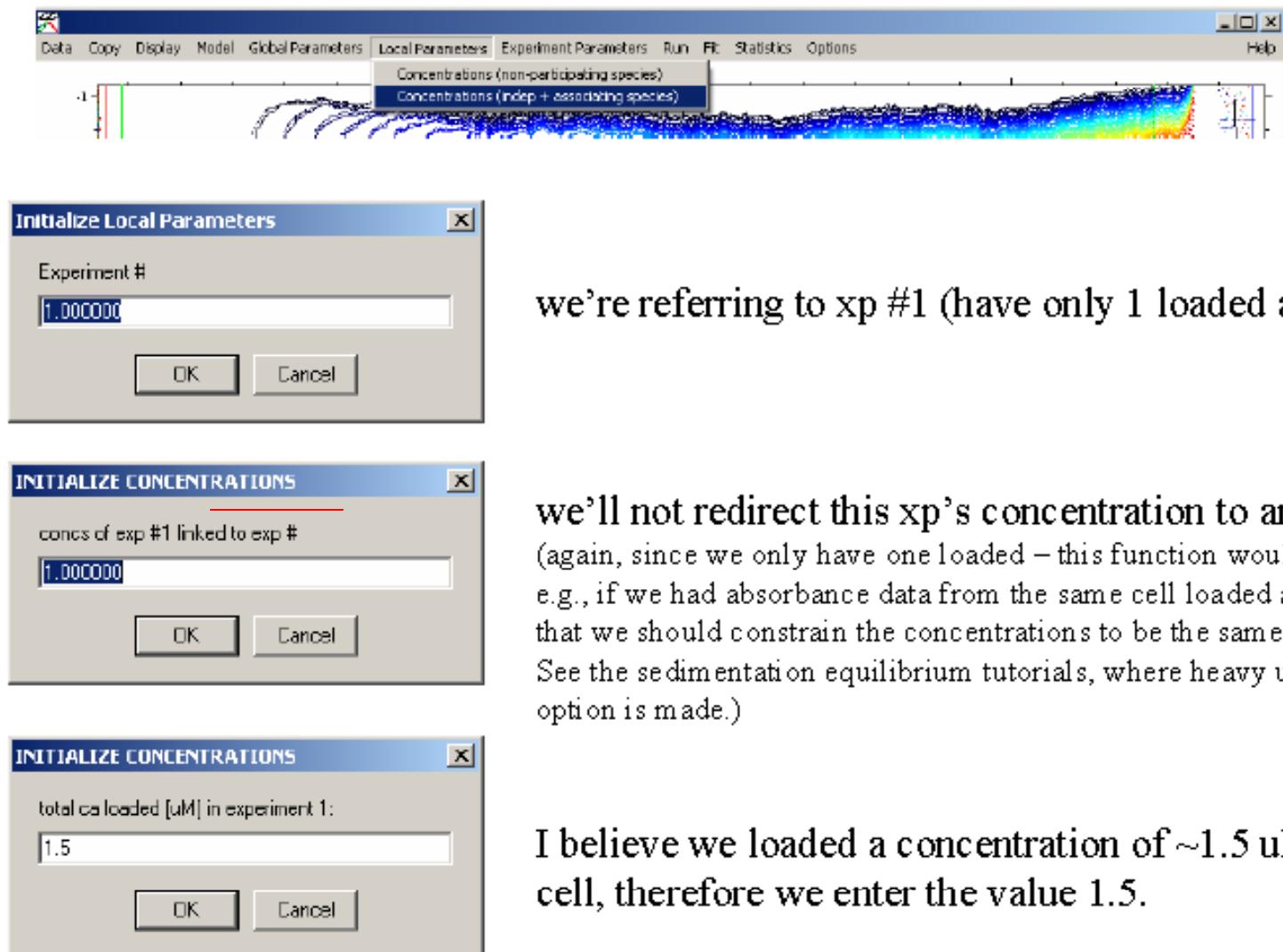
For the 'extinction coefficient A', enter a value of 275000. This is interference data, and based on the sensitivity for proteins of 3.3 fringes/mg/ml, we have the 'interference extinction coefficient' $2.75 \times M_w$.

Note we also have a pathlength of 1.2 cm, which will allow SEDPHAT to correct for the 12 mm thickness of the centerpiece used.

If this was absorbance data, we would enter here the conventional molar extinction coefficient (OD/M/cm), just as taken, for example, from SEDNTERP.

Click on OK

Next, choose the “concentrations (indep + associating species)” from the menu.

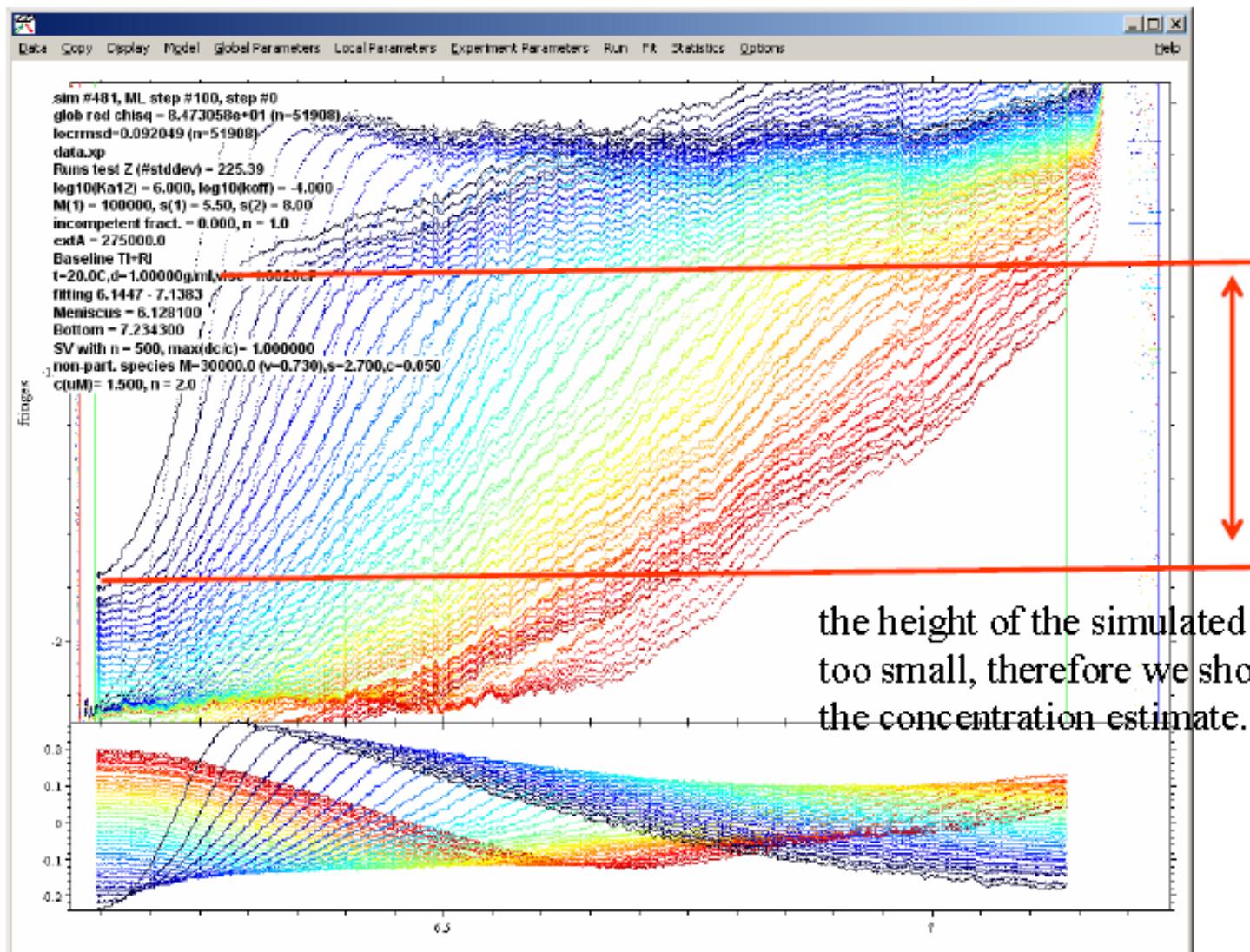


we're referring to xp #1 (have only 1 loaded anyway...)

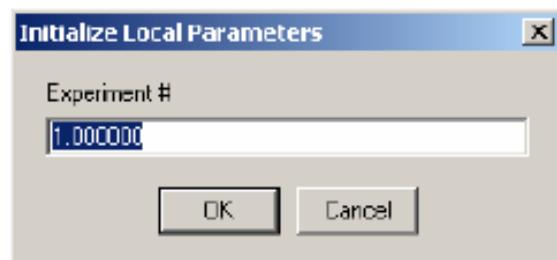
we'll not redirect this xp's concentration to another xp
(again, since we only have one loaded – this function would make sense,
e.g., if we had absorbance data from the same cell loaded as xp2 such
that we should constrain the concentrations to be the same – ‘redirected’.
See the sedimentation equilibrium tutorials, where heavy use of this
option is made.)

I believe we loaded a concentration of ~1.5 uM in this
cell, therefore we enter the value 1.5.

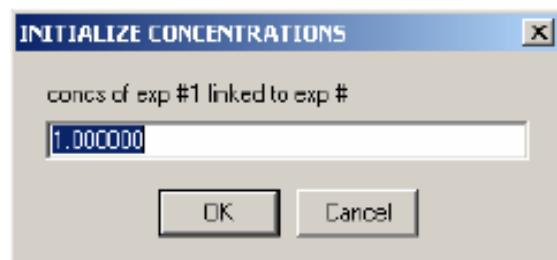
Do a RUN to test the starting guesses.



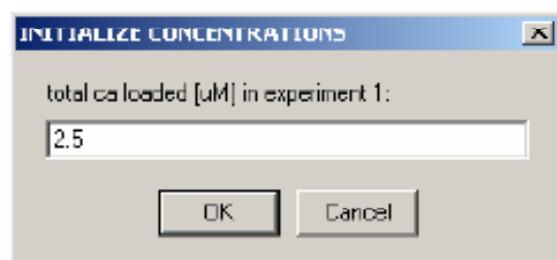
Again, choose the “concentrations (indep + associating species)” from the menu.



we're referring to xp #1 (have only 1 loaded anyway...)

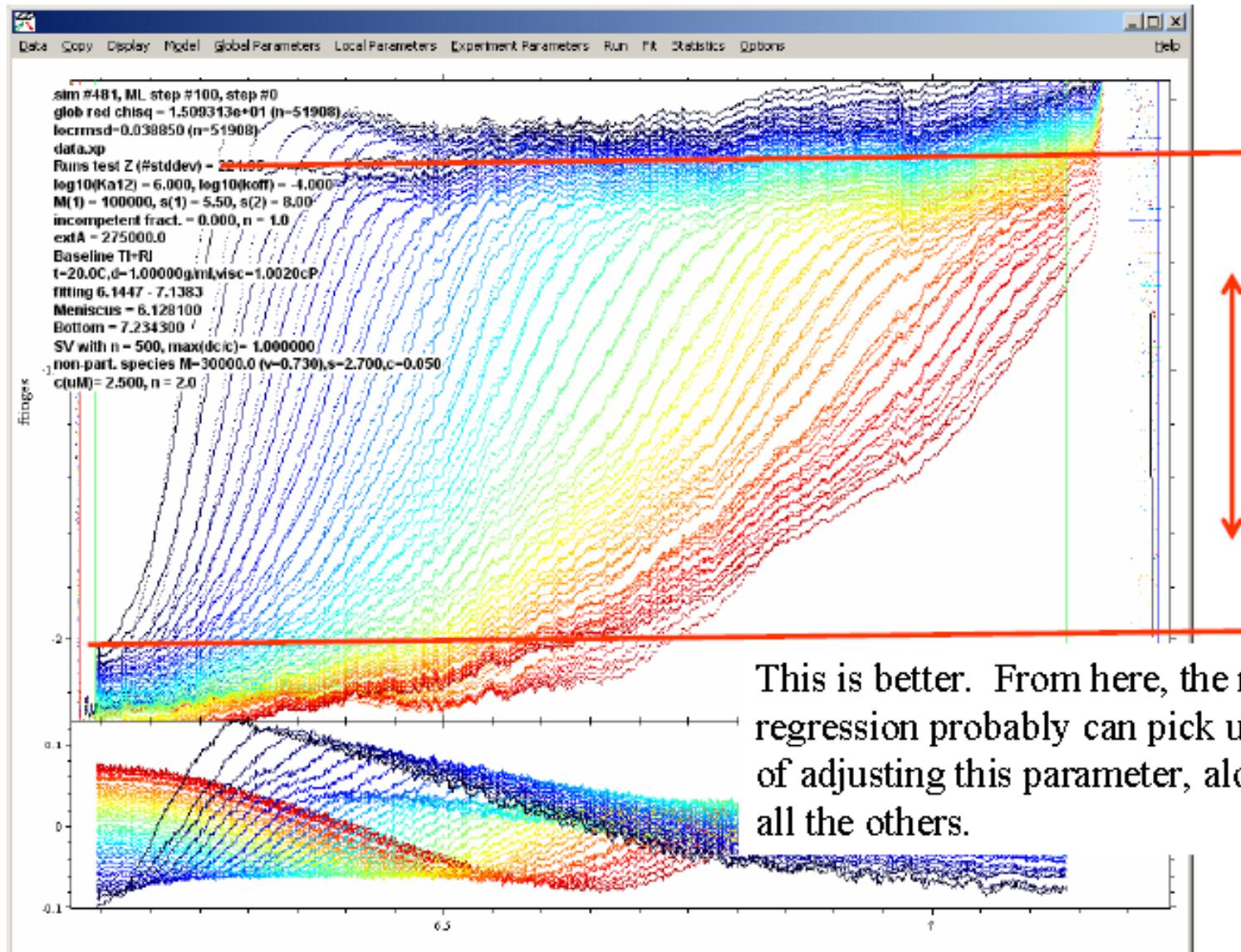


we'll not redirect this xp's concentration to another xp



Let's try a value of 2.5 uM.

Do a RUN to test the starting guesses.



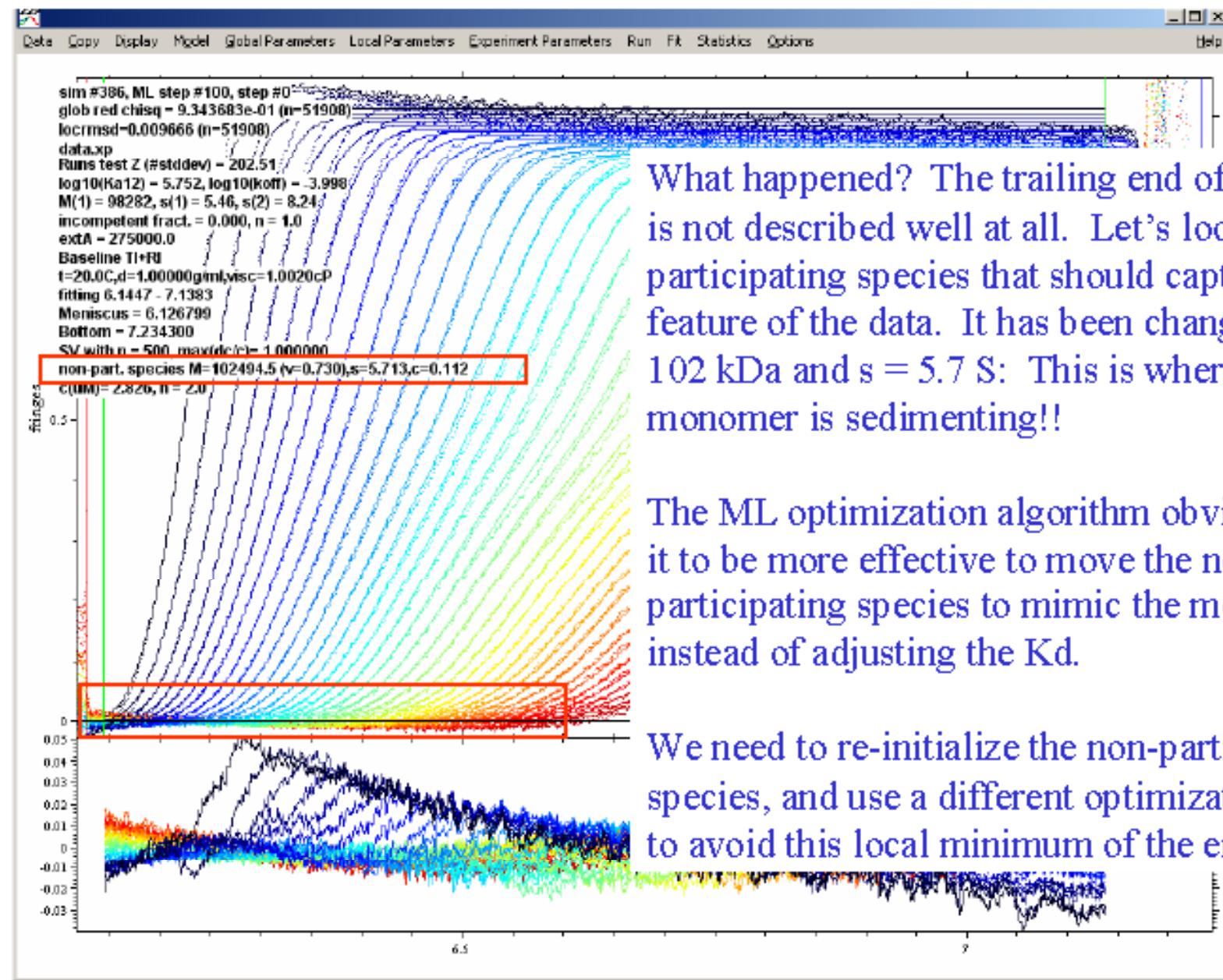
Now click on FIT. Note, the 'ML' in the first line indicates we are currently using Marquardt-Levenberg optimization, which will need to do a large number of simulations to initialize the gradient information and seems not to do much before updating the screen... nevertheless, this is a very efficient algorithm overall.

If your copy of SEDPHAT is configured to use the Simplex or the Simulated Annealing algorithm for minimization, the following sequence of events will be different. There are many different equally valid options how to do the fitting.

As it turns out, the following documentation will not be of a straightforward fit, but one where we have to avoid getting trapped into local minima. We use a combination of fitting algorithms to get to the best fit in a shallow minimum. (Remember the error surfaces of the banana function and the mountain range.) The details of the following are probably not reproducible, because the Simplex routine involves seed points generated randomly, and therefore it will slightly differ each time. However, it should be straightforward to apply the same strategy to your particular situation, as well as to other fits.

From here on, this is more an example on switching between fitting algorithms.

After the fit has stopped, hit control-N [subtract noise] and control-O [show best-fit info again]. This is where I ended up with after the fit. This does not look very good, yet.



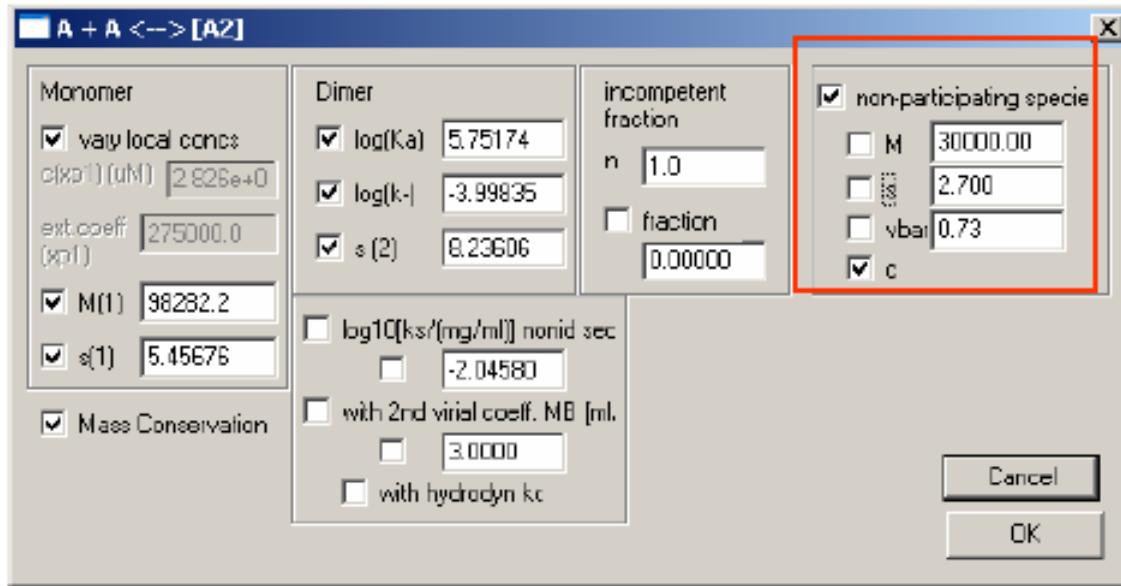
What happened? The trailing end of the boundary is not described well at all. Let's look at the non-participating species that should capture this feature of the data. It has been changed to $M = 102$ kDa and $s = 5.7$ S: This is where the monomer is sedimenting!!

The ML optimization algorithm obviously found it to be more effective to move the non-participating species to mimic the monomer, instead of adjusting the K_d .

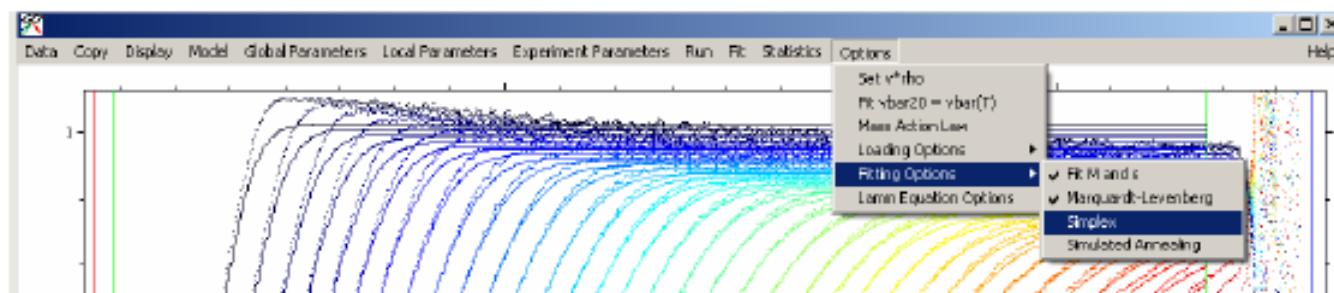
We need to re-initialize the non-participating species, and use a different optimization strategy to avoid this local minimum of the error surface.

in the Global parameters box, enter $M = 30000$ and $s = 2.7$

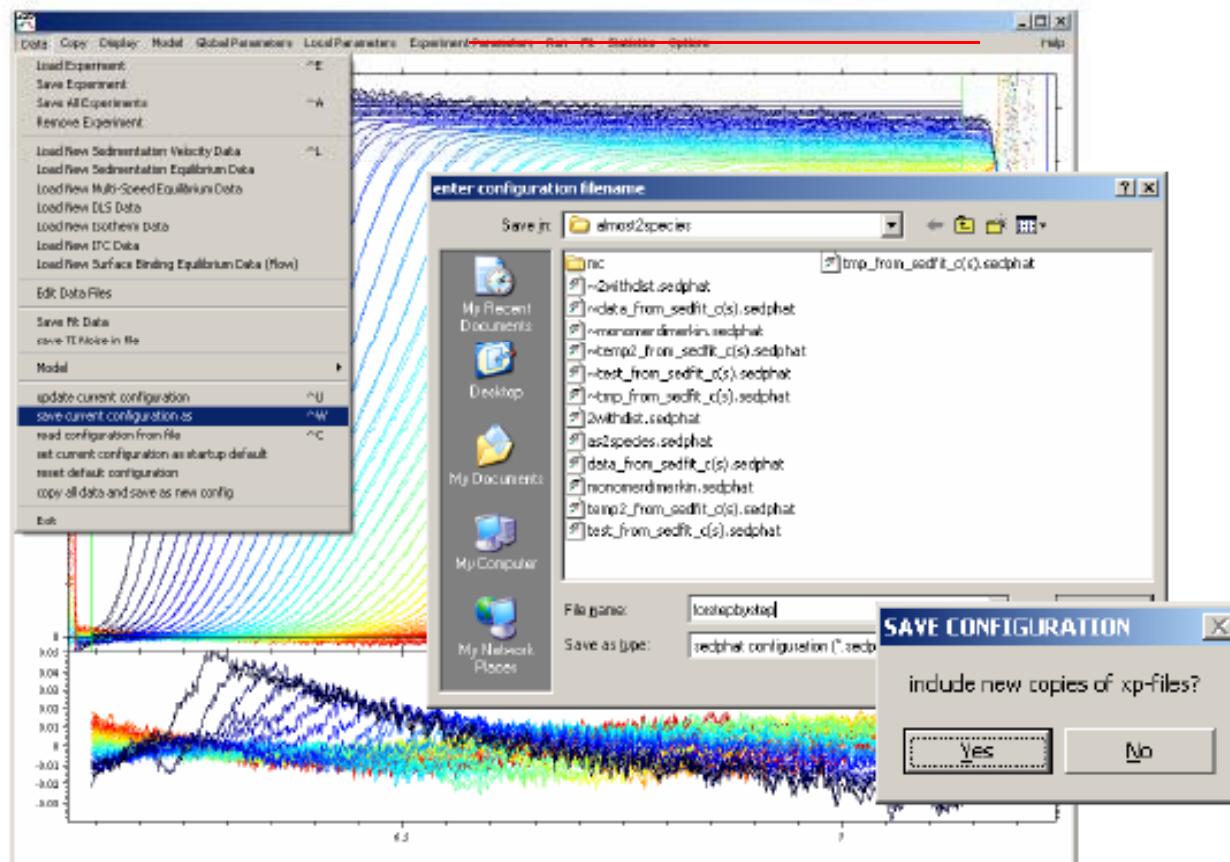
Also, we want to keep these parameters fixed now (we'll set them free later).



click on the “Simplex” menu function. This will toggle off Marquardt-Levenberg and switch on the Simplex algorithm

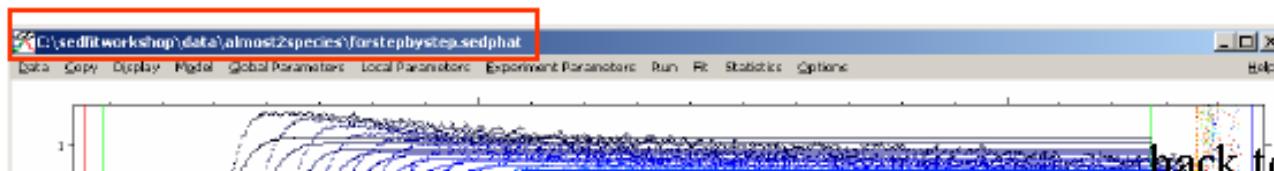


Also, just in case we need to go back, let us save this state of the analysis.
In the Data menu, click on “save current configuration as”, and enter a filename.



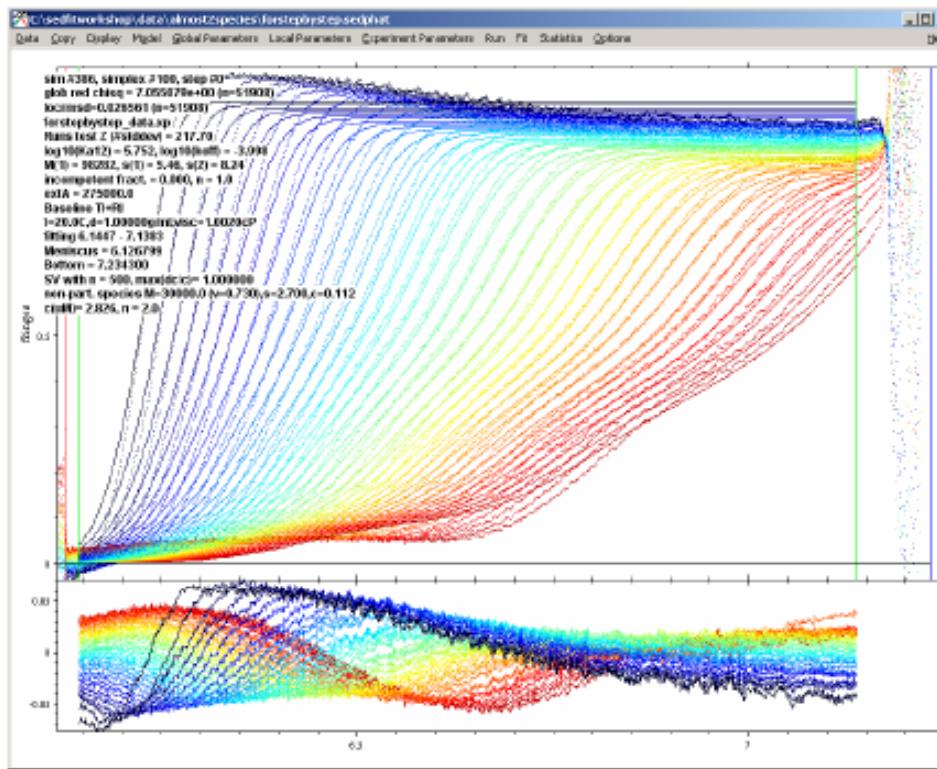
Yes, let's also
save the xp file.

This way, we can recreate this analysis later, and we'll be able to tell from the title bar the filename.



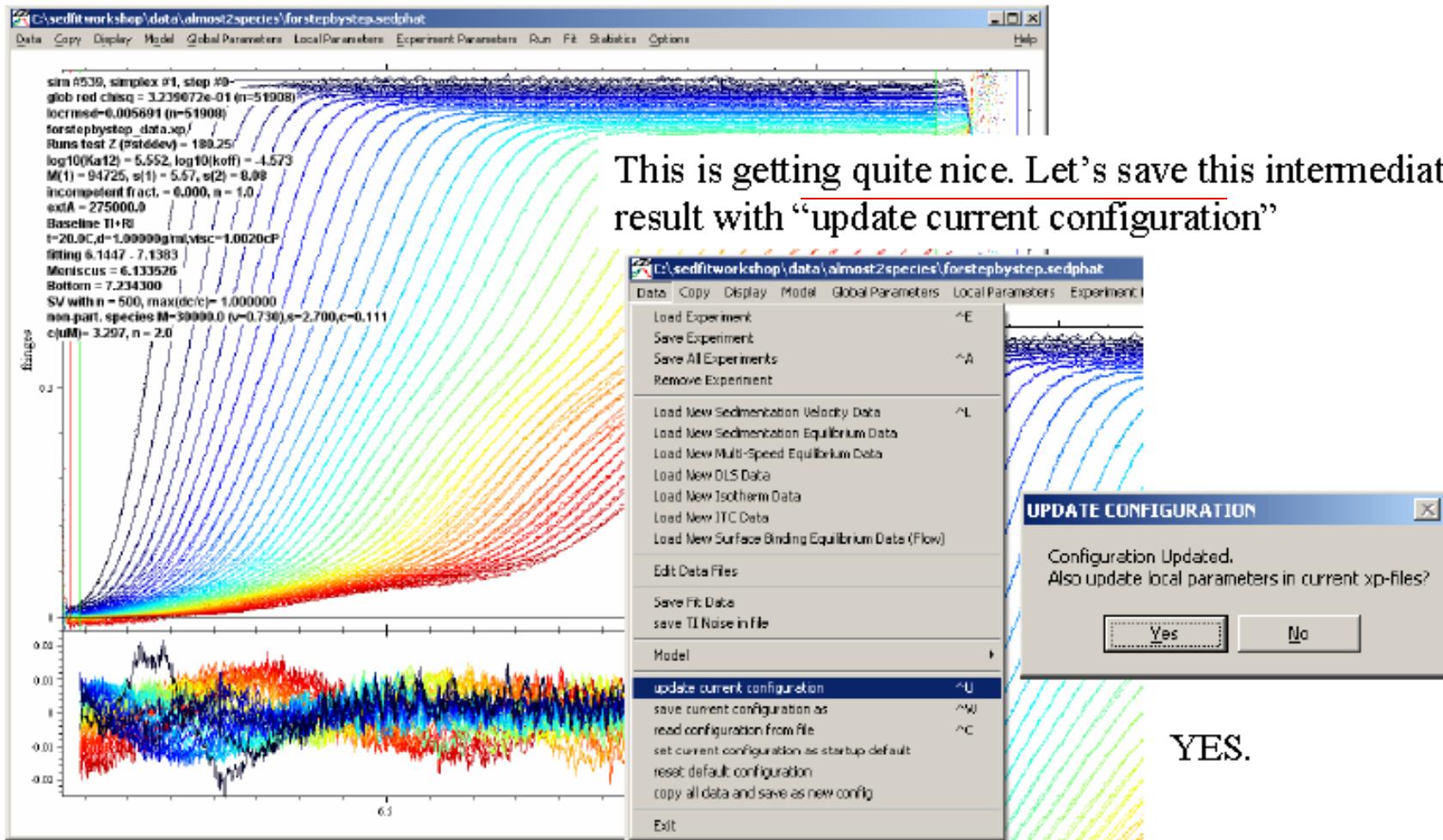
back to the fitting problem...

After just doing a RUN and control-N, we find this to be a good starting point. Not surprisingly, it seems to be worse then after the previous fit, but it will become hopefully much better again after we do a FIT.

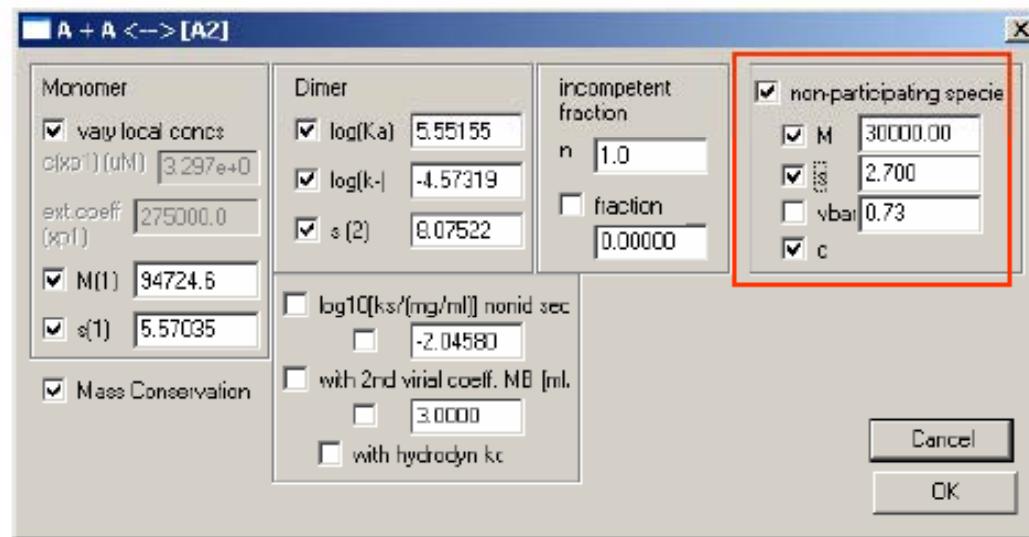


Click on FIT. You'll notice the different behavior of the display during a Simplex fit as compared to the previous Marquardt-Levenberg fit. Now, we see things changing and gradually improving earlier. However, at some point the convergence seems to become quite slow.

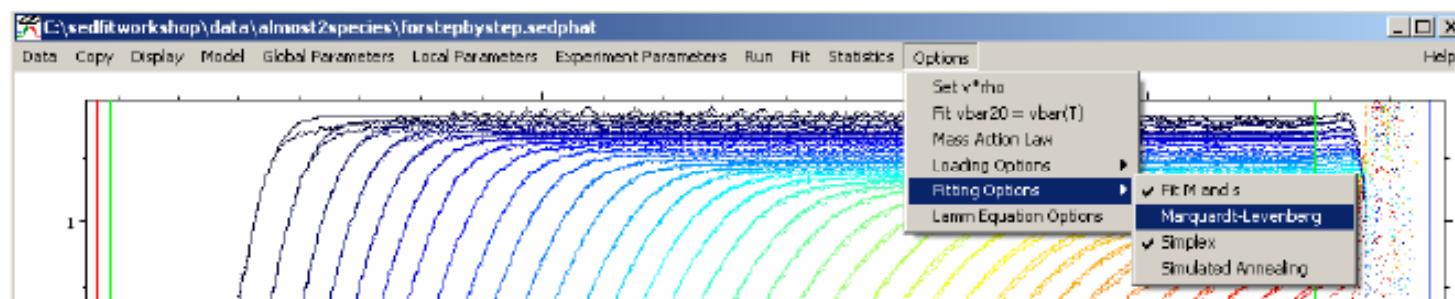
We could just let this go until it stops. I recommend that. However, being impatient, after the rmsd of the fit has dropped below 0.006, I interrupted this by hitting the space bar *once*. After control-N and control-O to get a good display, I find:



go back to the Global parameters box, and float M and s of the non-participating species

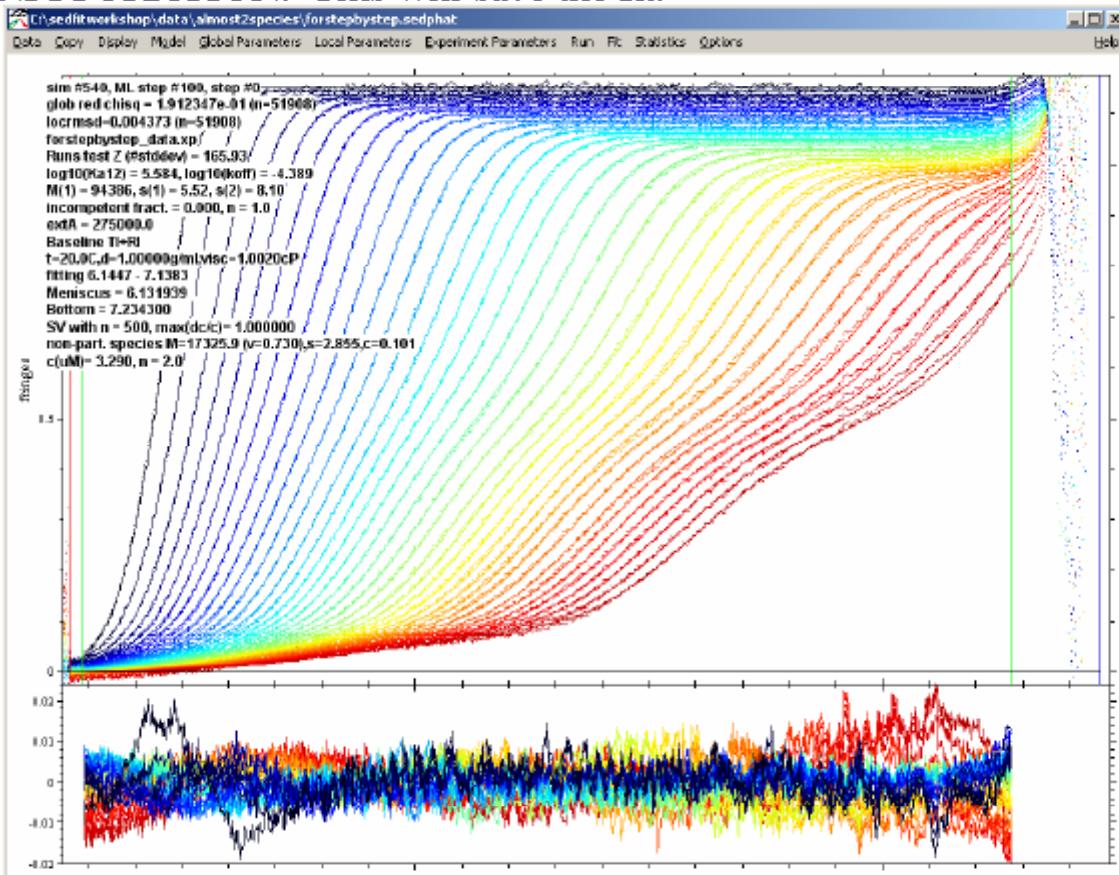


Also, we change back to Marquardt-Levenberg, to take advantage of its better ‘homing-in’ capability. Click on the “Marquardt-Levenberg” menu function to toggle.



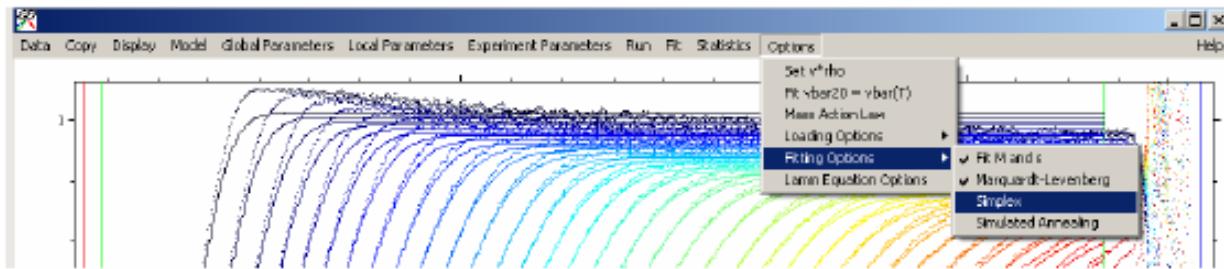
Then do a FIT

after a minute or two, we find the following best-fit (I used control-N and control-O to optimize display). use control-U or the menu function to UPDATE CURRENT CONFIGURATION. This will save the fit.



this fit does not look optimal, yet, as judged by the systematic residuals and misfit in the faster boundary. Therefore, we switch back to Simplex optimization and try to improve further.

click on the “Simplex” menu function. This will toggle off Marquardt-Levenberg and switch on the Simplex algorithm

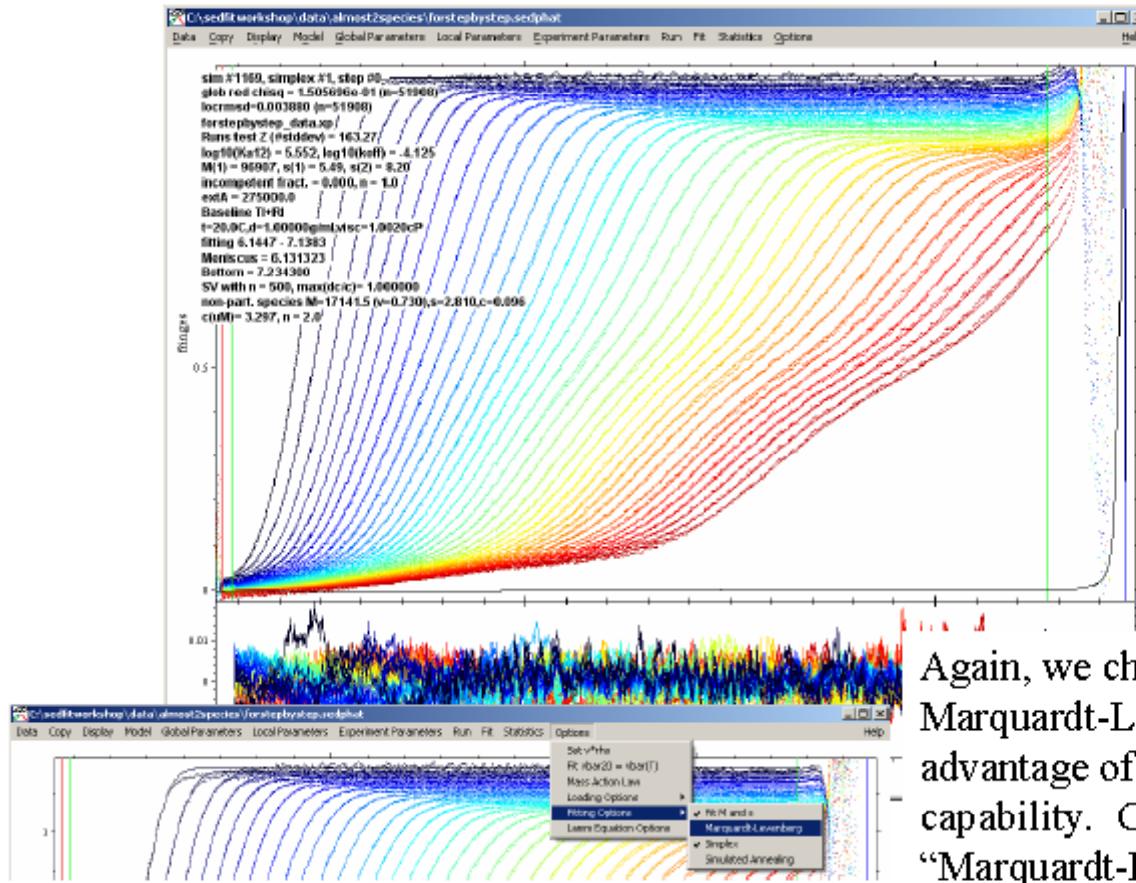


Then click on FIT

I saw initially a big improvement, indicating that the Marquardt-Levenberg again had been trapped in a small local minimum, from which the Simplex helped us escape. But after a while, convergence seemed to slow down very much, such that I again hit the space bar *once* to interrupt the fit.

[Hitting it once just stops the fit. It will then do a last simulation using the so-far best fitting parameters. If you inadvertently hit it twice, you are aborting this last simulation, and the fitted lines will be far off. However, you can rescue this by using a RUN command to re-establish the best-fit curves.]

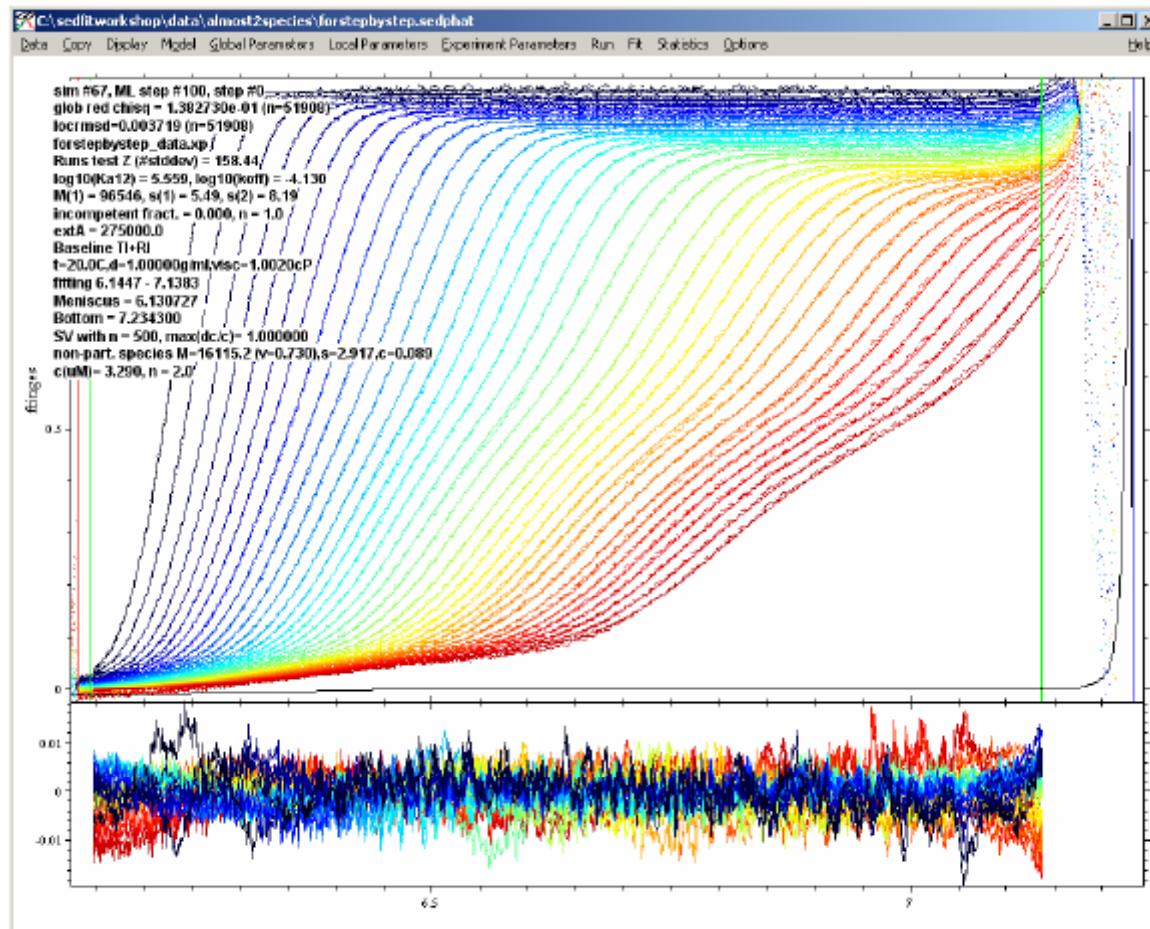
This is what I found. Use control-U or the menu function to UPDATE CURRENT CONFIGURATION. This will save the fit. Always also save the xp-file with it.



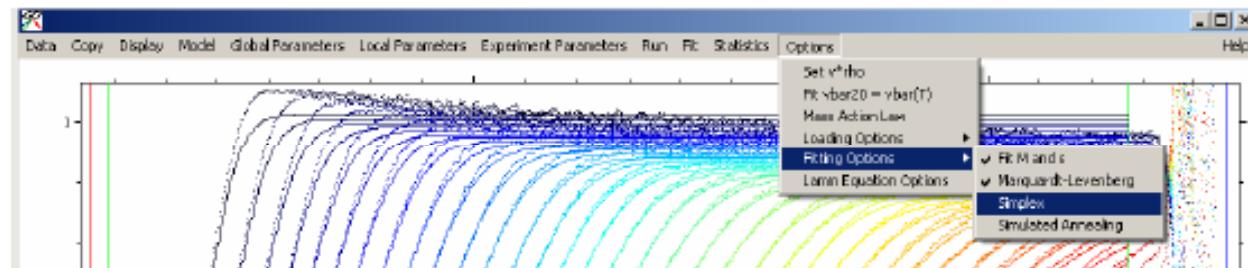
Again, we change back to Marquardt-Levenberg, to take advantage of its better ‘homing-in’ capability. Click on the “Marquardt-Levenberg” menu function to toggle.

Then do a FIT

This improved it slightly (judged by rmsd). Use control-U or the menu function to UPDATE CURRENT CONFIGURATION. This will save the fit.



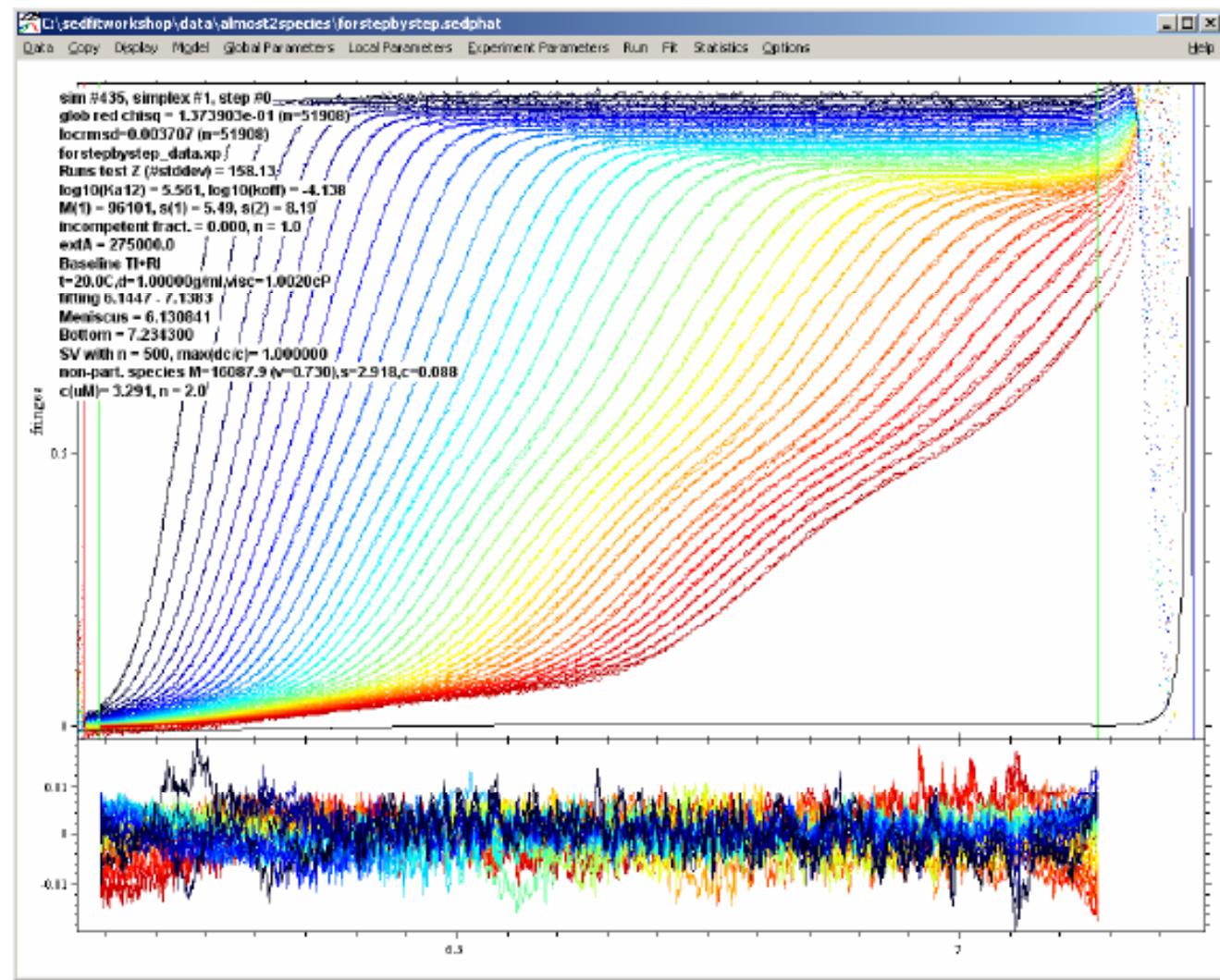
Again, click on the “Simplex” menu function. This will toggle off Marquardt-Levenberg and switch on the Simplex algorithm



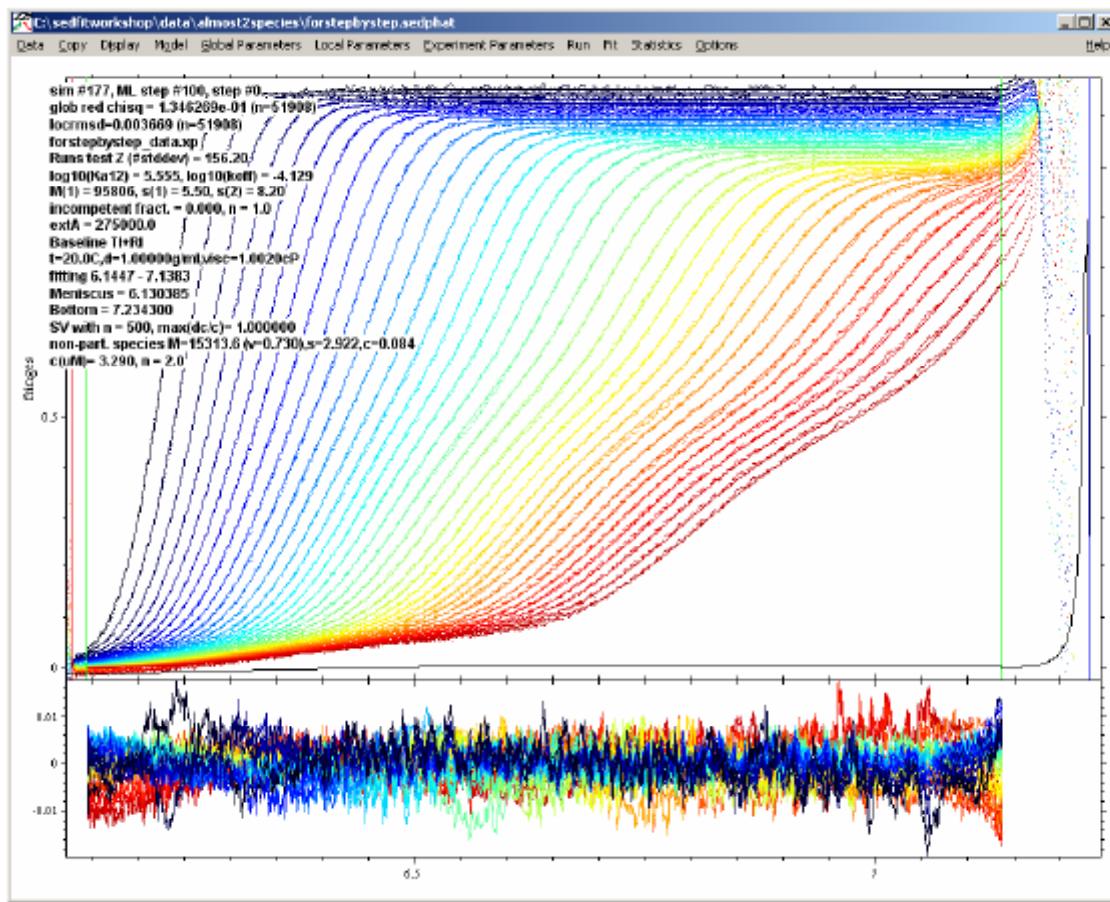
Then do a FIT

I find this still improves the fit slightly. After convergence gets very slow, I stopped it again and restarted the FIT. The reason for this is to take advantage of the random seeding of points each time you start a Simplex fit. This helped to get the rmsd further down. I repeat this procedure, but the second time it didn't seem to help.

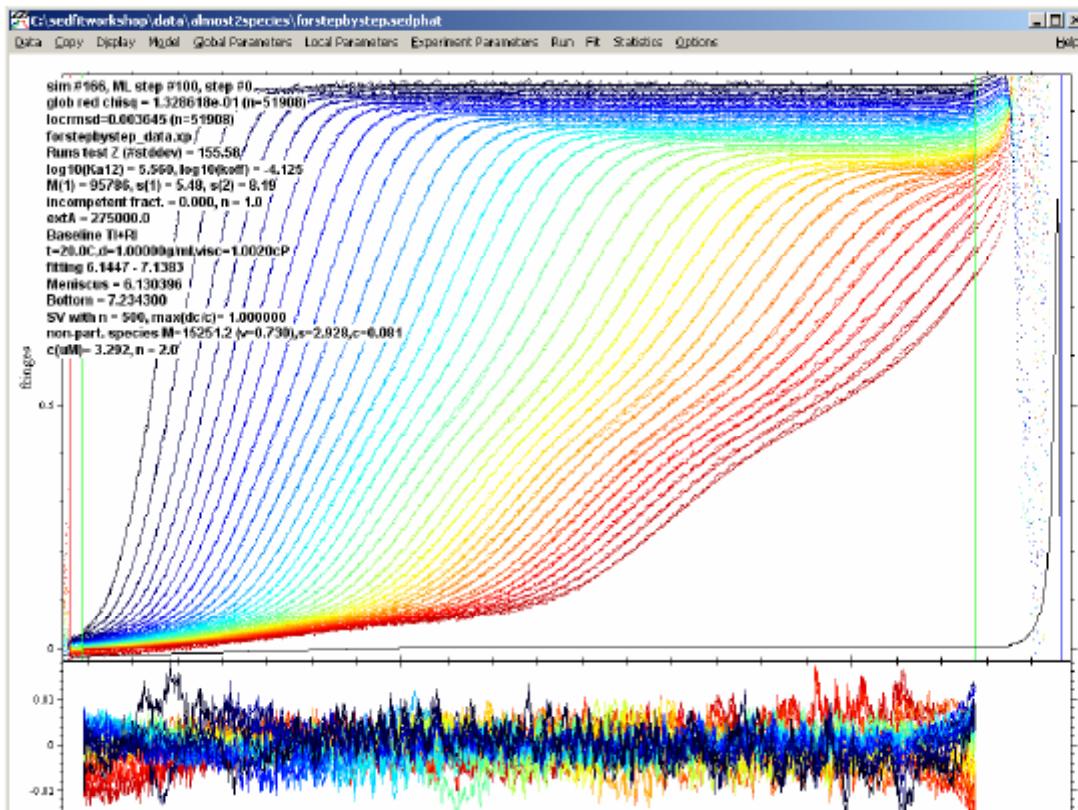
Save the configuration (control-U).



Then I switched back to the Marquardt-Levenberg method and tried again that way. This is the result: ***Further Improvement!!*** This is obviously a very badly shaped error surface, with a rather shallow valley that has many small ripples.



Same thing, again, after another round of alternating Simplex and Marquardt-Levenberg fits.



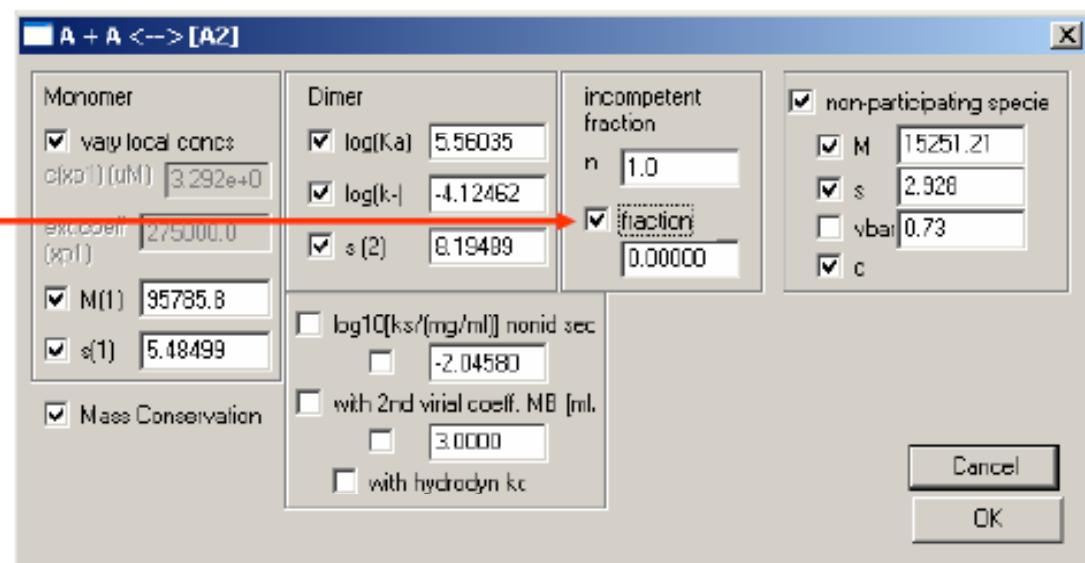
We could proceed in this way until we don't see any further improvement with any method. Another option would be to use the Simulated Annealing method, which I want to demonstrate in the following.

Simulated Annealing is better for error surfaces with many local minima, but more time-consuming. We could have used the Simulated Annealing method right from the beginning, but I didn't know yet that the error surface was so badly shaped.

At this point, let's float two more parameters that potentially could play a role in this fit:

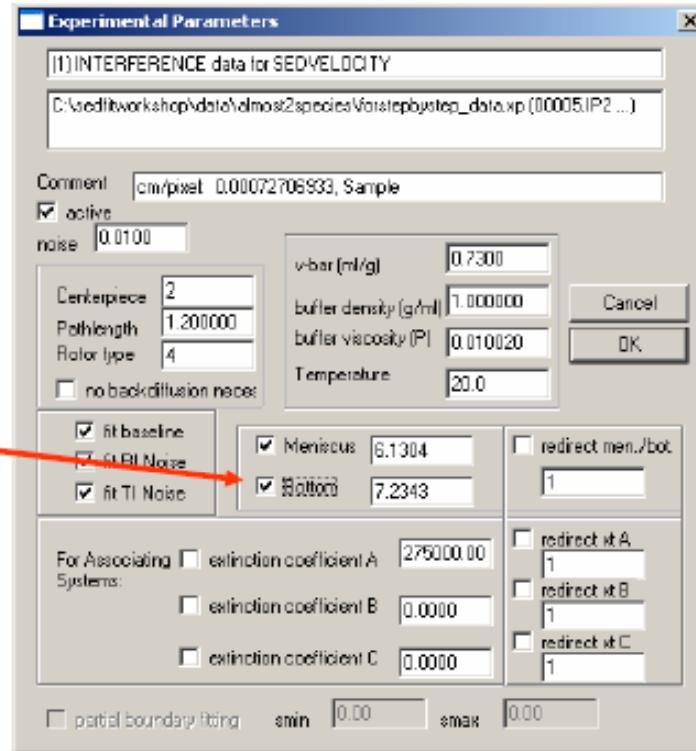
in the **global parameters**:
switch on the incompetent fraction

using only one experiment from a single concentration is not really very good to determine this parameter (because it likely will be correlated with the K_d), but the fine-points of the boundary shape could carry some information about it.

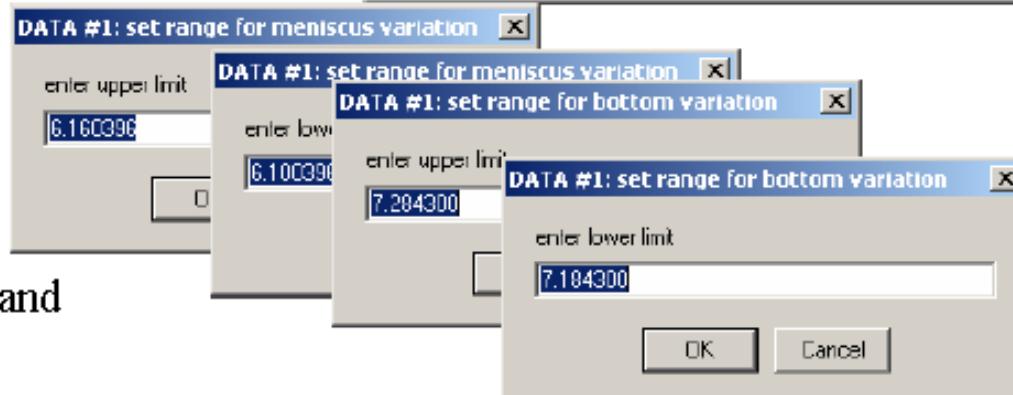


in the **experiment parameters**:
switch on floating of the
bottom.

even though we should not see back-diffusion from the main species, the contaminating smaller Mw species may experience back-diffusion slightly influencing our residuals close to the right fitting limit.

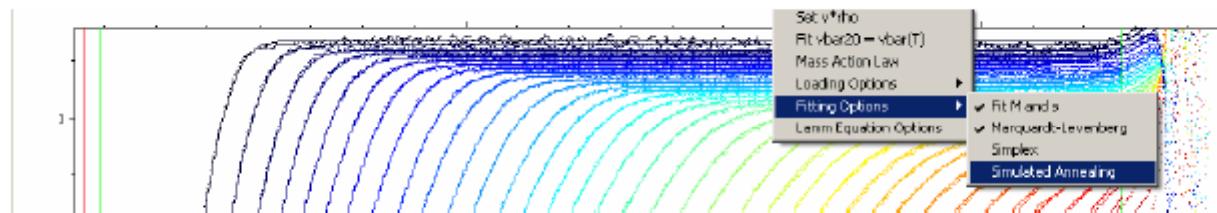


accept the default
limits for meniscus and
bottom variation.



before doing anything else, use **SAVE CURRENT CONFIGURATION AS**
and save under a different filename (I used 'forstepbystepSA').

Click on the ‘Simulated Annealing’ function. This will toggle it on (and the other algorithms off)



The following boxes will show up. Accept the default for all.

Fit with Simulated Annealing

starting 'temperature' (rel increase glob.chisq)
2.000000

OK Cancel

Fit with Simulated Annealing

of temperature steps
10.000000

OK Cancel

Fit with Simulated Annealing

of iterations per temperature
500.000000

OK Cancel

Fit with Simulated Annealing

final ML (1), Simplex (2), or none (0)
1.000000

OK Cancel

This will set the starting tolerance for accepting worse points during the procedure. A value of 2 means that points within a chisquare of twofold the starting chisquare will be tolerated.

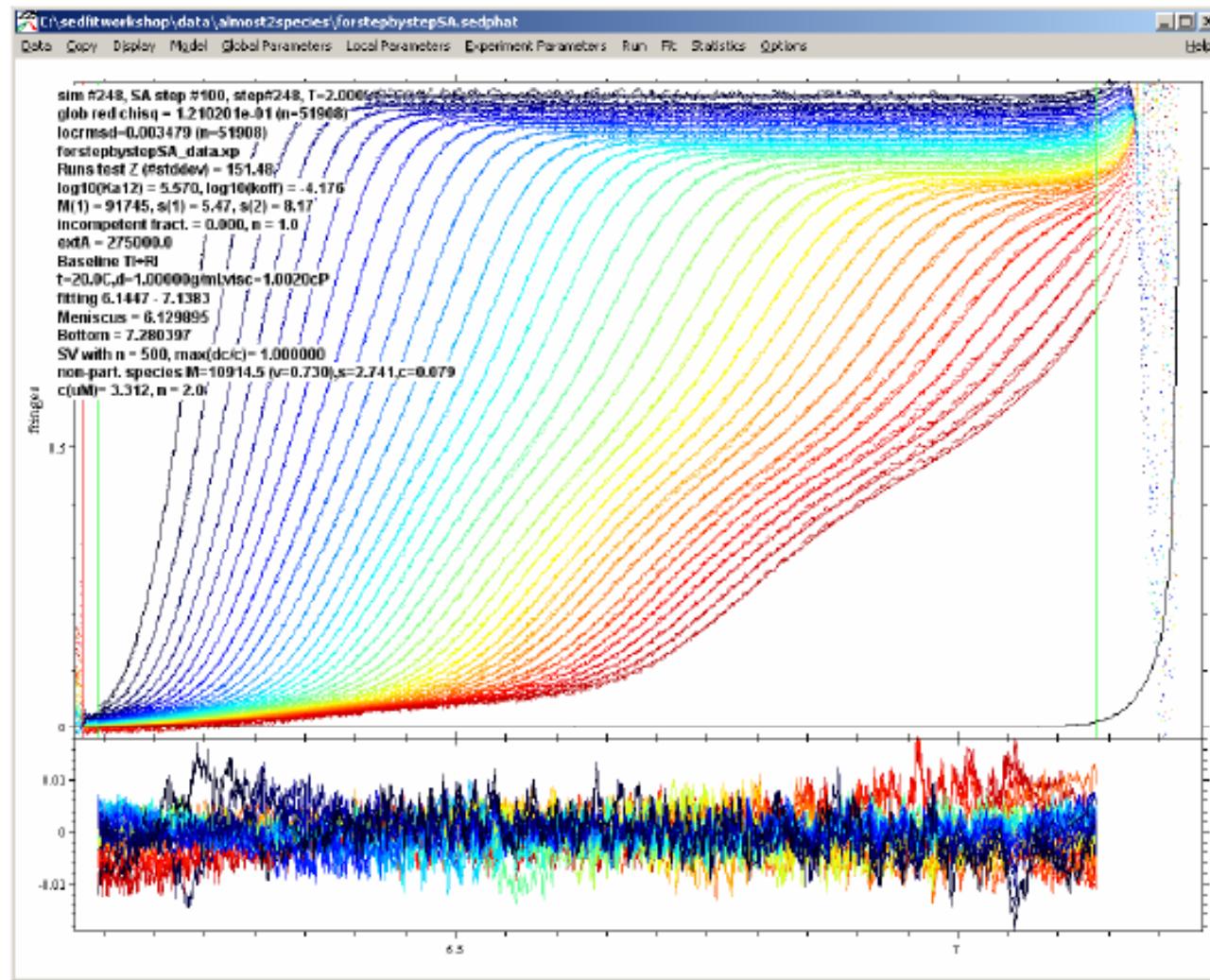
This determines the ‘temperature’ schedule, i.e. the number of times when we reduce the threshold chisquare for rejecting worse points.

Since we will not expect any form of convergence, we have to predefine the number of iterations to be done at each temperature. You can get a feel for a ball-park number by looking at how many iterations the Simplex algorithm used previously.

We know that we likely will need another algorithm to home-in. The combination with a final round of Marquardt-Levenberg optimization seems to be reasonable.

Then do a FIT
This will take a while.

This is where it ended up with SA – quite a bit better still. In fact, this is now better than the original c(s) fit we started with. We should continue searching with alternate Simplex, Simulated Annealing and/or Marquardt-Levenberg steps for better fits.



Finally, after a second trial with Simulated Annealing, no improvement was possible with Simplex and Marquardt-Levenberg. This is the best fit I found. Unfortunately, there's no guarantee that this is really the global optimum.

