

NOTE: This PDF file is for reference purposes only. This lab should be accessed directly from the web at <https://sassie-web.chem.utk.edu/docs/sassie-web-quick-start/quick-start.html>.

[Return to Main Documents Page](#)

## SASSIE-web: Quick Start

### Introduction

SASSIE-web is an online simulation and analysis tool for the modelling of biomolecular structures using small angle scattering data. It is based on the original, standalone, program SASSIE (Curtis *et al.* 2012) and retains all of its core features. This guide is designed to get you familiar with the basic features of the program as quickly as possible. The features covered will be:

- Monomer Monte Carlo simulation to create an ensemble of varied trial protein structures
- Calculation of theoretical scattering from protein models
- Comparison of theoretical and experimental scattering data
- Minimization of those structures that best fit the experimental scattering data

**Important Note:** Before you start this tutorial you will need to register for an account for and login to [SASSIE-web](#). Instructions on how to register can be found [here](#).

The only data needed to work through this tutorial are:

- A starting structure in PDB format: [gag\\_start.pdb](#)
- A PSF topology file: [gag\\_start.psf](#)
- An experimental scattering curve: [sans\\_data.sub](#)

You should download these files to your computer now. A good idea is to create a directory called *SASSIE-web-tutorial* and save all downloads from the tutorial in this location.

You will also need to familiarize yourself with a molecular viewing program that can display PDB and DCD files. We recommend [VMD](#) and provide a quick tutorial [here](#).

### SASSIE-web Interface

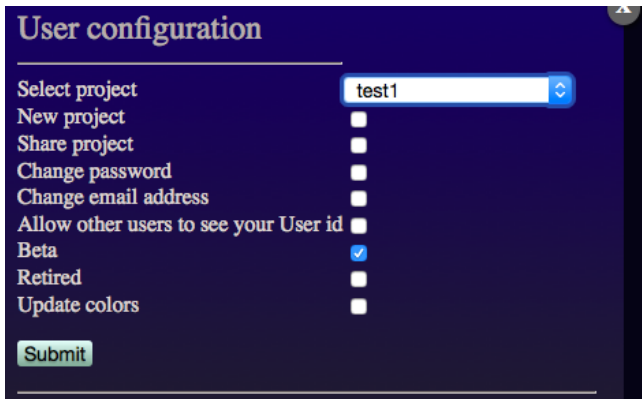
Once logged in to [SASSIE-web](#) the page should look something like the figure below.



- **Session Management:** Section at top right containing the following icons:
  - Filing Cabinet: File browser
  - Cogs: Job management
  - Head: User configuration
- **Main Menu Toggle:** Lines at top left:
  - Hides or shows the main menu
- **Main Menu:** Provides links to the different categories of modules available in SASSIE

During this tutorial, when instructed to select something from the **Main Menu** but no menu is visible on the left hand side of the page you must click on the **Main Menu toggle** to reveal it.

To choose a project name where your work will be stored and to access SASSIE modules that are still in Beta status, click on the Head icon. The User Configuration menu will appear.



User configuration

Select project

New project

Share project

Change password

Change email address

Allow other users to see your User id

Beta

Retired

Update colors

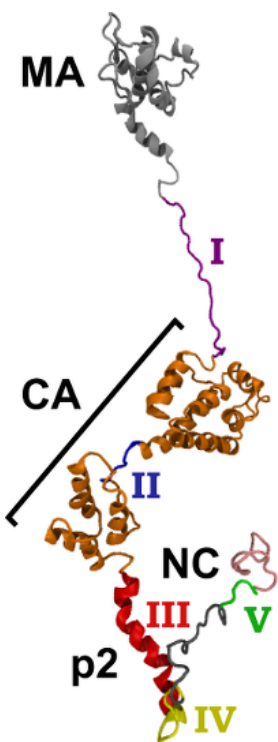
Submit

Choose an existing project or create a new project name. In addition, select the Beta checkbox. (You can also choose to select the Retired checkbox to access retired SASSIE modules. In addition, you can update the background and foreground screen colors by selecting the Update colors checkbox.) Click the 'Submit' button to connect to the project and to access the Beta modules. The project name should now appear at the top of the web page. Here, we have chosen project name 'test1'.



## Starting Structure

In this tutorial we will model the conformation of the HIV-1 Gag protein following the study of Datta *et al.* 2007. HIV Gag is a long polyprotein which is cleaved to form the functional proteins required by the virus. The viral proteins which form the domains in Gag are the matrix (MA), capsid (CA), p2 and nucleocapsid (NC). A structure stitched together from evidence based on crystal structures and models of the individual domains is shown below. A similar structure will be used as the starting structure for our simulations.



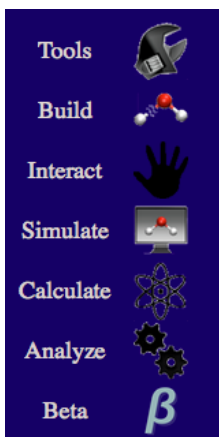
Domain	Flexible Region	Residues
MA		1 - 122
linker	I	123 - 144
CA		145 - 276
	II	277 - 282
		283 - 353
p2	III	354 - 377
linker	IV	378 - 389
NC		390 - 407
	V	408 - 412
		413 - 432

Datta *et al.* 2007 identified 5 flexible regions (labelled I-V), which are highlighted in the figure above. The table alongside the picture shows the residues which make up each region (we are going to need this information to select regions to be varied when we run Monte Carlo simulations).

## Data Interpolation

Data interpolation is necessary to create a new data file that is spaced on a uniform grid from the experimental data file. More information on the module is available in the [Data Interpolation documentation](#). Here we interpolate the SANS data from a HIV-1 Gag protein at a concentration of 1 mg/ml in a 100% D2O buffer.

Select the 'Tools' button at the top of the **Main Menu**. Click on **Main Menu Toggle** if necessary.



This will reveal a list of buttons for each tool running across the top of the page (just below the top bar with the **Session Management** and **Main Menu Toggle** icons).

Select the 'Data Interpolation' button from this menu.

You should now see a page like the one below. This page is used to enter all of the information needed to do the data interpolation.

### Data Interpolation

run name	<input type="text" value="run_0"/>	
experimental data file	<input type="button" value="Browse..."/> <input type="text" value="sans_data.sub"/>	or <input type="button" value="Browse server"/> Local: <input type="text" value="sans_data.sub"/>
output file name	<input type="text" value="sans_data.dat"/>	
I(0)	<input type="text" value="0.04"/>	<input type="button" value="↑"/> <input type="button" value="↓"/>
I(0) error	<input type="text" value="0.001"/>	<input type="button" value="↑"/> <input type="button" value="↓"/>
new delta q	<input type="text" value="0.02"/>	<input type="button" value="↑"/> <input type="button" value="↓"/>
number of new q-values	<input type="text" value="16"/>	<input type="button" value="↑"/> <input type="button" value="↓"/>

The figure shows the values for each field as required for data interpolation.

Edit the values on your screen to match the screenshot. An explanation of the field and how to edit it can be found below.

- run name:** user defined name of folder that will contain the results.
- experimental data file:** Name of input file with experimental data with at least three columns: q, I(q), and error in I(q). Here we use the *sans\_data.sub* file.
- output file name:** Name of file that will contain the interpolated data. Here we choose the name *sans\_data.dat*.
- I(0):** Experimentally determined value of scattering intensity at q = 0. Here we used the value of 0.04 that was derived from a Guinier fit to the data.
- I(0) error:** Experimentally determined value of the error of the scattering intensity at q = 0. Here we use the value of 0.001 that was obtained from the Guinier fit to the data.
- new delta q:** Desired spacing of q-values (1/Angstrom). This should be chosen so that your first interpolated data point falls within the q-range of the experimental data. For this tutorial, the value has been set to 0.02 since the first data point occurs at a value of ~0.013.
- number of new q-values:** Integer number of desired q-values. For this tutorial, the value has been set to 16 to that the maximum q value is 0.3.

Once you have understood the input fields and made sure that your values agree with the figure click on the 'Submit' button to start simulation.

As the run continues the progress bar beneath the submit button should update. A graph beneath this should will show the variation of the radius of gyration over the steps of the Monte Carlo simulation. Once complete the output should look similar to the figure below.

```

=====
DATA FROM RUN: Fri Jul  8 17:48:12 2016

Input file used : results/users/skrueger/test1/sans_data.sub

Interpolated data were written to ./run_0/data_interpolation/sans_data.dat

Interpolated data with S/N > 2 were written to ./run_0/data_interpolation
/stn_sans_data.dat

delta q = 0.020000 (1/Å)

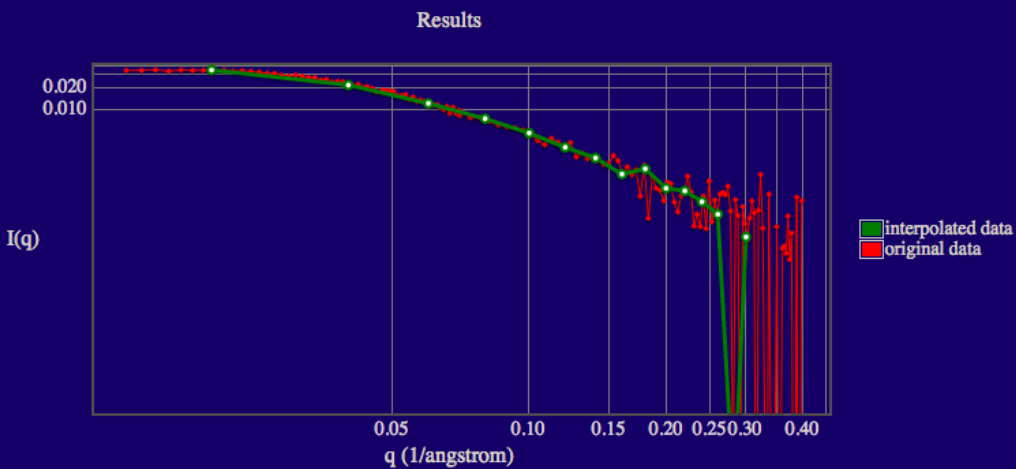
number of q-points = 16

q-range: 0 to 0.300000 (1/Å)
=====

```

progress:   
percent done: 100.0

original and interpolated data



The output will show a plot of the original and interpolated data, the name of the input file and the name of the interpolated data file as well as the directory in which it is located.

Note that roll-over help will indicate options to resize, zoom and reset the view of the plot.

### What have we generated:

test1/run\_0/data\_interpolation

- *sans\_data.dat*: text file containing the interpolated SANS data
- *stn\_sans\_data.dat*: text file containing the interpolated SANS data truncated at the q-value where the signal-to-noise drops below a value of 2

## PDB Scan

PDB Scan is used to assess whether an input PDB is ready for simulation and where possible to provide files enabling CHARMM forcefield parameterization. Information on missing atoms and residues and those not covered as standard by the CHARMM 27 forcefield are reported. PDB files do not need to have header information. At this time, only PDB files of proteins are supported. More information on the module is available in the [PDB Scan documentation](#). Here we examine the PDB file that describes the starting HIV-1 Gag protein structure.

Select the 'Build' button from the Main Menu of SASSIE-web and then click on the PDB Scan button.

You should now see a page like the one below. This page is used to enter all of the information needed to check the PDB file.

## PDB Scan

Only PDB files of proteins are supported

run name

pdb file input  or

**pdb file input:** The PDB file that we want to examine. Here we use the *gag\_start.pdb* file.

Once you have entered the file name, click on the 'Submit' button to start the file scan.

As the run continues the progress bar beneath the submit button should update. Once complete the output should look similar to the figure below.

```
##Structure Preparedness for Simulation
-----Chai
n Moltype      Missing      Non-CHARMM      Missing Heavy
Non-CHARMM     X      Protein      Residues      Residues      Atoms
Atoms
-----
0
-----
There are no critical issues with this PDB, it is ready to
simulate.
However, check the warnings to ensure the PDB contains the structure you are
interested in.

### Warnings
1. No information available about biological unit.
```

The text output region provides a brief summary of the PDB Scan report.

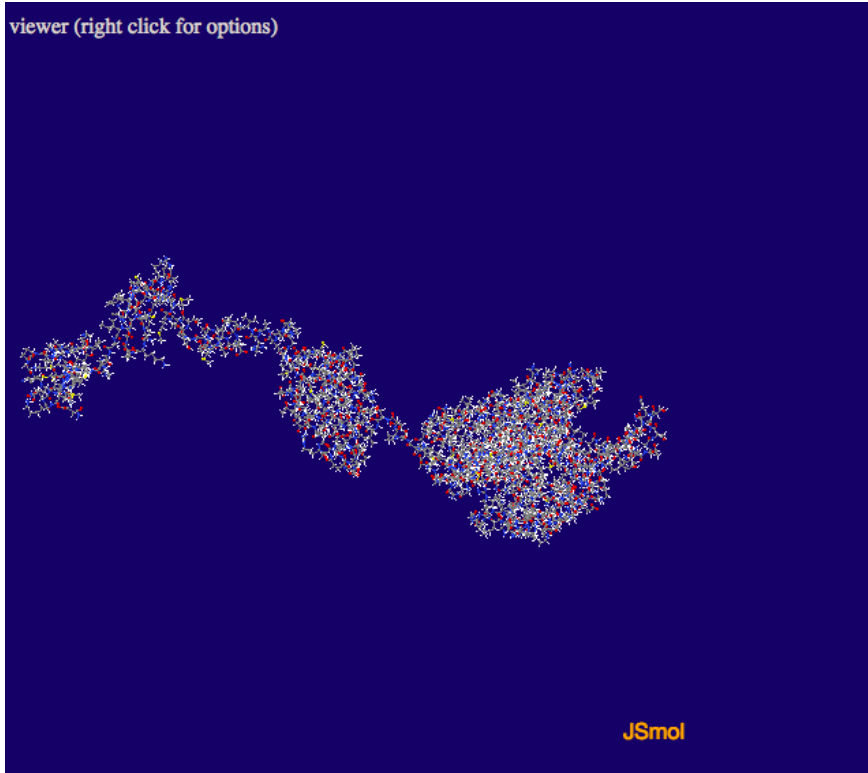
### What have we generated:

test1/run\_0/pdbscan

- *coor\_chain\_X\_gag\_start.pdb*: output PDB file; same as the input PDB file
- *gag\_start\_pdbscan.html*: HTML file containing the full PDB scan report

### Visualization

A JSmol visualization of the protein is produced and is shown below the text output region. Holding down the left mouse button and moving the cursor over the picture allows you to rotate the view, the scroll wheel facilitates zooming in and out. Right clicking on the image allows you to access all of the [JSmol](#) options.



The full PDB scan report can be found below the image of the structure.

## Structure Preparedness for Simulation

Chain	Moltype	Missing Residues	Non-CHARMM Residues	Missing Heavy Atoms	Non-CHARMM Atoms
X	Protein	0	0	0	0

There are no critical issues with this PDB, it is ready to simulate.

However, check the warnings to ensure the PDB contains the structure you are interested in.

### Warnings

- No information available about biological unit.

### Polymers

Chain ID	Type	No. Residues Found	Gaps Detected
X	Protein	431	No

### Heterogens

None

### System statistics

Calculated from PDB coordinates

#### Statistics:

Formula C2080H3378O632N616S24

Molecular weight 47897 Da : 47.90 kDa

Center of mass (22.84, -5.55, 1.55) Angstrom

(x,y,z)

Coordinate range Low: (-72.30, -36.67, -34.66), High: (81.54, 35.74, 44.40) Angstrom

(x,y,z)

Radius of gyration 44.47 Angstrom

Distance between N- & C-terminals 135.16 Angstrom

## Chain X

### Sequence:

```
123456789 | 123456789 | 123456789 | 123456789 | 123456789 |
0 GARASVLSGGELDKWEKIRLRPGGKKQYKPKHIVWASRELERFAVNPGLL
50 ETSEGCQRQLGQLQPSLQTGSEELRSLYNTIAVLYCVHQRIDVKDTKEAL
100DKIEEBEQNKSKKKAQQAADTGNNNSQVSQNYPIVQNLQGQMVHQAISPR
150LNAWVKVVEEKAFSPEVIMFSALEGATPQDLNMLNTVGGHQAAMQML
200KETINEEAAEWDVHPVHAGPIAPGQMRPRGSDIAGTTSITLQEQIGWMT
250NNPPIPVGGEIYKRWILLGLNKIVRMYSPSILDIRQGPKEPFRDYVDRFY
300KTLRAEQASQEVKNAATELLVQANPDCKTILKALGPAATLEEMMTACQ
350GVGGPGHKARVIAEAMSQVTSNATIMMQGNFRNRKTVKCFNCGKEGHI
400AKNCRAPRRKGCWKCCKGEGHQMCDTERQAN
```

### Gaps:

None

### Statistics:

Formula C<sub>2080</sub>H<sub>3378</sub>O<sub>632</sub>N<sub>616</sub>S<sub>24</sub>  
Molecular weight 47897 Da : 47.90 kDa  
Center of mass (22.84, -5.55, 1.55) Angstrom  
(x,y,z)  
Coordinate range Low: (-72.30, -36.67, -34.66), High: (81.54, 35.74, 44.40) Angstrom  
(x,y,z)  
Radius of gyration 44.47 Angstrom  
Distance between N- & C-terminals 135.16 Angstrom

This PDB file is ready for simulation so we can proceed to create an ensemble of structures for comparison to the SANS data.

## Structure Variation - Monte Carlo (Monomer)

The primary way to vary structures in SASSIE is via Monte Carlo simulations which rotate the backbone dihedral angles of flexible regions within proteins to sample a wide range of structures. More information can be found in the [Monomer Monte Carlo documentation](#). Here we setup and run such a simulation before visualizing the range of structures produced.

Select the 'Simulate' button from the Main Menu of SASSIE-web and then click on the 'Monomer Monte Carlo' button.

You should now see a page like the one below. This page is used to enter all of the information needed to run a Monte Carlo simulation.

### Monomer Monte Carlo

run name	<input type="text" value="run_0"/>	
reference pdb	<input type="button" value="Browse..."/> gag_start.pdb	or <input type="button" value="Browse server"/> Local: gag_start.pdb
output file name (dcd)	<input type="text" value="run_0.dcd"/>	
number of trial attempts	<input type="text" value="1000"/>	
return to previous structure	<input type="text" value="10"/>	
temperature (K)	<input type="text" value="300.0"/>	
molecule type	<input type="text" value="protein"/>	
number of flexible regions to vary	<input type="text" value="5"/>	
residue range for each flexible region	<input type="text" value="123-144,277-282,354-374,378-389,408-412"/>	
maximum angle(s)	<input type="text" value="30.0,30.0,30.0,30.0,30.0"/>	
structure alignment range	<input type="text" value="284-350"/>	
overlap basis	<input type="text" value="heavy atoms"/>	

---

Advanced Input

Check Box for Advanced Input

The figure shows the values for each field as required for our simulation. An explanation of some of the fields can be found below.

**reference pdb:** The starting structure for the simulation. Here we use the *gag\_start.pdb* file.

**number of trial attempts:** Number of times the simulation will try to vary the structure (some structures will be discarded by the Monte Carlo algorithm) For this tutorial set the value to 1000. For real studies tens of thousands of structures are needed.

**return to previous structure:** Number of discarded structures in a row that are considered before returning to a randomly-selected structure that was previously accepted

**number of flexible regions to vary:** single number

**residue range for each flexible region:** comma-separated list of the range of residues to vary for each flexible region

**maximum angle(s):** comma-separated list of the maximum angle sampled in a single Monte Carlo step for each flexible region

**structure alignment region:** a single range of residues for structural alignment of all the flexible segments. This makes it easy to make visual comparisons of each frame in the output trajectory.

**overlap basis:** Select either heavy atoms, all, backbone or enter atom name. The atom name option will spawn further inputs:

- **overlap basis:** Enter an atom name to check for overlap.
- **overlap cutoff (angstroms):** Overlap basis atoms closer than this distance defines an overlap condition.

Once you have understood the input fields and made sure that your values agree with the figure click on the 'Submit' button to start simulation.

As the run continues the progress bar beneath the submit button should update. A graph beneath this should will show the variation of the radius of gyration over the steps of the Monte Carlo simulation. Once complete the output should look similar to the figure below.

```

=====
DATA FROM RUN: Fri Jul 10 09:40:22 2015


Average accepted rg2 = 53.429966

Configurations and statistics saved in ./run_0/monomer_monte_carlo/ directory

lowest Rg = 37.832101   highest Rg = 67.812830
accepted 692 out of 1000 : 69.200000 percent
overlap check discarded 308 out of 1000 moves : 30.800000 percent
Rg cutoffs discarded 0 out of 1000 moves : 0.000000 percent

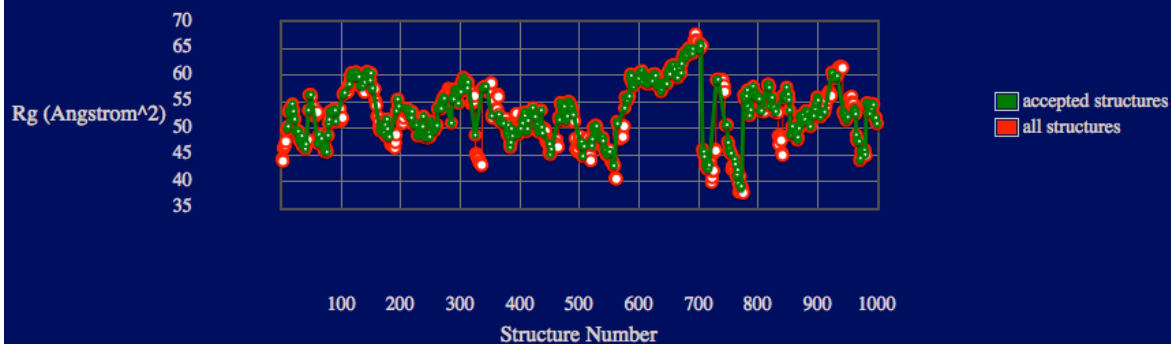
minimum x = -106.000069   maximum x = 144.011806 -> range: 250.011875 Angstroms
minimum y = -102.160418   maximum y = 95.319356 -> range: 197.479775 Angstroms
minimum z = -112.573038   maximum z = 78.222151 -> range: 190.795189 Angstroms
=====

```

progress:  100.0  
percent done:

all rg and accepted rg data

### Rg Results



## What have we generated:

test1/run\_0/monomer\_monte\_carlo

- *run\_0.dcd*: DCD file containing all of the structures accepted by our Monte Carlo simulation
- *run\_0.dcd.accepted\_rg\_results\_data.txt*: text file containing radius of gyration for all structures that made it into the DCD file
- *run\_0.dcd.all\_rg\_results\_data.txt*: text file containing radius of gyration for all structures generated
- *run\_0.dcd.stats*: text file containing statistics for our Monte Carlo run

## Visualization

You should now download the output trajectory using the file browser.

- Click on the filing cabinet icon in the **Session Management** area.

Note: the 'Configurations and statistics saved in!' line in the output gives a relative path under the project directory.

- Click on the triangle next to the 'test1' directory name to reveal the 'run\_0' directory created by our simulation.
- Check the box next to 'run\_0' to select it for download.
- Beneath the file tree is an option labelled 'Compression type'. Select an option suitable for your operating system from the list (for Windows select 'zipped' for Linux 'bzip2 tarball').

- Click the 'Download' button.

A progress bar will appear monitoring the upload of your files to the server. Once complete a link will appear beneath the download button.

- Click on this to download your data.

Once the download is complete uncompress the file in a location of your choice.

You should now load the PDB [gag\\_start.pdb](#) (you will find this in the run\_0 directory you just downloaded) and DCD (run\_0/monomer\_monte\_carlo/run\_0.dcd) into VMD to observe the variation produced even in our very short Monte Carlo simulation. Remember the DCD file contains coordinates alone, you need to load the PDB first so that the visualization software knows about the atoms they represent and how they are connected. You can also download and visualize the DCD file [run\\_0.dcd](#) generated from this particular Quick Start run for comparison.

## Initial SAS Curve Calculation

Next we calculate a theoretical scattering curve for each of the trial structures we have generated. The SASSIE workflow operates by calculating the scattering intensities at evenly spaced Q values and matching these against interpolated experimental values.

The file [sans\\_data.dat](#) contains our previously interpolated experimental data. In order to create the correct data points in our theoretical curves we need three pieces of information:

- Intensity and  $Q=0$ ,  $I(0)$ : 0.04
- Maximum value of Q: 0.3 (units are inverse Angstroms)
- Number of points in the curve: 16

A number of scattering calculators are available in SASSIE. Here we use SasCalc. More information can be found in the [SasCalc documentation](#). The starting structure **must** be a complete structure without missing residues or atoms (including hydrogen atoms) in order to obtain accurate scattering profiles. Atom and residue naming must be compatible with those defined in the CHARMM force field.

- Select 'Beta' from the Main Menu.
- Click the 'SasCalc' button.

Now you need to enter the information to run the scattering calculator. SasCalc can be used to calculate the scattering for SAXS and SANS and/or for several SANS contrasts at the same time.

The module is first run using the "converged number of golden vectors" option on just one structure. Choose this option from the **SasCalc method** menu in the Advanced Input section of the page.

Other than the values listed above you can keep the default values for this tutorial (see figure below).

### SasCalc

run name:

reference pdb:  gag\_start.pdb or  Local: gag\_start.pdb

trajectory file filename (dcd or pdb):  gag\_start.pdb or  Local: gag\_start.pdb

number of q values:

maximum q value:

---

**Neutron input**

number of contrast points:

D2O percentage [1]:  I(0) [1]:

number of exchangeable H regions:

exchangeable H region [1]:  fraction of exchangeable H [1]:

number of deuterated regions:

---

**X-ray input**

---

**Advanced Input**

SasCalc method:

tolerance of runtime average convergence:

check box to enable HyPred pRDF solvent model:

The single structure that we used to start the simulation is used as both the **reference pdb** and the **trajectory file filename** (PDB in this case) so it is already uploaded to the SASSIE-web server. Thus, you can either upload it again from your local computer or locate it on the server and read it from there.

To read the file from the server:

- Click on the 'Browse server' button next to the appropriate field.
- Navigate to the file test1/run\_0/monomer\_monte\_carlo/gag\_start.pdb
- Click 'OK'

Once you have understood the input fields and made sure that your values agree with the figure click on the 'Submit' button to start the calculation.

A scattering curve will be calculated for the starting structure (the progress bar should reach 100% and a message stating the run finished appear in the window beneath when the job has completed). Note that the files are written to a sub-directory of sascal/ that is named according to the D2O percentage in the solvent. This is useful when calculating the scattering curves for more than one contrast.

```
=====
DATA FROM RUN: Thu Jul  7 14:03:15 2016

Processed 1 DCD frame(s)

Data stored in directory: run_0/sascalc/neutron_D2Op_100

=====

progress:
percent done: 100.0
```

### What have we generated:

test1/run\_0/sascalc/neutron\_D2Op\_100

- *run\_0\_00001.iq*: files containing theoretical scattering data
- *run\_0\_00001.log*: log files containing information about the structure and calculation inputs
- *D2Op\_100.pdb*: input PDB file with element names including deuterium atoms that were added as a result of H-D exchange of exchangeable hydrogen atoms or deuteration of non-exchangeable hydrogen atoms.
- *HDexchange\_Info\_D2Op\_100.txt*: file describing which hydrogen atoms were replaced with deuterium as a result of H-D exchange of exchangeable hydrogen atoms.
- *Deuteration\_Info\_D2Op\_100.txt*: file describing which hydrogen atoms were replaced with deuterium atoms as a result of deuteration of non-exchangeable hydrogen atoms.

## SECOND SAS Curve Calculation

The run\_0\_00001.log file from the initial SAS Curve calculation indicates that 35 golden vectors were required for convergence to the desired tolerance (0.01 in this case).

```
#Structural Information:
1. PDB input file = /share/apps/genapp/sassie2/results/users/skrueger/test1/gag_start.pdb
2. Number of atoms = 6730; Mw = 48681.42
3. Dimensions = x: -72.30, 81.54, y: -36.67, 35.74, z: -34.66, 44.40
4. Maximum radial dimension Dmax = 187.52 A (contrast-independent)
5. Molecular center of mass = x: 22.80, y: -5.54, z: 1.55 (contrast-independent)
6. Molecular Rg = 44.465386 A (contrast-independent)
#Scattering Intensity:
7. Source = neutron
8. D2O % = 100.000
9. Non-exchangeable H deuteration file = run_0/sascalc/neutron_D2Op_100/HDexchange_Info_D2Op_100.txt
10. Exchangeable H deuteration file = run_0/sascalc/neutron_D2Op_100/Deuteration_Info_D2Op_100.txt
11. PDB output file with explicit D atoms = run_0/sascalc/neutron_D2Op_100/D2Op_100.pdb
12. I(q) .vs. q file = run_0/sascalc/neutron_D2Op_100/run_00_00001.iq
13: Convergence tolerance used = 0.010000
    Converged number of golden vectors (for complete scattering profile) = 35
14. Io = 0.040
15. Center of Mass = 24.03, y: -5.94, z: 1.25 (contrast-dependent)
16. Rg = 43.038362 A (contrast-dependent)
```

We now use this information to calculate the scattering curves for all of the generated structures using the "fixed number of golden vectors" option from the **SasCalc method** menu as shown below.

**SasCalc**

run name:

reference pdb:  gag\_start.pdb or  Local: gag\_start.pdb

trajectory file filename (dcd or pdb):  run\_0.dcd or  Local: run\_0.dcd

number of q values:

maximum q value:

---

**Neutron input**

number of contrast points:

D2O percentage [1]:  I(0) [1]:

number of exchangeable H regions:

exchangeable H region [1]:  fraction of exchangeable H [1]:

number of deuterated regions:

---

**X-ray input**

---

**Advanced Input**

SasCalc method:

number of golden vectors:

check box to enable HyPred pRDF:

solvent model:

The **reference pdb** is the starting structure and the **trajectory file filename** (DCD) comes from the result of the Monte Carlo simulation, so both are already on the SASSIE-web server. Thus, you can either upload them again from your local computer or locate them on the server and read them from there.

When all input fields are complete:

- Click 'Submit'

```
=====
DATA FROM RUN: Thu Jul  7 15:55:05 2016

Processed 692 DCD frame(s)

Data stored in directory: run_0/sascalc/neutron_D2Op_100
=====

progress:
percent done: 100.0
```

A scattering curve will be calculated for all of the structures generated by the Monte Carlo simulation (the progress bar should reach 100% and a message stating the run finished appear in the window beneath when the job has completed). Note that the files written during the initial SAS curve calculation will be **overwritten** since we chose the same run name (run\_0) in both cases. If you wish to save the files from the initial calculation, use a different run name.

### What have we generated:

test1/run\_0/sascalc/neutron\_D2Op\_100

- \*.iq: files containing theoretical scattering data for all frames in the DCD file (692 files in this case).
- \*.log: log files containing information about each structure and calculation inputs (692 files in this case).
- D2Op\_100.pdb: input PDB file with element names including deuterium atoms that were added as a result of H-D exchange of exchangeable hydrogen atoms or deuteration of non-exchangeable hydrogen atoms.
- HDexchange\_Info\_D2Op\_100.txt: file describing which hydrogen atoms were replaced with deuterium as a result of H-D exchange of exchangeable hydrogen atoms.
- Deuteration\_info\_D2Op\_100.txt: file describing which hydrogen atoms were replaced with deuterium atoms as a result of deuteration of non-exchangeable hydrogen atoms.

### Initial SAS Curve Comparison

Now we compare our theoretical curves to the experimental data to see which of our structures are plausible models of the real protein using Chi-Square Filter. More information can be found in the [Chi-Square Filter documentation](#).

- Select 'Analyze' from the Main Menu.
- Click the 'Chi-Square Filter' button.

We now need to select the path containing the theoretical scattering curves and the file containing the experimental data. In addition we need to input the value of  $I(0)$  to enable comparison of the two curves (see the picture below).

### Chi-Square Filter

run name

interpolated data file  sans\_data.dat or  Local: sans\_data.dat

I(0)

SAS type  ▾

SAS data path

chi-square type  ▾

number of weight files

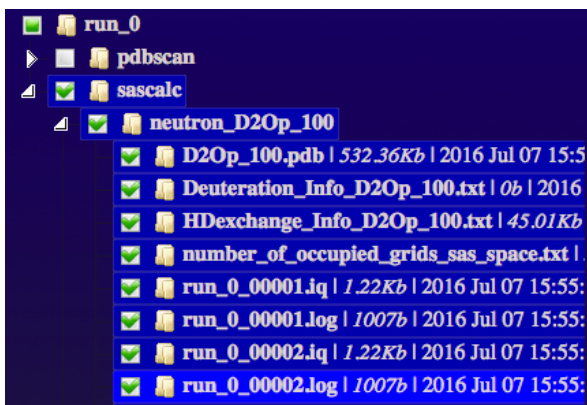
---

Advanced Input

Check Box for Advanced Input

To set the path to the scattering curves generated in the previous step:

- Click on the 'Browse server for a path' button
- Navigate to the test1 project folder and select the run\_0/sascalc/neutron\_D2Op\_100 folder (see picture below)



- Click 'OK'

#### interpolated data file

- Click on the 'Choose File' button
- Navigate to and select the [sans\\_data.dat](#) file on your local computer
- Click 'OK'

#### I(0)

- Enter the value 0.04

We eventually may want to create 'weight files' that record which frames meet criteria that make them successful models of our data. This means those with low chi square values. However, we don't know the range of chi square values we have at this stage. So, we set the 'number of weight files' to 0 at this time.

#### Sas type

- Choose SasCalc from the menu

#### number of weight files

- Set this to 0 using the down arrow associated with the input box. If you type '0' in the box, press the TAB button to make sure this value is accepted.

Note: There are list boxes that allow the selection of the format of the input theoretical curves and the metric used to compare the curves. Here we wish

to use the defaults of 'SasCalc' and 'reduced chi-square'.

Click 'Submit'.

Once complete you the run you should see outputs similar to those below.

```
=====
DATA FROM RUN: Thu Jul  7 19:49:16 2016

Data stored in directory: ./run_0/chi_square_filter/neutron_D2Op_100


PROCESSED 692 SAS FILES:

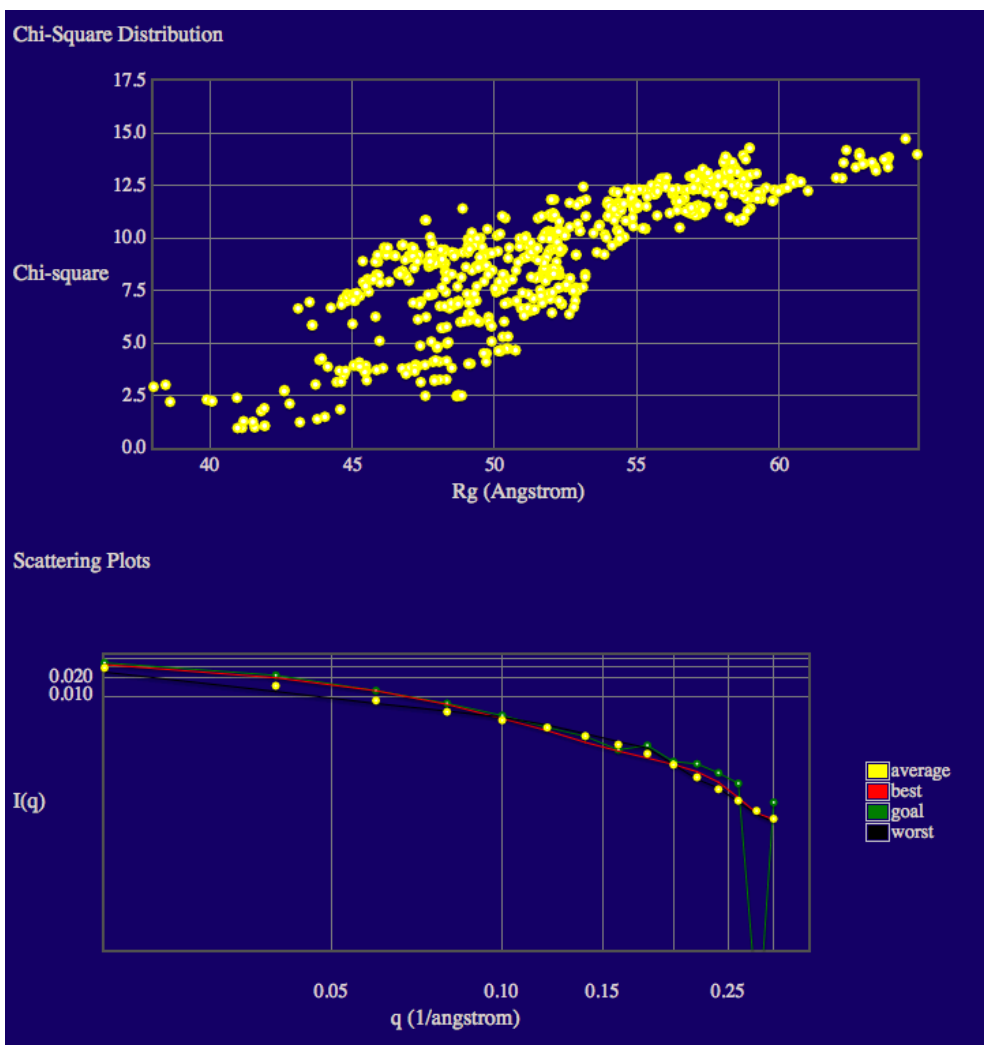
>> The BEST and WORST SAS spectra are in the file named : bestworstfile.txt
>> The AVERAGE SAS spectra is in the file named : averagefile.txt
>> Chi-square, Rg, and filename are in the file named : x2file.txt

BEST SINGLE STRUCTURE IS NUMBER 495 WTIH X2 = 0.946463 :      spectra:
run_0_00495

WORST SINGLE STRUCTURE IS NUMBER 486 WTIH X2 = 14.738848:      spectra:
run_0_00486

=====

progress:
percent done:  100.0
```



In the text output you will see the minimum chi square (X2) values is given.

The top plot shows the variation of chi squared (y-axis) with the radius of gyration (x-axis). Chi squared is a measure of the quality of fit of the theoretical curve to the experimental one. It is a percentage and the lower the value the better.

The bottom plot shows a direct comparison of the best, worst and average theoretical curves with experiment (goal).

### What have we generated:

test1/run\_0/chi\_square\_filter/neutron\_D2Op\_100

- *averagefile.txt*: average scattering curve for all structures
- *bestworstfile.txt*: best and worst scattering curves selected from all structures
- *sas\_spectra\_plot.txt*: goal, best, worst and average scattering curves
- *x2\_vs\_rg\_plot.txt*: chi squared against radius of gyration for all structures
- *x2file.txt*: chi squared for all structures

/spectra

- *spec\_\*.ciq*: scattering curves scaled to correct I(0) for each structure

## Second SAS Curve Comparison

Now that we know the range of chi square values that we have, we can compare the theoretical curves to the data a second time and create a weight file that flags all structures with chi square values below a certain number. Now, we set the 'number of weight files' to 1.

### Chi-Square Filter

run name

interpolated data file  sans\_data.dat or  Local: sans\_data.dat

I(0)

SAS type

SAS data path  Server: test1/run\_0/sascalc/neutron\_D2Op\_100

chi-square type

number of weight files

enter expression [1]

weight file name [1]

low Rg cutoff [1]

---

Advanced Input

Check Box for Advanced Input

**run name:**

- Since we already have a chi\_square\_filter folder in the run\_0 directory, set the run name to run\_1.

**number of weight files**

- Set this to 1 using the down arrow associated with the input box.

Weight files contain information on which frames in our simulation meet specific criteria provided in the expression box.

**enter expression**

- Enter the following expression:

x2 < 3.0

This selects all frames with a chi square less than 3.0. Adjust this value if necessary to suit the results from your simulation.

**weight file name**

- Enter x2\_lt\_3p0.txt

**low Rg cutoff**

- Enter a value if you wish to also restrict the Rg range to be above this value. The default value is 0 so that all Rg value are acceptable.
- Click 'Submit'.

Once complete you the run you should see outputs like those below.

```

=====
DATA FROM RUN: Thu Jul 7 20:25:48 2016

Data stored in directory: ./run_1/chi_square_filter/neutron_D2Op_100

PROCESSED 692 SAS FILES:

>> The BEST and WORST SAS spectra are in the file named : bestworstfile.txt
>> The AVERAGE SAS spectra is in the file named : averagefile.txt
>> Chi-square, Rg, and filename are in the file named : x2file.txt

BEST SINGLE STRUCTURE IS NUMBER 495 WITH X2 = 0.946463 :      spectra:
run_0_00495

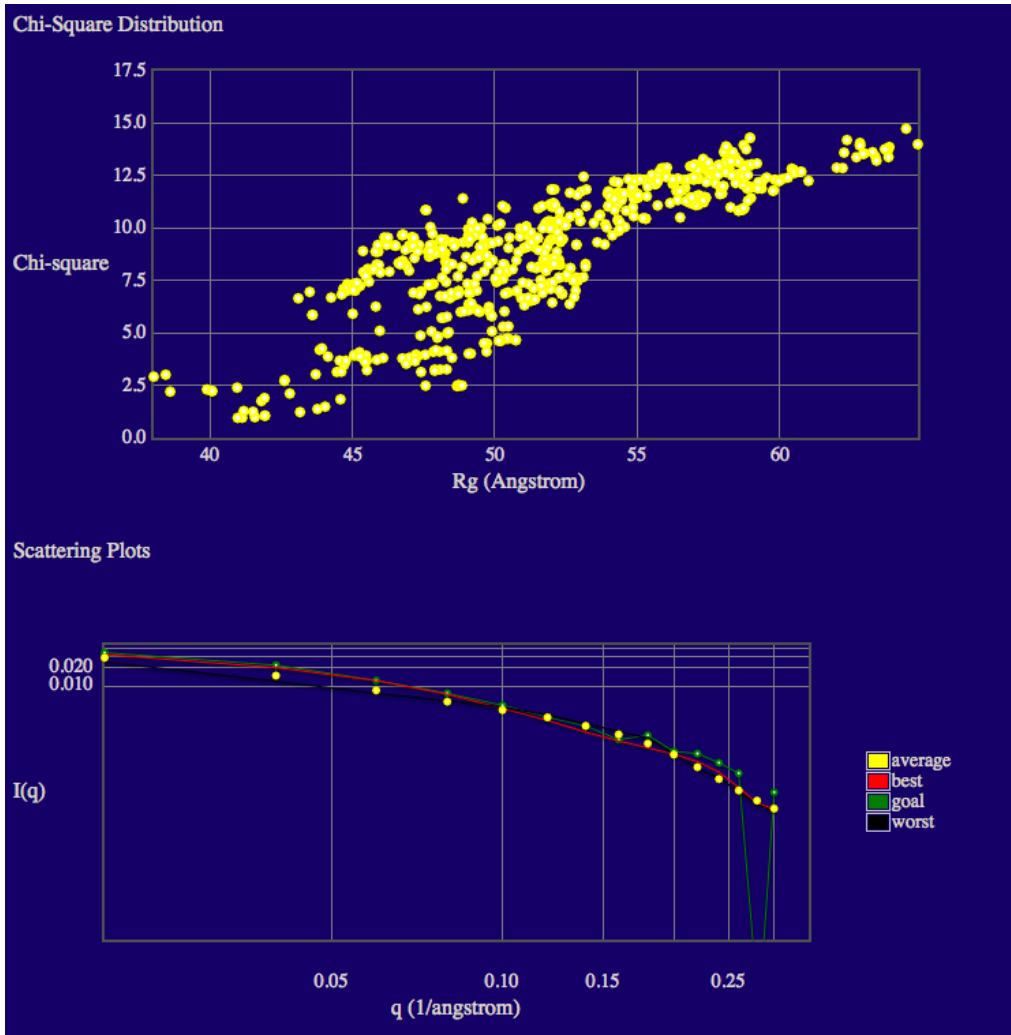
WORST SINGLE STRUCTURE IS NUMBER 486 WITH X2 = 14.738848:      spectra:
run_0_00486

=====

```

progress:

percent done: 100.0



These results are essentially the same as those from our first comparison above except that we have now generated a weight file.

**What have we generated:**

test1/run\_1/chi\_square\_filter/neutron\_D2Op\_100

- *averagefile.txt*: average scattering curve for all structures
- *bestworstfile.txt*: best and worst scattering curves selected from all structures
- *sas\_spectra\_plot.txt*: goal, best, worst and average scattering curves
- *x2\_vs\_rg\_plot.txt*: chi squared against radius of gyration for all structures
- *x2file.txt*: chi squared for all structures
- *x2\_lt\_3p0.txt*: weights file selecting only frames with chi squared < 3.0

/spectra

- *spec\_\*.ciq*: scattering curves scaled to correct I(0) for each structure

## Trajectory Filtering

Now we can filter out the best fit structures and visualize them using the Extract Utilities. More information can be found in the [Extract Utilities documentation](#).

- Select 'Tools' from the Main Menu.
- Click the 'Extract Utilities' button.

Extract Utilities

run name

extract trajectory

Trajectory Input

reference pdb   or  Local: gag\_start.pdb

trajectory file name   or  Local: run\_0.dcd

output file name (pdb or dcd)

extract SAS

select option

input name of weight file   or  Server: test1/run\_1

weight file

In this module we can select structures from the DCD we created from the Monte Carlo simulation using the weight files generated in the Chi-Square Filter module.

In this case, we chose to select the weight file from the server.

- Check the tick box labelled 'extract trajectory' (this will reveal the options shown in the screenshot)
- Select the usual 'reference pdb' and the DCD output from the Monte Carlo simulation
- Input 'best\_gag.dcd' as the 'output filename'
- Choose 'weight file' from the 'select option' listbox.
  - Where it says 'input name of weight file' select the 'x2\_lt\_3p0.txt' file generated in the last step that selects only the frames which had a chi squared value of less 3.0 (You may have chosen a different value and file name.)
- Click 'Submit'

When the process is finished your output should look like the one below.

```
=====
DATA FROM RUN: Thu Jul  7 20:37:17 2016

reading frames from results/users/skrueger/test1/run_0.dcd
writing frames to run_1/extract_utilities/best_gag.dcd
wrote 27 frames to run_1/extract_utilities/best_gag.dcd

=====

percent done: 100.0
```

In the event that none of your frames pass the filter then you can download these preprepared files and try the filtering process:

- DCD: [run\\_0.dcd](#)
- weights file: [x2\\_lt\\_3p0.txt](#)

### What have we generated:

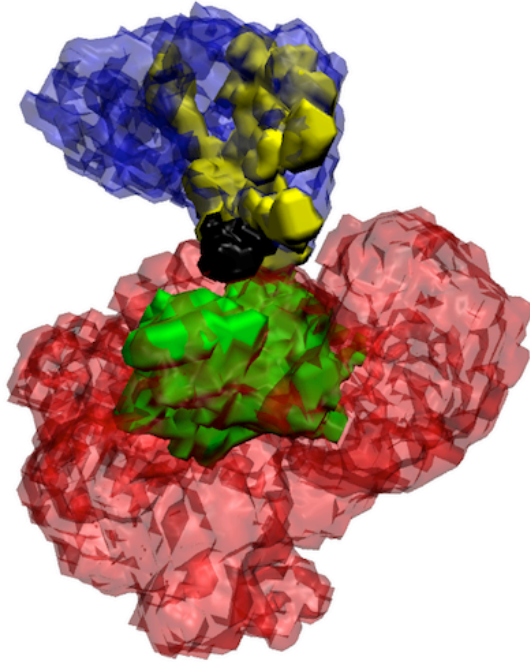
test1/run\_1/extract\_utilities

- *best\_gag.dcd*: DCD containing only the frames for which the theoretical scattering curve is a good match to experiment.

### Visualization

Download the 'best\_gag.dcd' file as you did the unfiltered DCD and then visualize the structure again in VMD (you will need to load a suitable PDB first as before). You should see that the filtered structures are all noticeably more compact than the starting structure and the majority of those in the unfiltered DCD.

Another way to visualize the structures sampled in the 'run\_0.dcd' and 'best\_gag.dcd' files for comparison is to use the Density Plot module. The density plot below shows the envelope sampled by all of the accepted structures as well as that sampled by only the best fit structures. The black region is the envelope represented by residues 283-353, which is approximately the alignment region that we defined in the Monomer Monte Carlo module. The blue and yellow regions represent the envelope sampled by residues 1-282 for all accepted structures (blue) and for the best fit structures (yellow). The red and green regions represent the envelope sampled by residues 354-431 for all accepted structures (red) and the best fit structures (green). This representation makes it easier to see that the envelope represented by the best fit structures is significantly smaller than that represented by all of the accepted structures. Remember that our sample has only 692 accepted structures and 27 best fit structures. For a real study, several thousand accepted structures would be needed to determine if this observation holds true.



More information on how to use the Density Plot module can be found in the [Density Plot documentation](#).

## Minimization of Best Structures

Now we can minimize the best fit gag structures using the Energy Minimization module. More information can be found in the [Energy Minimization documentation](#).

- Select 'Simulate' from the Main Menu.
- Click the 'Energy Minimization' button.

### Energy Minimization

run name:

reference pdb:  gag\_start.pdb or  Local: gag\_start.pdb

input filename (dcd or pdb):  No file selected. or  Server: test1/run\_1/extract\_utilities

PSF file name:  gag\_start.psf or  Local: gag\_start.psf

output file name (dcd):

number of processors:

keep run output files:

run type:

number of minimization steps:

---

### Advanced Input

Check Box for Advanced Input:


- **run name:** user defined name of folder that will contain the results.
- **reference pdb:** PDB file with naming information for coordinates that will be extracted. We are using the *gag\_start.pdb* file.
- **input filename (dcd or pdb):** file containing starting conformation(s) for simulation. The number of atoms must match that in the reference pdb. For files with multiple frames each one will be simulated. We are using the *run\_0.dcd* file.
- **PSF file name (dcd or pdb):** PSF file with topology information, must match the reference pdb and input dcd/pdb. Here we use the *gag\_start.psf* file.
- **output file name (dcd):** filename for the output DCD containing the final frames resulting from simulation
- **number of processors:** number of processors used to run the simulation (1-4). We are using use 2 processors.
- **keep run output files:** choice of whether to retain NAMD log files and other output from each simulation
- **run type:** select which of the four combinations of minimization and molecular dynamics to run
- **number of minimization steps:** number of steps of the conjugate gradient minimization to apply to each structure. We are using 1000 steps for each structure in the DCD file.
- Click 'Submit'
- **NOTE:** This will take awhile (~ 30 min).

When the process is finished your output should look like the one below.

```

=====
DATA FROM RUN: Thu Jul  7 20:42:31 2016
Total number of frames = 27
Minimized structures saved to : ./run_1/energy_minimization/
=====

```

progress:  100.0  
percent done: 100.0

## What have we generated:

test1/run\_1/energy\_minimization

- *min\_best\_gag.dcd*: DCD file containing the minimized best fit gag structures.
- *min\_best\_gag\_dcd.pdb*: reference PDB file that can be used to visualize the frames in the DCD file.

## Visualization

Download the 'min\_best\_gag.dcd' and 'min\_best\_gag.dcd.pdb' files and then visualize the structures in VMD. You can load 'best\_gag.dcd' (along with a suitable PDB file) again as well for comparison. Go through the two structures frame by frame. You should notice very little difference in the structures.

## Final SAS Curve Comparison

If desired, you can compare the minimized structures to the SANS data by calculating their theoretical SANS curves and comparing them to the SANS data again to see how different the best fit chi square values are after the minimization.

First, calculate the theoretical SANS curves using SasCalc.

The image shows a web-based form titled "SasCalc" with a dark blue background. The form is organized into several sections:

- run name:** A text input field containing "run\_2".
- reference pdb:** A "Browse..." button, a "No file selected." message, and an "or Browse server" button with the server path "Server: test1/run\_1/energy\_minimization /min\_best\_gag.dcd.pdb".
- trajectory file filename (dcd or pdb):** A "Browse..." button, a "No file selected." message, and an "or Browse server" button with the server path "Server: test1/run\_1/energy\_minimization /min\_best\_gag.dcd".
- number of q values:** A text input field containing "16".
- maximum q value:** A text input field containing "0.3".
- Neutron input:** A checked checkbox. Below it are several input fields:
  - number of contrast points:** A dropdown menu set to "1".
  - D2O percentage [1]:** A text input field containing "100.0".
  - I(0) [1]:** A text input field containing "0.04".
  - number of exchangeable H regions:** A dropdown menu set to "1".
  - exchangeable H region [1]:** A dropdown menu set to "moltype protein".
  - fraction of exchangeable H [1]:** A text input field containing "0.95".
  - number of deuterated regions:** A dropdown menu set to "0".
- X-ray input:** An unchecked checkbox.
- Advanced Input:**
  - SasCalc method:** A dropdown menu set to "fixed number of golden vectors".
  - number of golden vectors:** A text input field containing "35".
  - check box to enable HyPred pRDF solvent model:** An unchecked checkbox.
- Buttons:** "Submit" and "Reset to default values" buttons at the bottom left.

On the right side of the form, there is a vertical blue bar with the text "DOCS" and "FEEDBACK" in white.

run name:

Set the run name to run\_2. Once the inputs have been entered, click 'Submit'.

Once the run is complete, you should see outputs like those below.

```

=====
DATA FROM RUN: Fri Jul 8 15:50:48 2016

Processed 27 DCD frame(s)
Data stored in directory: run_2/sascalc/neutron_D2Op_100
=====

progress:
percent done: 100.0

```

### What have we generated:

test1/run\_2/sascalc/neutron\_D2Op\_100

- \*.iq: files containing theoretical scattering data for all frames in the DCD file (692 files in this case).
- \*.log: log files containing information about each structure and calculation inputs (692 files in this case).
- D2Op\_100.pdb: input PDB file with element names including deuterium atoms that were added as a result of H-D exchange of exchangeable hydrogen atoms or deuteration of non-exchangeable hydrogen atoms.
- HDexchange\_Info\_D2Op\_100.txt: file describing which hydrogen atoms were replaced with deuterium as a result of H-D exchange of exchangeable hydrogen atoms.
- Deuteration\_Info\_D2Op\_100.txt: file describing which hydrogen atoms were replaced with deuterium atoms as a result of deuteration of non-exchangeable hydrogen atoms.

Then, compare the theoretical SANS curves to the SANS data using Chi-Square Filter.

### Chi-Square Filter

run name	<input type="text" value="run_2"/>	
interpolated data file	<input type="button" value="Browse..."/> No file selected.	or <input type="button" value="Browse server"/> Server: test1/sans_data.dat
I(0)	<input type="text" value="0.04"/>	
SAS type	<input type="text" value="SasCalc"/>	
SAS data path	<input type="button" value="Browse server for a path"/> Server: test1/run_2/sascalc/neutron_D2Op_100	
chi-square type	<input type="text" value="reduced chi-square"/>	
number of weight files	<input type="text" value="0"/>	

---

Advanced Input

Check Box for Advanced Input

run name:

- Since we already have a chi\_square\_filter folder in the run\_1 directory, set the run name to run\_2.
- **number of weight files**
- Set the 'number of weight files' to 0 since we are already dealing with the best fit structures.
- Click 'Submit'.

Once the run is complete, you should see outputs like those below.

```

=====
DATA FROM RUN: Fri Jul 8 15:59:08 2016

Data stored in directory: ./run_2/chi_square_filter/neutron_D2Op_100

PROCESSED 27 SAS FILES:

>> The BEST and WORST SAS spectra are in the file named : bestworstfile.txt
>> The AVERAGE SAS spectra is in the file named : averagefile.txt
>> Chi-square, Rg, and filename are in the file named : x2file.txt

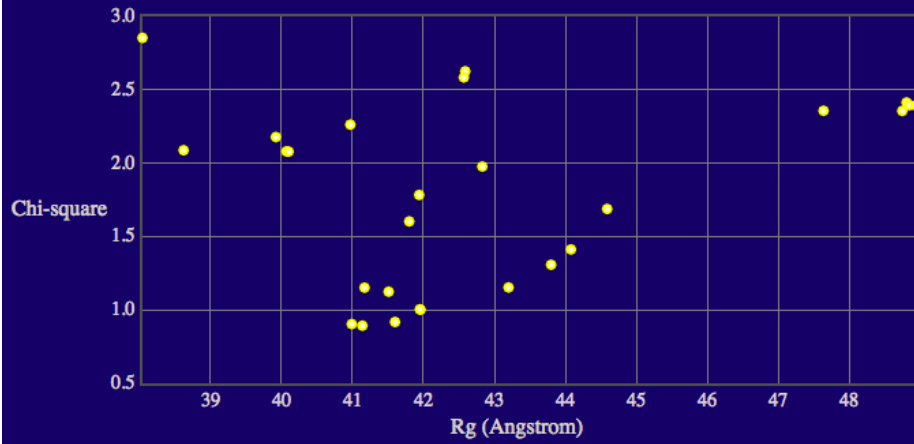
BEST SINGLE STRUCTURE IS NUMBER 15 WITH X2 = 0.892735 :          spectra:
run_2_00015

WORST SINGLE STRUCTURE IS NUMBER 27 WITH X2 = 2.854126:          spectra:
run_2_00027
=====

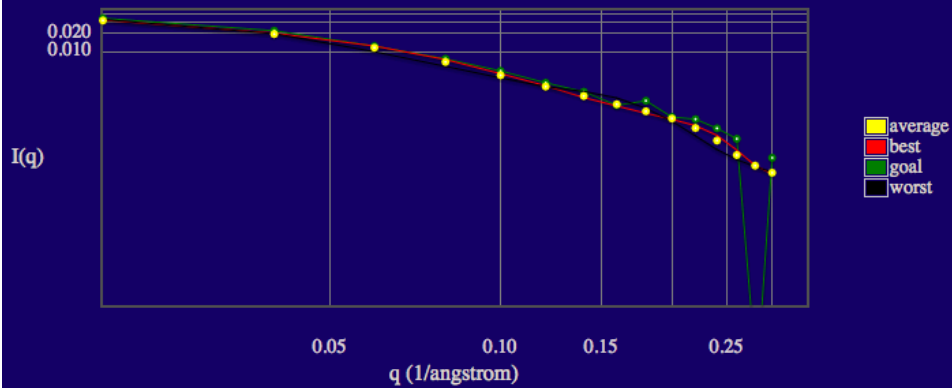
```

progress:   
percent done: 100.0

Chi-Square Distribution



Scattering Plots



**What have we generated:**

test1/run\_2/chi\_square\_filter

- *averagefile.txt*: average scattering curve for all structures
- *bestworstfile.txt*: best and worst scattering curves selected from all structures
- *sas\_spectra\_plot.txt*: goal, best, worst and average scattering curves
- *x2\_vs\_rg\_plot.txt*: chi squared against radius of gyration for all structures
- *x2file.txt*: chi squared for all structures

/spectra

- *spec\_\*.ciq*: scattering curves scaled to correct I(0) for each structure

## References

1. [Conformation of the HIV-1 Gag Protein in Solution](#) S. A. K. Datta, J. E. Curtis, W. Ratcliff, P. K. Clark, R. M. Crist, J. Lebowitz, S. Krueger, A. Rein, J. Mol. Biol. 365, 812-824 (2007). [BIBTeX](#), [Endnote](#), [Plain Text](#)
2. [SASSIE: A program to study intrinsically disordered biological molecules and macromolecular ensembles using experimental scattering restraints](#) J. E. Curtis, S. Raghunandan, H. Nanda, S. Krueger, Comp. Phys. Comm. 183, 382-389 (2012). [BIBTeX](#), [EndNote](#), [Plain Text](#)

[Return to Main Documents Page](#)

[Go to top](#)

Supported via CCP-SAS a joint EPSRC (EP/K039121/1) and NSF (CHE-1265821) grant